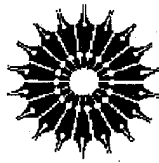




آشنایی با آنالیز عددی

کندال ای. اتکینسن

ترجمه دکتر علی دانایی



آشنایی با

آنالیز عددی

کندال ای. اتکینسن

ترجمه دکتر علی دانایی

مرکز نشر دانشگاهی



An Introduction to Numerical Analysis
Second Edition

Kendall E. Atkinson

John Wiley & Sons, 1989

آشنایی با آنالیز عددی
تألیف کندال ای. اتکینسن
ترجمه دکتر علی دانایی
ویراسته دکتر منوچهر وصال، دکتر محمدهادی شفیعیها
طراح جلد: شکوه پیله فروشها
نسخه پرداز: زهرا دلاوری
تهیه نمایه: فاطمه پیوندی
حروفچین: الهه عموحسن
ناظر چاپ: علی صادقی
مرکز نشر دانشگاهی
چاپ اول ۱۳۸۷
تعداد ۱۰۰۰
۹۵۰۰ تومان
لیتوگرافی: سعید
چاپ: دایره سفید
حق چاپ برای مرکز نشر دانشگاهی محفوظ است

فهرست نویسی پیش از انتشار کتابخانه ملی جمهوری اسلامی ایران

Atkinson, Kendall E.

سرشناسه: اتکینسون، کندال ای. / عنوان و نام پدیدآور: آشنایی با آنالیز عددی / کندال ای. اتکینسن؛ ترجمه علی دانایی؛ ویراسته منوچهر وصال، محمدهادی شفیعیها.

مشخصات نشر: تهران: مرکز نشر دانشگاهی، ۱۳۸۷.

مشخصات ظاهری: هشت، ۷۹۲ ص.

فروست: مرکز نشر دانشگاهی، ۱۲۹۳. ریاضی، آمار، و رایانه؛ ۱۵۵.

شابک: 978-964-01-1293-9

وضعیت فهرست نویسی: فیبا

یادداشت: عنوان اصلی: *An introduction to numerical analysis, 2nd ed, c1989.*

یادداشت: کتابنامه.

موضوع: آنالیز عددی.

شناسه افزوده: دانایی، علی، ۱۳۱۹-، مترجم.

شناسه افزوده: مرکز نشر دانشگاهی.

رده بندی کنگره: ۱۳۸۷ ۲۵ الف/ QA۲۹۷

رده بندی دیویی: ۵۱۸

شماره کتابشناسی ملی: ۱۱۲۳۹۲

بسم الله الرحمن الرحيم

فهرست

صفحه	عنوان
هفت	پیشگفتار مترجم
۱	پیشگفتار
۳	۱ خطا: منشأ، انتشار و تحلیل آن
۳	۱.۱ مقدمات ریاضی
۱۲	۲.۱ نمایش رایانه‌ی اعداد
۲۰	۳.۱ تعاریف و منشأهای خطا
۲۷	۴.۱ انتشار خطاها
۳۴	۵.۱ خطاهای مجموعیابی
۴۱	۶.۱ پایداری در آنالیز عددی
۵۰	مراجع
۵۲	مسائل
۶۰	۲ ریشه‌یابی معادله‌های غیرخطی
۶۴	۱.۲ روش نیم‌سازی (تصفیف)
۶۷	۲.۲ روش نیوتن
۷۵	۳.۲ روش خط قاطع
۸۳	۴.۲ روش مولر
۸۷	۵.۲ یک نظریه کلی برای روشهای بارستی تک نقطه‌ای
۹۶	۶.۲ برون‌یابی ایتکن برای دنباله‌های خطی - همگرای

۹۹	محاسبه عددی ریشه‌های چندگانه	۷.۲
۱۰۴	الگوریتم ریشه‌یابی برنت	۸.۲
۱۰۸	ریشه‌های چندجمله‌یها	۹.۲
۱۱۹	دستگاه معادلات غیرخطی	۱۰.۲
۱۲۵	روش نیوتن برای دستگاه‌های غیرخطی	۱۱.۲
۱۲۹	بهینه‌سازی نامقیّد	۱۲.۲
۱۳۴	مراجع	
۱۳۷	مسائل	

۱۴۹	نظریه درونیابی	۳
۱۵۰	نظریه درونیابی چندجمله‌یی	۱.۳
۱۵۷	تفاضلات منقسم نیوتن	۲.۳
۱۶۸	تفاضلات متناهی و جدول جهتدار فرمولهای درونیابی	۳.۳
۱۷۳	خطاها در داده‌ها و تفاضلات پیشرو	۴.۳
۱۷۶	نتایج دیگری درباره خطای درونیابی	۵.۳
۱۸۱	درونیابی ارمیت	۶.۳
۱۸۶	درونیابی چندجمله‌یی تکه‌یی	۷.۳
۲۰۲	درونیابی مثلثاتی	۸.۳
۲۱۱	مراجع	
۲۱۳	مسائل	

۲۲۳	تقریب توابع	۴
۲۲۴	قضیه وایرستراس و قضیه تیلر	۱.۴
۲۲۸	مسأله تقریب مینیمکس	۲.۴
۲۳۱	مسأله تقریب کمترین مربعات	۳.۴
۲۳۵	چندجمله‌یهای متعامد	۴.۴
۲۴۵	مسأله تقریب کمترین مربعات (ادامه)	۵.۴
۲۵۲	تقریبهای مینیمکس	۶.۴
۲۵۶	تقریبهای نزدیک مینیمکس	۷.۴
۲۷۰	مراجع	
۲۷۲	مسائل	

۲۸۰	انتگرالگیری عددی	۵
۲۸۳	قاعده ذوزنقه‌یی و قاعده سیمپسون	۱.۵
۲۹۵	فرمولهای انتگرالگیری نیوتن-کوتس	۲.۵
۳۰۳	انتگرالگیری گاوسی	۳.۵
۳۲۰	فرمولهای خطای مجانبی و کاربردهای آنها	۴.۵
۳۳۷	انتگرالگیری عددی خودکار	۵.۵
۳۴۵	انتگرالهای تکین	۶.۵
۳۵۶	مشتقگیری عددی	۷.۵
۳۶۳	مراجع	
۳۶۶	مسائل	

۳۷۳	روشهای عددی برای معادلات دیفرانسیل معمولی	۶
۳۷۶	وجود، یکتایی و نظریه پایداری	۱.۶
۳۸۲	روش اویلر	۲.۶
۴۰۱	روشهای چندگامی	۳.۶
۴۰۷	روش میانگاهی	۴.۶
۴۱۲	روش ذوزنقه‌یی	۵.۶
۴۲۰	الگوریتمی پیشگو-تصحیح‌کننده از مرتبه پایین	۶.۶
۴۲۹	پیدا کردن روشهای چندگامی از مراتب بالاتر	۷.۶
۴۴۴	نظریه همگرایی و پایداری برای روشهای چندگامی	۸.۶
۴۶۴	معادلات دیفرانسیل سرسخت و روش خطوط	۹.۶
۴۷۵	روشهای تک گامی و روشهای رونگه-کوتا	۱۰.۶
۴۹۳	مسائل مقدار مرزی	۱۱.۶
۵۰۹	مراجع	
۵۱۲	مسائل	

۵۲۳	جبر خطی	۷
۵۲۳	فضاهای برداری، ماتریسها و دستگاههای خطی	۱.۷
۵۳۳	ویژه بردارها و شکلهای متعارف (کانونی) ماتریسها	۲.۷
۵۴۵	نرمهای برداری و ماتریسی	۳.۷
۵۵۶	قضایای همگرایی و اختلال	۴.۷

۵۶۳	مراجع
۵۶۴	مسائل

۵۷۲	۸ حل عددی دستگاههای معادلات خطی
۵۷۳	۱.۸ حذف گاوسی
۵۸۲	۲.۸ محورگیری و مقیاس دهی در حذف گاوسی
۵۹۰	۳.۸ صورتهای دیگر حذف گاوسی
۵۹۹	۴.۸ تحلیل خطا
۶۱۲	۵.۸ روش تصحیح مانده
۶۱۷	۶.۸ روشهای بارستی
۶۲۶	۷.۸ پیشگویی خطا و شتاب
۶۳۱	۸.۸ حل عددی معادله پواسن
۶۳۸	۹.۸ روش گرادیان مزدوج
۶۴۹	مراجع
۶۵۲	مسائل

۶۶۳	۹ مسأله ویژه مقدار ماتریس
۶۶۴	۱.۹ جای ویژه مقدار، خطا، و قضایای پایداری
۶۸۲	۲.۹ روش توانی
۶۹۰	۳.۹ تبدیلات متعامد با استفاده از ماتریسهای هاوسهولدر
۷۰۲	۴.۹ ویژه مقدارهای یک ماتریس متقارن سه قطری
۷۰۷	۵.۹ روش QR
۷۱۴	۶.۹ محاسبه ویژه بردارها و بارست معکوس
۷۲۰	۷.۹ حل دستگاههای خطی به روش کمترین مربعات
۷۳۶	مراجع
۷۳۸	مسائل

۷۴۹	پیوستها
۷۵۶	پاسخهای تمرینهای انتخابی
۷۷۵	نمایه

پیشگفتار مترجم

کتابی که در اختیار شماست به چند دلیل برای ترجمه برگزیده شده است.

۱- کتابی است جامع در موضوعات مختلف آنالیز عددی که دانشجویان رشته‌های ریاضی، علوم کامپیوتر و مهندسی می‌توانند از آن استفاده کنند.

۲- کتابی است که برای پژوهشگران می‌تواند راهنمای بسیار مفیدی باشد.

۳- در موضوعی که وارد شده ابتدا به‌طور ساده آن را بیان کرده و دانشجویان یا پژوهشگر را به آخرین دستاوردهای آن موضوع رهبری می‌کند و این هدف را با معرفی مقالات و کتابهای سودمند که در پایان هر فصل آورده شده به نحو بسیار خوبی به انجام می‌رساند.

اساتید بزرگوار جناب آقای دکتر محمدهادی شفیعیها و جناب آقای دکتر منوچهر وصال به تناوب ویراستاری ادبی و علمی این کتاب را در مرکز نشر دانشگاهی با دقت و حوصله بسیار انجام داده‌اند و ضعف کار مترجم را برطرف نموده‌اند که از این بابت از این اساتید بزرگوار صمیمانه سپاسگزاری می‌شود و برای آنان طول عمر با برکت و توفیقات بیشتر آرزومندم. از کارکنان مرکز نشر دانشگاهی که با صبوری و حوصله در تایپ، غلطگیری و چاپ کتاب همکاری داشته‌اند تشکر می‌کنم.

از خوانندگانی که با راهنماییهای خود می‌توانند کاستیهای احتمالی را برطرف سازند پیشاپیش سپاسگزارم

علی دانایی

عضو هیأت علمی دانشگاه شیخ بهائی

پیشگفتار

ساختار اصلی ویراست دوم این کتاب، با چند بخش جدید که اضافه گردیده، اساساً مانند ویراست اول است. تمام بخشها بازنویسی شده‌اند، گاهی بسیار هم جدی این کار انجام گرفته و در بعضی مواقع مطالبی حذف شده‌اند. به علاوه، برخی از موضوعات جدید معرفی شده‌اند. مثلاً درونیابی مثلثاتی و FFT (بخش ۹.۶)، مشتقگیری عددی (بخش ۷.۵)، روش خطوط (بخش ۹.۶)، مسائل مقادیر مرزی (بخش ۱۱.۶)، روش گرادیان مزدوج (بخش ۹.۸) و حل دستگاههای معادلات خطی به روش کمترین مربعات (بخش ۷.۹). همچنین یک پیوست در مورد نرم‌افزار ریاضی اضافه شده است که بعضی بسته‌های نرم‌افزار شناخته‌تری را تشریح می‌کند. مراجع برای هر فصل به روز تبدیل شده است تا آخرین کتابها و آثار خواندنی پژوهشی را منعکس کند.

با عنایت به نرم‌افزار ریاضی، قویاً توصیه می‌کنم دانشجویان در حد امکان، برنامه‌های تجارتي با کیفیت بالای آنالیز عددی در یکی از کتابخانه‌های موجود را مورد استفاده قرار دهند. آنهایی که بیشتر به کار گرفته شده‌اند در IMSL و NAG می‌باشند که در پیوست تحت عنوان نرم‌افزار ریاضی در مورد آنها بحث شده است. در این کتاب تعدادی از بسته‌های نرم‌افزاری آنالیز عددی که برای اهداف خاص است توضیح داده شده‌اند. این نرم‌افزارها، وقتی کدهای اصلی مورد نیاز باشند، ارزشمندند، گرچه معمولاً برنامه‌های موجود در کتابخانه‌های تجارتي، برای برنامه‌سازی از مسائل، که دانشجویان با آن مواجه می‌شوند کافی‌اند. به علاوه کتابخانه‌های تجارتي معمولاً برای استفاده ساده‌ترند. همچنین بسیاری از بسته‌هایی که به منظور خاص تهیه شده و در این کتاب توضیح داده شده‌اند در کتابخانه‌های تجارتي مهم جای گرفته‌اند.

خطا: منشأ، انتشار و تحلیل آن

موضوع آنالیز عددی تهیه روشهای محاسباتی برای مطالعه و حل مسائل ریاضی است. در این کتاب روشهای عددی متداولترین مسائل ریاضی را مطالعه می‌کنیم و خطاهای موجود در این روشها را بررسی می‌نماییم. چون تقریباً تمام محاسبات امروزه با رایانه‌های رقمی انجام می‌گیرند از ضرورتهای این کار در انجام روشهای عددی نیز بحث خواهیم کرد.

مطالعه خطا مطلبی اساسی در آنالیز عددی است. بیشتر روشهای عددی جوابهایی به دست می‌دهند که فقط تقریبی از جواب درست مورد نظرند و درک، و در صورت امکان، توانایی برآورد یا کراندارکردن خطای حاصل حائز اهمیت است. در این فصل انواع مختلف خطاهایی که ممکن است در مسئله‌ای رخ دهند بررسی خواهند شد. نمایش اعداد در رایانه‌ها و درکنار آن خطا در حساب رایانه‌یی نیز بررسی می‌شود. قضایای کلی درباره انتشار خطاها در محاسبات، با بررسی مفصل خطا در روشهای مجموعیابی داده می‌شوند. بالاخره مفاهیم پایداری و مشروط کردن مسائل و روشهای عددی معرفی و تشریح خواهند شد. بخش اول شامل مقدمات ریاضی لازم برای این کار در فصول بعدی است.

۱.۱ مقدمات ریاضی

این بخش شامل مروری است بر قضایای حسابان که در این کتاب به‌کار برده می‌شوند. ابتدا بعضی

از قضایای مقدار میانگین را بیان می‌کنیم، سپس قضیهٔ تیلر برای توابع یک و دو متغیره را مورد بحث قرار می‌دهیم. این بخش را با نمادهایی که در فصول بعد کاربرد پیدا خواهند کرد پایان می‌دهیم.

قضیهٔ ۱.۱ (مقدار میانگین) بگیریم $f(x)$ بر بازهٔ متناهی $a \leq x \leq b$ پیوسته باشد، M و m را به صورت زیر تعریف می‌کنیم

$$m = \inf_{a \leq x \leq b} f(x), \quad M = \sup_{a \leq x \leq b} f(x)$$

پس برای هر عدد ζ در بازهٔ $[m, M]$ دست کم یک نقطهٔ ξ در $[a, b]$ وجود دارد به گونه‌ای که

$$f(\xi) = \zeta$$

به ویژه، نقاطی مانند \underline{x} و \bar{x} در $[a, b]$ هستند که برای آنها

$$m = f(\underline{x}), \quad M = f(\bar{x})$$

قضیهٔ ۲.۱ (مقدار میانگین) بگیریم $f(x)$ به ازای $a \leq x \leq b$ پیوسته و به ازای $a < x < b$ مشتقپذیر باشد. در این صورت دست کم یک نقطه چون ξ در (a, b) وجود دارد که برای آن

$$f(b) - f(a) = f'(\xi)(b - a)$$

قضیهٔ ۳.۱ (مقدار میانگین انتگرال) بگیریم $w(x)$ نامنفی و بر بازهٔ $[a, b]$ انتگرالپذیر باشد، و $f(x)$ در $[a, b]$ پیوسته باشد. در این صورت برای مقداری مانند $\xi \in [a, b]$ داریم

$$\int_a^b w(x)f(x)dx = f(\xi) \int_a^b w(x)dx$$

این قضایا در بیشتر کتابهای درسی مقدماتی حساب دیفرانسیل و انتگرال مورد بحث قرار گرفته‌اند و بنابراین ما از اثبات آنها صرف نظر می‌کنیم. بعضی از نتایج این قضایا در مسائل آخر این فصل بررسی شده‌اند.

یکی از مهمترین ابزارهای آنالیز عددی قضیهٔ تیلر و سری مربوطه به آن است که در سرتاسر این کتاب به کار برده شده است. این قضیه یک روش نسبتاً ساده‌ای برای تقریب توابع $f(x)$ به کمک چند جمله‌بیها، و بنابراین روشی برای محاسبهٔ $f(x)$ به دست می‌دهد.

قضیه ۴.۱ (قضیه تیلر) گیریم $f(x)$ ، $n + 1$ مشتق پیوسته، $n \geq 0$ ، بر $[a, b]$ ، داشته باشد و $x, x_0 \in [a, b]$ در این صورت

$$f(x) = p_n(x) + R_{n+1}(x) \quad (1.1.1)$$

$$p_n(x) = f(x_0) + \frac{(x - x_0)}{1!} f'(x_0) + \dots + \frac{(x - x_0)^n}{n!} f^{(n)}(x_0) \quad (2.1.1)$$

$$\begin{aligned} R_{n+1}(x) &= \frac{1}{n!} \int_{x_0}^x (x - t)^n f^{(n+1)}(t) dt \\ &= \frac{(x - x_0)^{n+1}}{(n + 1)!} f^{(n+1)}(\xi) \end{aligned} \quad (3.1.1)$$

که ξ مقداری است بین x و x_0 .

برهان نحوه به دست آوردن (۱.۱.۱) در بسیاری از کتابهای دیفرانسیل و انتگرال داده شده است. در به دست آوردن آن از انتگرالگیری جزء به جزء در اتحاد زیر که به دقت انتخاب شده، استفاده می شود

$$f(x) = f(x_0) + \int_{x_0}^x f'(t) dt$$

از n بار تکرار این عمل، روابط (۱.۱.۱) - (۳.۱.۱) با صورت انتگرالی باقیمانده $R_{n+1}(x)$ به دست می آید. صورت دوم $R_{n+1}(x)$ با استفاده از قضیه مقدار میانگین انتگرال، با انتخاب $w(t) = (x - t)^n$ حاصل می شود. ■

با استفاده از قضیه تیلر، فرمولهای استانده زیر را پیدا می کنیم

$$e^x = 1 + x + \frac{x^2}{2!} + \dots + \frac{x^n}{n!} + \frac{x^{n+1}}{(n + 1)!} e^{\xi_x} \quad (4.1.1)$$

$$\begin{aligned} \cos x &= 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \dots + (-1)^n \frac{x^{2n}}{(2n)!} \\ &+ (-1)^{n+1} \frac{x^{2n+2}}{(2n + 2)!} \cos(\xi_x) \end{aligned} \quad (5.1.1)$$

$$\begin{aligned} \sin x &= x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots + (-1)^{n-1} \frac{x^{2n-1}}{(2n - 1)!} \\ &+ (-1)^n \frac{x^{2n+1}}{(2n + 1)!} \cos(\xi_x) \end{aligned} \quad (6.1.1)$$

$$\begin{aligned} (1 + x)^\alpha &= 1 + \binom{\alpha}{1} x + \binom{\alpha}{2} x^2 + \dots + \binom{\alpha}{n} x^n \\ &+ \binom{\alpha}{n + 1} \frac{x^{n+1}}{(1 + \xi)^{n+1-\alpha}} \end{aligned} \quad (7.1.1)$$

که $\binom{\alpha}{k}$ به ازای هر عدد حقیقی α ، با تساوی زیر داده می‌شود

$$\binom{\alpha}{k} = \frac{\alpha(\alpha-1)\dots(\alpha-k+1)}{k!} \quad k = 1, 2, 3, \dots$$

در همه حالات نقطه مجهول ξ_x بین 0 و x واقع است.

یک حالت خاص مهم (۷.۱.۱) رابطه زیر است

$$\frac{1}{1-x} = 1 + x + x^2 + \dots + x^n + \frac{x^{n+1}}{1-x} \quad x \neq 1 \quad (۸.۱.۱)$$

در این حالت، $\alpha = -1$ و $-x$ جایگزین x شده است. باقیمانده صورتی ساده‌تر از (۷.۱.۱) دارد؛ و این مطلب به راحتی از ضرب طرفین رابطه (۸.۱.۱) در $1-x$ و سپس ساده کردن آن اثبات می‌شود. اگر ترتیب نوشتن رابطه (۸.۱.۱) را اندکی تغییر دهیم فرمول آشنای سری هندسی منتهای زیر به دست می‌آید

$$1 + x + x^2 + \dots + x^n = \frac{1-x^{n+1}}{1-x} \quad x \neq 1 \quad (۹.۱.۱)$$

نمایش سریهای نامتهای برای توابع طرف چپ روابط (۸.۱.۱) - (۴.۱.۱) را می‌توان با میل دادن $n \rightarrow \infty$ به دست آورد. سریهای نامتهای (۶.۱.۱) - (۴.۱.۱) به ازای همه مقادیر x همگرا هستند و سریهای (۸.۱.۱) - (۷.۱.۱) به ازای $|x| < 1$ همگرا هستند. از فرمول (۸.۱.۱) سری هندسی نامتهای معروف

$$\frac{1}{1-x} = \sum_{k=0}^{\infty} x^k \quad |x| < 1 \quad (۱۰.۱.۱)$$

به دست می‌آید.

سری تیلر برای هر تابع به قدر کافی مشتقپذیر $f(x)$ را می‌توان مستقیماً از تعریف (۲.۱.۱)، با هر تعداد جمله دلخواه محاسبه کرد، ولی به دلیل پیچیدگی مشتقگیری در بسیاری از توابع $f(x)$ ، اغلب بهتر است که تقریبهای چندجمله‌بیهای تیلر $p_n(x)$ یا سری تیلر آنها را با استفاده از فرمولهای قبلی (۸.۱.۱) - (۴.۱.۱) به دست آورد. ما سه مثال می‌آوریم که در آنها جملات خطا ساده‌تر از وقتی هستند که (۳.۱.۱) مستقیماً به کار برده می‌شود.

مثال ۱. گیریم $f(x) = e^{-x^2}$. در رابطه (۴.۱.۱) به جای x مقدار $-x^2$ را قرار می‌دهیم تا به دست آوریم

$$e^{-x^2} = 1 - x^2 + \frac{x^4}{2!} - \dots + (-1)^n \frac{x^{2n}}{n!} + (-1)^{n+1} \frac{x^{2x+2}}{(n+1)!} e^{\xi_x}$$

که در آن $-x^2 \leq \xi_x \leq 0$.

۲. گیریم $f(x) = \tan^{-1}(x)$. با قراردادن $x = -u^2$ در (۸.۱.۱) خواهیم داشت

$$\frac{1}{1+u^2} = 1 - u^2 + u^4 - \dots + (-1)^n u^{2n} + (-1)^{n+1} \frac{u^{2n+2}}{1+u^2}$$

و پس از انتگرالگیری بر بازه $[0, x]$ داریم

$$\tan^{-1}(x) = x - \frac{x^3}{3} + \frac{x^5}{5} - \dots + (-1)^n \frac{x^{2n+1}}{2n+1} + (-1)^{n+1} \int_0^x \frac{u^{2n+2}}{1+u^2} du \quad (۱۱.۱.۱)$$

با استفاده از قضیه مقدار میانگین در انتگرالها

$$\int_0^x \frac{u^{2n+2}}{1+u^2} du = \frac{x^{2n+3}}{2n+3} \cdot \frac{1}{1+\xi_x^2}$$

که در آن ξ_x بین 0 و x واقع است.

۳. گیریم $f(x) = \int_0^1 \sin(xt) dt$ با استفاده از (۶.۱.۱) داریم

$$\begin{aligned} f(x) &= \int_0^1 \left[xt - \frac{x^2 t^2}{2!} + \dots + (-1)^{n-1} \frac{(xt)^{2n-1}}{(2n-1)!} \right. \\ &\quad \left. + (-1)^n \frac{(xt)^{2n+1}}{(2n+1)!} \cos(\xi_{xt}) \right] dt \\ &= \sum_{j=1}^n (-1)^{j-1} \frac{x^j t^{j-1}}{(2j)!} + (-1)^n \frac{x^{2n+1}}{(2n+1)!} \int_0^1 t^{2n+1} \cdot \cos(\xi_{xt}) dt \end{aligned}$$

که در آن ξ_{xt} بین 0 و xt واقع است. انتگرال در باقیمانده به آسانی به وسیله $1/(2n+2)$ کراندار می‌شود؛ ولی می‌توانیم آن را به شکل ساده‌تری برگردانیم. اگرچه اثبات نکرده‌ایم، ولی می‌توان نشان داد که $\cos(\xi_{xt})$ تابعی است پیوسته از t . پس با استفاده از قضیه مقدار میانگین انتگرال داریم

$$\int_0^1 \sin(xt) dt = \sum_{j=1}^n (-1)^{j-1} \frac{x^j t^{j-1}}{(2j)!} + (-1)^n \frac{x^{2n+1}}{(2n+2)!} \cos(\zeta_x)$$

که ζ_x بین 0 و x قرار دارد.

قضیه تیلر در فضای دوبعدی بگیریم $f(x, y)$ تابع داده شده‌ای از دو متغیر مستقل x و y باشد. نشان خواهیم داد که چگونه قضیه قبلی تیلر برای بسط $f(x, y)$ حول نقطه داده شده (x_0, y_0) تعمیم داده می‌شود. این نتایج به‌سادگی برای توابع بیش از دو متغیر نیز قابل تعمیم‌اند. بگیریم نماد $L(x_0, y_0; x_1, y_1)$ معرف مجموعه همه نقاط واقع بر پاره خط واصل بین نقاط (x_0, y_0) و (x_1, y_1) باشد.

قضیه ۵.۱ بگیریم (x_0, y_0) و $(x_0 + \xi, y_0 + \eta)$ نقاط داده شده‌ای باشند، و فرض می‌کنیم $f(x, y)$ برای همه نقاط (x, y) در یک همسایگی $L(x_0, y_0; x_0 + \xi, y_0 + \eta)$ پیوسته و $n + 1$ بار مشتقپذیر باشد. در این صورت

$$\begin{aligned} f(x_0 + \xi, y_0 + \eta) &= f(x_0, y_0) + \sum_{j=1}^n \frac{1}{j!} \left[\xi \frac{\partial}{\partial x} + \eta \frac{\partial}{\partial y} \right]^j f(x, y) \Big|_{\substack{x=x_0 \\ y=y_0}} \\ &\quad + \frac{1}{(n+1)!} \left[\xi \frac{\partial}{\partial x} + \eta \frac{\partial}{\partial y} \right]^{n+1} f(x, y) \Big|_{\substack{x=x_0 + \theta\xi \\ y=y_0 + \theta\eta}} \end{aligned} \quad (12.1.1)$$

برای مقداری از $0 \leq \theta \leq 1$. نقطه $(x_0 + \theta\xi, y_0 + \theta\eta)$ نقطه‌ای بر خط $L(x_0, y_0; x_0 + \xi, y_0 + \eta)$ است.

برهان ابتدا به معنای نماد مشتق در (۱۲.۱.۱) توجه نمایید. به‌عنوان مثال

$$\left[\xi \frac{\partial}{\partial x} + \eta \frac{\partial}{\partial y} \right]^2 f(x, y) = \xi^2 \cdot \frac{\partial^2 f(x, y)}{\partial x^2} + 2\xi\eta \cdot \frac{\partial^2 f(x, y)}{\partial x \partial y} + \eta^2 \cdot \frac{\partial^2 f(x, y)}{\partial y^2}$$

معنی زیر نمایه‌های $x = x_0$ و $y = y_0$ آن است که مشتقات مختلف بایستی در (x_0, y_0) محاسبه شوند.

برهان (۱۲.۱.۱) بر پایه استفاده از قضیه اولیه تیلر برای

$$F(t) = f(x_0 + t\xi, y_0 + t\eta) \quad 0 \leq t \leq 1$$

استوار است.

با استفاده از قضیه ۴.۱

$$F(1) = F(0) + \frac{F'(0)}{1!} + \dots + \frac{F^{(n)}(0)}{n!} + \frac{F^{(n+1)}(\theta)}{(n+1)!}$$

به ازای $0 \leq \theta \leq 1$ ، روشن است که، $F(0) = f(x_0, y_0)$ و $F(1) = f(x_0 + \xi, y_0 + \eta)$ برای مشتق اول داریم

$$\begin{aligned} F'(t) &= \xi \frac{\partial f(x_0 + t\xi, y_0 + t\eta)}{\partial x} + \eta \frac{\partial f(x_0 + t\xi, y_0 + t\eta)}{\partial y} \\ &= \left[\xi \frac{\partial}{\partial x} + \eta \frac{\partial}{\partial y} \right] f(x, y) \Bigg|_{\substack{x=x_0+t\xi \\ y=y_0+t\eta}} \end{aligned}$$

مشتقات مراتب بالاتر به طریق مشابه محاسبه می‌شوند.

مثال به‌عنوان یک مثال ساده، بسط تابع $f(x, y) = x/y$ را حول نقطه $(\epsilon, 2) = (x_0, y_0)$ نظر می‌گیریم. $n = 1$ بگیریم. پس

$$\begin{aligned} \frac{x}{y} &= f(\epsilon, 2) + (x - \epsilon) \frac{\partial f(\epsilon, 2)}{\partial x} + (y - 2) \frac{\partial f(\epsilon, 2)}{\partial y} \\ &+ \frac{1}{2} \left[(x - \epsilon)^2 \cdot \frac{\partial^2 f(\delta, \gamma)}{\partial x^2} + 2(x - \epsilon)(y - 2) \frac{\partial^2 f(\delta, \gamma)}{\partial x \partial y} \right. \\ &\quad \left. + (y - 2)^2 \cdot \frac{\partial^2 f(\delta, \gamma)}{\partial y^2} \right] \\ &= 3 + \frac{1}{2}(x - \epsilon) - \frac{3}{4}(y - 2) \\ &+ \frac{1}{2} \left[(x - \epsilon)^2 \cdot 0 - 2(x - \epsilon)(y - 2) \frac{1}{\gamma^2} + (y - 2)^2 \cdot \frac{2\delta}{\gamma^3} \right] \end{aligned}$$

که در آن (δ, γ) نقطه‌ای است واقع بر خط $L(\epsilon, 2; x, y)$. برای (x, y) نزدیک به $(\epsilon, 2)$

$$\frac{x}{y} \doteq 3 + \frac{1}{2}x - \frac{3}{4}y$$

نمودار $z = 3 + \frac{1}{2}x - \frac{3}{4}y$ صفحه مماس بر نمودار $z = x/y$ در نقطه $(\epsilon, 2, 3) = (x, y, z)$ است.

بعضی نمادهای ریاضی از مفاهیم چندی در این کتاب استفاده شده است که شکل‌های ساده‌تر آنها در فصل‌های اولیه مورد نیازند. این مفاهیم شامل قضایایی دربارهٔ تفاضلهای منقسم توابع، فضاهای برداری و نرمهای بردار و ماتریس هستند. حداقل نمادهای لازم اکنون معرفی می‌شوند و بسط کامل آنها به جاهای مناسبتر کتاب موکول خواهد شد.

برای یک تابع داده شده $f(x)$ با فرض متمایز بودن x_2, x_1, x_0 تعاریف زیر را اختیار می‌کنیم

$$f[x_0, x_1] = \frac{f(x_1) - f(x_0)}{x_1 - x_0} \quad f[x_0, x_1, x_2] = \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0} \quad (13.1.1)$$

این نمادها را به ترتیب تفاضلهای منقسم مرتبه اول و مرتبه دوم $f(x)$ می‌نامند. روابط بین این نمادها و مشتقات $f(x)$ چنین‌اند

$$f[x_0, x_1] = f'(\xi) \quad f[x_0, x_1, x_2] = \frac{1}{2} f''(\zeta) \quad (14.1.1)$$

که ξ بین x_1, x_0 و ζ بین مینیم و ماکسیم x_2, x_1, x_0 واقع‌اند. تفاضلهای منقسم، برخلاف آنچه که ممکن است از رابطه (۱۳.۱.۱) انتظار داشت، مستقل از ترتیب شناسه‌های خود هستند. به عبارت دقیقتر

$$f[x_0, x_1] = f[x_1, x_0] \\ f[x_0, x_1, x_2] = f[x_i, x_j, x_k] \quad (15.1.1)$$

برای هر جایگشت (i, j, k) از اعداد $(0, 1, 2)$. اثبات این ویژگیها و سایر ویژگیها به صورت مسائل به‌خواننده واگذار می‌شوند. شرح کامل تفاضلهای منقسم در بخش ۲.۳ از فصل ۳ داده شده است. موضوعهای فضاهاى بردارى، ماتریسها و نرمهای بردار و ماتریس در فصل ۷، بلافاصله پس از فصلهای راجع به جبر خطی عددی مطرح شده‌اند. برخی از این مطالب را در اینجا ذکر می‌کنیم و اثبات آنها را برای فصل ۷ می‌گذاریم. دو فضای برداری هستند که کاربردهای بسیار دارند. این فضاها عبارت‌اند از:

$$\mathbf{R}^n = \left\{ x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \mid x_1, \dots, x_n \text{ اعداد حقیقی} \right\} \\ C[a, b] = \{ f(t) \mid f(t) \text{ پیوسته و حقیقی مقدار, } a \leq t \leq b \}$$

برای $x, y \in \mathbf{R}^n$ و عدد حقیقی α ، عبارت $x + y$ و αx چنین تعریف می‌شوند

$$x + y = \begin{bmatrix} x_1 + y_1 \\ \vdots \\ x_n + y_n \end{bmatrix} \quad \alpha x = \begin{bmatrix} \alpha x_1 \\ \vdots \\ \alpha x_n \end{bmatrix}$$

برای $f, g \in C[a, b]$ و عدد حقیقی α عبارت $f + g$ و αf چنین تعریف می‌شوند

$$(f + g)(t) = f(t) + g(t) \quad (\alpha f)(t) = \alpha f(t) \quad a \leq t \leq b$$

نرمهای بردارها برای اندازه‌گیری قدر بردار به‌کار برده می‌شوند. برای \mathbf{R}^n مقدماتاً دو نرم مختلف را تعریف می‌کنیم

$$\|x\|_{\infty} = \text{Max}_{1 \leq i \leq n} |x_i| \quad x \in \mathbf{R}^n \quad (16.1.1)$$

$$\|x\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2} \quad x \in \mathbf{R}^n \quad (17.1.1)$$

و برای $C[a, b]$ نرم را چنین تعریف می‌کنیم

$$\|f\|_{\infty} = \text{Max}_{a \leq t \leq b} |f(t)| \quad f \in C[a, b] \quad (18.1.1)$$

می‌توان نشان داد که تعاریف بالا در سه ویژگی خاص همهٔ نرمها صدق می‌کنند.

۱. $\|v\| \geq 0$; $\|v\| = 0$ اگر و فقط اگر $v = 0$ ، یعنی v بردار صفر باشد.

۲. $\|\alpha v\| = |\alpha| \cdot \|v\|$ ، به ازای همهٔ بردارهای v و اعداد حقیقی α .

۳. $\|v + w\| \leq \|v\| + \|w\|$ ، به ازای همهٔ بردارهای v و w .

رابطهٔ (۳) را معمولاً نابرابری مثلث می‌نامند. توضیح این نامگذاری و همچنین تعمیم بیشتر خواص نرمها برای $C[a, b]$ در فصل ۴ داده شده است.

نرمها را می‌توان برای ماتریسها نیز تعریف نمود. برای هر ماتریس $n \times n$

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & & & \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}$$

نرم را چنین تعریف می‌کنیم

$$\|A\|_{\infty} = \text{Max}_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \quad (19.1.1)$$

این تعریف، در ویژگیهای نرم بردار صدق می‌کند. علاوه بر آن می‌توان نشان داد که

$$\|AB\|_{\infty} \leq \|A\|_{\infty} \|B\|_{\infty} \quad (20.1.1)$$

$$\|Ax\|_{\infty} \leq \|A\|_{\infty} \|x\|_{\infty} \quad x \in \mathbf{R}^n \quad \text{برای هر} \quad (21.1.1)$$

که A و B ماتریسهای $n \times n$ دلخواه‌اند. اثبات این قضایا به‌عنوان مسأله واگذار می‌شود. این روابط به‌طور کلیتر در فصل ۷ نیز مورد بحث واقع شده‌اند.

مثال فضای برداری \mathbf{R}^2 و ماتریسهای 2×2 را در نظر می‌گیریم. در حالت خاص فرض می‌کنیم

$$A = \begin{bmatrix} 1 & -1 \\ 3 & 2 \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad y = Ax = \begin{bmatrix} -1 \\ 7 \end{bmatrix}$$

پس

$$\|A\|_{\infty} = 5 \quad \|x\|_{\infty} = 2 \quad \|y\|_{\infty} = 7$$

و برقراری (21.1.1) به‌سادگی دیده می‌شود. برای آنکه نشان دهیم نابرابری (21.1.1) را نمی‌توان بهتر کرد، بردار $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ را در نظر می‌گیریم. داریم

$$x = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad y = Ax = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$$

پس

$$\|y\|_{\infty} = 5 = \|A\|_{\infty} \|x\|_{\infty}$$

۲.۱ نمایش رایانه‌یی اعداد

رایانه‌های رقمی ابزار اصلی محاسبات در آنالیز عددی هستند و نتیجتاً درک شیوه کار آنها بسیار مهم است. در این بخش چگونگی نمایش اعداد را در رایانه‌ها ملاحظه می‌کنیم و در سایر بخشها نتایج این گونه نمایش اعداد و حساب با رایانه را خواهیم دید.

بیشتر رایانه‌ها یک وجه اعداد صحیح و یک وجه با ممیز شناور برای نمایش اعداد دارند. وجه اعداد صحیح فقط برای نمایش اعداد صحیح به کار برده می‌شود و چندان مورد علاقه ما نیست. وجه با ممیز شناور در نمایش اعداد حقیقی مورد استفاده قرار می‌گیرد. اعداد مجاز می‌توانند

دارای اندازه‌های بسیار متفاوت باشند، ولی محدودیتهایی چه از نظر اندازه و چه از لحاظ تعداد ارقام برای آنها وجود دارد. نمایش با ممیز شناور بستگی نزدیکی به آنچه که در بسیاری از کتابهای ریاضی دبیرستانی، نماد علمی خوانده می‌شود، دارد.

پایه اعدادی که در رایانه‌ها استفاده می‌شوند به‌ندرت اعشاری است. بیشتر رایانه‌های رقمی پایه دستگاه عددی ۲ (دودویی) یا شکلهای مختلفی از آن، چون پایه ۸ (هشت‌هستی) یا پایه ۱۶ (شانزده شانزدهی) را به‌کار می‌برند.

مثال (الف) در پایه ۲، ارقام ۰ و ۱ هستند. به‌عنوان مثال برای تبدیل یک عدد در پایه ۲ به اعشاری داریم

$$(1101101)_2 = 1 \times 2^6 + 1 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 = 27.25$$

هنگامی که از اعداد در پایه دیگری استفاده می‌کنیم معمولاً نماد $(x)_\beta$ را به‌کار می‌بریم تا نشان دهیم عدد x در دستگاه اعداد به پایه β بیان شده است.

(ب) در پایه ۱۶، ارقام عبارت‌اند از: ۰، ۱، ۲، ...، ۹، A، ...، F. مثلاً برای تبدیل چنین عددی به اعشاری داریم

$$(56C.F)_{16} = 5 \times 16^3 + 6 \times 16^2 + 12 \times 16^1 + 15 \times 16^0 = 13889375$$

تبدیل اعداد اعشاری به دودویی در ضمن مسائل بررسی شده است.

گیریم β عدد پایه‌ای باشد که در رایانه به‌کار گرفته شده است. پس یک عدد غیرصفر x در رایانه اساساً به شکل زیر ذخیره می‌شود

$$x = \sigma \cdot ({}^{\circ}a_1 a_2 \dots a_t)_\beta \cdot \beta^e \quad (1.2.1)$$

که در آن $\sigma = +1$ یا -1 و $0 \leq a_i \leq \beta - 1$ یک عدد صحیح و

$$({}^{\circ}a_1 a_2 \dots a_t)_\beta = \frac{a_1}{\beta^1} + \frac{a_2}{\beta^2} + \dots + \frac{a_t}{\beta^t}$$

حرف σ را علامت، e را نما و $({}^{\circ}a_1 a_2 \dots a_t)_\beta$ را جزء کسری عدد با ممیز شناور x خوانند. عدد β را همچنین پایه دستگاه عددی و نقطه قبل از a_1 در (۱.۲.۱) را ممیز پایه می‌نامند، برای

مثال ممیز اعشاری ($\beta = 10$)، ممیز دو دویی ($\beta = 2$). عدد صحیح t تعداد ارقام در پایه β را در نمایش فوق معین می‌کند. همیشه فرض می‌کنیم

$$a_1 \neq 0$$

که نمایش با ممیز شناور، به اصطلاح، نرمال شده را می‌دهد. همچنین فرض می‌کنیم که

$$L \leq e \leq U \quad (2.2.1)$$

که اندازه ممکن x را محدود می‌نماید. عدد $x = 0$ همیشه مجاز و مستلزم نمایش ویژه‌ای است. جدول ۱.۱ مقادیر β, t, L و U را برای شماری از رایانه‌های معمولی به دست می‌دهد. اساس استفاده از β, t, L و U برای معین کردن مشخصات حسابی از قراردادهای کتاب [فورسایت^۱ و همکاران، ۱۹۷۷، ص ۱۱] گرفته شده است. بعضی رایانه‌ها (مثل، CDC CYBER) جایگزاریهایی برای ممیز پایه به کار می‌برند. ما کارانه‌ای نما را به گونه‌ای تغییر داده‌ایم که حدود برای اندازه عدد با ممیز شناور، با توجه به نظریه این بخش، درست از آب درآید. نتایج نمایشهای با دقت مضاعف را در رایانه‌هایی که در سخت‌افزار آنها این درجه دقت گنجانده شده است نیز اضافه نموده‌ایم. در جدول ۱.۱ ستونهای بیشتری هستند که بعداً توضیح داده می‌شوند.

قطع کردن و گرد کردن بیشتر اعداد حقیقی x را نمی‌توان دقیقاً به وسیله نمایش با ممیز شناور، که قبلاً گفته شد، نشان داد، و بنابراین، در صورت امکان، بایستی با عددی نزدیک که قابل نمایش در رایانه است تقریب زده شوند. برای یک عدد دلخواه داده شده x ، $\text{fl}(x)$ را معرف رایانه‌یی آن، در صورت وجود، می‌گیریم. دو راه اساسی برای به دست آوردن $\text{fl}(x)$ از x ، وجود دارد: قطع کردن و گرد کردن.

گیریم عدد حقیقی x به شکل زیر نوشته شده باشد

$$x = \sigma \cdot ({}_r a_1 a_2 \dots a_t a_{t+1} \dots) \beta \cdot \beta^e \quad (3.2.1)$$

که $a_1 \neq 0$ ، e در (۲.۲.۱) صدق می‌کند. نمایش قطع شده رایانه‌یی x ، با عبارت زیر داده می‌شود

$$\text{fl}(x) = \sigma \cdot ({}_r a_1 \dots a_t) \beta \cdot \beta^e \quad (4.2.1)$$

جدول ۱.۱ نمایشهای با ممیز شناور در رایانه‌های گوناگون

ماشین	S/D	R/C	β	t	L	U	δ	M
CDC CYBER 170	S	R	۲	۴۸	-۹۷۶	۱۰۷۱	۳,۵۵E-۱۵	۲,۸۱E۱۴
CDC CYBER 205	S	C	۲	۴۷	-۲۸,۶۲۶	۲۸,۷۱۸	۱,۴۲E-۱۴	۱,۴۱E۱۴
CRAY-1	S	C	۲	۴۸	-۸۱۹۲	۸۱۹۱	۷,۱۱E-۱۵	۲,۸۱E۱۴
DEC VAX	S	R	۲	۲۴	-۱۲۷	۱۲۷	۵,۹۶E-۸	۱,۶۸E۷
DEC VAX	D	R	۲	۵۳	-۱۰۲۳	۱۰۲۳	۱,۱۱E-۱۶	۹,۰۱E۱۵
HP-11C,15C	S	R	۱۰	۱۰	-۹۹	۹۹	۵,۰۰E-۱۰	۱,۰۰E۱۰
IBM 3033	S	C	۱۶	۶	-۶۴	۶۳	۹,۵۴E-۷	۱,۶۸E۷
IBM 3033	D	C	۱۶	۱۴	-۶۴	۶۳	۲,۲۲E-۱۶	۷,۲۱E۱۶
Intel 8087	S	R	۲	۲۴	-۱۲۶	۱۲۷	۵,۹۶E-۸	۱,۶۸E۷
Intel 8087	D	R	۲	۵۳	-۱۰۲۲	۱۰۲۳	۱,۱۱E-۱۶	۹,۰۱E۱۵
PRIME 850	S	R	۲	۲۳	-۱۲۸	۱۲۷	۱,۱۹E-۷	۸,۳۹E۶
PRIME 850	S	C	۲	۲۳	-۱۲۸	۱۲۷	۱,۱۹E-۷	۸,۳۹E۶
PRIME 850	D	C	۲	۴۷	-۳۲,۸۹۶	۳۲,۶۳۹	۱,۴۲E-۱۴	۱,۴۱E۱۴

به این نمادها توجه کنید:

S/D: دقت معمولی یا مضاعف.

R/C: گردکردن یا قطع کردن.

β : عدد پایه.

t: رقمهای بعد از ممیز [(۱.۲.۱)] را ببینید

U, L: حدود نما [(۲.۲.۱)] را ببینید

δ : واحد گردکردن [(۱۲.۲.۱)] را ببینید

M: کران دقیق اعداد صحیح [(۱۶.۲.۱)] را ببینید

دلیل معرفی قطع کردن آن است که در بسیاری از رایانه‌ها به جای گردکردن در هر عمل حساب، از قطع کردن استفاده می‌کنند.

نمایش گردشده x به صورت زیر است:

$$fl(x) = \begin{cases} \sigma \cdot ({}^{\circ}r a_1 \dots a_t)_{\beta} \cdot \beta^e & 0 \leq a_{t+1} < \frac{\beta}{4} \quad (5.2.1) \\ \sigma \cdot [({}^{\circ}r a_1 \dots a_t)_{\beta} + ({}^{\circ}r \dots {}^{\circ}1)_{\beta}] \cdot \beta^e & \frac{\beta}{4} \leq a_{t+1} < \beta \quad (6.2.1) \end{cases}$$

در فرمول اخیر، $({}^{\circ}r \dots {}^{\circ}1)_{\beta}$ معرف β^{-t} است. اگرچه این تعریف $fl(x)$ تا حدی پیچیده است ولی همان تعریف استاندارد گردکردن را، که اغلب از دستگاه اعشاری آموخته‌ایم، به دست می‌دهد. گاهی این تعریف به گونه‌ای دیگر به کار برده می‌شود تا گردکردن ناریب به دست آید. در چنین حالتی، اگر

$$a_{t+1} = \frac{\beta}{4} \quad (1) \quad \text{و} \quad a_j = 0 \quad \text{برای} \quad j \geq t+2 \quad (2)$$

آنگاه گردکردن را اضافی یا نقصانی می‌گیرند بسته به اینکه به ترتیب، a_t فرد یا زوج باشد. این عمل به قاعده گردکردن نااریب می‌انجامد که بیشتر افراد برای گردکردن اعداد اعشاری آموخته‌اند، ولی ما از این پس تعریف ساده‌تر (۵.۲.۱) - (۶.۲.۱) را می‌پذیریم. برای بیشتر اعداد حقیقی x ، داریم $\text{fl}(x) \neq x$. با نگاه کردن به خطای نسبی (یا درصد)، می‌توان نشان داد که

$$\frac{x - \text{fl}(x)}{x} = -\varepsilon \quad (۷.۲.۱)$$

$$-\beta^{-t+1} \leq \varepsilon \leq 0 \quad \text{قطع شده} \text{fl}(x) \quad (۸.۲.۱)$$

$$-\frac{1}{4}\beta^{-t+1} \leq \varepsilon \leq \frac{1}{4}\beta^{-t+1} \quad \text{گرد شده} \text{fl}(x) \quad (۹.۲.۱)$$

ما نتیجه (۸.۲.۱) را برای قطع کردن اثبات می‌کنیم و اثبات نتیجه (۹.۲.۱) را برای گردکردن به عنوان یک مسئله واگذار می‌کنیم.

گیریم $\sigma = +1$ ، زیرا حالت $\sigma = -1$ علامت ε را تغییر نخواهد داد. از (۳.۲.۱) و (۴.۲.۱) به دست می‌آوریم

$$x - \text{fl}(x) = ({}^{\circ}r^{\circ} \dots {}^{\circ}a_{t+1} \dots)_{\beta} \cdot \beta^e$$

فرض می‌کنیم $\gamma = \beta - 1$

$$0 \leq x - \text{fl}(x) \leq ({}^{\circ}r^{\circ} \dots {}^{\circ}\gamma \gamma \dots)_{\beta} \cdot \beta^e$$

$$= \gamma [\beta^{-t-1} + \beta^{-t-2} + \dots] \cdot \beta^e$$

$$= \gamma \left[\frac{\beta^{-t-1}}{1 - \beta^{-1}} \right] \cdot \beta^e = \beta^{-t+e}$$

$$0 \leq \frac{x - \text{fl}(x)}{x} \leq \frac{\beta^{-t+e}}{({}^{\circ}ra_1a_2 \dots)_{\beta} \cdot \beta^e}$$

$$\leq \frac{\beta^{-t}}{({}^{\circ}r1^{\circ} \dots)_{\beta}} = \beta^{-t+1}$$

که حکم (۸.۲.۱) را ثابت می‌کند و اثبات (۹.۲.۱) به طریق مشابه صورت می‌گیرد. فرمول (۷.۲.۱) معمولاً به شکل هم‌ارز آن

$$\text{fl}(x) = (1 + \varepsilon)x \quad (۱۰.۲.۱)$$

نوشته می‌شود که ε توسط (۸.۲.۱) یا (۹.۲.۱) داده شده است. بنابراین $f(x)$ را می‌توان به‌عنوان یک اختلال کوچک نسبی x منظور کرد. این فرمول $f(x)$ ضمناً به ما امکان می‌دهد که دقیقاً نتایج خطاهای گردکردن/قطع‌کردن در عملیات حساب رایانه را بررسی کنیم. مثالهایی از این دست در بخشهای بعد آورده شده است. تعریف $f(x)$ و استفاده از (۱۰.۲.۱) اقتباس از [ویلکینسن، ۱۹۶۳] است، و در بسیاری از جاها در تحلیل نتایج خطای گردکردن مورد استفاده واقع شده است. دقت نمایش با ممیز شناور ما اکنون دو اندازه که تصور نسبتاً روشنی از دقت ممکن در نمایش با ممیز شناور را به دست می‌دهند، معرفی می‌کنیم. اندازه اول ارتباط نزدیکی با نتیجه خطای قبلی $f(x)$ که با (۷.۲.۱) تا (۹.۲.۱) داده شده، دارد.

واحد گردکردن یک رایانه عددی است مانند δ که در شرایط زیر صدق می‌کند: (۱) δ عدد مثبتی است با ممیز شناور، و (۲) کوچکترین عددی است که برای آن نامساوی زیر برقرار است.

$$f(1 + \delta) > 1 \quad (11.2.1)$$

بنابراین برای هر عدد $\delta < \hat{\delta}$ ، با ممیز شناور داریم $f(1 + \hat{\delta}) = 1$ که $\hat{\delta}$ و ۱ در حساب رایانه یکی هستند. این رابطه یک اندازه دقیقی به ما می‌دهد تا بدانیم چند رقم دقیق ممکن است در نمایش یک عدد وجود داشته باشد. در بسیاری از برنامه‌های رایانه‌های قابل حمل و با کیفیت بالا از این واحد گردکردن استفاده می‌کنند تا بیشترین دقتی را که ممکن است در رایانه به کار برد به دست آورند. واحد گردکردن δ به سادگی محاسبه و به صورت زیر بیان می‌شود

$$\delta = \begin{cases} \beta^{-t+1} & \text{طبق تعریف قطع‌کردن برای } f(x) \\ \frac{1}{4}\beta^{-t+1} & \text{طبق تعریف گردکردن برای } f(x) \end{cases} \quad (12.2.1)$$

ما این مطلب را برای حساب گردشده روی یک ماشین دودویی اثبات می‌کنیم. ابتدا باید نشان دهیم

$$f(1 + 2^{-t}) > 1 \quad (13.2.1)$$

می‌نویسیم

$$\begin{aligned} 1 + 2^{-t} &= [(0_r 100 \dots)_2 + (0_r 00 \dots 010 \dots)_2] \cdot 2^t \\ &\quad \text{در موضع } 1 + t \\ &= (0_r 10 \dots 010 \dots)_2 \cdot 2^t \end{aligned} \quad (14.2.1)$$

صورت $\text{fl}(1 + 2^{-t})$ را با توجه به اینکه در موضع $t + 1$ ام ارقام بعد از ممیز یک عدد ۱ وجود دارد تشکیل می‌دهیم. پس با توجه به (۶.۲.۱) داریم

$$\text{fl}(1 + 2^{-t}) = (\overset{\uparrow}{\text{در موضع } t} 0 \text{ ر } 10 \dots 0 1) 2 \cdot 2^1 = 1 + 2^{-t+1}$$

بنابراین (۱۳.۲.۱) برقرار است، اگرچه

$$\text{fl}(1 + 2^{-t}) \neq 1 + 2^{-t}$$

این واقعیت که δ نمی‌تواند از 2^{-t} کوچکتر باشد به سادگی از بازنگری (۱۴.۲.۱) حاصل می‌شود. اگر $\delta < \delta$ ، آنگاه $1 + \delta$ در موضع $t + 1$ ام ارقام بعد از ممیز در (۱۴.۲.۱) یک صفر خواهد داشت: و در این حال تعریف (۵.۲.۱) برای گردکردن ایجاب می‌کند که $\text{fl}(1 + \delta) = 1$.

یک اندازهٔ دوم برای بیشترین دقت ممکن در یک نمایش با ممیز شناور، پیدا کردن بزرگترین عدد صحیح M است به طوری که برای یک عدد صحیح m

$$0 \leq m \leq M \Rightarrow \text{fl}(m) = m \quad (15.2.1)$$

این رابطه هم ایجاب می‌کند که $\text{fl}(M + 1) \neq M + 1$. اثبات تساوی زیر به‌عنوان یک مسئله واگذار می‌شود

$$M = \beta^t \quad (16.2.1)$$

اعداد M و δ ، همراه با گردکردن (R) و قطع کردن (C) برای رایانه‌های مختلف، چه برای گردکردن (R) و چه برای قطع کردن (C)، در جدول ۱.۱ داده شده‌اند.

مثال برای رایانه‌های PRIME با دقت معمولی

$$M = 2^{23} = 8388608$$

بنابراین تمام اعداد صحیح با شش رقم اعشاری و اکثر اعداد با هفت رقم اعشاری را دقیقاً می‌توان در نمایش با ممیز شناور با دقت معمولی ذخیره کرد. برای واحد گردکردن

$$\delta = 2^{-23} \approx 19 \times 10^{-7} \quad \text{با حساب قطع شده}$$

$$\delta = 2^{-24} \approx 5.96 \times 10^{-8} \quad \text{با حساب گردشده}$$

کاربرهای PRIME برای دقت معمولی هم حساب قطع شده و هم حساب گردشده را در اختیار دارند.

تقریباً در همهٔ حالات، حساب گردشده بر حساب قطع شده ارجحیت بسیار دارد. این امر با جزئیات بیشتر در بخشهای بعدی بررسی شده است، ولی دلیل اصلی آن علامت اریبی ε در (۸.۲.۱)، در مقایسه با ناریبی آن در (۹.۲.۱) نهفته است.

پی‌ریز و سرریز وقتی که کرانه‌های نما در (۲.۲.۱) نقض شوند عدد متناظر x در (۱.۲.۱) نمی‌تواند در رایانه نشان داده شود. اکنون می‌بینیم که در این صورت در بارهٔ برد ممکن طول x چه می‌توانیم بگوییم. کوچکترین عدد مثبت با ممیز شناور عبارت است از

$$x_L = (0r10 \dots 0)_\beta \cdot \beta^L = \beta^{L-1}$$

با استفاده از $\gamma = \beta - 1$ ، بزرگترین عدد مثبت با ممیز شناور چنین خواهد شد

$$x_U = (0r\gamma \dots \gamma)_\beta \cdot \beta^U = (1 - \beta^{-t}) \cdot \beta^U$$

بنابراین تمام اعداد با ممیز شناور x باید در رابطهٔ زیر صدق کنند

$$x_L \leq |x| \leq x_U \quad (17.2.1)$$

در زبان فورترن اغلب رایانه‌ها، اگر یک عمل حساب به عددی چون x که $|x| > x_U$ ، منجر شود آنگاه خطای مخرب‌ی بروز خواهد کرد که موجب توقف برنامه می‌شود. این خطا را خطای سرریز می‌خوانند. برعکس، اگر

$$0 < |x| < x_L$$

آنگاه معمولاً $\text{fl}(x)$ صفر گرفته می‌شود و محاسبات ادامه می‌یابد. این خطا را خطای پی‌ریز می‌نامند.

مثال محاسبهٔ $x = s^{10}$ را در یک رایانهٔ بزرگ IBM در نظر بگیرید. در این صورت یک خطای پی‌ریز وجود دارد اگر

$$s^{10} < 16^{-65}$$

بنابراین اگر $|s| < 10^{-8} \times 1.49$ ، آنگاه $x = s^{10}$ صفر گرفته می‌شود. همچنین اگر نامساوی

$$s^{10} > 16^{63}$$

یا هم‌ارز آن، نامساوی

$$|s| > 3.9 \times 10^7$$

برقرار باشد یک خطای سرریز پیش خواهد آمد.

۳.۱ تعاریف و منشأهای خطا

اکنون به ذکر یک رده‌بندی کلی از راههای مهمی که خطاها در جواب یک مسأله وارد می‌شوند و بعضی راهها که بیرون از حیطهٔ ریاضیات قرار می‌گیرند، می‌پردازیم. مطلب را با چند تعریف ساده دربارهٔ خطا شروع می‌کنیم.

هنگام حل یک مسأله، ما دنبال جواب دقیق یا درستی هستیم که آن را با x_T نشان می‌دهیم. معمولاً تقریبهایی در حل مسأله پدید می‌آیند که یک جواب تقریبی x_A از آن نتیجه می‌شود. خطای x_A را با

$$\text{Error}(x_A) \equiv x_A \text{ خطای} = x_T - x_A$$

تعریف می‌کنیم. در بسیاری از مسائل ترجیح می‌دهیم که درصد یا خطای نسبی x_A را مطالعه کنیم

$$\text{Rel}(x_A) \equiv x_A \text{ نسبی خطای} = \frac{x_T - x_A}{x_T}$$

به شرطی که $x_T \neq 0$. به این مطلب قبلاً در اندازه‌گیری خطا در $\text{fl}(x)$ در (۷.۲.۱) اشاره کرده‌ایم.

$$x_A = \frac{1}{v} = 2.7142857\dots, \quad x_T = e = 2.7182818\dots \quad \text{مثال}$$

$$\text{Error}(x_A) = 0.003996\dots \quad \text{Rel}(x_A) = 0.00147\dots$$

به جای خطای نسبی، اغلب مفهوم ارقام با معنی را به‌کار می‌بریم. گوئیم x_A دارای m رقم (اعشاری) با معنی نسبت به x_T است اگر اندازهٔ خطای $x_T - x_A$ قدرمطلقاً نابزرگتر از 5 در $(m+1)$ امین رقم x_T داشته باشد، وقتی شمارش از اولین رقم غیرصفر x_T از چپ به راست صورت گیرد.

مثال (الف) $x_T = \frac{1}{e}$ ، $x_A = 0.333$ ، $|x_T - x_A| = 0.00033$. چون خطا در چهارمین رقم از چپ به راست، یعنی در اولین رقم غیرصفر x_T ، کوچکتر از 5 است، گوئیم x_A دارای سه رقم با معنی نسبت به x_T است.

$$(ب) \quad |x_T - x_A| = 0.002, \quad x_A = 23,494, \quad x_T = 23,496$$

عبارت x_A چهار رقم با معنی نسبت به x_T دارد، زیرا خطا در پنجمین رقم از چپ به راست اولین رقم غیر صفر در x_T ، کوچکتر از ۵ است. توجه داشته باشید که اگر x_A به چهار رقم گرد شود، یک خطای دیگر پیش می‌آید و x_A دیگر دارای چهار رقم با معنی نخواهد بود.

$$(ج) \quad |x_T - x_A| = 0.00006, \quad x_A = 0.2144, \quad x_T = 0.2138$$

عدد x_A دو رقم با معنی نسبت به x_T دارد، نه سه رقم.

آنچه در زیر می‌آید گاهی در اندازه‌گیری ارقام با معنی به‌کار می‌رود. اگر

$$\left| \frac{x_T - x_A}{x_T} \right| \leq 5 \times 10^{-m-1} \quad (1.3.1)$$

آنگاه x_A نسبت به x_T ، m رقم با معنی دارد. برای نشان دادن این موضوع حالتی را در نظر می‌گیریم که $1 < |x_T| \leq 10$ ، در این صورت از (۱.۳.۱) نتیجه می‌شود که

$$|x_T - x_A| \leq 5 \times 10^{-m-1} |x_T| < 0.5 \times 10^{-m}$$

چون $1 < |x_T| \leq 10$ ، در نتیجه x_A دارای m رقم با معنی است. اثبات برای حالت کلی x_T اساساً به همین گونه است، با استفاده از $x_T = \hat{x}_T \times 10^e$ و حالت $1 < |x_T| < 10$ ، e یک عدد صحیح است. توجه داشته باشید که (۱.۳.۱) شرط کافی برای داشتن m رقم با معنی x_A است و نه یک شرط لازم برای آن. مثالهای (الف) و (ب) که در بالا داده شدند، یک رقم با معنی بیشتر از آنچه که آزمون (۱.۳.۱) بیان می‌کند، دارند.

منشأهای خطا اکنون به ذکر یک رده‌بندی تقریبی از منشأهای مهم خطا می‌پردازیم.

(۱م) مدل‌سازی ریاضی از یک مسأله فیزیکی. منظور از تهیه مدل ریاضی برای یک وضعیت فیزیکی کوشش برای به‌دست آوردن روابط ریاضی بین کمیت‌هایی است که از لحاظ فیزیکی مورد توجه هستند. به علت پیچیدگی واقعیت‌های فیزیکی، برای ساختن یک مدل ریاضی که راحت‌تر بتوان به‌کاربرد، از فرض‌های ساده‌کننده گوناگونی استفاده می‌کنند. بر اثر این فرض‌ها، مدل حاصل از لحاظ دقت محدودیت‌هایی خواهد داشت که برحسب موارد استفاده از مدل، ممکن است این محدودیت‌ها در دسرافین باشند یا نباشند. در حالتی که مدل به اندازه کافی دقیق نباشد، حل عددی مدل نمی‌تواند این نبود اساسی دقت را مرتفع سازد.

مثال پرتابه‌ای به جرم m به هوا پرتاب شده است که مسیر حرکتش همیشه نزدیک به سطح زمین باقی

می‌ماند. دستگاه مختصات xyz در نظر می‌گیریم که مرکز آن بر سطح زمین و محور z عمود بر سطح زمین و جهت مثبت آن به طرف بالا انتخاب شده است. گیریم موقعیت پرتابه در لحظه t با استفاده از نمادگذاری استاندارد نظریه میدان برداری به صورت $\mathbf{r}(t) = x(t)\mathbf{i} + y(t)\mathbf{j} + z(t)\mathbf{k}$ باشد. یک مدل برای حرکت پرتابه توسط قانون دوم نیوتن به صورت زیر داده می‌شود

$$m \frac{d^2 \mathbf{r}(t)}{dt^2} = -mg\mathbf{k} - b \frac{d\mathbf{r}(t)}{dt} \quad (2.3.1)$$

که $b > 0$ مقداری ثابت و g شتاب حاصل از گرانش است. این معادله می‌گوید تنها نیروهایی که بر پرتابه اثر می‌کنند عبارت‌اند از: (۱) نیروی گرانش زمین، (۲) نیروی اصطکاک که مستقیماً با قدر مطلق سرعت $|\mathbf{v}(t)| = |d\mathbf{r}(t)/dt|$ متناسب است و در جهت عکس مسیر حرکت عمل می‌کند. در بعضی حالات، این یک مدل عالی است، و حتی لزومی به وارد کردن جمله اصطکاک ندارد. ولی این مدل، نیروهای مقاومت را که به صورت عمود بر صفحه حرکت، عمل می‌کنند در نظر نمی‌گیرد. مثلاً جهت عرضی باد را که اثر پدیده کوریولیس^۱ را خنثی می‌کند نادیده می‌گیرد. همچنین ممکن است نیروی اصطکاک در (۲.۳.۱) متناسب با $|\mathbf{v}(t)|^\alpha$ با $\alpha \neq 1$ باشد.

اگر مدلی برای منظوره‌های فیزیکی از لحاظ دقت مناسب باشد، می‌خواهیم یک طرح عددی برای حفظ این دقت به‌کار ببریم. ولی اگر مدل مناسب نباشد، آنالیز عددی نمی‌تواند دقت آن را بیشتر کند مگر به تصادف. از سوی دیگر پسندیده نیست که با وارد کردن مطالبی که در مقایسه با پدیده مورد مطالعه، نسبتاً کم اهمیت هستند، مدلی بسازیم که بیش از حد لازم پیچیده باشد. گاهی یک مدل پیچیده‌تر ممکن است مشکلات آنالیز عددی دیگری ایجاد کند، بی‌آنکه دقت بیشتر قابل توجهی عاید سازد. برای مطالعه کتابهایی که صراحتاً در زمینه مدلیابی ریاضی در علوم مطالبی دارند می‌توانید به بندر^۲ (۱۹۷۸)، لین و سگال^۳ (۱۹۷۴)، مکی و تامسن^۴ (۱۹۷۳) و رویینو^۵ (۱۹۷۵) مراجعه کنید.

(۲م) اشتباهات. تا پیش از پیدایش رایانه‌ها، خطاهای اتفاقی حساب همواره مسأله جدی بودند. طرحهایی برای رسیدگی، بعضاً خیلی پرزحمت، ابداع کرده بودند تا چنین خطاهایی را در صورت یافتن، پیش از آنکه از موضع خطا زیاد دور شوند، کشف و اصلاح کنند. برای مثالی از این نوع در طرحهای رسیدگی، که هنگام حل یک دستگاه معادلات خطی به‌کار رفته‌اند، به اثر فاده‌یوا^۶ مراجعه کنید.

1. Coriolis

2. Bender

3. Lin and Segal

4. Maki and Thompson

5. Rubinow

6. Fadeeva

با پیدایش رایانه‌های رقمی، نوع اشتباهات تغییر کرده است. خطاهای اتفاقی محاسباتی (مثلاً بدعمل کردن رایانه) اکنون نسبتاً نادر، و معمولاً خطاهای برنامه‌نویسی مشکل اساسی شده‌اند. اغلب یک خطای برنامه چندین بار ضمن اجرای برنامه تکرار می‌شود و وجود آن با خروجی عددی ناموجه آشکار می‌شود (اگرچه پیدا کردن منشأ خطا باز ممکن است دشوار باشد). ولی وقتی برنامه‌های رایانه‌ای رفته رفته پیچیده‌تر و طولانیتر می‌شوند وجود یک خطای کوچک برنامه و کشف و اصلاح آن ممکن است دشوار باشد، ولو اینکه این خطا ممکن است یک اختلاف ظریف ولی اساسی در نتایج به بار آورد. بدین جهت خطازدایی خوب برنامه خیلی مهم است، اگرچه در نگاه اول قابل توجه به نظر نیاید.

برای کشف خطاهای برنامه‌نویسی، داشتن راههایی که بتوان دقت خروجی برنامه را رسیدگی نمود حائز اهمیت است. وقتی برای اولین بار برنامه را اجرا می‌کنید باید، در صورت امکان، حالتی را که جواب درست آنها را می‌دانید به‌کار برید. برای یک برنامه پیچیده، آن را به زیر برنامه‌های کوچکتری خرد کنید تا بتوانید هر کدام را جداگانه مورد آزمایش قرار دهید. وقتی که کل برنامه رسیدگی و از درستی آن اطمینان حاصل شد خروجی را زیر نظر داشته باشید و ببینید که قابل قبول است یا نیست.

(م ۳) عدم قطعیت در داده‌های فیزیکی. بیشتر داده‌های حاصل از تجربیات فیزیکی متضمن یک نوع خطا یا عدم قطعیت هستند. این عدم قطعیت، بر دقت هر محاسبه‌ای که به این داده‌ها بستگی داشته باشد اثر می‌گذارد و در دقت جوابها محدودیتی ایجاد می‌کند. تکنیکهای تحلیل نتایج در محاسبات دیگر این خطا، کاملاً مشابه تکنیکهای تحلیل نتایج خطاهای گردکردن است، اگرچه خطا در داده‌ها معمولاً بسیار بزرگتر از خطاهای گردکردن است. در بخشهای بعد، مطالب این موضوع بیشتر مورد بحث قرار خواهند گرفت.

(م ۴) خطاهای ماشین. منظور از خطاهای ماشین خطاهایی هستند که ذاتاً در اثر به‌کار بردن شکل ممیز شناور اعداد ایجاد می‌شوند. بالاخص، منظور ما خطاهای قطع کردن/گردکردن و خطاهای پی‌ریز/سرریز هستند. خطاهای قطع کردن و گردکردن نتیجه محدودیت طول ماننسیس ممیز شناور هستند، و این خطاها در تمام عملیات حساب رایانه‌ای پیش می‌آیند. کلیه شکلهای این خطاهای ماشینی در بخش ۲.۱ مورد بحث قرار گرفته‌اند. در بخشهای آینده بعضی از پی‌آمدهای این خطاها مورد توجه قرار می‌گیرند. همچنین، برای سادگی نمادگذاری، از این پس هر جا که بتوانیم، از اصطلاح خطای گردکردن که شامل قطع کردن نیز خواهد بود استفاده خواهیم کرد.

(م ۵) خطای برشی ریاضی. این نام به خطای تقریب در حل عددی یک مسأله ریاضی اطلاق می‌شود، و خطایی است که عمدتاً با موضوع آنالیز عددی در ارتباط است. این خطا شامل

تقریب فرایندهای نامحدود از راه فرایندهای محدود، جایگزین کردن مسائل محاسبه ناپذیر با مسائل محاسبه پذیر است. برای روشن شدن بیشتر موضوع چند مثال می آوریم.

مثال (الف) با استفاده از دو جمله اول سری تیلر (۷.۱.۱)، داریم

$$\sqrt{1+x} \doteq 1 + \frac{1}{2}x \quad (۳.۳.۱)$$

که وقتی x کوچک باشد تقریب خوبی است. برای مطالعه مبحث کلی تقریب توابع، به فصل ۴ مراجعه کنید.

(ب) برای محاسبه یک انتگرال روی $[0, 1]$ از فرمول زیر استفاده می کنیم

$$\int_0^1 f(x) dx \doteq \frac{1}{n} \sum_{j=1}^n f\left(\frac{j-1}{n}\right) \quad n = 1, 2, 3, \dots \quad (۴.۳.۱)$$

این عمل، قاعده انتگرالگیری عددی نقطه میانی خوانده می شود. برای توضیح بیشتر به آخرین قسمت بخش ۲.۵ مراجعه کنید. موضوع کلی انتگرالگیری عددی در فصل ۵ آمده است.

(ج) برای مسأله معادله دیفرانسیل

$$Y'(t) = f(t, Y(t)) \quad Y(t_0) = Y_0 \quad (۵.۳.۱)$$

از تقریب مشتق

$$Y'(t) \doteq \frac{Y(t+h) - Y(t)}{h}$$

که در آن h مقدار کوچکی است استفاده می کنیم. گیریم $t_j = t_0 + jh$ ، $j \geq 0$ ، یک تابع جواب تقریبی $y(t_j)$ را با تساوی

$$\frac{y(t_{j+1}) - y(t_j)}{h} = f(t_j, y(t_j))$$

تعریف می کنیم. بنابراین داریم

$$y(t_{j+1}) = y(t_j) + hf(t_j, y(t_j)) \quad j \geq 0 \quad y(t_0) = Y_0$$

این روش همان روش اویلر برای حل مسأله مقدار اولیه در معادلات دیفرانسیل معمولی است. بحث جامع و تحلیل آن در ۲.۶ آمده است. در فصل ۶ بسط کامل روشهای عددی برای مسأله مقدار اولیه (۵.۳.۱) داده خواهد شد.

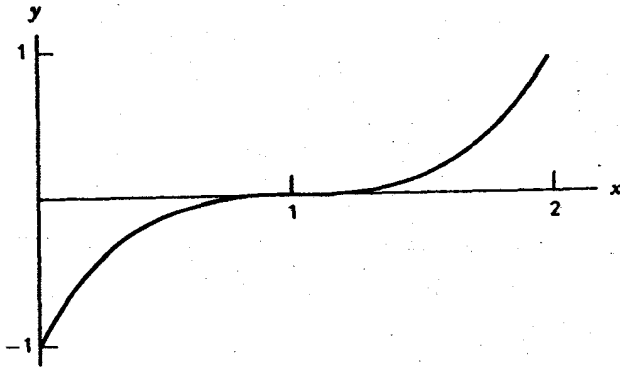
اغلب مسائل آنالیز عددی در فصلهای آینده عمدتاً شامل خطاهای برشی هستند. مورد استثنای عمده، حل دستگاههای معادلات خطی است که در آن خطاهای گردکردن منشأ اصلی خطا هستند.

نوفه در محاسبه تابع. یکی از نتایج بلافصل خطاهای گردکردن آن است که محاسبه تابع $f(x)$ با رایانه به یک تابع تقریبی $\hat{f}(x)$ که پیوسته نیست می‌انجامد. هر چند این امر تنها وقتی آشکار می‌شود که نمودار $\hat{f}(x)$ را روی مقیاسی به اندازه کافی کوچک بررسی کنیم. پس از هر عمل حساب که در محاسبه $f(x)$ به‌کار برده می‌شود معمولاً یک خطای گردکردن وجود دارد. وقتی نتیجه این خطاهای گردکردن در نظر گرفته شود، مقدار محاسبه شده $\hat{f}(x)$ به دست می‌آید که خطای آن، $f(x) - \hat{f}(x)$ ، وقتی x تغییر کند به صورت یک عدد کوچک تصادفی به نظر می‌رسد. این خطا در $\hat{f}(x)$ را نوفه نامند. هنگامی که نمودار $\hat{f}(x)$ روی مقیاسی به اندازه کافی کوچک بررسی شود، به صورت نوار تیره‌ای از نقاط دیده می‌شود که در آن مقادیر x تمام اعداد با ممیز شناور قابل قبول در ماشین را در بر می‌گیرد. این امر پیامدهایی برای بسیاری از برنامه‌های بعدی که در آنها از $\hat{f}(x)$ استفاده می‌کنند خواهد داشت. مثلاً محاسبه ریشه $f(x)$ با استفاده از $\hat{f}(x)$ به عدم قطعیت در محل ریشه منجر می‌شود، زیرا احتمالاً جای ریشه در محل تلاقی محور x ها و نمودار نوار تیره‌رنگ $\hat{f}(x)$ خواهد بود. مثال زیر نشان می‌دهد که این امر ممکن است موجب عدم قطعیت قابل ملاحظه‌ای در جای ریشه شود.

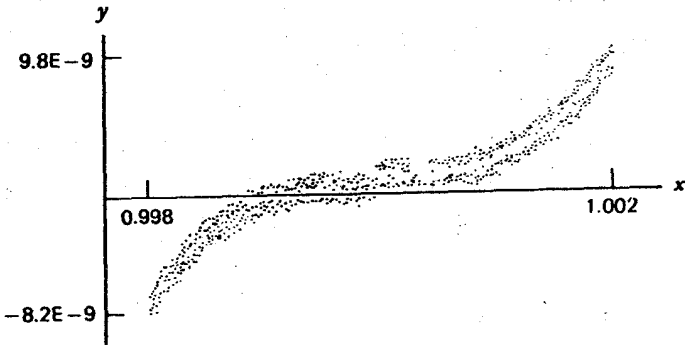
مثال گیریم

$$f(x) = x^3 - 3x^2 + 3x - 1 \quad (۶.۳.۱)$$

که همان $(x-1)^3$ است. ما (۶.۳.۱) را با BASIC با دقت معمولی در یک ریزرایانه معمولی، با استفاده از حساب گردشده دودویی و واحد گردکردن $10^{-8} \times 596 = 2^{-24} \doteq \delta$ محاسبه نموده‌ایم. نمودار $\hat{f}(x)$ روی بازه $[0, 2]$ در شکل ۱.۱ نشان داده شده است، و همان گونه که انتظار می‌رود یک منحنی پیوسته و هموار به نظر می‌آید. ولی نمودار روی یک بازه کوچکتر $[0.998, 1.002]$ ماهیت ناپیوسته $f(x)$ را، همان گونه که در شکل ۲.۱ دیده می‌شود، نشان می‌دهد. در این حالت اخیر، $\hat{f}(x)$ در 640 نقطه متساوی‌الفاصله برای مقادیر x در $[0.998, 1.002]$ محاسبه شده است که یک نوار تیره رنگ نمودار $f(x)$ را به وجود آورده است. از نمودار اخیر می‌توان دید که یک بازه بزرگ عدم قطعیت برای محل تقاطع $\hat{f}(x)$ با محور x ها وجود دارد. در بخش ۷.۲، فصل دوم، به این مبحث باز خواهیم گشت.



شکل ۱.۱ گراف (۶.۳.۱).



شکل ۲.۱ جزئیات گراف (۶.۳.۱).

خطاهای پی‌ریز/سرریز در محاسبات. پی‌آمد دیگری از خطاهای ماشین را در نظر می‌گیریم. حدود بالا و پایین اعداد با ممیز شناور، که در (۱۷.۲.۱) داده شده، ممکن است به خطاهایی در محاسبات بینجامد. گاهی این خطاها اجتناب‌ناپذیر، ولی بیشتر محصول طریقه محاسبات‌اند. برای نشان دادن این موضوع، محاسبه قدرمطلق یک عدد مختلط را در نظر بگیرید

$$|x + iy| = \sqrt{x^2 + y^2} \quad (۷.۳.۱)$$

این محاسبه ممکن است به پی‌ریز یا سرریز منجر شود، ولو اینکه قدرمطلق $|x + iy|$ در محدوده ماشین باشد. مثلاً اگر در رابطه (۱۷.۲.۱) $x = y = 10^{20}$ اختیار شود، آنگاه برای $x = y = 10^{20}$ رابطه (۷.۳.۱) سرریز می‌دهد ولو اینکه $|x + iy| \approx 1.4 \times 10^{20}$. برای احتراز از این مشکل

بزرگترین مقدار x یا y ، مثلاً x را معین می‌کنیم. رابطه (۷.۳.۱) را دوباره به شکل زیر می‌نویسیم

$$|x + iy| = |x| \cdot \sqrt{1 + \alpha^2}, \quad \alpha = \frac{y}{x} \quad (۸.۳.۱)$$

باید $\sqrt{1 + \alpha^2}$ را در حالت $0 \leq \alpha \leq 1$ ، محاسبه کنیم. این کار مشکل پی‌ریز و سرریز را برای مسائل (۷.۳.۱) از بین خواهد برد.

۴.۱ انتشار خطاها

در این بخش و بخشهای بعدی آثار محاسبات با اعدادی را که خطا دارند مد نظر قرار می‌دهیم. ابتدا اعمال اصلی حساب را در نظر می‌گیریم. فرض کنید ω یکی از اعمال جمع، تفریق، ضرب یا تقسیم حساب باشد؛ و گیریم $\hat{\omega}$ شکل رایانه‌یی همان عمل حساب، که معمولاً شامل خطای گردکردن است، باشد. گیریم x_A و y_A اعدادی باشند که برای محاسبه به‌کار رفته‌اند و فرض کنید که این اعداد دارای خطا باشند و مقدار درست آنها چنین باشد

$$x_T = x_A + \varepsilon \quad y_T = y_A + \eta$$

پس $x_A \hat{\omega} y_A$ عددی است که عملاً محاسبه شده است و برای خطای آن داریم

$$x_T \omega y_T - x_A \hat{\omega} y_A = [x_T \omega y_T - x_A \omega y_A] + [x_A \omega y_A - x_A \hat{\omega} y_A] \quad (۱.۴.۱)$$

کمیت کروشه اول را خطای منتشره گویند و کمیت کروشه دوم معمولاً خطای قطع‌کردن یا گردکردن است. برای این کمیت دوم، معمولاً داریم

$$x_A \hat{\omega} y_A = \text{fl}(x_A \omega y_A) \quad (۲.۴.۱)$$

که بدین معنی است که $x_A \omega y_A$ دقیقاً محاسبه و سپس گردشده است. از ترکیب (۱.۴.۱) و (۲.۴.۱)، داریم

$$|x_A \omega y_A - x_A \hat{\omega} y_A| \leq \frac{\beta}{4} |x_A \omega y_A| \beta^{-t} \quad (۳.۴.۱)$$

به شرطی که از گردکردن درست استفاده شده باشد.

برای خطای منتشره، حالت‌های خاص را بررسی می‌کنیم.

حالت الف. ضرب. برای خطا در $x_A y_A$

$$x_T y_T - x_A y_A = x_T y_T - (x_T - \varepsilon)(y_T - \eta) = x_T \eta + y_T \varepsilon - \varepsilon \eta$$

$$\begin{aligned} \text{Rel}(x_A y_A) &= \frac{x_T y_T - x_A y_A}{x_T y_T} = \frac{\eta}{y_T} + \frac{\varepsilon}{x_T} - \frac{\varepsilon}{x_T} \cdot \frac{\eta}{y_T} \\ &= \text{Rel}(x_A) + \text{Rel}(y_A) - \text{Rel}(x_A) \cdot \text{Rel}(y_A) \end{aligned}$$

برای $|\text{Rel}(x_A)| \ll 1$ و $|\text{Rel}(y_A)| \ll 1$

$$\text{Rel}(x_A y_A) \doteq \text{Rel}(x_A) + \text{Rel}(y_A) \quad (۴.۴.۱)$$

نماد « \ll » به معنای «بسیار کوچکتر است از» است.

حالت ب. تقسیم. با استدلالی مشابه

$$\text{Rel} \frac{x_A}{y_A} = \frac{\text{Rel}(x_A) - \text{Rel}(y_A)}{1 - \text{Rel}(y_A)} \quad (۵.۴.۱)$$

برای $|\text{Rel}(y_A)| \ll 1$:

$$\text{Rel} \left(\frac{x_A}{y_A} \right) \doteq \text{Rel}(x_A) - \text{Rel}(y_A) \quad (۶.۴.۱)$$

برای ضرب و تقسیم، خطاهای نسبی به سرعت افزایش نمی‌یابند.

حالت ج. جمع و تفریق

$$(x_T \pm y_T) - (x_A \pm y_A) = (x_T - x_A) \pm (y_T - y_A) = \varepsilon \pm \eta$$

$$\text{Error}(x_A \pm y_A) = \text{Error}(x_A) \pm \text{Error}(y_A) \quad (۷.۴.۱)$$

این امر کاملاً خوب و منطقی به نظر می‌رسد، ولی ممکن است گمراه کننده باشد. خطای نسبی $x_A \pm y_A$ در مقایسه با $\text{Rel}(x_A)$ و $\text{Rel}(y_A)$ ممکن است خیلی ناچیز باشد.

مثال گیریم $x_T = \pi$ و $x_A = ۳,۱۴۱۶$ و $y_T = \frac{22}{7}$ و $y_A = ۳,۱۴۲۹$ پس

$$x_T - x_A \doteq -۷,۳۵ \times 10^{-۶} \quad \text{Rel}(x_A) \doteq -۲,۳۴ \times 10^{-۶}$$

$$y_T - y_A \doteq -۴,۲۹ \times 10^{-۵} \quad \text{Rel}(y_A) \doteq -۱,۳۶ \times 10^{-۵}$$

$$(x_T - y_T) - (x_A - y_A) \doteq -۰,۰۰۱۲۶۴۵ - (-۰,۰۰۱۳) \doteq ۳,۵۵ \times 10^{-۵}$$

$$\text{Rel}(x_A - y_A) \doteq -۰,۰۲۸$$

اگرچه خطای $x_A - y_A$ کاملاً کوچک است ولی خطای نسبی $x_A - y_A$ بسیار بزرگتر از خطای نسبی x_A یا خطای نسبی y_A است.

خطاهای کاهش ارقام بامعنی. مثال اخیر نشان می‌دهد که از لحاظ خطای نسبی، از تفریق دو کمیت نزدیک به هم، ممکن است کاهش زیادی در دقت پیش آید. این عمل ممکن است از راههای مهم از دست دادن دقت هنگام انتشار خطا در یک محاسبه باشد.

اکنون چند مثال برای این پدیده می‌آوریم، توصیه‌هایی را برای چگونگی جلوگیری از این مشکلات در بعضی موارد، پیشنهاد می‌کنیم.

مثال حل معادله $ax^2 + bx + c = 0$ را وقتی $4ac$ در مقایسه با b^2 مقدار کوچکی است در نظر گرفته و از فرمولی معمول برای پیدا کردن ریشه‌ها استفاده می‌کنیم

$$r_T^{(1)} = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \quad r_T^{(2)} = \frac{-b - \sqrt{b^2 - 4ac}}{2a} \quad (۸.۴.۱)$$

برای روشن کردن مطلب، معادله $x^2 - ۲۶x + ۱ = 0$ را در نظر می‌گیریم. از فرمولهای (۸.۴.۱) نتیجه می‌گیریم

$$r_T^{(1)} = ۱۳ + \sqrt{۱۶۸} \quad r_T^{(2)} = ۱۳ - \sqrt{۱۶۸} \quad (۹.۴.۱)$$

اکنون فرض کنید از یک ماشین با پنج رقم ده‌دهی استفاده می‌کنیم. در این ماشین $\sqrt{۱۶۸} \doteq ۱۲٫۹۶۱$ بنابراین

$$r_A^{(1)} = ۱۳ + ۱۲٫۹۶۱ = ۲۵٫۹۶۱ \quad r_A^{(2)} = ۱۳ - ۱۲٫۹۶۱ = ۰٫۰۳۹ \quad (۱۰.۴.۱)$$

اگر جوابهای دقیق را بکار ببریم

$$\text{Rel}(r_A^{(1)}) \doteq ۱٫۸۵ \times ۱۰^{-۵} \quad \text{Rel}(r_A^{(2)}) \doteq ۱٫۲۵ \times ۱۰^{-۲} \quad (۱۱.۴.۱)$$

برای داده‌هایی که در محاسبه (۱۰.۴.۱) آمده‌اند، با استفاده از نمادگذاری (۷.۴.۱) داریم

$$x_T = x_A = ۱۳ \quad y_T = \sqrt{۱۶۸} \quad y_A = ۱۲٫۹۶۱$$

$$\text{Rel}(x_A) = 0 \quad \text{Rel}(y_A) \doteq ۳٫۷۱ \times ۱۰^{-۵}$$

دقت در $r_A^{(2)}$ بسیار کمتر از دقت در داده‌های x_A و y_A است که در محاسبه وارد شده‌اند. گوئیم ارقام با معنی در تفاضل $r_A^{(2)} = x_A - y_A$ از بین رفته‌اند، یا در محاسبه $r_A^{(2)}$ دچار خطای کاهش ارقام بامعنی شده‌ایم. در $r_A^{(1)}$ دقتی با پنج رقم با معنی داریم، حال آنکه در $r_A^{(2)}$ فقط دو رقم با معنی داریم.

برای رفع این نقیصه در این حالت خاص محاسبه $r_A^{(2)}$ ، (۹.۴.۱) را به شکل زیر می‌نویسیم

$$r_T^{(2)} = \frac{13 - \sqrt{168}}{1} \cdot \frac{13 + \sqrt{168}}{13 + \sqrt{168}} = \frac{1}{13 + \sqrt{168}}$$

سپس

$$\frac{1}{13 + \sqrt{168}} \doteq \frac{1}{25,961} \doteq 0.38519 \equiv r_A^{(2)} \quad (12.4.1)$$

در اینجا دو خطا وجود دارد، یکی مربوط به $\sqrt{168} \doteq 12,961$ و دیگری مربوط به تقسیم نهایی. ولی هر یک از آنها خطاهای نسبی کوچکی دارند [(۶.۴.۱) را ببینید] و مقدار جدید $r_A^{(2)}$ از مقدار قبلی بسیار دقیقتر است. با محاسبات دقیق، اکنون داریم

$$\text{Rel}(r_A^{(2)}) \doteq -1,0^3 \times 10^{-5},$$

که از (۱۱.۴.۱) خیلی بهتر است.

این محاسبه جدید $r_A^{(2)}$ نشان می‌دهد که خطای کاهش ارقام با معنی به شکل محاسبه بستگی دارد، نه به خطاهای موجود در داده‌های محاسبه. در این مثال به آسانی توانستیم راه دیگری برای محاسبات پیدا کنیم که خطای کاهش ارقام با معنی را مرتفع سازد، ولی این کار همیشه ممکن نخواهد بود. برای یک بحث کامل از روش عملی محاسبه ریشه‌های یک چند جمله‌یی درجه ۲، فورسایت (۱۹۶۹) را ببینید.

مثال در بسیاری از محاسبات کاهش ارقام بامعنی، تقریبهای چندجمله‌یی تیلر می‌توانند مشکل را حل کنند. ما این مطلب را با محاسبه

$$f(x) = \int_0^x e^{xt} dt = \frac{e^x - 1}{x} \quad x \neq 0 \quad (13.4.1)$$

نشان می‌دهیم. برای $x = 0$ ، $f(0) = 1$ ؛ و به سادگی می‌توان دید که $f(x)$ در $x = 0$ پیوسته است.

برای دیدن مسأله کاهش ارقام بامعنی، وقتی x کوچک است، $f(x)$ را با یک ماشین حساب دستی مناسب ده رقمی در $x = ۱۴ \times ۱۰^{-۹}$ محاسبه می‌کنیم. نتایج چنین‌اند

$$e^x \doteq ۱.۰۰۰۰۰۰۰۰۰۱$$

$$\frac{e^x - 1}{x} \doteq \frac{۱۰^{-۹}}{۱۴ \times ۱۰^{-۹}} \doteq ۰.۷۱۴ \quad (۱۴.۴.۱)$$

سمت راست تساویها، نتایج ماشین حساب را می‌دهد و جواب درست که با ده رقم گرد شده عبارت است از

$$f(x) = ۱.۰۰۰۰۰۰۰۰۰۱$$

در محاسبه (۱۴.۴.۱) در صورت کسر نه رقم اول ارقام بامعنی حذف شده است. برای اجتناب از خطای کاهش ارقام بامعنی از تقریب درجه دوم تیلر برای e^x استفاده و سپس $f(x)$ را ساده می‌کنیم

$$e^x = 1 + x + \frac{x^2}{2} + \frac{x^3}{6} e^\xi \quad 0 \leq \xi \leq x \leq 1$$

$$f(x) = 1 + \frac{x}{2} + \frac{x^2}{6} e^\xi \quad (۱۵.۴.۱)$$

بازای مقدار قبلی $x = ۱۴ \times ۱۰^{-۹}$ ، مقدار $f(x)$ با خطایی کمتر از $۱۰^{-۱۸}$ چنین خواهد شد

$$f(x) \doteq 1 + ۷ \times ۱۰^{-۱۰}$$

در حالت کلی (۱۵.۴.۱) را برای بازه‌ای چون $[\delta, 0]$ به‌کار می‌بریم و δ را طوری انتخاب می‌کنیم که مطمئن باشیم خطای

$$f(x) \doteq 1 + \frac{x}{2}$$

به اندازه کافی کوچک است. البته، می‌توانستیم از تقریب با درجه بالاتر برای e^x استفاده کنیم و مقدار δ را بزرگتر بگیریم.

به طور کلی، تقریبهای تیلر برای اجتناب از خطای کاهش ارقام بامعنی در محاسبات مفیدند. ولی در بعضی موارد خطای کاهش ارقام بامعنی مستلزم دقت بیشتر است.

مثال محاسبهٔ مجموعی به صورت

$$S = \sum_{j=1}^n x_j \quad (16.4.1)$$

را در نظر می‌گیریم با جمله‌های مثبت و منفی x_j ، که هر یک مقداری است تقریبی. بعلاوه فرض می‌کنیم که مجموع S خیلی کوچکتر از ماکسیم قدرمطلق x_j ها باشد. در محاسبهٔ چنین مجموعی با رایانه، احتمالاً دچار خطای کاهش ارقام بامعنی می‌شویم. ما این مطلب را با یک مثال نشان می‌دهیم.

با استفاده از فرمول تیلر (۴.۱.۱) برای e^x ، مقدار e^{-5} را محاسبه می‌کنیم:

$$e^{-5} \doteq 1 + \frac{(-5)}{1!} + \frac{(-5)^2}{2!} + \frac{(-5)^3}{3!} + \dots + \frac{(-5)^n}{n!} \quad (17.4.1)$$

تصور کنید که از رایانه‌ای با حساب ممیز شناور که با چهار رقم اعشار گرد می‌کند استفاده می‌کنیم، پس هر جمله از سری بالا را باید با چهار رقم با معنی گرد کرد. در جدول ۲.۱، مقادیر گردشدهٔ x_j همراه با مجموع دقیق جملات تا درجهٔ داده شده آمده است. مقدار درست e^{-5} تا چهار رقم بامعنی

جدول ۲.۱ محاسبهٔ (۱۷.۴.۱) با استفاده از حساب چهاررقم اعشاری

درجه	جمله	حاصل جمع	درجه	جمله	حاصل جمع
۰	۱٫۰۰	۱٫۰۰۰	۱۳	-۰٫۱۹۶۰	-۰٫۰۴۲۳۰
۱	-۵٫۰۰۰	-۴٫۰۰۰	۱۴	۰٫۷۰۰۱E-۱	۰٫۰۲۷۷۱
۲	۱۲٫۵۰	۸٫۵۰۰	۱۵	-۲۳۳۴E-۱	۰٫۰۰۴۳۷۰
۳	-۲۰٫۸۳	-۱۲٫۳۳	۱۶	۰٫۷۲۹۳E-۲	۰٫۱۱۶۶
۴	-۲۶٫۰۴	۱۳٫۷۱	۱۷	-۰٫۲۱۴۵E-۲	۰٫۰۰۹۵۱۸
۵	۲۶٫۰۴	-۱۲٫۳۳	۱۸	۰٫۵۹۵۸E-۳	۰٫۰۱۰۱۱
۶	۲۱٫۷۰	۹٫۳۷۰	۱۹	-۰٫۱۵۶۸E-۳	۰٫۰۰۹۹۵۷
۷	-۱۵٫۵۰	-۶٫۱۳۰	۲۰	۰٫۳۹۲۰E-۴	۰٫۰۰۹۹۹۶
۸	۹٫۶۸۸	۳٫۵۵۸	۲۱	-۰٫۹۳۳۳E-۵	۰٫۰۰۹۹۸۷
۹	-۵٫۳۸۲	-۱٫۸۲۴	۲۲	۰٫۲۱۲۱E-۵	۰٫۰۰۹۹۸۹
۱۰	۲٫۶۹۱	۰٫۸۶۷۰	۲۳	-۰٫۴۶۱۱E-۶	۰٫۰۰۹۹۸۹
۱۱	-۱٫۲۲۳	-۰٫۳۵۶۰	۲۴	۰٫۹۶۰۷E-۷	۰٫۰۰۹۹۸۹
۱۲	۰٫۵۰۹۷	۰٫۱۵۳۷	۲۵	-۰٫۱۹۲E-۷	۰٫۰۰۹۹۸۹

برابر 0.06738×10^6 است و با آخرین مجموع در جدول کاملاً متفاوت است. همچنین اگر (۱۷.۴.۱) را دقیقاً برای $n = 25$ محاسبه کنیم مقدار صحیح e^{-5} تا چهار رقم بامعنی به دست می‌آید. در این مثال، جمله‌های x_j نسبتاً بزرگ می‌شوند، ولی مجموع آنها عدد بسیار کوچکتر e^{-5} را تشکیل می‌دهد. این امر به معنای وجود خطاهای کاهش ارقام بامعنی در مجموعیابی است. احتراز از این مشکل در این حالت نسبتاً آسان است. یا می‌نویسیم

$$e^{-5} = \frac{1}{e^5}$$

e^5 را با یک سری که حذف جملات مثبت و منفی در آن وجود ندارد محاسبه می‌کنیم؛ یا $\frac{1}{e^5} = e^{-1}$ را مثلاً با استفاده از یک سری حساب می‌کنیم، و آن را در خودش چند بار ضرب می‌کنیم تا e^{-5} حاصل شود. در مورد سریهای دیگر، ممکن است که این‌گونه راه حل ساده‌ای وجود نداشته باشد.

انتشار خطا در محاسبه توابع. گیریم تابع $f(x)$ داده شده و $\hat{f}(x)$ معرف مقدار محاسبه‌شده $f(x)$ در رایانه باشد. اگر $f(x_T)$ معرف مقدار تابع مطلوب و $\hat{f}(x_A)$ مقداری است که در واقع حساب شده باشد. برای خطا، می‌نویسیم

$$f(x_T) - \hat{f}(x_A) = [f(x_T) - f(x_A)] + [f(x_A) - \hat{f}(x_A)] \quad (18.4.1)$$

کمیت اولین گروه را خطای منتشره می‌نامند، و کمیت دومین گروه، خطای محاسبه $f(x_A)$ در رایانه است. این خطای دوم معمولاً یک عدد تصادفی کوچک است که از خطاهای گردکردن در اعمال حساب تعریف‌کننده $f(x)$ به وجود می‌آید. ما قبلاً در بخش ۳.۱ به آن به عنوان نوبه در محاسبه $f(x)$ اشاره کرده‌ایم.

برای خطای منتشره، از قضیه مقدار میانگین داریم

$$f(x_T) - f(x_A) \doteq f'(x_T)(x_T - x_A) \quad (19.4.1)$$

در اینجا فرض می‌شود که x_A و x_T نسبتاً به یکدیگر نزدیک‌اند و $f'(x)$ وقتی x بین x_A و x_T است خیلی تغییر نمی‌کند.

مثال $0.00026 = \sin(\frac{\pi}{5}) - \sin(0.628) \doteq \cos(\frac{\pi}{5})[(\frac{\pi}{5}) - 0.628]$ که یک برآورد عالی برای خطاست.

با استفاده از قضیه تیلر (۱۲.۱.۱) برای توابع دو متغیره، می‌توانیم بحث اخیر را برای انتشار خطا در توابع دو متغیره تعمیم دهیم.

$$f(x_T, y_T) - f(x_A, y_A) \doteq f_x(x_T, y_T)(x_T - x_A) + f_y(x_T, y_T)(y_T - y_A) \quad (20.4.1)$$

که در آن $f_x \equiv \partial f / \partial x$ در اینجا فرض می‌کنیم که $f_x(x, y)$ و $f_y(x, y)$ برای (x, y) بین (x_T, y_T) و (x_A, y_A) خیلی تغییر نمی‌کنند.

مثال برای $f(x, y) = x^y$ داریم $f_x = yx^{y-1}$ ، $f_y = x^y \log(x)$ پس از (۲۰.۴.۱) به دست می‌آوریم

$$x_T^{y_T} - x_A^{y_A} \doteq \varepsilon \cdot y_T x_T^{y_T-1} + \eta [\log(x_T)] x_T^{y_T} \\ \text{Rel}(x_A^{y_A}) \doteq y_T [\text{Rel}(x_A) + \text{Rel}(y_A) \log(x_T)] \quad (21.4.1)$$

خطای نسبی در $x_A^{y_A}$ ممکن است بزرگ باشد، اگرچه $\text{Rel}(x_A)$ و $\text{Rel}(y_A)$ کوچک‌اند. به‌عنوان یک مثال دیگرگیریم $x_A = 1.2 \times 10^1$ ، $x_T = 1.2$ ، $y_T = y_A = 500$ در این صورت $x_T^{y_T} = 3.89604 \times 10^{39}$ و $x_A^{y_A} = 4.06179 \times 10^{39}$ و $\text{Rel}(x_A^{y_A}) = -0.0425$. این مقدار را با $\text{Rel}(x_A) = 8.3 \times 10^{-5}$ مقایسه کنید.

خطای داده‌ها. اگر داده‌های ورودی یک الگوریتم فقط r رقم دقت داشته باشند، گاهی توصیه می‌شود که در چنین حالت محاسباتی که در آنها این داده‌ها دخالت دارند، فقط حساب r رقمی را به‌کار برند. این کار معقولی نیست. البته قبول داریم که دقت محدود داده‌ها بر نتایج نهایی محاسبات الگوریتمی اثر می‌گذارد و جوابهایی به دست می‌دهد که دارای خطا هستند، مع‌هذا دلیلی وجود ندارد که وضعیت را با خطاهای گردکردنی که در اثر محاسبات با r رقم ناشی می‌شوند، بدتر کنیم. به‌عکس باید محاسبات را با دقت بیشتری انجام دهیم تا از تزلزل بیشتر در دقت نتایج الگوریتم جلوگیری شود. این امر موجب می‌شود که خطاهای گردکردن حسابی به مراتب ناچیزتر از خطای داده‌ها شوند و کمک می‌کند که دقت در حد دقت داده‌ها حفظ شود.

۵.۱ خطاهای مجموعیابی

در بسیاری از روشهای عددی، به‌ویژه در جبر خطی، مجموعیابی وجود دارد. در این بخش، وجوه متفاوتی از مجموعیابی را مد نظر قرار می‌دهیم، بویژه وقتی در حساب با ممیز شناور انجام گرفته باشند.

محاسبهٔ مجموع زیر را در نظر می‌گیریم

$$S = \sum_{j=1}^m x_j \quad (۱.۵.۱)$$

که x_1, \dots, x_m اعداد با ممیز شناورند. تعریف می‌کنیم

$$S_2 = fl(x_1 + x_2) = (x_1 + x_2)(1 + \varepsilon_2) \quad (۲.۵.۱)$$

که در آن از (۲.۴.۱) و (۱۰.۲.۱) استفاده کرده‌ایم. به‌طور بازگشتی

$$S_{r+1} = fl(S_r + x_{r+1}) \quad r = 2, \dots, m-1$$

تعریف می‌شود. پس

$$S_{r+1} = (S_r + x_{r+1})(1 + \varepsilon_{r+1}) \quad (۳.۵.۱)$$

کمیت‌های $\varepsilon_2, \dots, \varepsilon_m$ برحسب آنکه از قطع‌کردن یا گردکردن استفاده شده باشد، در (۸.۲.۱) یا (۹.۲.۱) صدق می‌کنند.

با بسط چند مجموع اولیه، به‌دست می‌آوریم:

$$S_2 - (x_1 + x_2) = \varepsilon_2(x_1 + x_2)$$

$$\begin{aligned} S_3 - (x_1 + x_2 + x_3) &= (x_1 + x_2)\varepsilon_2 + (x_1 + x_2)(1 + \varepsilon_2)\varepsilon_3 + x_3\varepsilon_3 \\ &\doteq (x_1 + x_2)\varepsilon_2 + (x_1 + x_2 + x_3)\varepsilon_3 \end{aligned}$$

$$\begin{aligned} S_4 - (x_1 + x_2 + x_3 + x_4) &\doteq (x_1 + x_2)\varepsilon_2 + (x_1 + x_2 + x_3)\varepsilon_3 \\ &\quad + (x_1 + x_2 + x_3 + x_4)\varepsilon_4 \end{aligned}$$

ما جملات حاصلضرب $\varepsilon_i x_j$ را نادیده گرفته‌ایم زیرا اندازه‌های بسیار کوچکتري دارند. به روش استقرای به‌دست می‌آوریم

$$\begin{aligned} S_m - \sum_1^m x_i &\doteq (x_1 + x_2)\varepsilon_2 + \dots + (x_1 + x_2 + \dots + x_m)\varepsilon_m \\ &= x_1(\varepsilon_2 + \varepsilon_3 + \dots + \varepsilon_m) + x_2(\varepsilon_2 + \varepsilon_3 + \dots + \varepsilon_m) \\ &\quad + x_3(\varepsilon_2 + \dots + \varepsilon_m) + \dots + x_m\varepsilon_m \end{aligned} \quad (۴.۵.۱)$$

از این فرمول نتیجه می‌گیریم که بهترین نحوه جمع‌کردن، اضافه‌کردن از کوچکترین عنصر به بزرگترین عنصر است. البته مثالهای نقض می‌توان ساخت، ولی برای تعداد زیادی مجموعیابی، قاعده فوق بایستی بهترین نحوه باشد. این، به‌ویژه وقتی اعداد x_i همگی دارای یک علامت باشند، درست است، زیرا در محاسبه مجموعهای میانی $x_1 + x_2 + \dots + x_m$ ، $m = 1, \dots, n$ ، هیچ‌گونه حذفی صورت نمی‌گیرد. در این حالت، اگر به جای گردکردن، قطع‌کردن به‌کار برده شود و اگر تمام x_i ها مثبت باشند آنگاه هیچ‌گونه حذفی در مجموعهای ε_i اتفاق نمی‌افتد. با نحوه جمع‌کردن از کوچکترین‌ها تا بزرگترین‌ها، اثر این خطاهای قطع‌کردن را به کمترین مقدار می‌رسانیم.

مثال جمله‌های x_j از مجموع S را به شکل زیر تعریف می‌کنیم: کسر $1/z$ را به کسر اعشاری برگردانده و آن را تا چهاررقم بامعنی گرد می‌کنیم و آن را با x_j نشان می‌دهیم. برای روش‌ترکردن خطای محاسبه S ، حساب اعشاری چهاررقمی یا ممیز شناور را به‌کار می‌بریم. جداول ۳.۱ و ۴.۱ نتایج محاسبه S را به چهار طریق متفاوت نشان می‌دهند. جمع‌کردن S از بزرگتر تا کوچکتر با LS و جمع‌کردن از کوچکتر تا بزرگتر با SL نشان داده شده است. در جدول ۳.۱ از حساب قطع شده با

$$-0.001 \leq \varepsilon_j \leq 0 \quad (5.5.1)$$

و در جدول ۴.۱ از حساب گرد شده با

$$-0.0005 \leq \varepsilon_j \leq 0.0005 \quad (6.5.1)$$

استفاده شده است. اعداد ε_j به (۴.۵.۱) اشاره دارند و کران آنها از (۸.۲.۱) و (۹.۲.۱) به‌دست می‌آیند.

جدول ۳.۱ محاسبه S در یک ماشین با استفاده از قطع کردن

n	مقدار واقعی	SL	خطا	LS	خطا
۱۰	۲,۹۲۹	۲,۹۲۸	۰.۰۰۱	۲,۹۲۷	۰.۰۰۲
۲۵	۳,۸۱۶	۳,۸۱۳	۰.۰۰۳	۳,۸۰۶	۰.۰۱۰
۵۰	۴,۴۹۹	۴,۴۹۱	۰.۰۰۸	۴,۴۷۹	۰.۰۲۰
۱۰۰	۵,۱۸۷	۵,۱۷۰	۰.۰۱۷	۵,۱۴۲	۰.۰۴۵
۲۰۰	۵,۸۷۸	۵,۸۴۱	۰.۰۳۷	۵,۷۸۶	۰.۰۹۲
۵۰۰	۶,۷۹۳	۶,۶۹۲	۰.۱۰۱	۶,۵۶۹	۰.۲۲۴
۱۰۰۰	۷,۴۸۶	۷,۲۸۴	۰.۲۰۲	۷,۰۶۹	۰.۴۱۷

جدول ۴.۱ محاسبه S در یک ماشین با استفاده از گردکردن

n	مقدار واقعی	SL	خطا	LS	خطا
۱۰	۲,۹۲۹	۲,۹۲۹	۰٫۰	۲,۹۲۹	۰٫۰
۲۵	۳,۸۱۶	۳,۸۱۶	۰٫۰	۳,۸۱۷	-۰٫۰۰۱
۵۰	۴,۴۹۹	۴,۵۰۰	-۰٫۰۰۱	۴,۴۹۸	۰٫۰۰۱
۱۰۰	۵,۱۸۷	۵,۱۸۷	۰٫۰	۵,۱۸۷	۰٫۰
۲۰۰	۵,۸۷۸	۵,۸۷۸	۰٫۰	۵,۸۷۶	۰٫۰۰۲
۵۰۰	۶,۷۹۳	۶,۷۹۴	-۰٫۰۰۱	۶,۷۸۳	۰٫۰۱۰
۱۰۰۰	۷,۴۸۶	۷,۴۸۶	۰٫۰	۷,۴۴۹	۰٫۰۳۷

در هر دو جدول فوق‌الذکر روشن است که استراتژی جمع‌کردن S از کوچکترین عنصر به بزرگترین بر جمع کردن S از بزرگترین عنصر به کوچکترین ارجحیت دارد. معهذاً آنچه بیشتر اهمیت دارد خطای بسیار کوچکتر گردکردن نسبت به قطع کردن است. تفاوت خیلی بیشتر از ۲ برابر است که از اندازه نسبی کرانها از (۵.۵.۱) و (۶.۵.۱) به دست می‌آید. ما بعداً تحلیلی از این موضوع می‌آوریم.

یک تحلیل آماری از انتشار خطا مجموع معمولی خطای

$$E = \sum_{j=1}^n \varepsilon_j \quad (۷.۵.۱)$$

از نوعی را که در مجموعیابی خطای (۴.۵.۱) پیش می‌آید در نظر می‌گیریم. یک کران ساده برای آن چنین است:

$$|E| \leq n\delta \quad (۸.۵.۱)$$

که δ کران $\varepsilon_1, \dots, \varepsilon_n$ است. پس، در مثال قبلی، برحسب آنکه قطع کردن یا گردکردن به کار برده شده باشد $\delta = ۰٫۰۰۰۵$ یا $\delta = ۰٫۰۰۰۱$ است. کران (۸.۵.۱) مربوط به بدترین حالت است که تمام خطاهای ε_j در حد ممکن بزرگ و همه دارای یک علامت‌اند.

هنگامی که گردکردن را به کار می‌بریم، رفتار متقارن علامت ε_j ، همان‌گونه که در (۹.۲.۱) نشان داده شد، تفاوت فاحشی در اندازه E ایجاد خواهد کرد. در چنین حالتی یک مدل بهتر این است که فرض کنیم خطاهای ε_j متغیرهای تصادفی هستند که در بازه $[-\delta, \delta]$ به شکل یکنواخت

توزیع شده و مستقل اند. پس

$$E = n \left(\frac{1}{n} \sum_1^n \varepsilon_j \right) = n\bar{\varepsilon}$$

میانگین نمونه‌ی $\bar{\varepsilon}$ یک متغیر تصادفی جدید است، که دارای توزیع احتمال با میانگین 0 و واریانس $\delta^2/3n$ است. برای محاسبه احتمال عبارتی که شامل $\bar{\varepsilon}$ هستند، لازم است توجه داشته باشیم که توزیع احتمال $\bar{\varepsilon}$ با توزیع نرمال با همان میانگین و واریانس، به خوبی تقریب زده می‌شود حتی برای مقادیر کوچکی چون $n \geq 10$. این امر از قضیه حد مرکزی در نظریه احتمال به دست می‌آید [مثلاً هاگ و کریگ (۱۹۷۸، فصل ۵) را ببینید]، با استفاده از توزیع نرمال تقریبی، احتمال برقراری

$$|\bar{\varepsilon}| \leq 0.398/\sqrt{n} \quad |E| \leq 0.398\sqrt{n}$$

برابر $1/2$ و احتمال برقراری

$$|\bar{\varepsilon}| \leq 1.498/\sqrt{n} \quad |E| \leq 1.498\sqrt{n} \quad (9.5.1)$$

برابر 0.99 است. نتیجه (۹.۵.۱) در مقایسه با (۸.۵.۱)، وقتی n تا حدی بزرگ باشد بهبود قابل توجهی است.

این تحلیل را می‌توان برای خطای قطع کردن نیز به کار برد. ولی در آن حالت، $-\delta \leq \varepsilon_j \leq 0$. میانگین نمونه‌ی $\bar{\varepsilon}$ اکنون میانگین $\delta/2$ دارد در حالی که واریانس همان $\delta^2/3n$ باقی می‌ماند. بنابراین احتمالی برابر 0.99 برای برقراری

$$\left(\frac{n}{2} - 1.498\sqrt{n} \right) \delta \leq -E \leq \left(\frac{n}{2} + 1.498\sqrt{n} \right) \delta \quad (10.5.1)$$

وجود دارد. برای مقدار بزرگ n ، این رابطه به ما اطمینان می‌دهد که E تقریباً برابر $n\delta/2$ است، که بسیار بزرگتر از (۹.۵.۱)، برای حالت گرد کردن، است.

هنگامی که این نتایج، یعنی (۹.۵.۱) و (۱۰.۵.۱) برای مجموعیایی معمولی (۴.۵.۱)، به کار برده شوند، علت احتمالی رفتار بسیار متفاوت خطای قطع کردن و گرد کردن در جدولهای ۳.۱ و ۴.۱ ملاحظه خواهد شد. به طور کلی حساب گرد کردن تقریباً همیشه بر حساب قطع کردن ارجحیت دارد.

گرچه تحلیل‌های آماری کرانه‌های واقع‌تری را به دست می‌دهند، ولی محاسبه آنها بسیار مشکل‌تر است. برای یک مثال پیچیده‌تر، صفحات ۵۹.۴۱ هنریچی^۱ (۱۹۶۲) را برای یک تحلیل آماری

خطای جواب عددی معادلات دیفرانسیل ببینید. یک مثال در جدول ۳.۶، فصل ۶ این کتاب، داده شده است.

حاصلضربهای داخلی. برای دو بردار $x, y \in R^m$.

$$x^T y = \sum_{j=1}^m x_j y_j \quad (11.5.1)$$

را حاصلضرب داخلی x و y می‌نامیم. (نماد x^T ترانهادهٔ ماتریس x را نشان می‌دهد). ویژگیهای حاصلضرب داخلی در فصل ۷ بررسی شده‌اند، ولی ما در اینجا اشاره می‌کنیم که

$$\|x\|_2 = \sqrt{x^T x} = \sqrt{\sum_{j=1}^m x_j^2} \quad (12.5.1)$$

$$\|x^T y\| \leq \|x\|_2 \|y\|_2 \quad (13.5.1)$$

نابرابری اخیر نابرابری کوشی - شوارتس خوانده می‌شود و حالت کلیتر آن در فصل ۴ اثبات شده است. مجموعهایی به شکل (۱۱.۵.۱) به کرات در مسائل جبر خطی پیش می‌آیند (برای مثال، ضرب ماتریسها). ما اکنون محاسبات عددی این گونه مجموعها را مورد بررسی قرار می‌دهیم.

فرض می‌کنیم x_i و y_i ، $i = 1, \dots, m$ ، اعدادی با ممیز شناورند. تعریف می‌کنیم

$$S_1 = \text{fl}(x_1 y_1)$$

$$S_{k+1} = \text{fl}(S_k + \text{fl}(x_k y_k)) \quad k = 1, 2, \dots, m-1 \quad (14.5.1)$$

پس مانند قبل، با استفاده از (۱۰.۲.۱)

$$S_1 = x_1 y_1 (1 + \varepsilon_1)$$

$$S_2 = [S_1 + x_2 y_2 (1 + \varepsilon_2)] (1 + \eta_2)$$

⋮

$$S_m = [S_{m-1} + x_m y_m (1 + \varepsilon_m)] (1 + \eta_m)$$

که عبارات ε_j و η_j ، برحسب آنکه از قطع کردن یا گرد کردن استفاده شده باشد به ترتیب در (۸.۲.۱)

و (۹.۲.۱) صدق می‌کنند. از ترکیب و مرتب‌کردن مجدد فرمولهای بالا، به دست می‌آوریم

$$S_m = \sum_{j=1}^m x_j y_j (1 + \gamma_j) \quad (15.5.1)$$

با

$$\begin{aligned} 1 + \gamma_j &= (1 + \varepsilon_j)(1 + \eta_j)(1 + \eta_{j+1}) \dots (1 + \eta_m) \quad \eta_1 = 0 \\ &\doteq 1 + \varepsilon_j + \eta_j + \eta_{j+1} + \dots + \eta_m \end{aligned} \quad (16.5.1)$$

این تقریب آخری بر این اساس است که از حاصلضربهای عبارات کوچک $\varepsilon_i \eta_k$ و $\eta_i \eta_k$ صرف‌نظر شده است. این امر ما را به تحلیل مشابهی که قبلاً برای مجموع (۱.۵.۱) انجام شد می‌رساند. تحلیل آماری خطا که بعد از (۷.۵.۱) آمده نیز معتبر است. برای یک کران مطمئن، می‌توان نشان داد که اگر $\delta < 10^{-m}$ آنگاه

$$|\gamma_j| \leq 10^{-1}(m+1-j)\delta \quad j = 1, \dots, m \quad (17.5.1)$$

که δ واحد گردکردن است که در (۱۲.۲.۱) داده شد [فورسایت و مولر^۱ (۱۹۶۷)، ص ۹۲ را ببینید] با به‌کارگرفتن این رابطه در (۱۵.۵.۱) و استفاده از (۱۳.۵.۱)،

$$\begin{aligned} |S - S_m| &\leq \sum_{j=1}^m |x_j y_j \gamma_j| \\ &\leq 10^{-1} m \cdot \delta \|x\|_2 \|y\|_2 \end{aligned} \quad (18.5.1)$$

این عبارت هیچ چیز در مورد خطای نسبی بیان نمی‌کند، زیرا $x^T y$ می‌تواند صفر باشد ولو اینکه تمام x_i ها و y_i ها مخالف صفر باشند.

این نتایج می‌گویند که خطای مطلق در $S_m \doteq x^T y$ خیلی به سرعت افزایش نمی‌یابد، به‌ویژه اگر از گردکردن واقعی استفاده شده باشد و تحلیل آماری قبلی (۷.۵.۱) را در نظر بگیریم. به‌رحال اغلب ممکن است که این خطا رابه سادگی و بدون صرف وقت زیاد به مقدار زیادی تقلیل داد، و این معمولاً در مسائل جبر خطی بسیار مهم است.

هر حاصلضرب $x_j y_j$ را با حسابی با دقت بالاتر محاسبه می‌کنیم و مجموعیابی را با همین حساب با دقت بالا انجام می‌دهیم. وقتی که محاسبهٔ مجموع کامل شد، نتیجه را با حساب با دقت معمولی گرد یا قطع می‌کنیم. برای مثال، وقتی، x_i و y_i با دقت معمولی حساب شده باشند، حاصلضربها و مجموعها را با دقت مضاعف محاسبه می‌کنیم. [در بیشتر رایانه‌ها دقت معمولی

مضاعف از نظر مدت زمان اجرا نسبتاً بهم نزدیک‌اند؛ مع‌هذا در سخت‌افزار بعضی رایانه‌ها دقت مضاعف گنجانده نشده بلکه فقط در نرم‌افزار آنها گنجانده شده است، که سرعت کمتری دارد. مجموع به دست آمده S_m در

$$S - S_m \doteq \delta S \quad (۱۹.۵.۱)$$

صدق می‌کند که بهبود قابل ملاحظه‌ای نسبت به (۱۸.۵.۱) یا (۱۵.۵.۱) است. این کار را می‌توان در بخشی از محاسبات با دقت معمولی به کار برد و دقت را تا حد قابل ملاحظه‌ای بهبود بخشید بدون آنکه تمام محاسبات را با دقت مضاعف انجام داد. در مسائل جبر خطی این امر ممکن است نیاز به حافظه را، در مقایسه با آنچه برای کل محاسبات با دقت مضاعف لازم است، نصف سازد.

۶.۱ پایداری در آنالیز عددی

شماری از مسائل ریاضی دارای جوابهایی هستند که در مقابل خطاهای محاسبه، مثلاً خطای گردکردن، کاملاً حساس‌اند. برای مواجهه با این پدیده مفاهیم پایداری و ضریب وضعیت^۱ را وارد می‌کنیم. ضریب وضعیت یک مسأله رابطه نزدیکی با ماکسیمم دقتی دارد که در حل آن مسأله با استفاده از اعداد با طول متناهی و حساب رایانه‌یی ممکن است به آن دست یافت. این مفاهیم را سپس برای روشهای عددی که در محاسبه جواب به کار می‌روند بسط می‌دهیم. به طور کلی علاقه‌مندیم از آن دسته روشهای عددی استفاده کنیم که در مقابل خطاهای کوچک، حساسیت بیشتری، نسبت به آنچه در خود مسأله ریاضی وجود دارد، نداشته باشند.

برای ساده کردن موضوع، بحث را به مسائلی که دارای شکل یک معادله

$$F(x, y) = 0 \quad (۱۶.۱)$$

هستند محدود می‌کنیم. متغیر x مجهولی است که باید پیدا شود و متغیر y داده‌ای است که جواب به آن بستگی دارد. این معادله می‌تواند معرّف انواع مختلف مسائل باشد. به عنوان مثال، $F(1)$ می‌تواند تابعی حقیقی - مقدار از متغیر حقیقی x باشد و y می‌تواند بردار ضرایب در تعریف F را نشان دهد؛ یا (۲) معادله می‌تواند یک معادله دیفرانسیل یا انتگرالی باشد، که در آن x تابع مجهول و y یک تابع داده شده یا مقادیر مرزی داده شده است.

مسأله (۱۶.۱) را پایدار گویند اگر جواب x به طور پیوسته به متغیر y بستگی داشته باشد. این بدان معناست که اگر $\{y_n\}$ دنباله مقادیری باشد که به نحوی به y میل کند، آنگاه مقادیر

۱. در ترجمه، اغلب کتب آنالیز عددی این واژه را با عدد شرط یا عدد وضعیت مترادف دانسته‌اند ولی ما ضریب وضعیت را مناسبتر تشخیص می‌دهیم (م).

جوابهای متناظر $\{x_n\}$ نیز باستی به طریقی به x میل نمایند. به عبارتی دیگر، اگر تغییرات جزئی در y ایجاد کنیم، این تغییرات باید تغییرات متناظر کوچکی در x ایجاد کنند. ملاک کوچک بودن این تغییرات به نرُمی بستگی دارد که برای اندازه‌گیری بزرگی بردارهای x و y استفاده شده است؛ انتخابهای مختلفی برای این نرُم‌ها وجود دارند که بسته به مسأله تغییر می‌کنند. مسائل پایدار را مسائل خوش حالت نیز می‌گویند، و ما از هر دو عبارت به‌عنوان مترادف یکدیگر استفاده می‌کنیم. اگر یک مسأله پایدار نباشد ناپایدار یا بدحالت خوانده می‌شود.

مثال (الف) جوابهای معادلهٔ زیر را در نظر می‌گیریم

$$ax^2 + bx + c = 0 \quad a \neq 0$$

هر جواب x یک عدد مختلط است. برای داده‌ها در این حالت، از بردار ضرایب $y = (a, b, c)$ استفاده می‌کنیم. از فرمول جواب معادلهٔ درجهٔ دوّم

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

به وضوح دیده می‌شود که هر دو جواب x با داده‌های $y = (a, b, c)$ به‌طور پیوسته تغییر می‌یابند.

(ب) مسألهٔ معادلهٔ انتگرالی زیر را در نظر می‌گیریم

$$\int_0^1 \frac{\cos x(t) dt}{1.25 - \cos(2\pi(s+t))} = y(s) \quad 0 \leq s \leq 1 \quad (2.6.1)$$

این یک مسألهٔ ناپایدار است. اختلالاتی چون $\delta_n(s) = y_n(s) - y(s)$ وجود دارند که برای آنها، وقتی که $n \rightarrow \infty$

$$\text{Max}_{0 \leq s \leq 1} |\delta_n(s)| \rightarrow 0 \quad n \rightarrow \infty \quad (3.6.1)$$

و جواب متناظر $x_n(s)$ ، به ازای همهٔ مقادیر $n \geq 1$ ، در رابطهٔ زیر صدق می‌کند

$$\text{Max}_{0 \leq s \leq 1} |x_n(s) - x(s)| = 1 \quad n \geq 1 \quad (4.6.1)$$

به‌ویژه، تعریف می‌کنیم $y_n(s) = y(s) + \delta_n(s)$

$$\delta_n(s) = \frac{1}{2^n} \cos(2n\pi s) \quad 0 \leq s \leq 1 \quad n \geq 1$$

در این صورت می‌توان نشان داد

$$x_n(s) - x(s) = \cos(2n\pi s)$$

که (۴.۶.۱) را ثابت می‌کند.

اگر مسأله (۱.۶.۱) ناپایدار باشد مشکلات جدی برای حل آن وجود دارد. معمولاً ممکن نیست چنین مسأله‌ای را، پیش از پی‌بردن بیشتر به ویژگیهای جواب، حل کرد، که این امر معمولاً از توجه به زمینه‌ای که مسأله ریاضی در آن طرح شده حاصل می‌شود. این موضوع اخیراً زمینه بسیار فعالی در تحقیقات ریاضیات کاربردی شده است. [برای مثال تیخونوف و ارسنین^۱ (۱۹۷۷) و وایا^۲ (۱۹۸۰) را ببینید].

مسائل زیادی وجود دارند که از لحاظ عملی به تعبیر فوق‌الذکر، پایدارند ولی، از نظر محاسبات آنالیز عددی، کماکان پردردسرنند. برای مواجهه با این مشکل، اندازه‌ای برای پایداری معرفی می‌کنیم که ضریب وضعیت نامیده می‌شود. این اندازه نشان می‌دهد که از نظر پایداری عملی بعضی از مسائل رفتار بهتری دارند.

ضریب وضعیت سعی دارد بدترین اثر ممکن روی x ، جواب (۱.۶.۱)، را وقتی اختلالات کوچکی در متغیر y ایجاد می‌شود اندازه‌گیری کند. گیریم δy اختلال y و δx جواب معادله اختلال‌یافته

$$F(x + \delta x, y + \delta y) = 0 \quad (۵.۶.۱)$$

باشد. تعریف می‌کنیم

$$K(x) = \sup_{\delta y} \frac{\| \delta x \| / \| x \|}{\| \delta y \| / \| y \|} \quad (۶.۶.۱)$$

برای اندازه بزرگی از نماد $\| \cdot \|$ استفاده کرده‌ایم. تعاریف (۱.۶.۱) - (۱۸.۱.۱) را برای بردارهایی از R^n و $C[a, b]$ یادآوری می‌نمائیم. در مثال (۲.۶.۱) برای اندازه‌گیری اختلالات هردوی x و y از نرم (۱۸.۱.۱) استفاده شده است. معمولاً x و y ممکن است متغیرهایی از انواع مختلف و بنابراین نرم‌های آنها متفاوت باشند. سوپریمم (۶.۶.۱)، روی تمام اختلالات کوچک δy ، که برای آنها هنوز مسأله (۵.۶.۱) معنی دارد، گرفته شده است. مسأله‌ای که ناپایدارند به $K(x) = \infty$ می‌انجامند.

عدد $K(x)$ ضریب وضعیّت (۱.۶.۱) خوانده می‌شود. این عدد اندازه حساسیت جواب x در مقابل تغییرات کوچک در داده‌های y است. اگر $K(x)$ خیلی بزرگ باشد، تغییرات نسبی کوچکی مانند δy برای y وجود دارند که تغییرات نسبی بزرگ δx را برای x موجب می‌شوند. ولی اگر $K(x)$ کوچک باشد، مثلاً $K(x) \leq 10^0$ ، آنگاه تغییرات نسبی کوچک در y همیشه تغییرات نسبی کوچکی در x ایجاد خواهند نمود. از آنجا که محاسبات عددی تقریباً همیشه دارای خطاهای کوچک محاسباتی مختلف می‌باشند، مسائل با ضریب وضعیّت بزرگ را نمی‌پسندیم. چنین مسائلی بدوضع خوانده می‌شوند و حل دقیق آنها معمولاً بسیار مشکل است.

مثال حل معادله

$$x - a^y = 0 \quad a > 0 \quad (7.6.1)$$

را در نظر می‌گیریم. با ایجاد اختلالی به اندازه δy در y ، داریم

$$\frac{\delta x}{x} = \frac{a^{y+\delta y} - a^y}{a^y} = a^{\delta y} - 1$$

برای ضریب وضعیّت (۷.۶.۱)

$$K(x) = \text{Sup}_{\delta y} \left| \frac{\delta x/x}{\delta y/y} \right| = \text{Sup}_{\delta y} \left| y \left(\frac{a^{\delta y} - 1}{\delta y} \right) \right|$$

اگر δy محدود به مقادیر کوچک باشد، داریم

$$K(x) \doteq |y \cdot \ln(a)| \quad (8.6.1)$$

بدون توجه به اینکه چگونه x را از (۷.۶.۱) محاسبه می‌کنیم، اگر $K(x)$ بزرگ باشد، تغییرات نسبی کوچک در y تغییرات نسبی بزرگی را در x ایجاد می‌نمایند. اگر $K(x) = 10^4$ ، و اگر مقدار y مورد استفاده دارای خطای نسبی 10^{-7} باشد، که به لحاظ حساب رایانه‌یی با طول متناهی و گردکردن ایجاد شده است، آنگاه احتمالاً مقدار به دست آمده x خطایی نسبی در حدود 10^{-3} خواهد داشت. این یک افت زیادی در دقت است و به ندرت می‌توان از آن احتراز کرد جز اینکه شاید تمام محاسبات را با حساب رایانه‌یی با دقت طولانیتر انجام داد، به شرطی که y نیز با دقت بیشتری به دست آید.

$$Y = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \cdots & \frac{1}{n} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \cdots & \frac{1}{n+1} \\ \vdots & & & & \vdots \\ \frac{1}{n} & \frac{1}{n+1} & \cdots & \cdots & \frac{1}{2n-1} \end{bmatrix} \quad (9.6.1)$$

را که ماتریس هیلبرت خوانده می‌شود در نظر می‌گیریم. مسأله محاسبه معکوس Y یا معادل آن، حل $YX = I$ که I ماتریس یکه است، یک مسأله خوش حالت است. جواب X با تعداد متناهی عملیات ساده حساب به دست می‌آید. ولی مسأله محاسبه X وقتی n بزرگ می‌شود، بتدریج بدوضع می‌شود.

بدوضع بودن معکوس عددی Y به طور عملی نشان داده خواهد شد. گیریم \hat{Y} نتیجه وارد کردن ماتریس Y در یک رایانه IBM-۳۷۰ باشد و درایه‌های ماتریس با شکل ممیز شناور با دقت معمولی ذخیره شده‌اند. درایه‌های کسری ماتریس در سیستم شانزده‌شانزدهمی (بایه ۱۶) بسط داده شده و سپس بعد از ۶ رقم شانزده‌شانزدهمی (در حدود ۷ رقم اعشاری) قطع شده‌اند. چون بیشتر درایه‌های Y دارای بسط متناهی شانزده‌شانزدهمی نیستند، خطای نسبی هر درایه \hat{Y} در حدود 10^{-6} خواهد بود. با استفاده از حساب با دقت بالاتر، می‌توانیم مقدار دقیق \hat{Y}^{-1} را محاسبه کنیم. وارون Y^{-1} به طور تحلیلی معلوم است، و می‌توانیم آن را با \hat{Y}^{-1} مقایسه نماییم. برای $n = 6$ ، بعضی درایه‌های \hat{Y}^{-1} با درایه‌های متناظر خود در Y^{-1} در اولین رقم غیر صفر خود تفاوت دارند. برای مثال، درایه‌های سطر ۶ ستون ۲ چنین‌اند

$$(Y^{-1})_{6,2} = 83160.00 \quad (\hat{Y}^{-1})_{6,2} = 73866.34$$

این امر، محاسبه Y^{-1} را به یک مسأله (بدوضع) بدل می‌کند، و این بدوضع با بزرگ شدن n بیشتر می‌شود. بر اثر دقت ضعیف \hat{Y}^{-1} در مقایسه با Y^{-1} ، ضریب وضعیّت در (۲.۶.۱) دست‌کم 10^6 خواهد شد. نباید این گونه تصور شود که این یک مثال نامعقول غیرعادی است که در عمل رخ نخواهد داد. این مثال خاص در نظریه تقریب کمترین مربعات پیش می‌آید (برای مثال بخش ۳.۴ را ببینید). زمینه کلی مسائل بدوضع در دستگاه‌های خطی و وارون ماتریسها با شرح بیشتری در فصل ۸ بررسی شده است.

پایداری الگوریتم‌های عددی. یک روش عددی، برای حل یک مسأله ریاضی هنگامی پایدار خوانده می‌شود که حساسیت جواب عددی نسبت به داده‌ها از آنچه در مسأله ریاضی اصلی وجود

دارد، بیشتر نباشد. ما این موضوع را، مجدداً با استفاده از (۱.۶.۱) به عنوان یک مدل برای این مسأله، دقیقتر بیان می‌کنیم. یک روش عددی برای حل (۱.۶.۱) معمولاً به یک دنباله از مسائل تقریبی

$$F_n(x_n, y_n) = 0 \quad (۱۰.۶.۱)$$

منجر می‌شود که به پارامتری چون n بستگی دارد. داده‌های y_n وقتی $n \rightarrow \infty$ به y میل می‌کنند؛ مقادیر تابع $F_n(z, w)$ برای تمام مقادیر (z, w) نزدیک به (x, y) ، با $n \rightarrow \infty$ به $F(z, w)$ میل خواهند کرد. برای مثال، (۱.۶.۱) را می‌توان یک معادلهٔ دیفرانسیل با مقدار اولیه مجسم نمود و (۱۰.۶.۱) ممکن است یک دنباله از تقریبات تفاضل - متناهی وابسته به $h = \frac{1}{n}$ را نشان دهد، همان‌گونه که در (۵.۳.۱) و به دنبال آن آمده است. حالت دیگر وقتی است که n نشان‌دهندهٔ تعداد ارقام استفاده‌شده در محاسبات باشد و خواهیم $F(x, y) = 0$ را با حساب، در حد ممکن با این دقت متناهی، حل کنیم.

برای هر یک از مسائل (۱۰.۶.۱) می‌توانیم یک ضریب وضعیت $K_n(x_n)$ درست مانند (۶.۶.۱) تعریف کنیم. با استفاده از این ضرایب وضعیت، تعریف می‌کنیم

$$\hat{K}(x) = \limsup_{n \rightarrow \infty} \sup_{k \geq n} K_k(x_k) \quad (۱۱.۶.۱)$$

گوییم روش عددی پایدار است اگر $\hat{K}(x)$ در حدود مقدار $K(x)$ در (۶.۶.۱) باشد، مثلاً اگر

$$\hat{K}(x) \leq 2K(x)$$

اگر این درست باشد، حساسیت (۱۰.۶.۱) در مقابل تغییرات داده‌ها در حدود حساسیت مسألهٔ اصلی (۱.۶.۱) خواهد بود.

بعضی مسائل و روشهای عددی به سادگی در چارچوب (۱.۶.۱)، (۶.۶.۱)، (۱۰.۶.۱) و (۱۱.۶.۱) قرار نمی‌گیرند، ولی یک صورت کلی از مسائل پایدار و ضرایب وضعیت را می‌توان تعریف کرد که معنی مشابهی به دست خواهد داد. کاربرد اصلی این مفاهیم در این کتاب، در زمینه‌های: (۱) ریشه‌یابی برای معادلات چندجمله‌یی (۲) حل معادلات دیفرانسیل؛ و (۳) مسائلی در جبر خطی عددی است. به‌طور کلی مسألهٔ کمی در مورد روشهای عددی ناپایدار در این کتاب وجود دارد. مشکل اصلی حل مسائل بدوضع خواهد بود.

مثال محاسبهٔ تابع بسیل

$$x = J_m(y) = \left(\frac{1}{2}y\right)^m \sum_{k=0}^{\infty} \frac{(-\frac{1}{4}y^2)^k}{k!(m+k)!} \quad m \geq 0 \quad (۱۲.۶.۱)$$

را در نظر می‌گیریم. این سری به سرعت همگراست، و محاسبه x به سادگی نشان می‌دهد که این مسأله در وابستگی که به y دارد خوش‌حالت است.

اکنون محاسبه $J_m(y)$ را با استفاده از رابطه بازگشتی سه تایی

$$J_{m+1}(y) = \frac{2m}{y} J_m(y) - J_{m-1}(y) \quad m \geq 1 \quad (۱۳.۶.۱)$$

با فرض معلوم بودن $J_0(y)$ و $J_1(y)$ ، در نظر می‌گیریم. اکنون به طور عددی نشان خواهیم داد که این روش برای محاسبه $J_m(y)$ ، حتی برای m های نه خیلی بزرگ، یک روش عددی ناپایدار است. y را مساوی ۱ می‌گیریم، که (۱۳.۶.۱) به صورت زیر درمی‌آید:

$$J_{m+1}(1) = 2m J_m(1) - J_{m-1}(1) \quad m \geq 1 \quad (۱۴.۶.۱)$$

از مقادیر $J_0(1)$ و $J_1(1)$ که تا ده رقم با معنی دقیق‌اند استفاده کرده‌ایم. مقادیر بعدی $J_m(1)$ از (۱۴.۶.۱) با حساب دقیق محاسبه و نتایج در جدول ۵.۱ داده شده‌اند. مقادیر واقعی برای مقایسه داده شده‌اند، و واگرایی سریع مقادیر تقریبی نسبت به این مقادیر واقعی را نشان می‌دهند. تنها خطاهایی که وجود داشتند مربوط به گردکردن مقادیر $J_0(1)$ و $J_1(1)$ هستند که اختلال بزرگ فزاینده‌ای در $J_m(1)$ وقتی m بزرگ می‌شود، ایجاد می‌کنند.

استفاده از رابطه بازگشتی سه‌جمله‌یی

$$f_{m+1}(x) = a_m(x)f_m(x) - b_m(x)f_{m-1}(x) \quad m \geq 1$$

جدول ۵.۱ مقادیر محاسبه‌شده $J_m(1)$

m	مقادیر واقعی	مقادیر محاسبه‌شده
۰	۰٫۷۶۵۱۹۷۶۸۶۶	۰٫۷۶۵۱۹۷۶۸۶۶
۱	۰٫۴۴۰۰۵۰۵۸۵۷	۰٫۴۴۰۰۵۰۵۸۵۷
۲	۰٫۱۱۴۹۰۳۴۸۴۸	۰٫۱۱۴۹۰۳۴۸۴۸
۳	۰٫۱۹۵۶۳۳۵۳۹۸E-۱	۰٫۱۹۵۶۳۳۵۳۵E-۱
۴	۰٫۲۴۷۶۶۳۸۹۶۴E-۲	۰٫۲۴۷۶۶۳۶۲E-۲
۵	۰٫۲۴۹۷۵۷۷۳۰۲E-۳	۰٫۲۴۹۷۳۶۱E-۳
۶	۰٫۲۰۹۳۸۳۳۸۰۰E-۴	۰٫۲۰۷۲۴۸E-۴
۷	۰٫۱۵۰۲۳۲۵۸۱۷E-۵	-۰٫۱۰۳۸۵E-۵

یک ابزار متداول در ریاضیات کاربردی و آنالیز عددی است. ولی همان گونه که قبلاً نشان داده شد، ممکن است به روشهای عددی ناپایدار منتهی شوند. برای یک تحلیل کلی از روابط بازگشتی سه‌جمله‌یی به گاؤچی^۱ (۱۹۶۷) نگاه کنید. در حالت‌های (۱۳.۶.۱) و (۱۴.۶.۱)، خطای از دست‌دادن تعداد زیادی از ارقام با معنی رخ می‌دهد.

بحث در آثار خواندنی

دانستن حساب رایانه‌یی برای برنامه‌نویسانی که با دقت عددی سروکار دارند حائز اهمیت است، به‌ویژه در نوشتن برنامه‌هایی که به طور گسترده مورد استفاده قرار می‌گیرند. همچنین در نوشتن برنامه‌هایی که با رایانه‌های متفاوت اجرا می‌شوند، ویژگیهای مختلف ممیز شناور در رایانه‌ها بایستی در نظر گرفته شوند. نحوه قدیمی برخورد با حساب ممیز شناور در کنوت^۲ (فصل ۴، ۱۹۸۱) و اشتربتس^۳ (۱۹۷۴) داده شده است.

بحث دقیق در موضوع انتشار خطا که سودمند باشد، به‌ویژه در مورد خطای گردکردن/قطع‌کردن، مشکل بوده است. بعضی مقالات قدیمی مهم در این مورد وجود دارند، ولی گرایشهای جاری به این موضوع به مقدار زیادی مدیون کارهای شادروان و بلیکینسن است. بیشتر کارهای او در جبر خطی عددی بوده است ولی سهم مهمی در بسیاری از زمینه‌های آنالیز عددی داشته است. برای یک آشنایی کلی با روشهای او در تحلیل انتشار خطا، همراه با کاربرد در چندین مسأله مهم به [بلیکینسن (۱۹۶۳)، (۱۹۶۵)، (۱۹۸۴)] نگاه کنید.

روش دیگر کنترل خطا، تحلیل بازه‌یی خوانده می‌شود. در این روش، به جای اینکه محاسبات را در یک عدد x_A ، دنبال کنیم آن را در طول یک بازه $[x_l, x_u]$ دنبال می‌کنیم و اعداد x_l و x_u حدود مقدار واقعی x_T را تضمین می‌کنند. مشکل این روش آن است که $x_u - x_l$ معمولاً بسیار بزرگتر از $|x_T - x_A|$ است، عمدتاً بدین‌علت که حذف خطاهای مختلف‌العلامه که ممکن است در محاسبه x_l و x_u پیش بیایند در نظر گرفته نمی‌شوند. برای آشنایی با این زمینه، و اینکه چگونه می‌توان این کرانه‌ها را در حالت‌های خاص بهبود بخشید به مور^۴ (۱۹۶۶) نگاه کنید. اخیراً، این زمینه با زمینه حساب رایانه‌یی بهم آمیخته شده است تا یک چارچوب نظری کلی برای عرضه الگوریتم‌هایی با کرانه‌های دقیق خطا را ممکن نماید. برای مثال‌هایی در این زمینه کتابهای درسی آلفلت و هرتسبرگر^۵ (۱۹۸۳) و کولیش و میرانکر^۶ (۱۹۸۱) و گزارشهای سمپوزیوم آلفلت و گریگوروف^۷ (۱۹۸۰) و برآوردهای مور (۱۹۷۹) را ببینید.

- | | | | |
|---------------------------|-------------------------|-------------|----------|
| 1. Gautschi | 2. Knuth | 3. Sterbenz | 4. Moore |
| 5. Alefeld and Herzberger | 6. Kulisch and Miranker | | |
| 7. Alefeld and Grigoreiff | | | |

به مسائل بدحالت در بخش ۶.۱ اشاره‌ای شد، ولی این موضوع در سالهای اخیر مورد توجه قرار گرفته است. مسائل زیادی در اندازه‌گیریهای غیرمستقیم فیزیکی وجود دارند که به مسائل بدحالت منجر می‌شوند که در این شکل، مسائل معکوس خوانده می‌شوند. کتاب لاورنتیف^۱ (۱۹۶۷) یک مقدمه کلی برای این مسائل به دست می‌دهد، اگر چه عمدتاً درباره (۱) ادامه تحلیلی توابع تحلیلی با یک متغیر مختلط، و (۲) مسائل معکوس در معادلات دیفرانسیل، بحث می‌کند. یکی از ابزارهای اصلی عددی که در مواجهه با مسائل بدحالت مورد استفاده قرار می‌گیرد، مرتب‌سازی خوانده می‌شود، و بسط جامعی از این روش در کتاب تیخونوف و آرسنین^۲ (۱۹۷۷) داده شده است. برای رایجترین نوشته‌های روشهای عددی در مسائل بدحالت مثلاً به گروچ^۳ (۱۹۸۴) و وایا (۱۹۸۰) نگاه کنید.

دو نوع جدید رایانه در دهه پانزده سال اخیر پیدا شده‌اند که اکنون نقش بیسابقه و فزاینده‌ای در آنالیز عددی پیدا کرده‌اند. اینها ریزرایانه‌ها و زبررایانه‌ها هستند. همه از ریزرایانه‌ها مطلع‌اند؛ دانشمندان و مهندسان به مقدار زیادی در محاسبات عددی خود از آنها استفاده می‌کنند. در ابتدا طرح حسابی در ریزرایانه‌ها نسبتاً ضعیف بود؛ و خطاهایی در عملیات حساب اصلی آنها وجود داشت. اخیراً استاندارده جدید بسیار خوبی برای حساب در ریزرایانه‌ها تولید شده است، که با آن می‌توان برنامه‌هایی عددی در سطح بالا و کارآمد نوشت. این استاندارد یعنی استاندارد برای حساب با ممیز شناور دودویی IEEE، در IEEE (۱۹۸۵) توضیح داده شده است. پیاده‌سازی این استاندارد در بیشتر ریزپردازها شروع شده است، برای مثال، پالمر و مورس^۴ (۱۹۸۴) را نگاه کنید.

نام زبررایانه به نوعی از ماشینهای طراحی شده اطلاق می‌شود که همگی قادرند محاسبات عددی را با سرعت بسیار زیاد مثلاً بیشتر از ۲۰ میلیون عمل با ممیز شناور در ثانیه را انجام دهند. این زمینه در حال گسترش است و خیلی سریع تغییر می‌کند، و لذا ما می‌توانیم فقط چند مرجع را که اشاره‌ای به اثر این ماشینها روی طراحی الگوریتم‌های عددی دارند نام ببریم. هاکنی و جسهوپ^۵ (۱۹۸۱) و کوین^۶ (۱۹۸۷) کتابهای عمومی در معماری زبررایانه‌ها و طراحی الگوریتم‌های عددی در آنها، هستند. پارتر^۷ (۱۹۸۴) یک گزارش سمپوزیوم است که کاربردهایی از زبررایانه‌ها را در مسائل گوناگون فیزیک به دست می‌دهد؛ و اورتگا و فویکت^۸ (۱۹۸۵) درباره زبررایانه‌ها وقتی برای حل معادلات دیفرانسیل جزئی به‌کار برده می‌شوند، بحث می‌کند. این ماشینها بتدریج در تمام زمینه‌های محاسبات، اهمیت فزاینده‌ای پیدا می‌کنند و معماری آنها احتمالاً بر روی رایانه‌های معمولی کوچکتر که در حال حاضر بسیار مورد استفاده قرار می‌گیرند، اثر خواهد گذاشت.

1. Lavrentiev

2. Tikhonov and Arsenin

3. Groetsch

4. Palmer & Morse

5. Hockney and Jesshope

6. Quinn

7. Parter

8. Ortega and Voigt

ریاضیات نمادی. زمینه‌ای است که به سرعت رشد می‌کند و با آن می‌توان به جای ریاضیات عددی از روش ریاضیات تحلیلی، مثلاً در پیدا کردن تابع اولیه دقیق، اگر وجود داشته باشد استفاده کرد. این زمینه تاکنون اثر قابل ملاحظه‌ای روی آنالیز عددی نداشته است، ولی به نظر می‌آید که این وضعیت در حال تغییر است، در بسیاری حالات، ریاضیات نمادی در بخشی از محاسبات به‌کار برده می‌شود و برای بقیه محاسبات از روشهای عددی استفاده می‌شود. یکی از پیچیده‌ترین زبانهای برنامه‌نویسی برای اجرای ریاضیات نمادی زبان MACSYMA است که در ژانویه (۱۹۸۴) توضیح داده شده است. برای یک بازنگری و بررسی تاریخی زبانهای برنامه‌نویسی در این زمینه به ون هولتسن و کالمت نگاه کنید.

ما این مبحث را با بحث درباره نرم‌افزار ریاضی خاتمه می‌دهیم. این زمینه به بحث درباره کارگیری الگوریتمهای عددی به صورت برنامه‌های کامپیوتری می‌پردازد، با توجهی دقیق به مسائل دقت، کارایی، انعطاف‌پذیری، قابلیت حمل، و سایر ویژگیهایی که سودمندی برنامه‌ها را تأمین می‌نماید. یک مجله مهم در این زمینه *ACM* در نرم‌افزار ریاضی است. برای یک مطالعه همه‌جانبه در این موضوع، از جمله نوشته‌های مهم برنامه‌ی که در سالهای اخیر توسعه یافته‌اند، به کاول (۱۹۸۴) نگاه کنید. در پیوست این کتاب درباره بعضی بسته‌های نرم‌افزاری که اخیراً برای آنالیز عددی در دسترس قرار گرفته‌اند، توضیح بیشتری داده‌ایم.

مراجع

- Alefeld, G., and R. Grigorieff, eds. (1980). *Fundamentals of Numerical Computation (Computer-oriented Numerical Analysis)*. Computing Supplementum 2, Springer-Verlag, Vienna.
- Alefeld, G., and J. Herzberger (1983). *Introduction to Interval Computations*. Academic Press, New York.
- Bender, E. (1978). *An Introduction to Mathematical Modelling*. Wiley, New York.
- Cowell, W., ed. (1984). *Sources and Development of Mathematical Software*. Prentice-Hall, Englewood Cliffs, N.J.
- Fadeeva, V. (1959). *Computational Methods of Linear Algebra*. Dover, New York.
- Forsythe, G. (1969). What is a satisfactory quadratic equation solver? In B. Dejon and P. Henrici (eds.), *Constructive Aspects of the Fundamental Theorem of Algebra*, pp. 53-61, Wiley, New York.
- Forsythe, G., and C. Moler (1967). *Computer Solution of Linear Algebraic Systems*. Prentice-Hall, Englewood Cliffs, N.J.
- Forsythe, G., M. Malcolm, and C. Moler (1977). *Computer Methods for Mathematical Computations*. Prentice-Hall, Englewood Cliffs, N.J.
- Gautschi, W. (1967). Computational aspects of three term recurrence relations.

SIAM Rev., 9, 24-82.

- Groetsch, C. (1984). *The Theory of Tikhonov Regularization for Fredholm Equations of the First Kind*. Pitman, Boston.
- Henrici, P. (1962). *Discrete Variable Methods in Ordinary Differential Equations*. Wiley, New York.
- Hockney, R., and C. Jesshope (1981). *Parallel Computers: Architecture, Programming, and Algorithms*. Adam Hilger, Bristol, England.
- Hogg, R., and A. Craig (1978). *Introduction to Mathematical Statistics*, 4th ed. Macmillan, New York.
- Institute of Electrical and Electronics Engineers (1985). Proposed standard for binary floating-point arithmetic. (IEEE Task P754), Draft 10.0. IEEE Society, New York.
- Knuth, D. (1981). *The Art of Computer Programming*, vol. 2, *Seminumerical Algorithms*, 2nd ed. Addison-Wesley, Reading, Mass.
- Kulisch, U., and W. Miranker (1981). *Computer Arithmetic in Theory and Practice*. Academic Press, New York.
- Lavrentiev, M. (1967). *Some Improperly Posed Problems of Mathematical Physics*. Springer-Verlag, New York.
- Lin, C., and L. Segal (1974). *Mathematics Applied to Deterministic Problems in the Natural Sciences*. Macmillan, New York.
- Maki, D., and M. Thompson (1973). *Mathematical Models and Applications*. Prentice-Hall, Englewood Cliffs, N.J.
- Moore, R. (1966). *Interval Analysis*. Prentice-Hall, Englewood Cliffs, N.J.
- Moore, R. (1979). *Methods and Applications of Interval Analysis*. Society for Industrial and Applied Mathematics, Philadelphia.
- Ortega, J., and R. Voigt (1985). Solution of partial differential equations on vector and parallel computers. *SIAM Rev.*, 27, 149-240.
- Parter, S., ed. (1984). *Large Scale Scientific Computation*. Academic Press, New York.
- Palmer, J., and S. Morse (1984). *The 8087 Primer*. Wiley, New York.
- Quinn, M. (1987). *Designing Efficient Algorithms for Parallel Computers*. McGraw-Hill, New York.
- Rand, R. (1984). *Computer Algebra in Applied Mathematics: An Introduction to MACSYMA*. Pitman, Boston.
- Rubinow, S. (1975). *Introduction to Mathematical Biology*. Wiley, New York.
- Sterbenz, P. (1974). *Floating-Point Computation*. Prentice-Hall, Englewood Cliffs, N.J.
- Tikhonov, A., and V. Arsenin (1977). *Solutions of Ill-posed Problems*. Wiley, New York.

- Van Hulzen, J., and J. Calmet (1983). Computer algebra systems. In B. Buchberger, G. Collins, R. Loos (eds.), *Computer Algebra: Symbolic and Algebraic Computation*, 2nd ed. Springer-Verlag, Vienna.
- Wahba, G. (1980). Ill-posed problems: Numerical and statistical methods for mildly, moderately, and severely ill-posed problems with noisy data. Tech. Rep. #595, Statistics Department, Univ. of Wisconsin, Madison. Prepared for the *Proc. Int. Symp. on Ill-Posed Problems*, Newark, Del., 1979.
- Wahba, G. (1984). Cross-validated spline methods for the estimation of multivariate functions from data on functionals, in *Statistics: An Appraisal*, H. A. David and H. T. David (eds.), Iowa State Univ. Press, Ames, pp. 205-235.
- Wilkinson, J. (1963). *Rounding Errors in Algebraic Processes*. Prentice-Hall, Englewood Cliffs, N.J.
- Wilkinson, J. (1965). *The Algebraic Eigenvalue Problem*. Oxford, England.
- Wilkinson, J. (1984). The perfidious polynomial. In Golub (ed.), *Studies in Numerical Analysis*, pp. 1-28, Mathematical Association of America.

مسائل

۱. (الف) گیریم $f(x)$ بر $a \leq x \leq b$ پیوسته باشد و میانگین

$$S = \frac{1}{n} \sum_{j=1}^n f(x_j)$$

را برای تمام نقاط x_j در بازه $[a, b]$ در نظر می‌گیریم. نشان دهید برای مقداری چون $\zeta \in [a, b]$

$$S = f(\zeta)$$

راهنمایی: از قضیه مقدار میانی استفاده کنید و برد مقادیر $f(x)$ و بنابراین برد S را در نظر بگیرید.

(ب) قسمت (الف) را برای مجموع

$$S = \sum_{j=1}^n w_j f(x_j)$$

به ازای تمام x_j ها در $[a, b]$ و تمام $w_j \geq 0$ ها تعمیم دهید.

۲. نابرابریهای زیر را نتیجه بگیرید

(الف) برای تمام مقادیر $x, z \geq 0$ $|e^x - e^z| \leq |x - z|$

(ب) برای $-\frac{\pi}{4} < x, z \leq \frac{\pi}{4}$ $|x - z| \leq |\tan(x) - \tan(z)|$

(ج) برای $0 \leq y \leq x$ و $p \geq 1$ $py^{p-1}(x-y) \leq x^p - y^p \leq px^{p-1}(x-y)$

۳. (الف) کران خطا در تقریب

$$\sin(x) \doteq x \quad |x| \leq \delta$$

را به دست آورید.

(ب) برای مقادیر کوچک δ ، خطای نسبی در $\sin(x) \doteq x$ را با استفاده از

$$\frac{\sin(x) - x}{\sin(x)} \doteq \frac{\sin(x) - x}{x} \quad x \neq 0$$

اندازه بگیرید. کران این خطای نسبی تعدیل یافته را برای $|x| \leq \delta$ به دست آورید. δ را به گونه‌ای پیدا کنید که این خطا کوچکتر از 10^{-6} شود که متناظر است با خطای یک درصد.

۴. فرض می‌کنیم $g \in C[a, b]$ ، نشان دهید برای مقداری از ζ در $[a, b]$

$$\int_a^b x^2 (b-x)^2 g(x) dx = \frac{h^5}{30} g(\zeta)$$

۵. برای توابع زیر یک سری تیلر بسازید و کران خطای برشی پس از n جمله را به دست آورید.

(الف) $\sin^{-1}(x), |x| < 1$ (ب) $\frac{1}{x} \int_0^{\infty} e^{-t} dt$

(ج) $\cos(x) + \sin(x)$ (د) $\frac{1}{x} \int_0^{\infty} \frac{\tan^{-1} t dt}{t}$

(ه) $\log \left[\frac{1+x}{1-x} \right], -1 < x < 1$ (و) $\log(1-x), -1 < x < 1$

۶. (الف) با استفاده از نتیجه (۱۱.۱.۱)، می‌توانیم نشان دهیم

$$\frac{\pi}{4} = \tan^{-1}(1) = \sum_{j=0}^{\infty} \frac{(-1)^{j+1}}{2j+1}$$

و می‌توانیم π را با ضرب آن در ۴ بدست آوریم. چرا این یک راه عملی برای محاسبه π نیست؟

(ب) با به‌کار بردن تقریب چندجمله‌یی تیلر یک راه عملی برای محاسبه π معرفی کنید.

۷. با به‌کار بردن قضیه تیلر برای تابع‌های دو متغیره، تقریبات خطی و درجه دوم را برای تابع‌های

$f(x, y)$ داده شده به ازای مقادیر کوچک x و y پیدا کنید. تابع صفحه مماس $z = p(x, y)$ را

که نمودار آن در نقطه $(0, 0, f(0, 0))$ بر $z = f(x, y)$ مماس است معین کنید

(الف) $\frac{(1+x)}{(1+y)}$ (ب) $\sqrt{1+4x-y}$

(ج) $x \cdot \cos(x-y)$ (د) $\cos(x + \sqrt{\pi^2 + y})$

۸. تفاضل منقسم مرتبه دوم $f[x_0, x_1, x_2]$ را که در (۱۳.۱.۱) تعریف شد در نظر می‌گیریم.

- (الف) ویژگی (۱۵.۱.۱) را که می‌گویند، مرتبهٔ شناسه‌های x_0, x_1, x_2 بر مقدار تفاضل منقسم اثری ندارد، ثابت کنید.
- (ب) فرمول (۱۴.۱.۱)

$$f[x_0, x_1, x_2] = \frac{1}{2} f''(\zeta)$$

- را که در آن ζ مقداری بین مینیمم و ماکسیمم x_0, x_1, x_2 است ثابت کنید.
- راهنمایی: در قسمت (الف) بدون کاستن از کلیت آن می‌توان فرض کرد $x_0 < x_1 < x_2$. از قضیهٔ تیلر برای تبدیل $f[x_0, x_1, x_2]$ که حول x_1 بسط داده شده و از قضیهٔ مقدار میانی برای ساده‌کردن جملهٔ خطا استفاده کنید.
- (ج) با فرض آنکه $f(x)$ دویار پیوسته - مشتق‌پذیر باشد، نشان دهید که $f[x_0, x_1, x_2]$ می‌تواند به طور پیوسته به حالتی تعمیم یابد که بعضی یا همهٔ نقاط x_0, x_1, x_2 بر هم منطبق شوند. برای مثال نشان دهید که

$$f[x_0, x_1, x_0] \equiv \lim_{x_2 \rightarrow x_0} f[x_0, x_1, x_2]$$

- موجود است و فرمولی برای محاسبهٔ آن پیدا کنید.
۹. (الف) نشان دهید که نرم‌های برداری (۱۶.۱.۱) و (۱۸.۱.۱) در سه ویژگی عمومی نرمها که به دنبال (۱۸.۱.۱) بیان شده‌اند صدق می‌کنند.
- (ب) نشان دهید $\|x\|_2$ در (۱۷.۱.۱) یک نرم برداری است، خود را به $n = 2$ محدود کنید.
- (ج) نشان دهید نرم ماتریسی (۱۹.۱.۱) در روابط (۲۰.۱.۱) و (۲۱.۱.۱) صدق می‌کند. برای سادگی فقط ماتریسهای 2×2 را در نظر بگیرید.
۱۰. اعداد زیر را به هم‌ارزهای اعشاری آنها برگردانید.

- (الف) $(101010101)_2$ (ب) $(2A3 \cdot FF)_{16}$
- (ج) $(0101010101 \dots)_2$ (د) $(.AAAA \dots)_{16}$
- (ا) $(00001100110011 \dots)_2$
- (و) $(11, \dots)_2$ که n تا یک در پراوتر باشد.

۱۱. برای تبدیل یک عدد صحیح اعشاری مثبت x به هم‌ارز دودویی آن،

$$x = (a_n a_{n-1} \dots a_1 a_0)_2$$

چنین می‌نویسیم

$$x = a_n \cdot 2^n + a_{n-1} \cdot 2^{n-1} + \dots + a_1 \cdot 2^1 + a_0 \cdot 2^0$$

بر این اساس الگوریتم زیر را به‌کار برید

$$(i) \quad x := x_0; \quad j := 0$$

(ii) تا وقتی $x_j \neq 0$ ، عملیات زیر را انجام دهید:

a_j را برابر باقیمانده عدد صحیح $x_j/2$ قرار دهید

x_{j+1} را برابر خارج قسمت عدد صحیح $x_j/2$ قرار دهید

j را برابر $j + 1$ قرار دهید

پایان حلقه تا وقتی

زبان الگوریتم باید نیاز به توضیح نداشته باشد. از این الگوریتم برای تبدیل اعداد صحیح زیر

به هم‌ارز دودویی آنها استفاده کنید.

(الف) ۴۹ (ب) ۱۲۷ (ج) ۱۲۹

۱۲. برای تبدیل یک کسر مثبت اعشاری $x < 1$ به هم‌ارز دودویی آن

$$x = (0.a_1a_2a_3\dots)_2$$

آن را به شکل زیر بنویسید

$$x = a_1 \cdot 2^{-1} + a_2 \cdot 2^{-2} + a_3 \cdot 2^{-3} + \dots$$

بر این اساس الگوریتم زیر را به‌کار برید

$$(i) \quad x := x_0; \quad j := 1$$

(ii) تا وقتی $x_j \neq 0$ ، عملیات زیر را انجام دهید:

a_j را برابر قسمت صحیح $2 \cdot x_j$ قرار دهید

x_{j+1} را برابر جزء کسری $2 \cdot x_j$ قرار دهید

j را برابر $j + 1$ قرار دهید

پایان حلقه تا وقتی

این الگوریتم را برای تبدیل اعداد اعشاری زیر به هم‌ارز دودویی آنها به‌کار برید

(الف) ۸۱۲۵ (ب) ۰۶۲۵ (ج) ۰۱ (د) ۰۲

(ه) ۰۴ (و) $\frac{1}{\sqrt{2}} = 0.۱۴۲۸۵۷۱۴۲۸۵۷\dots$

۱۳. مسائل ۱۱ و ۱۲ را برای تبدیل یک عدد صحیح اعشاری به هم‌ارز شانزده‌شانزده‌ی آنها به‌کار برید.

۱۴. خروجی قطعه برنامه فورتن زیر را با یک رایانهٔ دودویی با استفاده از حساب قطع‌شده، پیش‌بینی کنید

```

I = 0
X = 0.0
H = .1
10  I=I+1
    X = X + H
    PRINT*, I, X
    IF(X .LT. 1.0) Go To 10
    
```

اگر عبارت " $X = X + H$ " با عبارت " $X = I * X$ " عوض شود آیا خروجی تغییر می‌کند؟

۱۵. کرانه‌های (۹.۲.۱) را برای خطای نسبی در تمایش گرد شده با ممیز شناور (۵.۲.۱) - (۶.۲.۱) به‌دست آورید.

۱۶. نتیجهٔ کران بالای $M = \beta^t$ در (۱۶.۲.۱) را به‌دست آورید.

۱۷. (الف) برنامه‌ای بنویسید که یک خطای سرریز در رایانهٔ شما تولید کند. برای مثال، عدد x را به‌عنوان ورودی بگیرید و آن را چندبار به توان دو برسانید.

(ب) برنامه‌ای بنویسید که به‌طور تجربی بزرگترین عدد با ممیز شناور مجاز را معین کند.

۱۸. (الف) یک مدل ساده برای افزایش جمعیت چنین است

$$\frac{dN}{dt} = kN \quad t \geq t_0, \quad N(t_0) = N_0.$$

که $N(t)$ جمعیت در زمان t است و $k > 0$. نشان دهید که این رابطه مستلزم یک نرخ هندسی افزایش در جمعیت است:

$$N(t+1) = CN(t) \quad t \geq t_0.$$

فرمولی برای C پیدا کنید.

(ب) یک مدل پیچیده‌تر برای رشد جمعیت چنین است

$$\frac{dN}{dt} = kN[1 - bN] \quad t \geq t_0, \quad N(t_0) = N_0.$$

که $b, k > 0$ و $1 - bN_0 > 0$. جواب این معادله دیفرانسیل را پیدا کنید. این جواب را با جواب قسمت (الف) مقایسه کنید. تفاوت‌هایی را که در افزایش جمعیت در دو مدل پیش می‌آید برای هر دو حالت مقادیر بزرگ و کوچک t شرح دهید.

۱۹. دو تابع زیر را با رایانه خود محاسبه کنید

$$f(x) = x^3 - 3x^2 + 3x - 1 \quad (\text{الف})$$

$$f(x) = x^3 + 2x^2 - x - 2 \quad (\text{ب})$$

آنها را برای نمونه بزرگی از مقادیر x حول ۱ محاسبه کنید و سپس نوع رفتار آنها را آن‌گونه که در شکل ۲.۱ نشان داده شده به دست آورید. نتایج به دست آمده برای دو تابع را با هم مقایسه کنید.

۲۰. برنامه‌ای بنویسید که به طور تجربی حد زیر را محاسبه کند

$$\lim_{p \rightarrow \infty} (x^p + y^p)^{1/p}$$

در اینجا x و y اعداد مثبت‌اند. اولاً محاسبات را در شکلی که هم‌اکنون نشان داده شد انجام دهید. ثانیاً محاسبات را به روشی که در (۸.۳.۱) به کار رفته تکرار کنید. برنامه را برای انواع مقادیر بزرگ و کوچک x و y ، مثلاً $x = y = 10^{10}$ و $x = y = 10^{-10}$ اجرا کنید.

۲۱. برای اعداد x_A و x_T در زیر، چند رقم با معنی در x_A با توجه به x_T وجود دارد؟

$$x_T = ۴۵۱,۰۰۱ \quad , x_A = ۴۵۱,۰۲۳ \quad (\text{الف})$$

$$x_T = -۰,۰۴۵۱۸ \quad , x_A = -۰,۰۴۵۱۱۳ \quad (\text{ب})$$

$$x_T = ۲۳,۴۶۰۴ \quad , x_A = ۲۳,۴۲۱۳ \quad (\text{ج})$$

۲۲. گیریم تمام اعداد زیر به ارقامی که در اعداد نشان داده شده دقیقاً گرد شده باشند

$$(\text{الف}) ۰,۹۴۷ + ۱,۱۰۶۲, \quad (\text{ب}) ۱۲,۷۵۳ - ۲۳,۴۶, \quad (\text{ج}) (۶,۸۳) \times (۲,۷۴۷)$$

(د) $۰,۰۶۴ / ۸,۴۷۳$. برای هر محاسبه، کوچکترین بازه‌ای را معین کنید که جواب در آن قرار بگیرد، وقتی از مقادیر درست به جای مقادیر تقریبی استفاده شود.

۲۳. فرمول (۵.۴.۱) را برای $\text{Rel}(x_A/y_A)$ اثبات کنید.

۲۴. معادله $x^2 - 40x + 1 = 0$ داده شده است، ریشه‌های آن را تا ۵ رقم با معنی پیدا کنید.

$\sqrt{399} \approx 19,975$ را که تا ۵ رقم با معنی گرد شده به کار برید.

۲۵. راه‌های درست جلوگیری از خطای از دست دادن ارقام با معنی را در محاسبات زیر مشخص کنید

$$\log(1+x) - \log(x) \quad \text{مقادیر بزرگ } x \quad (\text{الف})$$

$$x \doteq y \quad \sin(x) - \sin(y) \quad (\text{ب})$$

$$x \doteq y \quad \tan(x) - \tan(y) \quad (\text{ج})$$

$$x \doteq 0 \quad \frac{1 - \cos(x)}{x^2} \quad (\text{د})$$

$$x \doteq 0 \quad \sqrt{1+x} - 1 \quad (\text{ه})$$

۲۶. برای احتراز از خطای از دست دادن ارقام بامعنی در محاسبات زیر از تقریبهای تیلر استفاده کنید.

$$f(x) = \frac{e^x - e^{-x}}{2x} \quad (\text{الف})$$

$$f(x) = \frac{\log(1-x) + xe^{x/2}}{x^3} \quad (\text{ب})$$

در هر دو حالت مقدار $\lim_{x \rightarrow 0} f(x)$ چیست؟

۲۷. ارزیابی $\cos(x)$ را برای مقادیر بزرگ x با استفاده از تقریب تیلر (۵.۱.۱)

$$\cos(x) \doteq 1 - \frac{x^2}{2!} + \dots + (-1)^n \frac{x^{2n}}{(2n)!}$$

در نظر می‌گیریم. برای مشاهده مشکلاتی که در استفاده از این تقریب پیش می‌آید، آن را برای محاسبه $\cos(2\pi) = 1$ به‌کار برید. n را به گونه‌ای بیابید که خطای تقریب تیلر کمتر از 5×10^{-6} باشد. نوع محاسبه‌ای را که در (۱۷.۴.۱) و جدول ۲.۱ به‌کار رفته تکرار کنید. چگونه $\cos(x)$ برای مقادیر بزرگ x باید ارزیابی شود؟

۲۸. فرض کنید می‌خواهید مقادیر:

$$\cos(1473) \quad (\text{الف}), \quad \cos(1471) \quad (\text{ج})$$

$$\tan^{-1}(2621) \quad (\text{ب}), \quad e^{2653} \quad (\text{د})$$

را محاسبه کنید. در هر حالت، فرض می‌کنیم فقط جدولی از مقادیر تابع برحسب شناسه x که با نمونه‌های 10^{-6} داده شده‌اند در دست باشد. مقداری از جدول را که شناسه آن نزدیکترین مقدار را به شناسه‌های شما داشته باشد انتخاب کنید. خطای نتیجه را تخمین بزنید.

۲۹. فرض کنید $x_A = 937$ سه رقم صحیح با معنی با توجه به x_T دارد. کران خطای نسبی x_A را به دست آورید. برای $f(x) = \sqrt{1-x}$ ، کران خطا و خطای نسبی در $f(x_A)$ را با توجه به $f(x_T)$ پیدا کنید.

۳۰. اعداد زیر تا ارقامی که نشان داده شده‌اند به‌طور دقیق گرد شده‌اند. خطاهای مقادیر توابع را برحسب خطاهای شناسه‌های آنها تخمین بزنید. خطاهای نسبی را تخمین بزنید. خطاهای نسبی را کراندار کنید.

$$\sin[(2685)(314)] \quad (\text{الف}), \quad \ln(1712) \quad (\text{ب}), \quad (156)^{344} \quad (\text{ج})$$

۳۱. یک برنامه رایانه‌ی برای مجموع

$$S = \sum_1^n a_i$$

به سه طریق بنویسید (۱) از کوچکترین مقدار به بزرگترین مقدار، (۲) از بزرگترین مقدار به کوچکترین مقدار (۳) با دقت مضاعف، با دقت معمولی خطای گرد کردن/قطع کردن در نتیجه مجموعیابی. نتیجه دقت مضاعف را برای پیدا کردن خطا در دو نتیجه با دقت معمولی فوق به کار برید. نتایج را چاپ کنید. همچنین یک برنامه کلی بنویسید که سری زیر را با دقت معمولی تولید کند، و برنامه‌ای را که هم‌اکنون داده شد برای محاسبه مجموع سریهای زیر به کار برید.

راهنمایی: در نوشتن برنامه، به خاطر سادگی فرض کنید که جملات سری از بزرگ به کوچک مرتب شده‌اند.

$$\sum_1^n \frac{(-1)^j}{j} \quad (\text{د}) \quad , \quad \sum_1^n \frac{1}{j^2} \quad (\text{ج}) \quad , \quad \sum_1^n \frac{1}{j^3} \quad (\text{ب}) \quad , \quad \sum_1^n \frac{1}{j} \quad (\text{الف})$$

۳۲. حاصلضرب $a_0 a_1 \dots a_m$ را در نظر بگیرید که در آن $a_0, a_1, \dots, a_m, (m+1)$ عدد ذخیره شده در رایانه‌ای هستند که از حساب n رقمی در پایه β استفاده می‌کند. تعریف می‌کنیم $p_1 = \text{fl}(a_0 a_1), p_2 = \text{fl}(a_2 p_1), p_3 = \text{fl}(a_3 p_2), \dots, p_m = \text{fl}(a_m p_{m-1})$. اگر بنویسیم $p_m = a_0 a_1 \dots a_m (1+w)$ ، به ازای $i = 0, 1, \dots, m$ یک کران دقیق برای w چیست؟ برآورد آماری برای اندازه w چیست؟

ریشه‌یابی معادله‌های غیرخطی

پیدا کردن یک یا چند ریشهٔ یک معادلهٔ

$$f(x) = 0 \quad (۱.۰.۲)$$

یکی از متداولترین مسائلی است که در ریاضیات کاربردی پیش می‌آید. در بیشتر موارد جوابهای صریح در دست نیست و باید به همین قانع باشیم که بتوانیم یک ریشه را با هر درجه از دقت مشخصی حساب کنیم. روشهای عددی در یافتن ریشه‌ها، روشهای بارستی (تکراری) خوانده می‌شوند که موضوع اصلی این فصل هستند.

کار را با روشهای بارستی برای حل (۱.۰.۲)، هنگامی که $f(x)$ یک تابع پیوسته مشتقپذیر حقیقی-مقدار از متغیر حقیقی x باشد، آغاز می‌کنیم. روشهای بارستی در حل این رده از معادلات کاملاً کلی نیاز به یک یا چند حدس اولیهٔ x برای ریشهٔ مطلوب α از $f(x)$ دارند. حدس اولیهٔ x را می‌توان معمولاً، با به‌کارگیری شرایطی که اولین بار مسأله در آن ظاهر شده، پیدا کرد؛ در غیر این صورت یک نمودار سادهٔ $y = f(x)$ اغلب برای برآورد x کافی خواهد بود.

مسألهٔ مهم دیگری که در این فصل مورد بحث قرار می‌گیرد یافتن یک یا چند ریشه از معادلهٔ

چند جمله‌ای زیر است

$$p(x) \equiv a_0 + a_1x + \dots + a_nx^n = 0 \quad a_n \neq 0 \quad (2.0.2)$$

روشهای مسأله اول اغلب اختصاص به ۲.۰.۲ یافته است و روش ما نیز همین خواهد بود. ولی نوشته‌های زیادی در مورد روشهایی وجود دارند که به ویژه برای معادله‌های چند جمله‌ای با استفاده از ویژگیهای خاص آنها از راهی طبیعی به وجود آمده‌اند. این روشها مهمترین روشهایی هستند که در ایجاد برنامه‌های خودکار رایانه‌ای برای حل ۲.۰.۲ به کار برده می‌شوند و ما در اینجا به بعضی از آنها اشاره خواهیم کرد.

سومین رده مسائلی که در این فصل مورد بحث قرار می‌گیرند حل دستگاه معادله‌های غیرخطی هستند. این دستگاهها به شکلهای گوناگون هستند و آنالیز عددی مربوط به آنها هم مفصل و هم پیچیده است. ما فقط اشاره‌ای به این موضوع خواهیم کرد و بعضی از روشهای موفق را که از نظر کاربرد نسبتاً کلی هستند ذکر خواهیم کرد. مطالعه مفصل این موضوع نیاز به دانش کافی در جبر خطی نظری و عددی دارد که تا فصلهای ۷ تا ۹ به این مطالب پرداخته شده است.

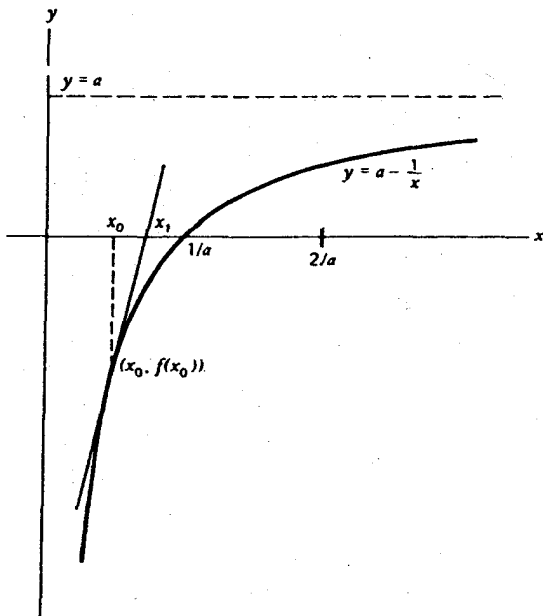
آخرین رده مسائلی مورد بحث در این فصل مسائل بهینه‌سازی هستند. در این حالت در پی ماکسیمم‌سازی یا مینیمم‌سازی یک تابع حقیقی-مقدار $f(x_1, \dots, x_n)$ و پیدا کردن نقطه (x_1, \dots, x_n) هستیم که در آن نقطه مقدار بهینه به دست می‌آید. این مسائل اغلب به یک دستگاه معادله‌های غیرخطی بدل می‌شوند، ولی معمولاً بهتر است که روشهای ویژه‌ای مستقیماً برای بهینه‌سازی تهیه شود. موضوع بهینه‌سازی به نحوی گسترده پیشرفت کرده و تکامل یافته است، ما فقط به اختصار آن را معرفی و مطالعه می‌کنیم.

برای روشن کردن مفهوم روش بارستی در پیدا کردن ریشه ۱.۰.۲، مطلب را با یک مثال آغاز می‌کنیم.

حل معادله

$$f(x) \equiv a - \frac{1}{x} = 0 \quad (3.0.2)$$

را به ازای مقدار داده شده $a > 0$ در نظر می‌گیریم. این مسأله یک کاربرد عملی در رایانه‌های بدون عمل تقسیم دارد. این امر در بعضی رایانه‌های اولیه واقعیت داشت، و بعضی رایانه‌های جدید نیز از الگوریتمی که ذیلاً داده می‌شود برای عمل تقسیم استفاده می‌کنند.



شکل ۱.۲ حل بارستی $a - (1/x) = 0$

گیریم x_0 یک جواب تقریبی معادله باشد. در نقطه $(x_0, f(x_0))$ خط مماس بر نمودار $y = f(x)$ را رسم می‌کنیم (شکل ۱.۲). بگیریم x_1 نقطه‌ای باشد که خط مماس محور x ها را قطع می‌کند. این نقطه باید یک تقریب بهتری برای ریشه α باشد. برای به دست آوردن یک معادله برای x_1 ، شیب خط مماس و مشتق $f(x)$ در x_0 را برابر قرار

می‌دهیم

$$f'(x_0) = \frac{f(x_0) - 0}{x_0 - x_1}$$

با توجه به ۳.۰.۲ و اندکی عملیات به دست می‌آوریم

$$x_1 = x_0(2 - ax_0)$$

فرمول بارستی کلی از تکرار این روند با قراردادن x_1 به جای x_0 و ادامه این کار تا بینهایت، حاصل می‌شود

$$x_{n+1} = x_n(2 - ax_n) \quad n \geq 0 \quad (4.0.2)$$

یک شکل مناسبتر برای منظوره‌های نظری با وارد کردن شکل مقیاس شده مانده

$$r_n = 1 - ax_n \quad (5.0.2)$$

به دست می‌آید. با استفاده از آن خواهیم داشت:

$$x_{n+1} = x_n(1 + r_n) \quad n \geq 0 \quad (۶.۰.۲)$$

پس خطا برابر است با

$$e_n = \frac{1}{a} - x_n = \frac{r_n}{a} \quad (۷.۰.۲)$$

اکنون به تحلیل همگرایی این روش، سرعت و وابستگی آن به x_0 می‌پردازیم. ابتدا

$$r_{n+1} = 1 - ax_{n+1} = 1 - ax_n(1 + r_n) = 1 - (1 - r_n)(1 + r_n)$$

$$r_{n+1} = r_n^2 \quad (۸.۰.۲)$$

از راه استقرار

$$r_n = r_0^{2^n} \quad n \geq 0 \quad (۹.۰.۲)$$

با توجه به ۷.۰.۲، خطای e_n وقتی $n \rightarrow \infty$ ، به صفر می‌گراید اگر و تنها اگر r_n به صفر میل کند. با توجه به (۹.۰.۲)، r_n به صفر می‌گراید اگر و تنها اگر $|r_0| < 1$ ، یا هم‌ارز با آن

$$-1 < 1 - ax_0 < 1$$

$$0 < x_0 < \frac{2}{a} \quad (۱۰.۰.۲)$$

برای آنکه x_n به $1/a$ میل کند لازم و کافی است که x_0 طوری انتخاب شود که در (۱۰.۰.۲) صدق کند.

برای بررسی سرعت همگرایی، وقتی که (۱۰.۰.۲) برقرار باشد، فرمولهایی برای خطا و خطای نسبی به دست می‌آوریم. برای بررسی سرعت همگرایی در حالتی که (۱۰.۰.۲) برقرار باشد

$$e_{n+1} = \frac{r_{n+1}}{a} = \frac{r_n^2}{a} = \frac{e_n^2 a^2}{a}$$

$$e_{n+1} = a e_n^2 \quad (۱۱.۰.۲)$$

$$\frac{e_{n+1}}{1/a} = e_n^2 a^2 = \left[\frac{e_n}{1/a} \right]^2$$

$$\text{Rel}(x_{n+1}) = \text{Rel}(x_n)^2 \quad n \geq 0 \quad (۱۲.۰.۲)$$

نماد $\text{Rel}(x_n)$ معرّف خطای نسبی در x_n است. با توجه به معادله (۱۱.۰.۲)، می‌گوییم e_n به طور مربعی به صفر همگرا است. برای نشان‌دادن سرعت کاهش خطا، فرض می‌کنیم $\text{Rel}(x_0) = 0.1$ ، در این صورت $\text{Rel}(x_4) = 10^{-16}$. با هر بارست شمار ارقام با معنی دو برابر می‌شود.

این مثال نحوه ساختن یک روش بارستی را برای حل یک معادله نشان می‌دهد؛ و یک تحلیل کامل همگرایی در آن داده شده است. این تحلیل شامل اثبات همگرایی، تعیین بازه همگرایی برای انتخاب x_0 ، و تعیین سرعت همگرایی است. این مفاهیم با استفاده از روشهای کلّیتری در حل (۱۰.۰.۲) در بخشهای بعد بررسی شده‌اند.

تعریف یک دنباله از بارستهای $\{x_n \mid n \geq 0\}$ را به نقطه α همگرا از مرتبه $p \geq 1$ خوانند اگر به ازای مقداری مانند $c > 0$

$$|\alpha - x_{n+1}| \leq c |\alpha - x_n|^p \quad n \geq 0 \quad (13.0.2)$$

اگر $p = 1$ ، دنباله را همگرای خطی به α خوانند. در این حالت باید $c < 1$. ثابت c نرخ همگرایی خطی x_n به α نامیده می‌شود.

با استفاده از این تعریف، مثال قبلی (۶.۰.۲) (۵.۰.۲) دارای همگرایی مرتبه ۲ است که همگرایی مربعی نیز نامیده می‌شود. این تعریف مرتبه، همیشه برای بعضی روشهای بارستی همگرای خطی مناسب نیست. با استفاده از استقرأ در (۱۳.۰.۲) با $p = 1$ ، به دست می‌آوریم

$$|\alpha - x_n| \leq c^n |\alpha - x_0| \quad n \geq 0 \quad (14.0.2)$$

این رابطه مستقیماً همگرایی x_n به α را نشان می‌دهد. برای بعضی روشهای بارستی، می‌توانیم (۱۴.۰.۲) را مستقیماً ثابت کنیم، حال آنکه ممکن است (۱۳.۰.۲) به ازای هیچ مقدار $c < 1$ درست نباشد. در چنین حالتی، این روش بازهم همگرای خطی با نرخ c خوانده می‌شود.

۱.۲ روش نیم‌سازی (تنصیف)

فرض کنید $f(x)$ بر بازه داده شده $[a, b]$ پیوسته باشد و در رابطه زیر نیز صدق کند

$$f(a)f(b) < 0 \quad (1.1.2)$$

با استفاده از قضیه مقدار میانی ۱.۱ از فصل ۱، تابع $f(x)$ باید حداقل یک ریشه در $[a, b]$ داشته

باشد. معمولاً $[a, b]$ به گونه‌ای انتخاب می‌شود که تنها یک ریشه α را شامل شود، ولی الگوریتم زیر، برای روش نیم‌سازی به موجب (۱.۱.۲) همیشه به ریشه‌ای چون α در $[a, b]$ همگراست.

الگوریتم نیم‌ساز $(f, a, b, \text{ریشه}, \varepsilon)$

$$1. \text{تعریف کنید } c := (a + b)/2$$

۲. اگر $b - c \leq \varepsilon$ ، آنگاه $c :=$ ریشه را بپذیرید و خارج شوید.

۳. اگر $f(c) \leq 0$ (علامت $f(c)$). (علامت $f(b)$)، آنگاه $c := a$ ؛ در غیر این صورت $c := b$.

۴. به گام اول برگردید.

بازه $[a, b]$ در هر گذر از الگوریتم نصف می‌شود. به موجب گام ۳، $[a, b]$ همیشه یک ریشه $f(x)$ را در بر دارد. چون یک ریشه α در $[a, b]$ است، بایستی یا در $[a, c]$ باشد یا در $[c, b]$ و در نتیجه

$$|c - \alpha| \leq b - c = c - a$$

این توجیه آزمون گام ۲ است. پس از تکمیل الگوریتم، c یک مقدار تقریبی ریشه با

$$|c - \alpha| \leq \varepsilon$$

خواهد بود.

مثال بزرگترین ریشه حقیقی α از معادله

$$f(x) \equiv x^6 - x - 1 = 0 \quad (2.1.2)$$

را پیدا کنید. به سادگی می‌توان نشان داد که $1 < \alpha < 2$ ، و ما این بازه را به‌عنوان بازه اولیه $[a, b]$ به‌کار می‌بریم. الگوریتم نیم‌ساز با $\varepsilon = 0.00005$ به‌کار گرفته شده‌است. نتایج در جدول ۱.۲ نشان داده شده‌اند. دوبارست اول بازه اولیه شامل α را به‌دست می‌دهند و سایر مقادیر $n \geq 1, c_n$ ، معرّف نقاط میانی متوالی هستند که با استفاده از نیم‌سازی پیدا شده‌اند. آخرین مقدار $c_{15} = 1.13474$ به‌عنوان تقریبی برای α با $0.00004 \leq |\alpha - c_{15}|$ پذیرفته شده‌است. جواب درست چنین است:

$$\alpha \approx 1.13472413840152 \quad (3.1.2)$$

خطای واقعی در c_{15} عبارت است از

$$\alpha - c_{15} = -0.000016$$

جدول ۱.۲ مثال روش نیم‌سازی

n	c_n	$f(c_n)$	n	c_n	$f(c_n)$
	$2.0 = b$	6.10	۸	1.3672	0.02062
	$1.0 = a$	-1.0	۹	1.3477	0.00043
۱	1.5	8.89063	۱۰	1.3379	-0.000960
۲	1.25	1.56470	۱۱	1.3428	-0.000459
۳	1.125	-0.09771	۱۲	1.3452	-0.000208
۴	1.1875	0.61665	۱۳	1.3464	-0.000083
۵	1.15625	0.23327	۱۴	1.3470	-0.000020
۶	1.14063	0.06158	۱۵	1.3474	0.000016
۷	1.13281	-0.1958			

این خطا بسیار کوچکتر از کران خطای پیش‌بینی شده‌است. ممکن است چنین به نظر آید که می‌توانستیم با توقف در یکی از بارستهای قبلی قدری از محاسبات بکاهیم. ولی راهی برای پیش‌بینی دقت بهتر احتمالی در بارست قبلی نیست، و بنابراین راهی وجود ندارد که بتوانیم تشخیص دهیم که این بارست به اندازه کافی دقیق هست یا نیست. برای مثال c_9 به اندازه کافی دقیق است، ولی هیچ راهی که به ما این حقیقت را در حین محاسبات بگوید وجود ندارد. برای بررسی سرعت همگرایی، گیریم c_n معرّف n امین مقدار c در الگوریتم باشد. لذا به سادگی می‌توان دید که

$$\alpha = \lim_{n \rightarrow \infty} c_n$$

$$|\alpha - c_n| \leq \left[\frac{1}{2}\right]^n (b - a) \quad (4.1.2)$$

که $b - a$ معرّف طول بازه اولیه‌ای است که در نیم‌سازی داده شده است. با استفاده از صورت (۴.۱.۲) برای تعریف همگرایی خطی، می‌گوییم روش نیم‌سازی همگرایی خطی با نرخ $1/2$ است. ممکن است خطای واقعی در هر گام با نرخ $1/2$ کاهش نیابد، ولی بر مبنای (۴.۱.۲) نرخ متوسط کاهش $1/2$ است. مثال قبلی نتیجه (۴.۱.۲) را نشان می‌دهد.

الگوریتم نیم‌ساز چند عیب دارد. اول، حدود دقت ماشین را آن‌گونه که در بخش ۱-۲ فصل یک توضیح داده شد، رعایت نمی‌کند. یک برنامه عملی واحد گرد کردن ماشین را در نظر می‌گیرد و اگر لازم باشد ϵ داده شده را تعدیل می‌نماید [(۱۲.۲.۱)]. مشکل عمده دیگر در روش نیم‌سازی آن است که در مقایسه با سایر روشهایی که در بخشهای بعدی تعریف می‌شوند، به کندی همگرا می‌شود. امتیازهای

عمده روش نیم‌سازی عبارت‌اند از: (۱) همگرایی آن تضمین شده است (به شرطی که f روی $[a, b]$ پیوسته و (۱.۱.۲) برقرار باشد) و (۲) یک کران خطای قابل قبول در دسترس است. روشهایی که در هر گام کرانه‌های بالایی و پایینی ریشه α را به دست می‌دهند روشهای حصارکشی (خصر) خوانده می‌شوند. در بخش ۸.۲ ما به تشریح یک الگوریتم خصر می‌پردازیم که مزایای روش نیم‌سازی مذکور در قبل و همگرایی سریعتر روش خط قاطع را (که در بخش ۲-۳ توضیح داده می‌شود) توأمأ دارد.

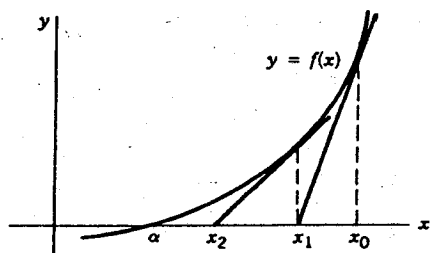
۲.۲ روش نیوتن

فرض کنید که یک برآورد اولیه x_0 برای α ، ریشه مطلوب $f(x) = 0$ ، معلوم باشد. روش نیوتن دنباله‌ای از بارستهای $\{x_n : n \geq 1\}$ را تولید می‌کند، که امیدواریم به α همگرا شود. چون x_0 نزدیک به α فرض شده است، نمودار $y = f(x)$ را در مجاورت ریشه α با رسم خط مماس بر نمودار در $(x_0, f(x_0))$ تقریب می‌زنیم. سپس ریشه این خط مماس^۱ را برای تقریب زدن α به کار می‌بریم؛ این تقریب جدید را x_1 می‌نامیم. این روند را به دفعات بسیار زیاد تکرار می‌کنیم تا دنباله‌ای از بارستهای x_n به دست آید. مانند مثال (۳.۰.۲) در ابتدای این فصل، این کار ما را به فرمول بارستی زیر می‌رساند

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \quad n \geq 0 \quad (1.2.2)$$

این روند در شکل ۲.۲ برای بارستهای x_1 و x_2 نشان داده شده است.

روش نیوتن معروفترین شیوه برای یافتن ریشه‌های یک معادله است. این روش به راههای گوناگونی برای حل مسائل غیرخطی مشکلتر، مثلاً برای دستگاههای معادلات غیرخطی و معادلات انتگرالی و دیفرانسیل غیرخطی، تعمیم داده شده است. این روش همیشه بهترین روش برای یک مسأله داده شده نیست، ولی سادگی صوری و سرعت زیاد آن اغلب موجب می‌شود که اولین روشی باشد که برای حل یک مسأله غیرخطی به کار می‌رود.



شکل ۲.۲ روش نیوتن

۱. یعنی محل تقاطع خط مماس با محور x ها

به‌عنوان یک روش دیگر دستیابی به (۱.۲.۲)، از بسط سری تیلر استفاده می‌کنیم. تابع $f(x)$ را حول x_n بسط می‌دهیم،

$$f(x) = f(x_n) + (x - x_n)f'(x_n) + \frac{(x - x_n)^2}{2} f''(\xi)$$

که ξ بین x و x_n قرار دارد. x را مساوی α قرار می‌دهیم و با استفاده از $f(\alpha) = 0$ این عبارت را نسبت به α حل می‌کنیم تا به دست آوریم

$$\alpha = x_n - \frac{f(x_n)}{f'(x_n)} - \frac{(\alpha - x_n)^2}{2} \cdot \frac{f''(\xi_n)}{f'(x_n)}$$

که ξ_n بین x_n و α واقع است. می‌توانیم جملهٔ خطا (آخرین جمله) را حذف کنیم تا یک تقریب بهتر از x_n برای α به دست آوریم، و ما این تقریب را به صورت x_{n+1} در (۱.۲.۲) داریم. در این صورت

$$\alpha - x_{n+1} = -(\alpha - x_n)^2 \cdot \frac{f''(\xi_n)}{2f'(x_n)} \quad n \geq 0 \quad (2.2.2)$$

از این فرمول استفاده می‌کنیم تا نشان دهیم که روش نیوتن همگرایی مرتبهٔ دوم، $p = 2$ در (۱۳.۰.۲)، دارد.

مثال باز هم معادلهٔ

$$f(x) \equiv x^6 - x - 1 = 0$$

را برای پیدا کردن بزرگترین ریشه حل می‌کنیم. روش نیوتن (۱.۲.۲) به‌کار برده شده و نتایج در جدول ۲.۲ نشان داده شده‌اند. محاسبات با حساب تقریبی با ۱۶ رقم ممیز شناور انجام شده

جدول ۲.۲ مثال روش نیوتن

n	x_n	$f(x_n)$	$\alpha - x_n$	$x_{n+1} - x_n$
۰	۲.۰	۶۱.۰	-۸.۶۵۳E-۱	
۱	۱.۶۸۰۶۲۸۲۷۳	۱۹.۸۵	-۵.۴۵۹E-۱	-۰.۴۹۹E-۱
۲	۱.۴۳۰۷۳۸۹۸۹	۶.۱۴۷	-۲.۹۶۰E-۱	-۱.۷۵۸E-۱
۳	۱.۲۵۴۹۷۰۹۵۷	۱.۶۵۲	-۱.۲۰۲E-۱	-۹.۳۴۳E-۲
۴	۱.۱۶۱۵۳۸۴۳۳	۲.۹۴۳E-۱	-۲.۶۸۱E-۲	-۲.۵۱۹E-۲
۵	۱.۱۳۶۳۵۳۲۷۴	۱.۶۸۳E-۲	-۱.۶۲۹E-۳	-۱.۶۲۳E-۳
۶	۱.۱۳۴۷۳۰۵۲۸	۶.۵۷۴E-۵	-۶.۳۹۰E-۶	-۶.۳۹۰E-۶
۷	۱.۱۳۴۷۲۴۱۳۹	۱.۰۱۵E-۹	-۹.۸۷۰E-۱۱	-۹.۸۷۰E-۱۱

و جدول بارستها ازگردکردن این مقادیر دقیقتر به دست آمده است. آخرین ستون، $x_{n+1} - x_n$ برآوردی برای $\alpha - x_n$ است؛ این مطلب بعداً در این بخش مورد بحث واقع شده است.

هرگاه یک بارست تا حدی به α نزدیک باشد، روش نیوتن به سرعت همگرا می شود. این موضوع در بارستهای x_4, x_5, x_6, x_7 نمایش داده شده است. بارستهای x_1, x_2, x_3 همگرایی اولیه آهسته را نشان می دهند که با حدس اولیه ضعیف x_0 امکان دارد. اگر حدس اولیه $x_0 = 1$ انتخاب شده بود، x_4 تا x_7 رقم با معنی و x_5 تا x_4 رقم با معنی دقیق می شدند. این نتایج اگر با نتایجی که از روش نیم سازی در جدول ۱.۲ داده شده مقایسه شوند، سرعت زیادتر روش نیوتن بلافاصله آشکار می شود.

تحلیل همگرایی. یک قضیه همگرایی داده خواهد شد که سرعت و همچنین بازه ای را که از آن حدسهای اولیه را می توان انتخاب کرد نشان می دهد.

قضیه ۱.۲ فرض کنید $f(x)$ ، $f'(x)$ و $f''(x)$ به ازای جمیع مقادیر x در یک همسایگی α پیوسته باشند و فرض می کنیم $f(\alpha) = 0$ ، $f'(\alpha) \neq 0$. اگر x_0 به اندازه کافی نزدیک به α انتخاب شود، بارستهای x_n ، $n \geq 0$ در (۱.۲.۲) به سمت α همگرا می شوند. به علاوه

$$\lim_{n \rightarrow \infty} \frac{\alpha - x_{n+1}}{(\alpha - x_n)^2} = -\frac{f''(\alpha)}{2f'(\alpha)} \quad (3.2.2)$$

که اثبات می کند بارستها همگرایی مرتبه $p = 2$ دارند.

برهان یک بازه به اندازه کافی کوچک $I = [\alpha - \varepsilon, \alpha + \varepsilon]$ را که در آن $f'(x) \neq 0$ ، برمی گزینیم [با توجه به پیوستگی $f'(x)$ ، این بازه وجود دارد]، و سپس فرض می کنیم

$$M = \frac{\max_{x \in I} |f''(x)|}{2 \min_{x \in I} |f'(x)|}$$

با توجه به (۲.۲.۲) داریم

$$|\alpha - x_1| \leq M |\alpha - x_0|^2$$

$$M |\alpha - x_1| \leq (M |\alpha - x_0|)^2$$

می گیریم $|\alpha - x_0| \leq \varepsilon$ و $|\alpha - x_0| < 1$ و $M |\alpha - x_0| < 1$ پس $M |\alpha - x_1| < 1$ و $M |\alpha - x_1| \leq M |\alpha - x_0|$ که در نتیجه $|\alpha - x_1| \leq \varepsilon$. می توانیم از راه استقرا

همین استدلال را برای x_1, x_2, \dots به‌کار برده نشان دهیم که به‌ازای جمیع مقادیر $n \geq 1$ روابط

$$|\alpha - x_n| \leq \varepsilon \quad \text{و} \quad |\alpha - x_n| < 1 \quad \text{برقرارند.}$$

برای نشان دادن همگرایی، (۲.۲.۲) را به‌کار می‌بریم تا به‌دست آوریم

$$|\alpha - x_{n+1}| \leq M |\alpha - x_n|^2$$

$$M |\alpha - x_{n+1}| \leq (M |\alpha - x_n|)^2 \quad (۴.۲.۲)$$

با استفاده از روش استقراء

$$M |\alpha - x_n| \leq (M |\alpha - x_0|)^{2^n}$$

$$|\alpha - x_n| \leq \frac{1}{M} (M |\alpha - x_0|)^{2^n} \quad (۵.۲.۲)$$

برقراری رابطه $|\alpha - x_n| < 1$ نشان می‌دهد که وقتی $n \rightarrow \infty$ ، آنگاه $x_n \rightarrow \alpha$.

در فرمول (۲.۲.۲) نقطه مجهول ξ_n بین x_n و α واقع است و در نتیجه وقتی $n \rightarrow \infty$

آنگاه $\xi_n \rightarrow \alpha$. بنابراین

$$\lim_{n \rightarrow \infty} \frac{\alpha - x_{n+1}}{(\alpha - x_n)^2} = -\lim_{n \rightarrow \infty} \frac{f''(\xi_n)}{2f'(x_n)} = \frac{-f''(\alpha)}{2f'(\alpha)}$$

ستون خطا در جدول ۲.۲ را می‌توان برای نمایش (۳.۲.۲) به‌کار برد. به‌ویژه در این مثال:

$$-\frac{f''(\alpha)}{\alpha f'(\alpha)} = -۲,۴۱۷ \quad \frac{\alpha - x_6}{(\alpha - x_5)^2} = -۲,۴۱$$

گیریم M معرّف حد طرف راست (۲.۲.۲) باشد. پس اگر x_n نزدیک به α باشد، (۲.۲.۲)

نتیجه می‌دهد که

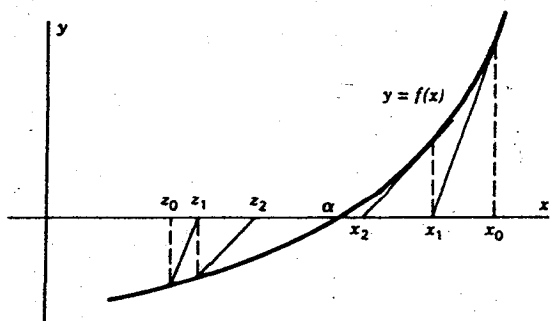
$$M(\alpha - x_{n+1}) \doteq [M(\alpha - x_n)]^2$$

برای اینکه x_n به α همگرا شود، این عبارت می‌گوید که احتمالاً باید داشته‌باشیم

$$|\alpha - x_n| < \frac{1}{M} \quad (۶.۲.۲)$$

بنابراین M میزانی است برای آنکه بدانیم چه اندازه x_n باید به α نزدیک باشد تا همگرایی به α

تضمین شود. چند مثال با مقادیر بزرگ M در مسائل آخر این فصل داده شده‌اند.



شکل ۳.۲ روش نیوتن - فوریه

یک راه دیگر برای تحلیل خطای روش نیوتن استفاده از ترسیم و قضیه زیر است. گیریم $f(x)$ روی بازه $[a, b]$ که شامل α است دوبار پیوسته مشتقپذیر باشد. به علاوه فرض می‌کنیم $f(b) > 0, f(a) < 0$ و همچنین

$$f'(x) > 0 \quad f''(x) > 0 \quad \text{برای } a \leq x \leq b \quad (۷.۲.۲)$$

در این صورت $f(x)$ در $[a, b]$ اکیداً صعودی است و یک ریشه یکتای α در $[a, b]$ وجود دارد. همچنین به ازای $a \leq x < \alpha, f(x) < 0$ ، و به ازای $\alpha < x \leq b, f(x) > 0$.

گیریم $x_0 = b$ و بارستهای نیوتن x_n را به صورت (۱-۲-۲) تعریف می‌کنیم. سپس یک دنباله جدید از بارستها را به شکل زیر تعریف می‌کنیم.

$$z_{n+1} = z_n - \frac{f(z_n)}{f'(x_n)} \quad n \geq 0 \quad (۸.۲.۲)$$

با $z_0 = a$ بارستهای حاصله در شکل ۳.۲ نمایش داده شده‌اند. با استفاده از $\{z_n\}$ کرانه‌های بالا و پایین عالی برای α به دست می‌آوریم. استفاده از (۸.۲.۲) با روش نیوتن، روش نیوتن-فوریه نامیده می‌شود.

قضیه ۲.۲ مانند گذشته، فرض می‌کنیم $f(x)$ بر بازه $[a, b]$ دوبار پیوسته مشتقپذیر باشد و $f(b) > 0, f(a) < 0$ و شرطهای (۷-۲-۲) برقرار باشند. در این صورت بارستهای x_n اکیداً به α کاهش می‌یابند و بارستهای z_n اکیداً به α افزایش پیدا می‌کنند. به علاوه،

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - z_{n+1}}{(x_n - z_n)^2} = \frac{f''(\alpha)}{2f'(\alpha)} \quad (۹.۲.۲)$$

که نشان می‌دهد فاصله بین x_n و z_n با افزایش n کاهش درجه ۲ دارد.

برهان ابتدا نشان می‌دهیم که

$$z_0 < z_1 < \alpha < x_1 < x_0. \quad (۱۰.۲.۲)$$

با توجه به تعاریف (۱.۲.۲) و (۸.۲.۲)

$$x_1 - x_0 = \frac{-f(x_0)}{f'(x_0)} < 0$$

$$z_1 - z_0 = \frac{-f(z_0)}{f'(x_0)} > 0$$

از فرمول خطای

$$\alpha - x_1 = -(\alpha - x_0)^2 \frac{f''(\xi_0)}{2f'(x_0)} < 0$$

بالاخره

$$\begin{aligned} \alpha - z_1 &= \alpha - z_0 + \frac{f(z_0)}{f'(x_0)} = \alpha - z_0 + \frac{f(z_0) - f(\alpha)}{f'(x_0)} \\ &= \alpha - z_0 - \frac{f'(\zeta_0)(\alpha - z_0)}{f'(x_0)} \quad z_0 < \zeta_0 < \alpha \text{ چون } \\ &= (\alpha - z_0) \left[\frac{f'(x_0) - f'(\zeta_0)}{f'(x_0)} \right] > 0 \end{aligned}$$

زیرا $f'(x)$ در $[a, b]$ یک تابع صعودی است. از ترکیب این نتایج اثبات (۱۰.۲.۲) حاصل می‌شود. این برهان را می‌توان به‌طور استقرایی تکرار کرد تا ثابت شود که

$$z_n < z_{n+1} < \alpha < x_{n+1} < x_n \quad n \geq 0 \quad (۱۱.۲.۲)$$

دنباله $\{x_n\}$ توسط α از پایین کراندار است، پس یک اینفیمم \bar{x} دارد؛ همچنین دنباله $\{z_n\}$

یک سوپرمم \bar{z} دارد:

$$\lim_{n \rightarrow \infty} x_n = \bar{x} \geq \alpha \quad \lim_{n \rightarrow \infty} z_n = \bar{z} \leq \alpha$$

با حد گرفتن در (۱.۲.۲) و (۸.۲.۲) به‌دست می‌آوریم

$$\bar{x} = \bar{x} - \frac{f(\bar{x})}{f'(\bar{x})} \quad \bar{z} = \bar{z} - \frac{f(\bar{z})}{f'(\bar{z})}$$

که ما را به $f(\bar{x}) = 0 = f(\bar{z})$ می‌رساند. چون α ریشه یکتای $f(x)$ در $[a, b]$ است، این ثابت می‌کند که $\{x_n\}$ و $\{z_n\}$ به α همگرا هستند.

برهان (۹.۲.۲) پیچیده‌تر است و خوانندگان را به اثر آستروفسکی (۱۹۷۳، ص ۷۰) ارجاع می‌دهیم. با توجه به قضیه ۱.۲ و فرمول (۳.۲.۲)، دنباله $\{x_n\}$ همگرایی درجه ۲ به α است. قضیه (۹.۲.۲) نشان می‌دهد که

$$|\alpha - x_n| \leq |z_n - x_n|$$

یک کران خطاست که کاهش درجه ۲ دارد.

فرضهای قضیه ۲.۲ بدل می‌شود به اینکه $f(x)$ دو بار پیوسته مشتقپذیر در یک همسایگی α باشد و

$$f'(\alpha)f''(\alpha) \neq 0 \quad (۱۲.۲.۲)$$

پس، یک بازه $[a, b]$ حول α با $f'(x)$ و $f''(x)$ مخالف صفر در این بازه وجود دارد. لذا مسأله ریشه‌یابی $f(x) = 0$ در (۷-۲-۲) صدق می‌کند یا به آسانی به یک مسأله هم‌ارزی که در آن صدق می‌کند بدل می‌شود. برای مثال، اگر $f'(\alpha) < 0$ و $f''(\alpha) > 0$ ، مسأله ریشه‌یابی $g(x) = 0$ را با شرط $g(x) \equiv f(-x)$ در نظر می‌گیریم. ریشه g برابر $-\alpha$ خواهد بود و $g(x)$ در شرایط (۷.۲.۲) در بازه‌ای حول $-\alpha$ صدق می‌کند. مثال عددی قضیه ۲.۲ به بخش مسائل واگذار خواهد شد.

برآورد خطا. روند قبلی، کرانهای بالا و پایین ریشه را به دست می‌دهد با فاصله $x_n - z_n$ که کاهش درجه ۲ دارد. ولی در بیشتر کاربردها، روش نیوتن به تنهایی، بدون (۸.۲.۲)، به کار برده می‌شود. در آنجا، از روش زیر استفاده می‌کنیم.

با استفاده از قضیه مقدار میانگین،

$$f(x_n) = f(x_n) - f(\alpha) = f'(\xi_n)(x_n - \alpha)$$

$$\alpha - x_n = \frac{-f(x_n)}{f'(\xi_n)}$$

که ξ_n بین x_n و α است. اگر $f'(x)$ در بین x_n و α به سرعت تغییر نکند، داریم $f'(\xi_n) \doteq f'(x_n)$ ، و

$$\alpha - x_n \doteq \frac{-f(x_n)}{f'(x_n)} = x_{n+1} - x_n$$

که تساوی آخری از تعریف روش نیوتن به دست می‌آید. برای روش نیوتن برآورد خطای استاندارد عبارت است از

$$\alpha - x_n \doteq x_{n+1} - x_n \quad (۱۳.۲.۲)$$

و این نکته در جدول ۲.۲ نشان داده شده است. برای خطای نسبی می‌نویسیم

$$\frac{\alpha - x_n}{\alpha} \doteq \frac{x_{n+1} - x_n}{x_{n+1}}$$

الگوریتم نیوتن. با به‌کارگیری فرمول (۱.۲.۲) و برآورد خطای (۱۳.۲.۲)، الگوریتم زیر را می‌دهیم.

الگوریتم نیوتن (f, df, x_0, ε , ریشه، $itmax$)

۱. توجه: df تابع مشتق $f'(x)$ است، $itmax$ ماکسیمم تعداد بارستهاست که بایستی محاسبه شود و ier یک اعلام خطا برای استفاده کننده است.

$$itnum := ۱.۲$$

$$denom := df(x_0).۳$$

۴. اگر $denom = 0$ آنگاه $ier := ۲$ و خارج شوید.

$$x_1 := x_0 - f(x_0)/denom.۵$$

۶. اگر $|x_1 - x_0| \leq \varepsilon$ آنگاه قرار دهید $ier := 0$ ، ریشه، و خارج شوید.

۷. اگر $itnum = itmax$ قرار دهید $ier := ۱$ و خارج شوید.

۸. در غیر این صورت $itnum := itnum + ۱$ ، $x_0 := x_1$ ، و به گام ۳ برگردید.

همانند الگوریتم نیم‌ساز قبلی، هیچ محدودیتی برای حساب رایانه در نظر گرفته نشده است، مع‌هذا یک برنامه عملی برای انجام آن مورد نیاز است. همچنین روش نیوتن تضمینی برای همگرایی نمی‌کند، و لذا یک آزمون روی تعداد بارستها (گام ۷) لازم است.

وقتی روش نیوتن همگرا باشد، معمولاً این همگرایی خیلی سریع خواهد بود، که خود مزیتی است نسبت به روش نیم‌سازی. ولی ممکن است همگرا نباشد. یک منبع دیگر اشکال در بعضی موارد لزوم دانستن $f'(x)$ به‌طور صریح است. در بعضی مسائل ریشه‌یابی این امکان‌پذیر نیست. در روش بخش بعد، این نقص، به بهای اندک کاهش سرعت همگرایی، برطرف خواهد شد.

۳.۲ روش خط قاطع

همانند روش نیوتن، نمودار $y = f(x)$ در مجاورت ریشه α با یک خط مستقیم تقریب زده می‌شود. در این حالت، فرض می‌کنیم که x_0 و x_1 دو برآورد اولیه از ریشه α باشند. نمودار $y = f(x)$ را با خط قاطعی که توسط $(x_0, f(x_0))$ و $(x_1, f(x_1))$ تعیین می‌شود تقریب می‌زنیم. ریشه آن را با x_2 نشان می‌دهیم که امیدواریم یک تقریب بهتری برای α باشد. این مطلب در شکل ۴.۲ نمایش داده شده است. با استفاده از فرمول شیب خط قاطع داریم

$$\frac{f(x_1) - f(x_0)}{x_1 - x_0} = \frac{f(x_1) - 0}{x_1 - x_2}$$

از حل آن برحسب x_2 خواهیم داشت:

$$x_2 = x_1 - f(x_1) \cdot \frac{x_1 - x_0}{f(x_1) - f(x_0)}$$

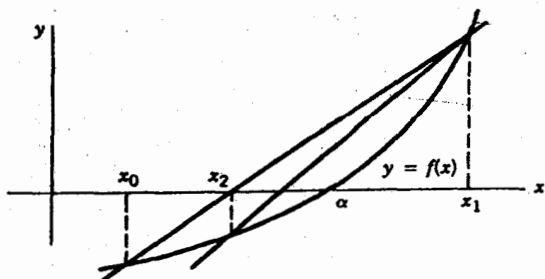
با استفاده از x_1 و x_2 این فرآیند را تکرار می‌کنیم تا x_3 و غیره به دست آیند. فرمول عمومی بر این پایه چنین است

$$x_{n+1} = x_n - f(x_n) \cdot \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} \quad n \geq 1 \quad (۱.۳.۲)$$

این روش خط قاطع است. همانند روش نیوتن همگرایی این روش تضمین نمی‌شود، ولی وقتی همگرا شود، سرعت آن معمولاً بیشتر از سرعت روش نیم‌سازی است.

مثال باز پیدا کردن بزرگترین ریشه

$$f(x) \equiv x^6 - x - 1 = 0$$



شکل ۴.۲ روش خط قاطع

جدول ۳.۲ مثال روش خط قاطع

n	x_n	$f(x_n)$	$\alpha - x_n$	$x_n - x_{n-1}$
۰	۲٫۰	۶۱٫۰	$۸٫۶۵E - ۱$	
۱	۱٫۰	-۱٫۰	$۱٫۳۵E - ۱$	$۱٫۶۱E - ۲$
۲	۱٫۰۱۶۱۲۹۰۳۲	$-۹٫۱۵۴E - ۱$	$۱٫۱۹E - ۱$	$۱٫۷۴E - ۱$
۳	۱٫۱۹۰۵۷۷۷۶۹	$۶٫۵۷۵E - ۱$	$-۵٫۵۹E - ۲$	$-۷٫۲۹E - ۲$
۴	۱٫۱۱۷۶۵۵۸۳۱	$-۱٫۶۸۵E - ۱$	$-۱٫۷۱E - ۲$	$۱٫۴۹E - ۲$
۵	۱٫۱۳۲۵۳۱۵۵۰	$-۲٫۲۴۴E - ۲$	$۲٫۱۹E - ۳$	$۲٫۲۹E - ۳$
۶	۱٫۱۳۴۸۱۶۸۰۸	$۹٫۵۳۶E - ۴$	$-۹٫۲۷E - ۵$	$-۹٫۳۲E - ۵$
۷	۱٫۱۳۴۷۲۳۶۴۶	$-۵٫۰۶۶E - ۶$	$۴٫۹۲E - ۷$	$۴٫۹۲E - ۷$
۹	۱٫۱۳۴۷۲۴۱۳۸	$-۱٫۱۳۵E - ۹$	$۱٫۱۰E - ۱۰$	$۱٫۱۰E - ۱۰$

را در نظر می‌گیریم. روش خط قاطع (۱.۳.۲) به‌کار برده شده و بارست ادامه یافته‌است تا جایی که تفاضل‌های متوالی $x_n - x_{n-1}$ به اندازه کافی کوچک شده‌اند. نتایج عددی در جدول ۳.۲ داده شده‌اند. محاسبات با یک ماشین دودویی تقریباً با دقت ۱۶ رقم اعشاری انجام گرفته است و نتایج جدول از گردکردن نتایج رایانه به‌دست آمده‌اند.

همگرایی با افزایش n به‌طور فزاینده‌ای سریع می‌شود. یک راه اندازه‌گیری آن محاسبه نسبت

زیر است

$$\frac{\alpha - x_{n+1}}{\alpha - x_n} \quad n \geq 0$$

برای روشهای خطی، این عبارت وقتی که x_n به α همگرا می‌شود، معمولاً ثابت است. ولی در این مثال، این نسبتها با افزایش n کوچکتر می‌شوند. یک توضیح شهودی این است که خط مستقیمی که نقاط $(x_{n-1}, f(x_{n-1}))$ و $(x_n, f(x_n))$ را به هم وصل می‌کند، در مجاورت $x = \alpha$ به تدریج تقریب بسیار دقیقی برای نمودار $y = f(x)$ می‌شود، و در نتیجه ریشه x_{n+1} خط مستقیم، یک برآورد بسیار بهبودیافته‌ای برای α خواهد شد. همچنین توجه کنید که بارستهای x_n با ازدیاد n ، به نحوی ظاهراً تصادفی نسبت به α بالا و پایین می‌روند. توضیح این مطلب از فرمول خطای (۳.۳.۲) داده می‌شود که در پایین آمده است.

تحلیل خطا. اگر طرفین تساوی (۱.۳.۲) را در -۱ ضرب و سپس α را به طرفین آن اضافه کنیم، به دست می‌آوریم

$$\alpha - x_{n+1} = \alpha - x_n + f(x_n) \cdot \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})}$$

با اندکی دستکاری جبری در طرف دست راست فرمول زیر به دست می‌آید

$$\alpha - x_{n+1} = -(\alpha - x_{n-1})(\alpha - x_n) \frac{f[x_{n-1}, x_n, \alpha]}{f[x_{n-1}, x_n]} \quad (۲.۳.۲)$$

کمیت‌های $f[x_{n-1}, x_n]$ و $f[x_{n-1}, x_n, \alpha]$ تفاضلات منقسم اول و دوم نیوتن هستند که در فصل ۱ تعریف شدند. خواننده بایستی با قراردادن این مقادیر از (۱۳.۱.۱) و سپس ساده کردن آن، صحت (۲.۳.۲) را کنترل کند. با استفاده از (۱۴.۱.۱)، فرمول (۲.۳.۲) چنین خواهد شد

$$\alpha - x_{n+1} = -(\alpha - x_{n-1})(\alpha - x_n) \cdot \frac{f''(\zeta_n)}{2f'(\xi_n)} \quad (۳.۳.۲)$$

که ξ_n بین x_n و x_{n-1} قرار دارد و ζ_n بین x_{n-1} و α یا بین x_n و α واقع است. با استفاده از این فرمول خطا، می‌توانیم همگرایی روش خط قاطع را امتحان کنیم.

قضیه ۳.۲ فرض می‌کنیم $f(x)$ ، $f'(x)$ و $f''(x)$ به ازای جميع مقادیر x در بازه‌ای شامل α پیوسته باشند، و $f'(\alpha) \neq 0$. در این صورت اگر حدسه‌های اولیه x_0 و x_1 به اندازه کافی نزدیک به α انتخاب شوند، بارسته‌های x_n از (۱-۳-۲) به α همگرا می‌شوند. مرتبه همگرایی چنین خواهد شد: $p = (1 + \sqrt{5})/2 \approx 1,62$.

برهان در همسایگی $I = [\alpha - \varepsilon, \alpha + \varepsilon]$ با مقداری از $\varepsilon > 0$ ، در همه جا روی I ، $f'(x)$ مخالف صفر است. پس تعریف می‌کنیم

$$M = \frac{\max_{x \in I} |f''(x)|}{2 \min_{x \in I} |f'(x)|}$$

در این صورت به ازای جميع مقادیر x_0 و x_1 متعلق به $[\alpha - \varepsilon, \alpha + \varepsilon]$ ، با استفاده از (۳.۳.۲)

$$|e_2| \leq |e_1| \cdot |e_0| \cdot M,$$

$$M |e_2| \leq M |e_1| \cdot M |e_0|$$

به علاوه فرض می‌کنیم که x_0 و x_1 به گونه‌ای انتخاب شده باشند که

$$\delta \equiv \text{Max}\{|M| e_0|, |M| e_1|\} < 1 \quad (۴.۳.۲)$$

در این صورت $|M| e_2| < 1$ زیرا

$$|M| e_2| \leq \delta^2$$

همچنین $|M| e_2| \leq \delta^2 < \delta$ ایجاب می‌کند که

$$|e_2| < \frac{\delta}{M} = \text{Max}\{|e_1|, |e_0|\} \leq \varepsilon$$

و بنابراین $x_2 \in [\alpha - \varepsilon, \alpha + \varepsilon]$ می‌توانیم این استدلال را به‌طور استقرایی به‌کار ببریم تا نشان دهیم که $x_n \in [\alpha - \varepsilon, \alpha + \varepsilon]$ و به ازای $n \geq r$ $|M| e_n| \leq \delta$ برای اثبات همگرایی و به دست آوردن مرتبه همگرایی، استفاده از (۳.۳.۲) را ادامه می‌دهیم تا به دست آوریم

$$|M| e_3| \leq |M| e_2| \cdot |M| e_1| \leq \delta^2 \cdot \delta = \delta^3$$

$$|M| e_4| \leq |M| e_3| \cdot |M| e_2| \leq \delta^5$$

برای

$$|M| e_n| \leq \delta^{q_n} \quad (۵.۳.۲)$$

$$|M| e_{n+1}| \leq |M| e_n| \cdot |M| e_{n-1}| \leq \delta^{q_n + q_{n-1}} = \delta^{q_{n+1}}$$

بنابراین

$$q_{n+1} = q_n + q_{n-1} \quad n \geq 1 \quad (۶.۳.۲)$$

با $q_0 = q_1 = 1$. این دنباله فیبوناچی اعداد است و یک فرمول صریح برای آن می‌توان داد:

$$q_n = \frac{1}{\sqrt{5}} [r_0^{n+1} - r_1^{n+1}] \quad n \geq 0 \quad (۷.۳.۲)$$

$$r_0 = \frac{1 + \sqrt{5}}{2} \doteq ۱.۶۱۸ \quad r_1 = \frac{1 - \sqrt{5}}{2} \doteq -۰.۶۱۸$$

بنابراین

$$q_n \doteq \frac{1}{\sqrt{5}} (۱.۶۱۸)^{n+1} \quad \text{به ازای مقدار بزرگ } n \quad (۸.۳.۲)$$

برای مثال، $q_0 = 8$ ، و فرمول (۸.۳.۲) عدد ۲۵×۸ را می‌دهد. به فرمول (۵.۳.۲) برمی‌گردیم،

کران خطای

$$|e_n| \leq \frac{1}{M} \delta^{q_n} \quad n \geq 0 \quad (۹.۳.۲)$$

را به دست می‌آوریم که q_n به وسیله (۷.۳.۲) داده شده است. چون وقتی $n \rightarrow \infty$ ، $q_n \rightarrow \infty$ خواهیم داشت $x_n \rightarrow \alpha$.

با یک استخراج دقیقتر، در واقع عملاً می‌توانیم نشان دهیم که مرتبه همگرایی چنین است: $p = (1 + \sqrt{5})/2$. برای ساده کردن بیان به جای آن نشان می‌دهیم که این نرخی است که کران (۹.۳.۲) با آن کاهش می‌یابد. گیریم B_n کران بالای (۹.۳.۲) را نشان دهد. پس

$$\frac{B_{n+1}}{B_n^{r_0}} = \frac{\frac{1}{M} \delta^{q_{n+1}}}{\left[\frac{1}{M}\right]^{r_0} \cdot \delta^{r_0 q_n}} = M^{r_0-1} \delta^{q_{n+1}-r_0 q_n} \leq \delta^{-1} M^{r_0-1} \equiv c$$

زیرا $q_{n+1} - r_0 q_n = r_1^{n+1} > -1$ پس

$$B_{n+1} \leq c B_n^{r_0}$$

که مستلزم داشتن مرتبه همگرایی $p = r_0 = (1 + \sqrt{5})/2$ است. یک نتیجه مشابه برای خطاهای واقعی e_n برقرار است؛ به علاوه

$$\lim_{n \rightarrow \infty} \frac{|e_{n+1}|}{|e_n|^{r_0}} = \left| \frac{f''(\alpha)}{2f'(\alpha)} \right|^{(\sqrt{5}-1)/2} \quad (۱۰.۳.۲)$$

با استفاده از فرمول خطای (۳.۳.۲) می‌توان رفتار نوسانی بارستهای x_n حول α در آخرین مثال را توضیح داد. برای x_n و x_{n-1} نزدیک به α (۳.۳.۲) ایجاب می‌کند که داشته باشیم

$$\alpha - x_{n+1} \doteq -(\alpha - x_n)(\alpha - x_{n-1}) \cdot \frac{f''(\alpha)}{2f'(\alpha)} \quad (۱۱.۳.۲)$$

علامت $\alpha - x_{n+1}$ از روی علامت دو خطای قبلی، همراه با علامت $f''(\alpha)/f'(\alpha)$ مشخص می‌شود.

شرط (۴.۳.۲) به ما اطلاع می‌دهد که مقادیر اولیه x_1 و x_0 چقدر باید به α نزدیک باشند تا همگرایی داشته باشیم. اگر کمیت M خیلی بزرگ باشد یا بالاخص، اگر

$$\left| \frac{f''(\alpha)}{2f'(\alpha)} \right|$$

خیلی بزرگ باشد، $\alpha - x_n$ و $\alpha - x_{n+1}$ باید به همین قیاس کوچکتر باشند. همگرایی ممکن است بدون (۴.۳.۲) اتفاق افتد، ولی احتمالاً چنین حالتی کاملاً اتفاقی است. به‌عنوان یک آزمون خطا، همین برآورد خطای (۱۳.۲.۲) را که برای روش نیوتن استفاده شده به‌کار می‌بریم، یعنی

$$\alpha - x_n \doteq x_{n+1} - x_n$$

استفاده از آن در جدول ۳.۲، در آخرین مثال نشان داده شده است. چون روش خط قاطع ممکن است همگرا نباشد، در برنامه‌های اجرای آن باید تعداد بارستها یک کران بالا داشته باشد، همانند الگوریتم نیوتن در بخش اخیر. یک مسأله دیگر که ممکن است با روش خط قاطع حل شود، محاسبه مشتق تقریبی

$$a_n = \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}} \quad (۱۲.۳.۲)$$

است که در این صورت روش خط قاطع (۱.۳.۲) چنین نوشته می‌شود

$$x_{n+1} = x_n - \frac{f(x_n)}{a_n} \quad (۱۳.۳.۲)$$

محاسبه a_n متضمن کاهش ارقام با معنی هم در صورت و هم در مخرج است. بنابراین وقتی $a_n, x_n \rightarrow \alpha$ با دقت کمتری تقریب مشتق f خواهد بود. با وجود این به‌دست آوردن بهبود در دقت x_n را ادامه می‌دهیم تا اینکه به ترازوفه $f(x)$ به ازای مقدار x نزدیک به α برسیم. در این موقع، a_n ممکن است با $f'(\alpha)$ بسیار تفاوت داشته باشد، و x_{n+1} ممکن است از ریشه به سرعت دور شود. بدین دلیل، دنیس و اشناپل (۱۹۸۳ صص ۳۱ و ۳۲) استفاده از (۱۲.۳.۲) را توصیه می‌کنند تا اینکه $x_n - x_{n-1}$ به اندازه کافی کوچک شود. سپس اینان تقریب دیگری را برای $f'(x)$ توصیه می‌کنند:

$$f'(x_n) \doteq a_n = \frac{f(x_n + h) - f(x_n)}{h}$$

با مقدار ثابت h . توصیه آنها برای h چنین است:

$$h = \sqrt{\delta} \cdot T_\alpha$$

که T_α یک تقریب غیر صفر قابل قبولی برای α ، مثلاً x_n ، است و δ واحد گردکردن رایانه است [به (۱۲-۲-۱)] نگاه کنید]. توصیه آنها استفاده از h است وقتی که $|x_n - x_{n-1}|$ کوچکتر از h باشد.

هزینه روش خط قاطع در محاسبه‌های تابعی کمی افزایش می‌یابد ولی احتمالاً نه بیش از هزینه یک یا دوبارست.

روش خط قاطع به‌عنوان یک روند ریشه‌یابی کارا و آسان در عمل، برای انواع بسیاری از مسائل کاملاً توصیه می‌شود. این مزیت را هم دارد که برعکس روش نیوتن نیازی به دانستن $f'(x)$ ندارد. در بخش ۸.۲، روش خط قاطع یک جزء مهم از یک الگوریتم دیگر ریشه‌یابی را تشکیل می‌دهد که همگرایی آن تضمین شده است.

مقایسه روش نیوتن و روش خط قاطع. روش نیوتن و روش خط قاطع خیلی بهم نزدیک‌اند. اگر از تقریب

$$f'(x_n) \doteq \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$$

در فرمول (۱.۲.۲) نیوتن استفاده کنیم، فرمول خط قاطع (۱.۳.۲) را به‌دست آوریم. شرایط همگرایی تقریباً یکی هستند [برای مثال، (۶.۲.۲) و (۴.۳.۲)] را برای شرایط خطای اولیه ببینید، و فرمولهای خطا متشابه‌اند [(۲.۲.۲) و (۳.۳.۲)] را ببینید. با این حال، دو تفاوت عمده وجود دارد. روش نیوتن در هر بارست به دو محاسبه تابعی نیاز دارد، که $f(x_n)$ و $f'(x_n)$ هستند، در حالی که روش خط قاطع فقط یک محاسبه تابعی در هر بارست می‌طلبد که $f(x_n)$ است [به شرطی که مقدار تابع $f(x_{n-1})$ مورد نیاز از آخرین بارست حفظ‌شده باشد]. روش نیوتن در هر بارست پرخرجتر است. از طرف دیگر روش نیوتن سریعتر همگرایی می‌شود [مرتبه همگرایی در روش نیوتن $p = 2$ و در روش خط قاطع $p \doteq 1.62$]، و در نتیجه برای یک دقت مطلوب نیاز به بارستهای کمتری دارد. یک تحلیل از تأثیر این دو تفاوت در روشهای خط قاطع و نیوتن در زیر داده شده است.

اکنون زمان لازم برای دستیابی به ریشه مطلوب α در یک بازه تحمل ε را در نظر می‌گیریم. برای ساده کردن تحلیل خود، فرض می‌کنیم که حدسهای اولیه به ریشه مطلوب کاملاً نزدیک باشند. تعریف می‌کنیم

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \quad n \geq 0$$

$$\bar{x}_{n+1} = \bar{x}_n - f(\bar{x}_n) \cdot \frac{\bar{x}_n - \bar{x}_{n-1}}{f(\bar{x}_n) - f(\bar{x}_{n-1})} \quad n \geq 1$$

و می‌گیریم $\bar{x}_0 = x_0$. \bar{x}_1 را بر پایه فرمول همگرایی زیر تعریف می‌کنیم. از (۳-۲-۲) و (۱۰-۳-۲)،

به ترتیب به دست می‌آوریم:

$$\begin{aligned} |\alpha - x_{n+1}| &\doteq c |\alpha - x_n|^r & n \geq 0, \quad c &= \left| \frac{f''(\alpha)}{2f'(\alpha)} \right| \\ |\alpha - \bar{x}_{n+1}| &\doteq d |\alpha - \bar{x}_n|^r & r &= \frac{1 + \sqrt{\delta}}{2}, \quad d = c^{r-1} \end{aligned}$$

برای خطا در بارستهای نیوتن به طریق استقرا

$$\begin{aligned} c |\alpha - x_{n+1}| &\doteq (c |\alpha - x_n|)^r \\ c |\alpha - x_n| &\doteq (c |\alpha - x_{n-1}|)^{r^n} \\ |\alpha - x_n| &\doteq \frac{1}{c} (c |\alpha - x_0|)^{r^n} \quad n \geq 0 \end{aligned}$$

همچنین برای بارستهای روش قاطع داریم:

$$\begin{aligned} |\alpha - \bar{x}_n| &\doteq d |\alpha - \bar{x}_{n-1}|^r \\ &\doteq d^{1+r+\dots+r^{n-1}} |\alpha - \bar{x}_0|^{r^n} \end{aligned}$$

با استفاده از فرمول (۹-۱-۱) برای یک سری متناهی هندسی، چنین به دست می‌آوریم

$$d^{1+r+\dots+r^{n-1}} = d^{(r^n-1)/(r-1)} = c^{r^n-1}$$

و بنابراین

$$|\alpha - \bar{x}_n| \doteq c^{r^n-1} |\alpha - \bar{x}_0|^{r^n} = \frac{1}{c} [c |\alpha - x_0|]^{r^n}$$

برای اینکه نابرابری $|\alpha - x_n| \leq \varepsilon$ برقرار باشد، برای بارستهای نیوتن، باید داشته باشیم

$$(c |\alpha - x_0|)^{r^n} \leq c\varepsilon$$

$$n \geq \frac{K}{\log 2} \quad K = \log \left[\frac{\log \varepsilon c}{\log c |\alpha - x_0|} \right]$$

گیریم m زمان محاسبه $f(x)$ باشد، و s زمان محاسبه $f'(x)$. در این صورت حداقل زمان لازم برای به دست آوردن دقت مورد نظر با روش نیوتن عبارت است از

$$T_N = (m + ms)n = \frac{(1+s)mK}{\log 2} \quad (۱۴.۳.۲)$$

برای روش قاطع، یک محاسبه مشابه نشان می‌دهد که $|\alpha - \bar{x}_n| \leq \varepsilon$ ، اگر

$$n \geq \frac{K}{\log r}$$

بنابراین کمترین زمان لازم برای به دست آوردن دقت مورد نظر عبارتست از

$$T_S = mn = \frac{mK}{\log r} \quad (۱۵.۳.۲)$$

برای مقایسه زمانهای روش خط قاطع و روش نیوتن داریم

$$\frac{T_S}{T_N} = \frac{\log 2}{(1+s)\log r}$$

روش خط قاطع سریعتر از روش نیوتن است اگر این کسر کوچکتر از واحد باشد،

$$\frac{T_S}{T_N} < 1$$

$$s > \frac{\log 2}{\log r} - 1 \approx 0.44 \quad (۱۶.۳.۲)$$

اگر زمان برای محاسبه $f'(x)$ بیش از ۴۴ درصد زمان لازم برای محاسبه $f(x)$ باشد، روش خط قاطع کاراتر است. در عمل، عوامل خیلی زیادتری در هزینه‌های نسبی دو روش اثر می‌گذارند، لذا باید با احتیاط از عامل ۰.۴۴ استفاده شود.

استدلال قبلی در نشان دادن اینکه سرعت ریاضی همگرایی تصویر کامل مسأله نیست، مفید است. زمان کلی محاسبات، کاربرد آسان یک الگوریتم، پایداری، و عوامل دیگر نیز موجب برتری نسبی یک الگوریتم بر دیگری می‌شود.

۴.۲ روش مولر

روش مولر برای پیدا کردن ریشه‌های حقیقی و ریشه‌های مختلط یک تابع مفید، و از نظر اجرای رایانه‌ای، معقول و ساده است. آن را به دست می‌آوریم، در همگرایی آن بحث می‌کنیم و چند مثال عددی می‌آوریم.

روش مولر تعمیم راهی است که به روش خط قاطع منجر شده است. سه نقطه x_0 ، x_1 و x_2 داده شده‌اند، چند جمله‌ای درجه دومی می‌سازیم که از سه نقطه $(x_i, f(x_i))$ ، $i = 0, 1, 2$ بگذرد؛ یکی از ریشه‌های این چند جمله‌ای به عنوان یک برآورد بهبود یافته برای ریشه α ی $f(x)$ به کار برده می‌شود.

این چند جمله‌ای درجهٔ دوم با رابطهٔ زیر داده می‌شود.

$$p(x) = f(x_2) + (x - x_2)f[x_2, x_1] + (x - x_2)(x - x_1)f[x_2, x_1, x_0] \quad (۱.۴.۲)$$

تفاضلات منقسم $f[x_2, x_1]$ و $f[x_2, x_1, x_0]$ در (۱۳.۱.۱) در فصل ۱ تعریف شده‌اند. برای کنترل برقراری

$$p(x_i) = f(x_i) \quad i = 0, 1, 2$$

تنها x_i را در (۱.۴.۲) قرار می‌دهیم و عبارت نتیجه را با استفاده از (۱۳.۱.۱) خلاصه می‌کنیم. فرمولهای دیگری برای $p(x)$ وجود دارند که در فصل ۳ داده شده‌اند، ولی شکلی که در (۱.۴.۲) نشان داده شده مناسبترین شکل برای تعریف روش مولر است. فرمول (۱.۴.۲) را شکل تفاضلات منقسم نیوتن در درون‌یابی چند جمله‌ای گویند و در بخش ۳-۲ی فصل ۳ به‌طور کلی بحث شده است. برای پیدا کردن صفرهای (۱.۴.۲) ما آن را به شکل مناسبتری بازنویسی می‌کنیم

$$y = f(x_2) + w(x - x_2) + f[x_2, x_1, x_0](x - x_2)^2$$

$$w = f[x_2, x_1] + (x_2 - x_1)f[x_2, x_1, x_0]$$

$$= f[x_2, x_1] + f[x_2, x_0] - f[x_0, x_1]$$

می‌خواهیم کوچکترین مقدار $x - x_2$ را پیدا کنیم که در معادلهٔ $y = 0$ صدق کند، بنابراین ریشه‌ای از (۱.۴.۲) را پیدا می‌کنیم که نزدیکترین مقدار به x_2 باشد. جواب چنین است

$$x - x_2 = \frac{-w \pm \sqrt{w^2 - 4f(x_2)f[x_2, x_1, x_0]}}{2f[x_2, x_1, x_0]}$$

با انتخاب علامتی که صورت را در حدممکن کوچک نماید. به عدت خطاهای کاهش ارقام بامعنی که به‌طور ضمنی در این فرمول وجود دارد، صورت را گویا می‌کنیم، تا فرمول با رستی تازه‌ای به‌دست آوریم

$$x_3 = x_2 - \frac{2f(x_2)}{w \pm \sqrt{w^2 - 4f(x_2)f[x_2, x_1, x_0]}} \quad (۲.۴.۲)$$

با انتخاب علامتی که مقدار مخرج را حداکثر نماید.

برای تعریف دنبالهٔ بارسته‌های $\{x_n; n \leq 0\}$ ، (۲.۴.۲) را بطور بازگشتی تکرار می‌کنیم.

اگر این بارسته‌ها به یک نقطهٔ α همگرا شوند و اگر $f'(\alpha) \neq 0$ ، آنگاه α یک ریشهٔ $f(x)$

خواهد بود. برای مشاهده این امر از (۱۴.۱.۱) فصل ۱ و (۲.۴.۲) استفاده می‌کنیم تا به دست آوریم

$$w \rightarrow f'(\alpha) \quad n \rightarrow \infty \quad \text{وقتی}$$

$$\alpha = \alpha - \frac{2f(\alpha)}{f'(\alpha) \pm \sqrt{[f'(\alpha)]^2 - 2f(\alpha)f''(\alpha)}}$$

که نشان می‌دهد کسر طرف راست باید صفر باشد. چون بنا به فرض $f'(\alpha) \neq 0$ ، روش انتخاب علامت در مخرج مستلزم این است که مخرج مخالف صفر باشد. پس صورت باید صفر باشد، که نشان می‌دهد $f(\alpha) = 0$. فرض $f'(\alpha) \neq 0$ می‌گوید که α یک ریشه ساده است (بخش ۷.۲ را برای بحث درباره ریشه‌های ساده و چندگانه ببینید).

با استدلالی مشابه استدلالی که در روش خط قاطع به کار برده شد، می‌توان نشان داد که

$$\lim_{n \rightarrow \infty} \frac{|\alpha - x_{n+1}|}{|\alpha - x_n|^p} = \left| \frac{f^{(p)}(\alpha)}{p! f'(\alpha)} \right|^{(p-1)/2} \quad p \neq 1, 84 \quad (3.4.2)$$

به شرطی که $f(x)$ در یک همسایگی α سه بار پیوسته مشتقپذیر باشد و $f'(\alpha) \neq 0$. مرتبه p ریشه مثبت

$$x^3 - x^2 - x - 1 = 0$$

است.

با روش خط قاطع، انتخاب مقادیر حقیقی برای x_0 و x_1 به مقدار حقیقی x_2 می‌انجامد. ولی با روش مولر، انتخاب مقادیر حقیقی x_0, x_1, x_2 ، می‌تواند و باید به ریشه‌های مختلط $f(x)$ منتهی شود. این یک جنبه مهم روش مولر است و یک دلیل برای استفاده از آن.

مثالهای ذیل با استفاده از یک برنامه بازگانی محاسبه شده‌اند که به خودی خود یک پیاده‌سازی روش مولر را به دست می‌دهد. بدون اینکه حدسه‌های اولیه داده شده باشند، با این برنامه ریشه‌های $f(x)$ را تقریباً با یک ترتیب افزایشی پیدا می‌کند. پس از آنکه تقریبهای z_1, z_2, \dots به عنوان ریشه‌ها پیدا می‌شوند، از تابع

$$g(x) = \frac{f(x)}{(x - z_1) \dots (x - z_r)} \quad (4.4.2)$$

برای پیدا کردن بقیه ریشه‌های $f(x)$ استفاده می‌شود. [برای بحث در خطاهای این استفاده از $g(x)$ ، به پیترو و ویلکینسن (۱۹۷۱) مراجعه کنید]. برای آنکه ریشه تقریبی z در برنامه قابل قبول باشد، باید در یکی از شرایط زیر (که استفاده کننده مشخص می‌کند) صدق کند:

جدول ۴.۲ روش مولر، مثال ۲

IT	ریشه	f(ریشه)
۹	$۱٫۱۵۷۲۲۱۱۷۳۶E - ۱$	$۵٫۹۶E - ۸$
۱۰	$۶٫۱۱۷۵۷۴۸۴۵۲E - ۱ + ۹٫۰۱E - ۲۰i$	$-۲٫۹۸E - ۷ + ۹٫۰۶E - ۱۱i$
۱۴	$۲٫۸۳۳۷۵۱۳۳۷۷E۰ - ۵٫۰۵E - ۱۷i$	$۲٫۵۵E - ۵ - ۴٫۷۸E - ۸i$ $۷٫۱۳E - ۵ + ۹٫۳۷E - ۶i$
۱۳	$۴٫۵۹۹۲۲۷۶۳۹۴E۰ - ۵٫۹۵E - ۱۵i$	$۷٫۱۳E - ۵ + ۹٫۳۷E - ۶i$
۸	$۱٫۵۱۲۶۱۰۲۶۹۸E۰ + ۲٫۹۸E - ۱۶i$	$۳٫۳۴E - ۶ - ۲٫۳۵E - ۷i$
۱۹	$۱٫۳۰۰۰۶۰۵۴۹۹eE۱ + ۹٫۰۴E - ۱۸i$	$۲٫۳۲E - ۱ + ۴٫۱۵E - ۷i$
۱۶	$۹٫۶۲۱۳۱۶۸۴۲۵E۰ - ۴٫۹۷E - ۱۷i$	$-۳٫۶۶E - ۲ - ۵٫۳۸E - ۷i$
۱۴	$۱٫۷۱۱۶۸۵۵۱۸۷E۱ - ۸٫۴۸E - ۱۷i$	$-۱٫۶۸E + ۰ + ۲٫۴۰E - ۵i$
۱۳	$۲٫۲۱۵۱۰۹۰۳۷۹E۱ + ۹٫۳۵E - ۱۸i$	$۸٫۶۱E - ۱ + ۲٫۶۰E - ۵i$
۷	$۶٫۸۴۴۵۲۵۴۵۳۱E۰ - ۳٫۴۳E - ۲۸i$	$-۴٫۴۹E - ۳ - ۱٫۲۲E - ۱۸i$
۴	$۲٫۸۴۸۷۹۶۷۲۵۱E۱ + ۵٫۷۷E - ۲۵i$	$-۶٫۳۴E + ۱ - ۲٫۹۶E - ۱۱i$
۴	$۳٫۷۰۹۹۱۲۱۰۴۴E۱ + ۲٫۸۰E - ۲۴i$	$۲٫۱۲E + ۳ + ۷٫۷۲E - ۹i$

جدول ۵.۲ روش مولر، مثال ۳

IT	ریشه	f(ریشه)
۴۱	$۲٫۹۹۸۷۵۲۶ - ۶٫۹۸E - ۴i$	$-۳٫۳۳E - ۱۱ + ۶٫۷۰E - ۱۱i$
۱۷	$۲٫۹۹۹۷۵۹۱ - ۲٫۶۸E - ۴i$	$۵٫۶۸E - ۱۴ - ۶٫۴۸E - ۱۴i$
۳۱	$۳٫۰۰۰۳۰۹۵ - ۳٫۱۷E - ۴i$	$-۳٫۴۱E - ۱۳ + ۳٫۲۲E - ۱۴i$
۱۰	$۳٫۰۰۰۳۰۴۶ + ۳٫۱۴E - ۴i$	$۳٫۹۸E - ۱۳ - ۳٫۸۳E - ۱۴i$
۶	$۵٫۹۷E - ۱۵ + ۳٫۰۰۰۰۰۰۰۰۰۰۰۰i$	$۴٫۳۸E - ۱۱ - ۱٫۱۹E - ۱۱i$
۳	$۵٫۹۷E - ۱۵ - ۳٫۰۰۰۰۰۰۰۰۰۰۰۰i$	$۴٫۳۸E - ۱۱ + ۱٫۱۹E - ۱۱i$

$$|f(x)| \leq ۱۰^{-۱۰} \cdot ۱$$

۲. دقت z تا ۸ رقم بامعنی باشد.

در جدولهای ۴.۲ و ۵.۲، ریشه‌ها به همان ترتیبی که پیدا شدند داده شده‌اند. ستون IT تعداد بارستها را که برای هر ریشه محاسبه شده‌اند می‌دهد. مثالها همگی برای یک چند جمله‌ای $f(x)$ داده شده‌اند، ولی برنامه برای حالت کلی $f(x)$ طرح‌ریزی شده و x می‌تواند مختلط باشد.

یک نظریه کلی برای روشهای بارستی تک نقطه‌ای ۸۷

مثال ۱. $f(x) = x^{20} - 1$. همهٔ 20 ریشه با دقت ده رقم بامعنی یابیشتر پیدا شده‌اند. در همهٔ حالات، ریشه تقریبی z در $10^{-10} < |f(z)|$ و به‌طور کلی بسیار کمتر صدق می‌کرده است. تعداد بارستها از ۱ تا ۱۸ تغییر کرده و دارای میانگینی برابر ۸٫۵ بوده است.

۲. $f(x) =$ چند جمله‌ای درجهٔ ۱۲ لاگیر^۱. اجزای حقیقی ریشه‌ها، همان‌گونه که نشان داده شده است، به تعداد ارقامی که گرد شده‌اند، درست بوده‌اند، ولی اجزای انگاری باید صفر می‌بودند. نتایج عددی در جدول ۴٫۲ داده شده‌اند. توجه نمایید که برای بسیاری از ریشه‌های تقریبی، $f(x)$ کاملاً بزرگ است.

۳.

$$f(x) = x^6 - 12x^5 + 63x^4 - 216x^3 + 567x^2 - 972x + 729 \\ = (x^2 + 9)(x - 3)^4$$

نتایج عددی در جدول ۵٫۲ داده شده‌اند. به عدم دقت در اولین چهارریشه توجه نمایید، که به لحاظ نوبت $f(x)$ مربوط به مکرر بودن ریشهٔ $\alpha = 3$ ذاتی است. بخش ۷٫۲ را برای یک بحث کامل در مسألهٔ محاسبهٔ ریشه‌های مکرر ببینید.

دو مثال اخیر نشان می‌دهند که چرا دو آزمون خطا لازم‌اند، و اشاره می‌کنند که چرا در این روش کار برای محاسبهٔ هر ریشه به یک تعداد حداکثر مجاز برای تعداد بارستها نیاز است. شکل (۲٫۴٫۲) روش مولر مربوط به تراوب^۲ (۱۹۶۴، صص ۲۱۰ تا ۲۱۳) است. برای بحث محاسباتی، ویتلی^۳ (۱۹۶۸) را ببینید.

۵٫۲ یک نظریه کلی برای روشهای بارستی تک نقطه‌ای

اکنون حل یک معادلهٔ $x = g(x)$ را برای پیدا کردن یک ریشهٔ α از راه بارستی

$$x_{n+1} = g(x_n) \quad n \geq 0 \quad (1.5.2)$$

در نظر می‌گیریم که در آن x_0 یک حدس اولیه برای α است. روش نیوتن برای این الگو چنین می‌شود

$$g(x) \equiv x - \frac{f(x)}{f'(x)} \quad (2.5.2)$$

هر جواب $x = g(x)$ یک نقطهٔ ثابت g خوانده می‌شود. اگرچه ما معمولاً علاقه‌مند به حل معادلهٔ $f(x) = 0$ هستیم، ولی راههای گوناگونی برای نوشتن این معادله به شکل یک مسألهٔ نقطه ثابت وجود دارد. در اینجا این فرایند بازنویسی را تنها با چند مثال نشان خواهیم داد.

جدول ۶.۲ مثالهای بارستی برای $x^2 - 3 = 0$

	حالت (۱)	حالت (۲)	حالت (۳)
n	x_n	x_n	x_n
۰	۲٫۰	۲٫۰	۲٫۰
۱	۳٫۰	۱٫۵	۱٫۷۵
۲	۹٫۰	۲٫۰	۱٫۷۳۲۱۴۳
۳	۸۷٫۰	۱٫۵	۱٫۷۳۲۰۵۱

مثال حل معادله $x^2 - a = 0$ را برای $a > 0$ در نظر می‌گیریم.

(i) $x = x^2 + x - a$ ، یا به صورت کلیتر، به ازای مقداری چون $c \neq 0$ ، $x = x + c(x^2 - a)$.

(ii) $x = a/x$

(iii)

$$x = \frac{1}{2} \left(x + \frac{a}{x} \right) \quad (۳.۵.۲)$$

یک مثال عددی با $a = 3$ ، $x_0 = 2$ و $x_\infty = \sqrt{3} = ۱٫۷۳۲۰۵۱$ می‌آوریم. با $x_0 = 2$ ، نتایج عددی برای (۳.۵.۲) در این حالتها در جدول ۶.۲ داده شده است.

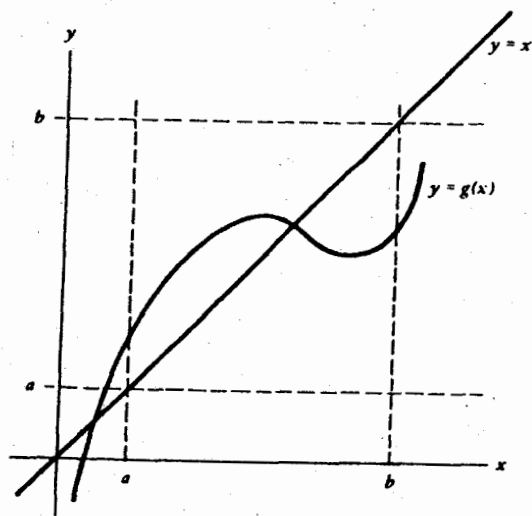
طبیعی است که سؤال شود چه چیزی موجب می‌شود که رفتار برنامه‌های بارستی مختلف به گونه‌ای باشد که در این مثال دیده می‌شود. یک نظریه کلی به دست می‌دهیم که این رفتار را توضیح دهد و به تحلیل روشهای بارستی جدید کمک نماید.

لم ۴.۲. گیریم $g(x)$ در بازه $a \leq x \leq b$ پیوسته باشد و فرض می‌کنیم برای هر $a \leq x \leq b$ داشته باشیم $a \leq g(x) \leq b$. (گوییم g ، $[a, b]$ را به خود $[a, b]$ می‌فرستد و آن را با $g([a, b]) \subset [a, b]$ نشان می‌دهیم.) در این صورت $x = g(x)$ حداقل یک جواب در $[a, b]$ دارد.

برهان تابع پیوسته $g(x) - x$ را در نظر می‌گیریم. در $x = a$ این تابع مثبت و در $x = b$ منفی است. پس طبق قضیه مقدار میانی، باید حداقل یک ریشه در بازه $[a, b]$ داشته باشد. در شکل ۵.۲، ریشه‌ها نقاط برخورد $y = x$ با $y = g(x)$ هستند. ■

لم ۵.۲. گیریم $g(x)$ روی $[a, b]$ پیوسته باشد و فرض می‌کنیم $g([a, b]) \subset [a, b]$. به علاوه فرض می‌کنیم ثابتی چون $0 < \lambda < 1$ وجود دارد که

$$|g(x) - g(y)| \leq \lambda |x - y| \quad x, y \in [a, b]$$



شکل ۵.۲ مثال لم ۴.۲

پس $x = g(x)$ یک جواب یکتای α در $[a, b]$ دارد. همچنین بارستهای

$$x_n = g(x_{n-1}) \quad n \geq 1$$

برای هر انتخابی از x_0 در $[a, b]$ به α همگراست و

$$|\alpha - x_n| \leq \frac{\lambda^n}{1 - \lambda} |x_1 - x_0| \quad (5.5.2)$$

برهان فرض می‌کنیم $x - g(x)$ دو جواب α و β در $[a, b]$ دارد. پس

$$|\alpha - \beta| = |g(\alpha) - g(\beta)| \leq \lambda |\alpha - \beta|$$

$$(1 - \lambda) |\alpha - \beta| \leq 0$$

از آنجا که $0 < \lambda < 1$ ، نتیجه می‌شود که $\alpha = \beta$. همچنین با توجه به لم قبلی می‌دانیم که حداقل یک ریشه α در $[a, b]$ موجود است.

برای آزمون همگرایی بارستهای x_n ، ابتدا توجه نمایید که آنها همگی در $[a, b]$ باقی می‌مانند. برای مشاهده این امر، توجه داشته باشید که می‌توان با استفاده از این که با استقراء ریاضی اثبات

نمود که به ازای همهٔ مقادیر n ، $x_n \in [a, b]$ برای همگرایی،

$$|\alpha - x_{n+1}| = |g(\alpha) - g(x_n)| \leq \lambda |\alpha - x_n| \quad (۶.۵.۲)$$

و با استقرا

$$|\alpha - x_n| \leq \lambda^n |\alpha - x_0| \quad n \geq 0 \quad (۷.۵.۲)$$

که وقتی $n \rightarrow \infty$ ، $\lambda^n \rightarrow 0$ ؛ بنابراین $x_n \rightarrow \alpha$.

برای اثبات کران (۵.۵.۲)، با

$$|\alpha - x_0| \leq |\alpha - x_1| + |x_1 - x_0| \leq \lambda |\alpha - x_0| + |x_1 - x_0|$$

$x_{n+1} = g(x_n) \in [a, b]$ ایجاب می‌کند

آغاز می‌کنیم که در آخرین مرحله از (۶.۵.۲) استفاده شده است. سپس آن را نسبت به $|\alpha - x_0|$ حل می‌کنیم، داریم

$$|\alpha - x_0| \leq \frac{1}{1-\lambda} |x_1 - x_0| \quad (۸.۵.۲)$$

از ترکیب این رابطه با (۷.۵.۲) اثبات کامل می‌شود.

کران (۶.۵.۲) نشان می‌دهد که دنبالهٔ $\{x_n\}$ به‌طور خطی همگراست، و بر پایهٔ تعریف (۱۳.۰.۲) نرخ همگرایی با کران λ محدود می‌شود. همچنین از برهان فوق، می‌توانیم یک کران خطا شاید دقیقتر از (۵.۵.۲) پیدا کنیم. از تکرار برهانی که به (۸.۵.۲) انجامید، به‌دست می‌آوریم

$$|\alpha - x_n| \leq \frac{1}{1-\lambda} |x_{n+1} - x_n|$$

وانگهی، استفاده از (۶.۵.۲) کران زیر را به‌دست می‌دهد

$$|\alpha - x_{n+1}| \leq \frac{\lambda}{1-\lambda} |x_{n+1} - x_n| \quad (۹.۵.۲)$$

وقتی λ قابل محاسبه باشد، این رابطه یک کران عملی در بسیاری از حالات ارائه می‌دهد. کرانه‌های دیگر خطا و تخمینها در بخش بعدی مورد بحث قرار گرفته‌اند.

اگر $g(x)$ در $[a, b]$ مشتقپذیر باشد، آنگاه برای تمام $x, y \in [a, b]$

$$g(x) - g(y) = g'(\xi)(x - y) \quad \xi \text{ بین } y, x$$

تعریف می‌کنیم

$$\lambda = \text{Max}_{a \leq x \leq b} |g'(x)|$$

پس برای تمام مقادیر $x, y \in [a, b]$

$$|g(x) - g(y)| \leq \lambda |x - y|$$

قضیه ۶.۲ فرض کنید که $g(x)$ بر $[a, b]$ به طور پیوسته مشتقپذیر باشد، و $g([a, b]) \subset [a, b]$ ، و اینکه

$$\lambda = \text{Max}_{a \leq x \leq b} |g'(x)| < 1 \quad (۱۰.۵.۲)$$

در این صورت

(i) $x = g(x)$ یک جواب یکتای α در $[a, b]$ دارد.

(ii) برای هر انتخاب x_0 در $[a, b]$ با $x_{n+1} = g(x_n)$ ، $n \geq 0$

$$\lim_{n \rightarrow \infty} x_n = \alpha$$

$$|\alpha - x_n| \leq \lambda^n |\alpha - x_0| \leq \frac{\lambda^n}{1 - \lambda} |x_1 - x_0| \quad (\text{iii})$$

$$\lim_{n \rightarrow \infty} \frac{\alpha - x_{n+1}}{\alpha - x_n} = g'(\alpha) \quad (۱۱.۵.۲)$$

برهان همه نتیجه‌ها از لم‌های قبلی به دست می‌آیند مگر برای نرخ همگرایی (۱۱.۵.۲). برای آن، رابطه زیر را به کار می‌بریم

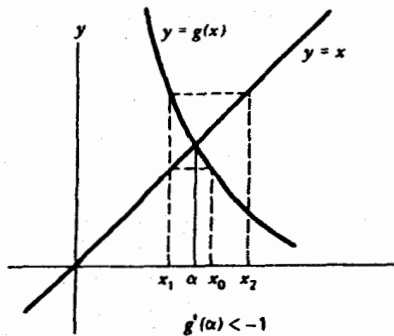
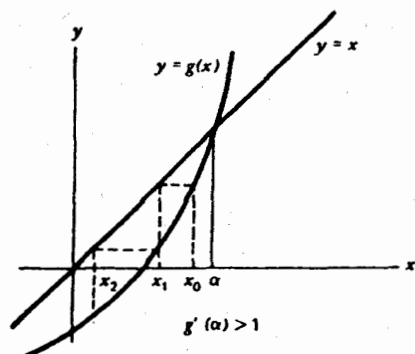
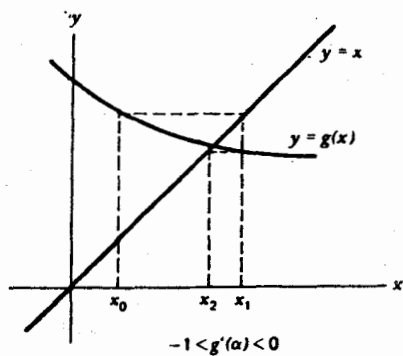
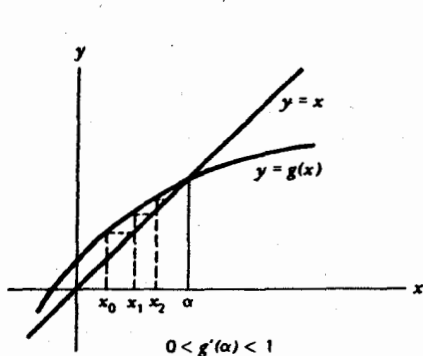
$$\alpha - x_{n+1} = g(\alpha) - g(x_n) = g'(\xi_n)(\alpha - x_n) \quad n \geq 0 \quad (۱۲.۵.۲)$$

که ξ_n نقطه مجهولی بین x_n و α است. چون $x_n \rightarrow \alpha$ ، باید داشته باشیم $\xi_n \rightarrow \alpha$ ، و بنابراین

$$\lim_{n \rightarrow \infty} \frac{\alpha - x_{n+1}}{\alpha - x_n} = \lim_{n \rightarrow \infty} g'(\xi_n) = g'(\alpha)$$

اگر $g'(\alpha) \neq 0$ ، آنگاه دنباله $\{x_n\}$ با مرتبه دقیقاً $p=1$ ، یعنی همگرایی خطی، به α همگرا می‌شود. ■

این قضیه برای دستگاه‌های m معادله غیر خطی با m مجهول تعمیم می‌یابد. کافی است x را یک عنصر \mathbf{R}^m ، $g(x)$ را تابعی از \mathbf{R}^m به \mathbf{R}^m بگیریم و به جای قدر مطلقها نرمهای برداری و ماتریسی و به جای $g'(x)$ ماتریس ژاکوبی $g'(x)$ را قرار دهیم. فرض $g([a, b]) \subset [a, b]$ باید با فرض قویتری جایگزین شود و برای انتخاب ناحیه‌ای که $[a, b]$ را تعمیم دهد باید به دقت عمل شود. لم‌ها تعمیم می‌یابند ولی اثبات آنها آسان نخواهد بود. در این مورد در بخش ۲-۱۰ بیشتر توضیح داده شده است. برای ملاحظه اهمیت فرض (۱۰.۵.۲) در مورد اندازه $g'(x)$ ، فرض کنید $|g'(\alpha)| > 1$. در این صورت اگر دنباله‌ای از بارستهای $x_{n+1} = g(x_n)$ می‌داشتیم و $\alpha = g(\alpha)$ یک ریشه بود، (۱۲.۵.۲) را می‌داشتیم. اگر x_n به اندازه کافی به α نزدیک شود، آنگاه $|g'(\xi_n)| > 1$ و خطای $|\alpha - x_{n+1}|$ از $|\alpha - x_n|$ بزرگتر می‌شود. پس اگر $|g'(\alpha)| > 1$ ، همگرایی میسر نخواهد بود. ما محاسبه بارستها را در چهار حالت با نمودار نشان داده‌ایم (شکل ۶.۲ را ملاحظه کنید). برای ساده کردن کاربرد قضیه قبلی، قضیه زیر را می‌آوریم.



شکل ۶.۲ مثالهایی از دنباله‌های همگرا و ناهمگرای $x_{n+1} = g(x_n)$

یک نظریه کلی برای روشهای بارستی تک نقطه‌ای ۹۳

قضیه ۷.۲.۲. α یک جواب $x = g(x)$ باشد و فرض می‌کنیم $g(x)$ در بازه‌ای از همسایگی α به طور پیوسته مشتقپذیر با $1 < |g'(\alpha)|$ باشد در این صورت نتایج قضیه ۶.۲ کماکان درست‌اند به شرطی که x به اندازه کافی به α نزدیک انتخاب شود.

برهان یک عدد λ که $1 < \lambda < |g'(\alpha)|$ را برقرار کند انتخاب می‌کنیم، سپس یک بازه $I = [\alpha - \varepsilon, \alpha + \varepsilon]$

$$\text{Max}_{x \in I} |g'(x)| \leq \lambda < 1$$

را در نظر می‌گیریم. داریم $g(I) \subset I$ ، زیرا $|\alpha - x| \leq \varepsilon$ ایجاب می‌کند که داشته باشیم

$$|\alpha - g(x)| = |g(\alpha) - g(x)| = |g'(\xi)| \cdot |\alpha - x| \leq \lambda |\alpha - x| \leq \varepsilon$$

■ اکنون قضیه قبلی را با استفاده از $[a, b] = [\alpha - \varepsilon, \alpha + \varepsilon]$ به کار می‌بریم.

مثال با مراجعه به مثال پیشین در این بخش، $g'(\alpha)$ را محاسبه می‌کنیم.

$$g(x) = x^2 + x - 3 \quad g'(\alpha) = g'(\sqrt{3}) = 2\sqrt{3} + 1 > 1 \quad (\text{i})$$

$$g(x) = \frac{3}{x} \quad g'(\sqrt{3}) = \frac{-3}{(\sqrt{3})^2} = -1 \quad (\text{ii})$$

$$g(x) = \frac{1}{4}\left(x + \frac{3}{x}\right) \quad g'(x) = \frac{1}{4}\left(1 - \frac{3}{x^2}\right) \quad g'(\sqrt{3}) = 0 \quad (\text{iii})$$

مثال برای $x = x + c(x^2 - 3)$ مقدار c را به گونه‌ای انتخاب می‌کنیم که همگرایی تضمین شود.

چون جواب $\alpha = \sqrt{3}$ است و چون $g'(x) = 1 + 2cx$ باید طوری باشد که

$$-1 < 1 + 2c\sqrt{3} < 1$$

برای یک نرخ خوب همگرایی، c را طوری می‌گیریم که

$$1 + 2c\sqrt{3} = 0$$

و این انتخاب

$$c = \frac{-1}{2\sqrt{3}}$$

جدول ۷.۲ مثال عددی از بارستی (۱۳.۵.۲)

n	x_n	$\alpha - x_n$	نسبت
۰	۲٫۰	$-۲٫۶۸E - ۱$	
۱	۱٫۷۵	$-۱٫۷۹E - ۲$	۰٫۰۶۶۸
۲	۱٫۷۳۴۳۷۵۰	$-۲٫۳۲E - ۳$	۰٫۱۳۰
۳	۱٫۷۳۲۳۶۰۸	$-۳٫۱۰E - ۴$	۰٫۱۳۴
۴	۱٫۷۳۲۰۹۲۳	$-۴٫۱۵E - ۵$	۰٫۱۳۴
۵	۱٫۷۳۲۰۵۶۴	$-۵٫۵۶E - ۶$	۰٫۱۳۴
۶	۱٫۷۳۲۰۵۱۶	$-۷٫۴۵E - ۷$	۰٫۱۳۴
۷	۱٫۷۳۲۰۵۰۹	$-۱٫۰۰E - ۷$	۰٫۱۳۴

را به دست می‌دهد. از $\frac{1}{4} = c$ استفاده می‌کنیم. در این صورت $g'(\sqrt{3}) = 1$ که برنامه بارستی زیر را می‌دهد

$$x_{n+1} = x_n - \frac{1}{4}(x_n^2 - 3) \quad n \geq 0 \quad (13.5.2)$$

نتایج عددی در جدول ۷.۲ داده شده‌اند. ستون نسبت‌ها مقادیر

$$\frac{\alpha - x_n}{\alpha - x_{n-1}} \quad n \geq 1$$

را به دست می‌دهد. نتایج با مقدار نظری $g'(\sqrt{3})$ به خوبی توافق دارد.

روشهای تک نقطه‌ای مرتبه بالاتر. مبحث نظریه روشهای بارستی تک نقطه‌ای را با در نظر گرفتن روشهایی که دارای مرتبه همگرایی بزرگتر از یک باشند، مثلاً روش نیوتن، کامل می‌کنیم.

قضیه ۸.۲. اگریم α یک ریشه $x = g(x)$ باشد و $g(x)$ به ازای تمام x های نزدیک به α و $p, p \geq 2$ بار پیوسته مشتقپذیر باشد. به علاوه گیریم

$$g'(\alpha) = \dots = g^{(p-1)}(\alpha) = 0 \quad (14.5.2)$$

در این صورت اگر حدس اولیه x_0 به اندازه کافی نزدیک به α انتخاب شده باشد، بارستی

$$x_{n+1} = g(x_n) \quad n \geq 0$$

دارای همگرایی مرتبه p بوده و

$$\lim_{n \rightarrow \infty} \frac{\alpha - x_{n+1}}{(\alpha - x_n)^p} = (-1)^{p-1} \cdot \frac{g^{(p)}(\alpha)}{p!}$$

برهان $g(x)$ را حول α بسط می‌دهیم؛ برای یک مقدار از ξ_n بین x_n و α داریم:

$$x_{n+1} = g(x_n) = g(\alpha) + (x_n - \alpha)g'(\alpha) + \dots + \frac{(x_n - \alpha)^{p-1}}{(p-1)!} g^{(p-1)}(\alpha) + \frac{(x_n - \alpha)^p}{p!} g^{(p)}(\xi_n)$$

با استفاده از (۱۴.۵.۲) و $\alpha = g(\alpha)$

$$\alpha - x_{n+1} = -\frac{(x_n - \alpha)^p}{p!} g^{(p)}(\xi_n)$$

با استفاده از قضیه ۷.۲ و $x_n \rightarrow \alpha$ برهان کامل می‌شود.

روش نیوتن را می‌توان با این قضیه تجزیه و تحلیل نمود

$$g(x) = x - \frac{f(x)}{f'(x)} \quad g'(x) = \frac{f(x)f''(x)}{[f'(x)]^2}$$

$$g'(\alpha) = 0 \quad g''(\alpha) = \frac{f''(\alpha)}{f'(\alpha)}$$

این رابطه و روابط (۱۴.۵.۲) همان نتیجه همگرایی (۳.۲.۲) قبلی را برای روش نیوتن به دست می‌دهند. برای کاربرد قضیه ۸.۲ در مثالهای دیگر، مسائل آخر این فصل را ملاحظه کنید.

نظریه این بخش تنها برای روشهای بارستی تک نقطه‌ای است، بنابراین روش خط قاطع و روش مولر مورد توجه قرار نگرفته‌اند. یک نظریه نقطه ثابت متناظر برای روشهای چند مرحله‌ای نقطه ثابت نیز وجود دارد که می‌توان در تراوب^۱ (۱۹۶۲) یافت. ما در اینجا از آن صرف نظر کرده‌ایم، اساساً به این دلیل که در فصلهای بعدی فقط نظریه بارستی نقطه ثابت تک نقطه‌ای مورد نیازند.

۶.۲ برون‌یابی ایتکن^۱ برای دنباله‌های خطی - همگرا

از (۱۱.۵.۲)، قضیه ۶.۲ برای یک بارست همگرایی $x_{n+1} = g(x_n)$ ، $n \geq 0$ داریم

$$\lim_{n \rightarrow \infty} \frac{\alpha - x_{n+1}}{\alpha - x_n} = g'(\alpha) \quad (۱.۶.۲)$$

در این فصل، توجه خود را فقط به حالت همگرایی خطی معطوف می‌سازیم. بنابراین فرض می‌کنیم

$$0 < |g'(\alpha)| < 1 \quad (۲.۶.۲)$$

برآورد خطا را در بارستها آزمایش می‌کنیم و راهی برای شتاب دادن به همگرایی دنباله $\{x_n\}$ بیان می‌کنیم.

مسئله را با در نظر گرفتن نسبت‌های زیر آغاز می‌کنیم

$$\lambda_n = \frac{x_n - x_{n-1}}{x_{n-1} - x_{n-2}} \quad n \geq 2 \quad (۳.۶.۲)$$

ادعا می‌کنیم که:

$$\lim_{n \rightarrow \infty} \lambda_n = g'(\alpha) \quad (۴.۶.۲)$$

برای اثبات آن، می‌نویسیم

$$\lambda_n = \frac{(\alpha - x_{n-1}) - (\alpha - x_n)}{(\alpha - x_{n-2}) - (\alpha - x_{n-1})}$$

با استفاده از (۱۲.۵.۲)

$$\lambda_n = \frac{(\alpha - x_{n-1}) - g'(\xi_{n-1})(\alpha - x_{n-1})}{(\alpha - x_{n-1})/[g'(\xi_{n-2})] - (\alpha - x_{n-1})} = \frac{1 - g'(\xi_{n-1})}{1/[g'(\xi_{n-2})] - 1}$$

$$\lim_{n \rightarrow \infty} \lambda_n = \frac{1 - g'(\alpha)}{1/[g'(\alpha)] - 1} = g'(\alpha)$$

کمیت λ_n قابل محاسبه بوده و هرگاه عملاً به مقداری چون λ همگرا شود، می‌گیریم $\lambda = g'(\alpha)$.

از $\lambda_n = g'(\alpha)$ برای برآورد خطا در بارستهای x_n استفاده می‌کنیم. گیریم

$$\alpha - x_n = \lambda_n(\alpha - x_{n-1})$$

$$\begin{aligned} \alpha - x_n &= (\alpha - x_{n-1}) + (x_{n-1} - x_n) \\ &\doteq \frac{1}{\lambda_n}(\alpha - x_n) + (x_{n-1} - x_n) \\ \alpha - x_n &\doteq \frac{\lambda_n}{1 - \lambda_n}(x_n - x_{n-1}) \end{aligned} \quad (5.6.2)$$

این فرمول خطای ایتنن برای x_n است، و دقت آن با همگرایی $\{\lambda_n\}$ به $g'(\alpha)$ افزایش می‌یابد.

مثال بارستی زیر را در نظر می‌گیریم.

$$x_{n+1} = 6.28 + \sin(x_n) \quad n \geq 0 \quad (6.6.2)$$

ریشه دقیق آن $\alpha \doteq 6.01550307297$ است. نتایج بارستی، همراه با مقادیر λ_n ، $\alpha - x_n$ ، $x_n - x_{n-1}$ و برآورد خطای (6.6.2)، در جدول ۸.۲ داده شده است. مقادیر λ_n به

$$g'(\alpha) = \cos(\alpha) \doteq 0.9644$$

همگراست و برآورد (6.6.2)، یک نشانگر دقیق از خطای واقعی است. اندازه $g'(\alpha)$ نیز نشان می‌دهد که بارستها خیلی به کندی همگرا می‌شوند، و در این حالت، $x_{n+1} - x_n$ نشانگر دقیقی از $\alpha - x_n$ نخواهد بود.

جدول ۸.۲ بارست (6.6.2)

n	x_n	$x_n - x_{n-1}$	λ_n	$\alpha - x_n$	نسبت
۰	۶.۰۰۰۰۰۰۰			۱.۵۵E-۲	
۱	۶.۰۰۰۵۸۴۵	۵.۸۴۵E-۴		۱.۴۹E-۲	
۲	۶.۰۰۱۱۴۵۸	۵.۶۱۳E-۴	۰.۹۶۰۳	۱.۴۴E-۲	۱.۳۶E-۲
۳	۶.۰۰۱۶۸۴۸	۵.۳۹۰E-۴	۰.۹۶۰۴	۱.۳۸E-۲	۱.۳۱E-۲
۴	۶.۰۰۲۲۰۲۶	۵.۱۷۸E-۴	۰.۹۶۰۶	۱.۳۳E-۲	۱.۲۶E-۲
۵	۶.۰۰۲۷۰۰۱	۴.۹۷۴E-۴	۰.۹۶۰۷	۱.۲۸E-۲	۱.۲۲E-۲
۶	۶.۰۰۳۱۷۸۰	۴.۷۸۰E-۴	۰.۹۶۰۹	۱.۲۳E-۲	۱.۱۷E-۲
۷	۶.۰۰۳۶۳۷۴	۴.۵۹۳E-۴	۰.۹۶۱۰	۱.۱۸E-۲	۱.۱۳E-۲

(5.6.2)

فرمول برون‌یابی ایتنکن. همان (۵.۶.۲) است که به صورت برآوردی از α نوشته شده است:

$$\alpha \doteq x_n + \frac{\lambda_n}{1 - \lambda_n}(x_n - x_{n-1}) \quad (۷.۶.۲)$$

طرف راست را با \hat{x}_n برای $n \geq 2$ نشان می‌دهیم. با گذاردن (۳.۶.۲) در (۷.۶.۲) فرمول برای \hat{x}_n را می‌توان به شکل زیر نوشت

$$\hat{x}_n = x_n - \frac{(x_n - x_{n-1})^2}{(x_n - x_{n-1}) - (x_{n-1} - x_{n-2})} \quad n \geq 2 \quad (۸.۶.۲)$$

همان فرمولی که در بسیاری از کتابهای درسی داده شده است.

مثال نتایج جدول ۸.۲ را برای فرمول بارستی (۶.۶.۲) به‌کار می‌بریم. با $n = 7$ از (۷.۶.۲) یا از (۸.۶.۲) استفاده می‌کنیم،

$$\hat{x}_7 = ۶٫۰۱۴۹۵۱۸ \quad \alpha - \hat{x}_7 = ۵٫۵۱E - ۴$$

بنابراین برون‌یاب \hat{x}_7 به‌طور قابل ملاحظه‌ای از x_7 دقیقتر است. اکنون بارست خطی و برون‌یابی ایتنکن را در یک الگوریتم ساده ترکیب می‌نمائیم.

الگوریتم ایتنکن ($g, x_0, \varepsilon, \text{root}$)

۱. توجه: فرض شده است که $|g'(\alpha)| < 1$ و بارست خطی معمولی با استفاده از x_0 به α همگراست.

۲. $x_2 := g(x_1)$ و $x_1 := g(x_0)$.

۳. $\hat{x}_2 := x_2 - \frac{(x_2 - x_1)^2}{(x_2 - x_1) - (x_1 - x_0)}$.

۴. اگر $|\hat{x}_2 - x_2| \leq \varepsilon$ ، آنگاه $\hat{x}_2 := \text{root}$ و خارج شوید.

۵. $\hat{x}_2 := x_0$. قرار داده و به مرحله ۲ بروید.

این الگوریتم، به شرطی که مفروضات مرحله اول برقرار باشند، معمولاً همگراست.

مثال برای نمایاندن الگوریتم ایتنکن، مثال قبلی (۶.۶.۲) را تکرار می‌کنیم. نتایج عددی در جدول ۹.۲ داده شده‌اند. مقادیر x_3, x_6, x_8 مقادیر برون‌یاب ایتنکن هستند که با (۷.۶.۲) تعریف شده‌اند. مقادیر λ_n فقط برای حالت‌های $n = 2, 5, 8$ داده شده‌اند، زیرا فقط در این صورت، خطاهای $\alpha - x_n$ ، $\alpha - x_{n-2}$ و $\alpha - x_{n-1}$ بطور خطی کاهش می‌یابند همان‌گونه که برای $\lambda_n \doteq g'(\alpha)$ لازم است.

جدول ۹.۲ الگوریتم ایتمکن برای (۶.۶.۲)

n	x_n	λ_n	$\alpha - x_n$
۰	۶٫۰۰۰۰۰۰۰۰۰۰۰۰۰۰		۱٫۵۵E - ۲
۱	۶٫۰۰۰۰۵۸۴۵۰۱۸۰۱		۱٫۴۹E - ۲
۲	۶٫۰۰۰۱۱۴۵۷۷۰۷۶۱	۰٫۹۶۰۲۵	۱٫۴۴E - ۲
۳	۶٫۰۱۴۷۰۵۱۴۷۵۴۳		۷٫۹۸E - ۴
۴	۶٫۰۱۴۷۳۳۶۴۸۷۲۰		۷٫۶۹E - ۴
۵	۶٫۰۱۴۷۶۱۱۲۸۹۵۵	۰٫۹۶۴۱۸	۷٫۴۲E - ۴
۶	۶٫۰۱۵۵۰۰۸۰۲۰۶۰		۲٫۲۷E - ۶
۷	۶٫۰۱۵۵۰۰۸۸۲۹۳۵		۲٫۱۹E - ۶
۸	۶٫۰۱۵۵۰۰۹۶۰۹۳۱	۰٫۹۶۴۳۹	۲٫۱۱E - ۶
۹	۶٫۰۱۵۵۰۳۰۷۲۹۴۷		۲٫۰۵E - ۱۱

برون‌یابی اغلب با روشهای بارستی خطی کند همگرا در حل دستگاههای بزرگ معادلات خطی به‌کار برده می‌شود. روشهای واقعی که به‌کار می‌روند با روشی که در بالا توضیح داده شد، تفاوت دارند، ولی آنها نیز بر پایه ایده کلی یافتن رفتار کیفی خطا، مانند (۱.۶.۲)، و سپس به‌کار بردن آن برای برآورد بهتری برای جواب، استوارند. این ایده همچنین در بسط روشهای عددی انتگرالگیری، حل معادلات دیفرانسیل و سایر مسائل ریاضی دنبال می‌شود.

۷.۲ محاسبه عددی ریشه‌های چندگانه

گوئیم $f(x)$ یک ریشه α از مرتبه $p > 1$ دارد اگر

$$f(x) = (x - \alpha)^p h(x) \quad (1.7.2)$$

که در آن $h(\alpha) \neq 0$ و $h(x)$ در $x = \alpha$ پیوسته است. p را یک عدد صحیح مثبت فرض می‌کنیم، اگرچه بعضی از مطالب زیر برای مقادیر غیر صحیح p نیز معتبرند. اگر $h(x)$ در $x = \alpha$ به اندازه کافی مشتق‌پذیر باشد، آنگاه (۱-۷-۲) هم‌ارز است با

$$f(\alpha) = f'(\alpha) = \dots = f^{(p-1)}(\alpha) = 0 \quad f^{(p)}(\alpha) \neq 0 \quad (2.7.2)$$

وقتی ریشه تابعی را با رایانه پیدا می‌کنیم، همیشه یک بازه عدم قطعیت حول ریشه وجود دارد، و اگر ریشه چندگانه باشد، این وضع بدتر خواهد شد. برای آنکه این موضوع بهتر روشن شود، محاسبه دو تابع $f_1(x) = x^2 - 3$ و $f_2(x) = 9 + x^2(x^2 - 6)$ را در نظر می‌گیریم. بنابراین

$\alpha = \sqrt{3}$ ریشهٔ یگانه تابع f_1 و ریشهٔ دوگانهٔ تابع f_2 است. اگر حساب اعشاری چهاررقمی را به‌کار ببریم، برای $x \leq ۱۷۳۱$ داریم $f_1(x) < 0$ و $f_1(۱۷۳۲) = 0$ و برای $x > ۱۷۳۳$ داریم $f_1(x) > 0$. ولی برای $۱۷۳۸ \leq x \leq ۱۷۲۶$ داریم $f_2(x) = 0$. پس مقدار دقت در یافتن ریشهٔ $f_2(x)$ محدود شده است. مثال دومی برای اثر نوفه در محاسبهٔ یک ریشهٔ چندگانه برای تابع $f(x) = (x-1)^3$ در شکل‌های ۱.۱ و ۲.۱ از بخش ۳.۱ فصل یک نشان داده شده است و برای یک مثال نهایی مثال زیر را در نظر می‌گیریم.

مثال تابع

$$f(x) = (x - ۱۷۱)^3(x - ۲۱۱) \\ = x^4 - ۵۴x^3 + ۱۰۷۵۶x^2 - ۸۹۵۴x + ۲۷۹۵۱ \quad (۳.۷.۲)$$

را روی یک ریزکامپیوتر IBM با استفاده از حساب با دقت مضاعف (در BASIC) محاسبه می‌کنیم. ضرایب دارای بسط دودویی محدود نیستند بنابراین دقیقاً در محاسبه وارد نمی‌شوند (غیر از ضریب جملهٔ x^4). محاسبهٔ چند جمله‌ای $f(x)$ به شکل بسط داده شدهٔ (۳.۷.۲) انجام شده و طرح ضرب تودرتو به‌کار گرفته شده است:

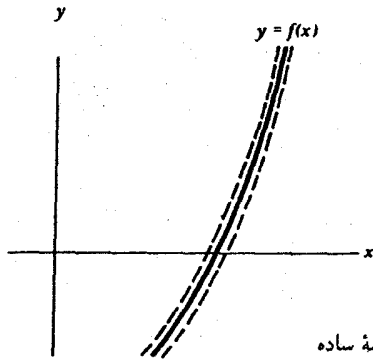
$$f(x) = ۲۷۹۵۱ + x(-۸۹۵۴ + x(۱۰۷۵۶ + x(-۵۴ + x))) \quad (۴.۷.۲)$$

نتایج عددی در جدول ۱۰.۲ داده شده‌اند. توجه داشته باشید حساب به‌کار برده شده در نمایش با ممیز شناور در حدود ۱۶ رقم اعشاری دارد. پس همان گونه که از نتایج عددی جدول مشهود است، نمی‌توان دقتی بیش از ۶ رقم اعشار در محاسبهٔ ریشهٔ $f(x)$ یعنی $\alpha = ۱۷۱$ ، انتظار داشت. همچنین به نتیجهٔ کاربرد دو نمایش مختلف (۳.۷.۲) و (۴.۷.۲) توجه کنید.

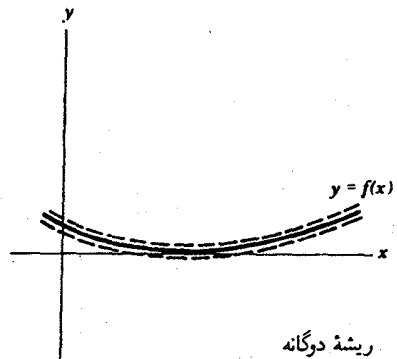
در محاسبهٔ هر تابع $f(x)$ به لحاظ استفاده از حساب با دقت محدود در اثر گردکردن یا قطع کردن، یک عدم قطعیت پیدا می‌شود. این موضوع در بخش ۳ فصل یک تحت نام نوفه در محاسبهٔ تابع بحث شده است. برای ریشه‌های چندگانه، این امر به عدم قطعیت قابل ملاحظه‌ای برای محل ریشه می‌انجامد. در شکل ۷.۲، خط پر نمودار $y = f(x)$ و خط‌چین‌ها ناحیهٔ عدم قطعیت در محاسبهٔ $f(x)$ را که در اثر خطای گردکردن و حساب با ارقام محدود به‌وجود می‌آید، نشان می‌دهند. بازهٔ عدم قطعیت در پیدا کردن ریشهٔ $f(x)$ ، با تقاطع نوار عدم قطعیت و محور x داده می‌شود. روشن است که این بازه برای ریشه چندگانه بزرگتر و برای ریشهٔ ساده کوچکتر است، اگرچه پهنای قائم نوارها حول $y = f(x)$ یکی است.

جدول ۱۰.۲ ارزیابی $f(x) = (x - ۱.۱)^۲(x - ۲.۱)$

x	$f(x) : (۳.۷.۲)$	$f(x) : (۴.۷.۲)$
۱.۰۹۹۹۹۲	۳.۸۶E - ۱۶	۵.۵۵E - ۱۶
۱.۰۹۹۹۹۴	۳.۸۶E - ۱۶	۲.۷۶E - ۱۶
۱.۰۹۹۹۹۶	۲.۷۶E - ۱۶	۰
۱.۰۹۹۹۹۸	-۵.۵۵E - ۱۷	۱.۱۱E - ۱۶
۱.۱۰۰۰۰۰	۵.۵۵E - ۱۷	۰
۱.۱۰۰۰۰۰۲	۵.۵۵E - ۱۷	۵.۵۵E - ۱۷
۱.۱۰۰۰۰۰۴	-۵.۵۵E - ۱۷	۰
۱.۱۰۰۰۰۰۶	-۱.۶۷E - ۱۶	-۱.۶۷E - ۱۶
۱.۱۰۰۰۰۰۸	-۶.۱۱E - ۱۶	-۵.۰۰E - ۱۶



ریشه ساده



ریشه دوگانه

شکل ۷.۲ نوار عدم قطعیت در محاسبه یک تابع

روش نیوتن و ریشه‌های چندگانه. مسأله دیگر در ریشه‌های چندگانه این است که روشهای ریشه‌یابی قبلی، وقتی که ریشه مورد جستجو چندگانه باشد، به خوبی ریشه ساده انجام نمی‌گیرند. اکنون این را برای روش نیوتن بررسی می‌کنیم.

روش نیوتن را به‌عنوان یک روش نقطه ثابت، مانند (۲.۵.۲)، با $f(x)$ که در (۱.۷.۲) صدق می‌کند، در نظر می‌گیریم:

$$x_{n+1} = g(x_n) \quad g(x) = x - \frac{f(x)}{f'(x)} \quad x \neq \alpha$$

پیش از به‌دست آوردن $g'(\alpha)$ ، با استفاده از (۱.۷.۲)، $g(x)$ را ساده می‌کنیم:

$$f'(x) = (x - \alpha)^p h'(x) + p(x - \alpha)^{p-1} h(x)$$

$$g(x) = x - \frac{(x - \alpha)h(x)}{ph(x) + (x - \alpha)h'(x)}$$

با مشتگیری

$$g'(x) = 1 - \frac{h(x)}{ph(x) + (x - \alpha)h'(x)} - (x - \alpha) \cdot \frac{d}{dx} \left[\frac{h(x)}{ph(x) + (x - \alpha)h'(x)} \right]$$

و

$$g'(\alpha) = 1 - \frac{1}{p} \neq 0 \quad p > 1 \quad \text{برای} \quad (5.7.2)$$

بنابراین روش نیوتن یک روش خطی با نرخ همگرایی $(p - 1)/p$ است.

مثال کوچکترین ریشه معادله زیر را با استفاده از روش نیوتن پیدا کنید.

$$f(x) = -4,689,999 + x(9,1389 + x(-5,56 + x)) \quad (6.7.2)$$

نتایج عددی در جدول ۱۱.۲ نشان داده شده‌اند. محاسبات با یک ریز کامپیوتر IBM PC با دقت مضاعف (با BASIC) انجام شده‌اند. فقط قسمتهایی از نتایج آورده شده‌اند، که حالت کلی محاسبات را نشان می‌دهند. ستونی که با نسبت مشخص شده‌است، نرخ همگرایی خطی را که با λ_n در (۳.۶.۲) اندازه‌گیری می‌شود، نشان می‌دهد.

جدول ۱۱.۲ روش نیوتن برای (۶.۷.۲)

n	x_n	$f(x_n)$	$\alpha - x_n$	نسبت
۰	۱٫۲۲	-۱٫۸۸E-۴	۱٫۰۰E-۲	
۱	۱٫۲۲۴۹۸۶۷۳۷۴	-۴٫۷۱E-۵	۵٫۰۱E-۳	
۲	۱٫۲۲۷۴۹۰۰۲۲۲	-۱٫۱۸E-۵	۲٫۵۱E-۳	۰٫۵۰۲
۳	۱٫۲۲۸۱۷۴۴۱۷۰۵	-۲٫۹۵E-۶	۱٫۲۶E-۳	۰٫۵۰۱
۴	۱٫۲۲۹۳۷۱۸۷۴۶	-۷٫۳۸E-۷	۶٫۲۸E-۴	۰٫۵۰۱
۵	۱٫۲۲۹۶۸۵۸۸۴۶	-۱٫۸۵E-۷	۳٫۱۴E-۴	۰٫۵۰۰
۱۸	۱٫۲۲۹۹۹۹۹۶۲۱	-۲٫۸۹E-۱۵	۳٫۸۰E-۸	۰٫۵۰۵
۱۹	۱٫۲۲۹۹۹۹۹۸۲۳	-۶٫۶۶E-۱۶	۱٫۷۷E-۸	۰٫۵۲۵
۲۰	۱٫۲۲۹۹۹۹۹۹۲۴	-۱٫۱۱E-۱۶	۷٫۵۸E-۹	۰٫۴۹۶
۲۱	۱٫۲۲۹۹۹۹۹۹۶۳	۰٫۰	۳٫۶۶E-۹	۰٫۳۸۳

واضح است که در این حالت روش نیوتن برای پیدا کردن ریشه (۶.۷.۲)، خطی با نرخ $g'(x) = \frac{1}{4}$ با (۵.۷.۲) سازگار است، زیرا $\alpha = ۱۲۳$ یک ریشه با چندگانی $p = ۲$ است. بارستهای پایانی در جدول تحت تأثیر نوبه محاسبه رایانه‌ای $f(x)$ واقع شده‌اند. اگرچه نمایش با ممیز شناور تقریباً شامل ۱۶ رقم است، ولی در این حالت دقت فقط در حدود هشت رقم است. برای بهبود روش نیوتن، علاقه‌مندیم $g(x)$ به گونه‌ای باشد که $g'(\alpha) = 0$. با توجه به (۵.۷.۲)، تعریف می‌کنیم

$$g(x) = x - p \frac{f(x)}{f'(x)}$$

در این صورت به‌سادگی ملاحظه می‌شود که $g'(\alpha) = 0$ ؛ بنابراین

$$\begin{aligned} \alpha - x_{n+1} &= g(\alpha) - g(x_n) \\ &= -g'(\alpha)(x_n - \alpha) - \frac{1}{4}(x_n - \alpha)^2 g''(\xi_n) \end{aligned}$$

که ξ_n بین x_n و α است. بنابراین

$$\alpha - x_{n+1} = -\frac{1}{4}(\alpha - x_n)^2 g''(\xi_n)$$

که نشان می‌دهد روش

$$x_{n+1} = x_n - p \frac{f(x_n)}{f'(x_n)} \quad n = 0, 1, 2, \dots \quad (۷.۷.۲)$$

همانند روش اصلی نیوتن برای ریشه‌های ساده دارای همگرایی مرتبه ۲ است.

مثال رابطه (۷.۷.۲) را برای مثال قبلی (۶.۷.۲)، با $P = ۲$ برای یک ریشه مضاعف به‌کار می‌بریم. نتایج در جدول ۱۲.۲، با استفاده از همان رایانه قبلی، نشان داده شده‌اند. بارستها به سرعت همگرا می‌شوند، و سپس حول ریشه نوسان می‌کنند. دقت (یا عدم آن)، نوبه در $f(x)$ و چندگانی ریشه را منعکس می‌کند.

روش نیوتن را می‌توان برای تعیین چندگانی p به‌کار برد، همانند جدول ۱۱.۲ که با (۵.۷.۲) ترکیب شده باشد، و سپس (۷.۷.۲) را می‌توان به‌کار برد و سرعت همگرایی را افزایش داد. ولی عدم قطعیت ذاتی در ریشه که از چندگانی ریشه نتیجه می‌شود کماکان باقی می‌ماند. این نقص

جدول ۱۲.۲ روش اصلاح‌شده (۷.۷.۲) نیوتن، در حل (۶.۷.۲)

n	x_n	$f(x_n)$	$\alpha - x_n$
۰	۱٫۲۲	-۱٫۸۸E-۴	۱٫۰۰E-۲
۱	۱٫۲۲۹۹۷۳۴۷۴۸	-۱٫۳۱E-۹	۲٫۶۵E-۵
۲	۱٫۲۲۹۹۹۹۹۹۹۹۸	-۱٫۱۱E-۱۶	۱٫۸۵E-۱۰
۳	۱٫۲۳۰۰۰۰۰۳۲۰۸	-۱٫۹۲E-۱۳	-۳٫۲۱E-۷
۴	۱٫۲۳۰۰۰۰۰۰۰۰۱	-۱٫۱۱E-۱۶	-۸٫۵۴E-۱۱
۵	۱٫۲۲۹۹۹۹۹۳۰۴۶	-۹٫۰۴E-۱۳	۶٫۹۵E-۷

را می‌توان فقط با تغییر صورت مسأله ریشه‌یابی از راه تحلیلی به صورتی که در آن α ریشه ساده باشد رفع کرد. ساده‌ترین راه انجام آن، تشکیل مشتق مرتبه $(p-1)$ ام $f(x)$ و سپس حل

$$f^{(p-1)}(x) = 0 \quad (۸.۷.۲)$$

است.

مثال مثال قبلی دارای ریشه با چندگانی $p = ۲$ است، که یک ریشه ساده برای

$$f'(x) = 3x^2 - 11.12x + 9.1389$$

است. با استفاده از آخرین بارست در جدول ۱۱.۲ به عنوان حدس اولیه، و به‌کار بردن روش نیوتن برای پیدا کردن ریشه $f'(x)$ که هم‌اکنون داده شد، برای محاسبه α با دقت کامل رایانه، فقط یک بارست لازم است.

۸.۲ الگوریتم ریشه‌یابی برنت

الگوریتمی را توضیح می‌دهیم که هم مزایای روش نیم‌سازی و هم روش خط قاطع را دارد، و در عین حال نقاط ضعف هیچ یک از آنها را ندارد. این الگوریتم از (فصل ۴، برنت^۱ (۱۹۷۳)) گرفته شده و تکمیل یافته الگوریتم دیگری است که (دکتر^۲ (۱۹۶۹)) ارائه کرده است. الگوریتم به بازه کوچکی می‌انجامد که ریشه را در بردارد. چنانچه تابع حول ریشه ξ ، به اندازه کافی هموار باشد، آنگاه مرتبه همگرایی، همانند روش خط قاطع، زیر خطی است.

در توضیح الگوریتم، نمادهای (برنت (۱۹۷۳، ص ۴۷)) را به‌کار می‌بریم. با دو مقدار a و b ، به برنامه وارد می‌شویم که برای این دو مقدار، (۱) حداقل یک ریشه ξ از $f(x)$ بین a و b .

وجود دارد و (۲) $f(a) \cdot f(b) \leq 0$. در برنامه همچنین یک تحمل مطلوب t داده می‌شود که از آن تحمل حد توقف δ ناشی می‌شود:

$$\delta = t + 2\epsilon |b| \quad (۱.۸.۲)$$

که ϵ واحد گردکردن رایانه است [۱-۲-۱۱] از فصل ۱ را ببینید.

در یک مرحله نمونه این الگوریتم، b بهترین برآورد جاری ریشه ζ ، a مقدار قبلی b ، c بارست گذشته است که طوری انتخاب شده است که ریشه ζ بین b و c قرارگیرد (در ابتدا $c = a$). تعریف می‌کنیم $m = \frac{1}{4}(c - b)$.

اگر (۱) $f(b) = 0$ ، یا (۲) $|m| \leq \delta$ الگوریتم را متوقف نمایید. در هر یک از دو حالت، مقدار تقریبی ریشه را $b = \hat{\zeta}$ قرار دهید. برای حالت (۲)، چون b معمولاً از راه روش خط قاطع به دست آمده است، ریشه ζ معمولاً به b نزدیکتر از c است. از این رو معمولاً

$$|\zeta - \hat{\zeta}| \leq \delta$$

گرچه تنها چیزی که می‌تواند تضمین شود رابطه زیر است:

$$|\zeta - \hat{\zeta}| \leq 2\delta$$

اگر آزمون خطا برقرار نشد، قرار دهید

$$i = b - f(b) \frac{b - a}{f(b) - f(a)} \quad (۲.۸.۲)$$

سپس قرار دهید

$$b'' = \begin{cases} i & \text{اگر } i \text{ بین } b \text{ و } b + m = \frac{b+c}{4} \text{ واقع باشد} \\ b + m & \text{در غیر این صورت (که روش خط قاطع است)} \end{cases}$$

در حالتی که a ، b و c متمایز باشند، روش خط قاطع در تعریف i ، با یک روش عکس درون‌یابی درجه دو، جایگزین می‌شود. این عمل به یک همگرایی اندکی سریعتر برای کل الگوریتم می‌انجامد. به دنبال تعریف b'' ، تعریف می‌کنیم

$$b' = \begin{cases} b'' & \text{اگر } |b - b''| > \delta \\ b + \delta \cdot \text{sign}(m) & \text{اگر } |b - b''| \leq \delta \end{cases} \quad (۳.۸.۲)$$

اگر از ریشه فاصله داشته باشید، آنگاه $b' = b''$. با این انتخاب، روش چنین است: (۱) درون‌یابی خطی (یا درجه دو) است، یا: (۲) روش نیم‌سازی است؛ معمولاً برای توابع هموار $f(x)$ ، حالت (۱) است. این حالت معمولاً به مقداری از m می‌رسد که کوچک نمی‌شود. برای به‌دست آوردن بازه کوچکی که شامل ریشه ζ باشد، وقتی که به آن نزدیک هستیم، از $b' := b + \delta \cdot \text{sign}(m)$ استفاده می‌کنیم، که مرحله‌ای از δ در جهت c است. به علت راهی که برای تعیین مقدار جدید c انتخاب کردیم، اکنون معمولاً بازه کوچکی جدیدی حول ζ به‌دست می‌آید. پرنهت قبل از انتخاب b جدید، معمولاً b' که اینجا گفته شد، یک قدم مهم دیگر که تکنیکی است برمی‌دارد.

وقتی که b' به‌دست آمد، b را مساوی b' و a را مساوی مقدار قدیمی b قرار می‌دهیم. اگر علامت $f(b)$ با استفاده از مقدار جدید b ، با علامت مقدار قدیمی b یکی باشد، مقدار c عوض نمی‌شود؛ در غیر این صورت، c را برابر مقدار قدیمی b قرار می‌دهیم، که بازه کوچکتری حول ζ به‌دست می‌دهد. دقت مقدار b ، همانگونه که قبلاً توضیح داده شد، اکنون امتحان می‌شود. پرنهت، خیلی مواظب بوده‌است که از مشکلات پی‌ریز و سرریز در روش خود احتراز کند و نتیجتاً تا حدودی خواندن برنامه پیچیده شده است.

مثال هر یک از حالات زیر با یک IBM PC، با هم‌پردازنده حساب ۸۰۸۷ و حساب با دقت ساده که در استاندارد IEEE برای حساب با ممیز شناور صدق می‌کند، حل شده است. تحمل $t = 10^{-5}$ بوده و چون در تمام حالات ریشه ζ برابر با ۱ بوده است،

$$\delta = 10^{-5} + 2 |b| \times (5.96 \times 10^{-8}) \approx 1.01 \times 10^{-5}$$

توابع به همان شکلی محاسبه شده‌اند که در اینجا داده شده‌اند و بازه اولیه $[a, b] = [0, 3]$ بوده است. در جدول مقادیر b و c تا ۷ رقم اعشاری گرد شده‌اند.

حالت ۱. $f(x) = (x-1)[1+(x-1)^2]$. نتایج عددی در جدول ۱۳.۲ داده شده‌اند. لزوم به‌کار بردن $b' = b + \delta(x-1)^2 \cdot \text{sign}(m)$ برای به‌دست آوردن یک بازه کوچکی که شامل ζ باشد این مثال را روشن می‌کند.

حالت ۲. $f(x) = x^2 - 1$. نتایج عددی در جدول ۱۴.۲ داده شده‌اند.

حالت ۳. $f(x) = -1 + x(3+x(-3+x))$. ریشه $\zeta = 1$ سه‌گانی است، 50° بارست برای همگرایی به ریشه تقریبی 1.000001 به‌کار رفته است. با مقادیر اولیه $a = 0$ ، $b = 3$ ، روش نیم‌سازی فقط ۱۹ بارست برای همین دقت را می‌طلبد. اگر روش پرنهت با روش نیم‌سازی،

جدول ۱۳.۲ مثال ۱ روش برنت

b	$f(b)$	c
۰°	-۲,۰۰E + ۰	۳°
۰ر۵	-۶,۲۵E - ۱	۳°
۰ر۷۱۳۹۰۳۸	-۳,۱۰E - ۱	۳°
۰ر۹۱۵۴۵۰۷	-۸,۵۲E - ۲	۳°
۰ر۹۹۰۱۷۷۹	-۹,۸۲E - ۳	۳°
۰ر۹۹۹۸۵۶۷	-۱,۴۳E - ۴	۳°
۰ر۹۹۹۹۹۹۹	-۱,۱۹E - ۷	۳°
۰ر۹۹۹۹۹۹۹	-۱,۱۹E - ۷	۱ر۰۰۰۰۰۱۰

جدول ۱۴.۲ مثال ۲ روش برنت

b	$f(b)$	c
۰°	-۱,۰°	۳°
۰ر۳۳۳۳۳۳۳	-۸,۸۹E - ۱	۳°
۰ر۳۳۳۳۳۳۳	-۸,۸۹E - ۱	۱ر۶۶۶۶۶۶۷
۰ر۷۷۷۷۷۷۸	-۴,۰۰E - ۱	۱ر۶۶۶۶۶۶۷
۱ر۰۶۸۶۸۷	۱,۴۲E - ۱	۰ر۷۷۷۷۷۷۸
۰ر۹۹۱۷۳۳۶	-۱,۶۵E - ۲	۱ر۰۶۸۶۸۷
۰ر۹۹۹۷۲۴۴	-۵,۵۱E - ۴	۱ر۰۶۸۶۸۷
۱ر۰۰۰۰۰۰۰	۲,۳۸E - ۷	۰ر۹۹۹۷۲۴۴
۱ر۰۰۰۰۰۰۰	۲,۳۸E - ۷	۰ر۹۹۹۹۹۰۰

در رده همه توابع پیوسته مقایسه شود، آنگاه تعداد بارستها برای تحمل خطای δ تقریباً چنین است:

$$\log_2 \left(\frac{b-1}{\delta} \right) \quad \text{برای روش نیم‌سازی}$$

$$\left[\log_2 \left(\frac{b-a}{\delta} \right) \right]^2 \quad \text{برای الگوریتم برنت}$$

پس همان‌گونه که مثال ما نشان می‌دهد، حالت‌هایی وجود دارند که روش نیم‌سازی بهتر است. ولی برای توابع به اندازه کافی هموار با $f'(\alpha) \neq 0$ ، الگوریتم برنت تقریباً همیشه سریعتر است.

حالت ۴. $f(x) = (x-1) \exp[-1/(x-1)^2]$ ریشه $x = 1$ چندگانی بی‌نهایت دارد، زیرا برای همه مقادیر $r \geq 0$ ، $f^{(r)}(1) = 0$. نتایج عددی در جدول ۱۵.۲ داده شده‌اند.

جدول ۱۵.۲ مثال ۴ روش برنت

b	$f(b)$	c
۰°	$-۳,۶۸E - ۱$	۳°
$۰,۵۷۳۱۷۵۴$	$-۱,۷۶E - ۳$	۳°
$۰,۵۹۵۹۳۳۱$	$-۱,۶۳E - ۳$	۳°
$۰,۶۰۹۸۴۴۳$	$-۵,۴۷E - ۴$	۳°
$۰,۶۱۳۶۳۵۴$	$-۴,۷۶E - ۴$	$۱,۸۰۴۹۲۲$
$۰,۶۳۸۹۲۵۸$	$-۱,۶۸E - ۴$	$۱,۸۰۴۹۲۲$
$۱,۲۲۱۹۲۴$	$۳,۳۷E - ۱۰$	$۰,۶۳۸۹۲۵۸$
$۱,۲۲۱۹۱۴$	$۳,۳۷E - ۱۰$	$۰,۶۳۸۹۲۵۸$
$۱,۲۱۶۵۸۵$	$۱,۲۰E - ۱۰$	$۰,۶۳۸۹۲۵۸$
$۰,۹۲۷۷۵۵۳$	۰°	$۱,۲۱۶۵۸۵$

توجه کنید که برنامه یک ریشه دقیق برای شکل ماشینی $f(x)$ یافته است. این به خاطر عدم دقت ذاتی در محاسبه تابع $f(x)$ است؛ بخش قبلی درباره ریشه‌های چندگانه را ببینید. این ریشه البته خیلی نادقیق است، ولی این چیزی نیست که برنامه بتواند آن را اصلاح کند.

برنامه اصلی برنت، که در سال ۱۹۷۳ منتشر شده هنوز عمومیّت دارد و زیاد از آن استفاده می‌شود. مع‌ذالک اصلاح و توسعه آن هنوز ادامه دارد. برای یکی از آنها، و برای مرور بقیه، ل' (۱۹۸۳) را ببینید.

۹.۲ ریشه‌های چندجمله‌یها

ما اکنون حل معادله چندجمله‌ای

$$p(x) \equiv a_0 + a_1x + \dots + a_nx^n = 0 \quad a_n \neq 0 \quad (۱.۹.۲)$$

را در نظر می‌گیریم. این مسأله از راه‌های گوناگونی ظاهر می‌شود و نوشته‌های زیادی برای پرداختن به آن، پدید آمده است. گاهی یک ریشه خاص خواسته می‌شود و یک حدس اولیه خوب در دست است. در چنین حالتی، بهترین راه، اصلاح یکی از روشهای بارستی برای بهره‌گیری از شکل خاص چندجمله‌ای است. در حالات دیگر، ممکن است اطلاع کمی درباره جای ریشه‌ها داشته باشیم، پس روشهای دیگری باید به‌کار برده شوند، که انواع آنها زیاد است. در این بخش، ما فقط یک گشت مختصری، بدون ادعا برکامل بودن آن، در قسمت ریشه‌یابی چندجمله‌یها، خواهیم زد. بر اصلاح جزئی روشهای بخشهای

گذشته تکیه خواهیم کرد و مسائل مربوط به پایداری عددی را مورد ملاحظه قرار خواهیم داد. با مروری بر بعضی قضایای معین محدوده ریشه یا جایابی تقریبی ریشه‌های (۱.۹.۲) بحث را شروع خواهیم کرد.

قضیه‌های جایابی. چون $p(x)$ یک چندجمله‌ای است، قضیه‌های زیادی درباره ریشه‌های آن می‌توان بیان کرد، که درباره توابع دیگر درست نخواهند بود. معروفترین آنها، قضیه اساسی جبر است، که به ما اجازه می‌دهد $p(x)$ را به صورت حاصلضربی یکتا (جز از لحاظ ترتیب) که شامل ریشه‌ها باشد بنویسیم،

$$p(x) = a_n(x - z_1) \dots (x - z_n) \quad (2.9.2)$$

z_1, \dots, z_n ریشه‌های $p(x)$ اند که به تعداد درجه چندگانی خود تکرار می‌شوند. اکنون چند قضیه معروف جایابی تقریبی و محدود کردن این ریشه‌ها را ذکر می‌کنیم. قاعده علامت دکارت. که به‌کار برده می‌شود تا ریشه‌های حقیقی و مثبت $p(x)$ را محدود نماید، به شرطی که ضرایب a_0, \dots, a_n تماماً حقیقی باشند.

گیریم ν تعداد تغییر علامتها در ضرایب $p(x)$ در (۱.۹.۲) باشد وقتی جملات صفرناده گرفته شوند. گیریم k تعداد ریشه‌های حقیقی مثبت $p(x)$ باشد که چندگانی آنها به حساب آمده‌اند. در این صورت $k \leq \nu$ و $\nu - k$ زوج است.

یک برهان این قضیه در هنریچی^۱ (۱۹۷۴، ص ۴۴۲) و هاسهولدر^۲ (۱۹۷۰، ص ۸۲) داده شده‌است.

مثال عبارت $p(x) = x^2 - x - 1$ دارای $\nu = 1$ تغییر علامت در ضرایب است. بنابراین $k = 1$ ؛ زیرا در غیر این صورت $k = 0$ و $\nu - k = 1$ یک عدد صحیح زوج نیست که یک تناقض است.

قاعده تغییر علامت دکارت برای محدود کردن تعداد ریشه‌های منفی $p(x)$ نیز به‌کار برده می‌شود. آن را برای چندجمله‌ای

$$q(x) = p(-x)$$

به‌کار می‌بریم. ریشه‌های مثبت این چندجمله‌ای، همان ریشه‌های منفی $p(x)$ هستند. این قاعده را برای مثال اخیر در $q(x) = x^2 + x - 1$ به‌کار می‌بریم. باز هم یک ریشه حقیقی مثبت برای $q(x)$ وجود دارد، پس $p(x)$ دارای یک ریشه حقیقی منفی است.

یک کران بالا برای تمام ریشه‌های $p(x)$ با رابطه زیر داده می‌شود

$$|z_i| \leq R \equiv 1 + \max_{0 \leq i \leq n-1} \left| \frac{a_i}{a_n} \right| \quad (3.9.2)$$

این قضیه از آگوستن کوشی^۱ (۱۸۲۹) است و اثبات آن در هاوسهولدر (۱۹۷۰، ص ۷۱) داده شده است. قضیه مشابه دیگری از کوشی هست که بر مبنای در نظر گرفتن چند جمله‌بیهای زیر قرار دارد

$$|a_n| x^n + |a_{n-1}| x^{n-1} + \dots + |a_1| x - |a_0| = 0 \quad (4.9.2)$$

$$|a_n| x^n - |a_{n-1}| x^{n-1} - \dots - |a_1| x - |a_0| = 0 \quad (5.9.2)$$

فرض می‌کنیم $a_0 \neq 0$ ، که هم‌ارز است با آنکه فرض کنیم $x = 0$ یک ریشه $p(x)$ نیست. در این صورت با استفاده از قاعده علامت دکارت، هر یک از چند جمله‌بیهای فوق، دقیقاً یک ریشه مثبت دارد؛ آنها را به ترتیب ρ_1 و ρ_2 می‌نامیم. پس تمام ریشه‌های z_j از $p(x)$ در رابطه زیر صدق می‌کنند.

$$\rho_1 \leq |z_j| \leq \rho_2 \quad (6.9.2)$$

اثبات کران بالا در هنریچی (۱۹۷۴، ص ۴۵۸) و هاوسهولدر (ص ۷۰، ۱۹۷۰) داده شده است. اثبات کران پایین را می‌توان بر اساس روش زیر، که برای ساختن کران پایین (۳.۹.۲) نیز به‌کار می‌رود، بنانهاد. چند جمله‌ای

$$q(x) = x^n p\left(\frac{1}{x}\right) = a_n + a_{n-1}x + \dots + a_1x^{n-1} + a_0x^n \quad a_0 \neq 0 \quad (7.9.2)$$

را در نظر می‌گیریم. پس ریشه‌های $q(x)$ برابر $1/z$ هستند، که z یک ریشه $p(x)$ است. اگر قضیه کران بالای (۶.۹.۲) برای (۷.۹.۲) به‌کار برده شود کران پایین (۶.۹.۲) به‌دست می‌آید. این کاربرد را در ضمن یک مسأله نشان خواهیم داد.

چون هر یک از چند جمله‌بیهای (۴.۹.۲) و (۵.۹.۲) یک ریشه یگانه ساده مثبت دارد، روش نیوتن را می‌توان به‌سادگی برای ترسیم R_1 و R_2 به‌کار برد. به‌عنوان یک حدس اولیه، از کران بالای (۳.۹.۲) استفاده می‌کنیم یا با حدسهای اولیه مثبت کوچکتری تجربه می‌کنیم. مثالهای این قضیه‌ها را در مسائل مطرح می‌کنیم.

قضیه‌های بسیار دیگری از این نوع وجود دارند، که هم هنریچی (۱۹۷۴، فصل ۶) و هم هاوسهولدر (۱۹۷۰) مقالات عالی در این موضوع نوشته‌اند.

ضرب تودرتو. یک راه بسیار کارآمد برای محاسبهٔ چندجمله‌ای $p(x)$ که در (۱.۹.۲) داده شده، استفاده از ضرب تودرتو است:

$$p(x) = a_0 + x(a_1 + x(a_2 + \cdots + x(a_{n-1} + a_n x) \cdots)) \quad (۸.۹.۲)$$

در فرمول (۱.۹.۲)، n جمع و $2n - 1$ ضرب و در فرمول (۸.۹.۲)، n جمع و n ضرب وجود دارد که یک صرفه‌جویی قابل ملاحظه‌ای است.

برای کارهای بعدی، مناسبتر است که ضرایب کمکی زیر را معرفی کنیم؛ گیریم $b_n = a_n$

$$b_k = a_k + z b_{k+1}, \quad k = n - 1, n - 2, \dots, 0 \quad (۹.۹.۲)$$

با ملاحظهٔ (۸.۹.۲) به‌سادگی ملاحظه می‌شود که

$$p(z) = b_0. \quad (۱۰.۹.۲)$$

چندجمله‌ای زیر را معرفی می‌کنیم

$$q(x) = b_1 + b_2 x + \cdots + b_n x^{n-1} \quad (۱۱.۹.۲)$$

در این صورت

$$b_0 + (x - z)q(x) = b_0 + (x - z)[b_1 + b_2 x + \cdots + b_n x^{n-1}]$$

$$= (b_0 - b_1 z) + (b_1 - b_2 z)x + \cdots$$

$$+ (b_{n-1} - b_n z)x^{n-1} + b_n x^n$$

$$= a_0 + a_1 x + \cdots + a_n x^n = p(x)$$

$$p(x) = b_0 + (x - z)q(x) \quad (۱۲.۹.۲)$$

که $q(x)$ خارج‌قسمت و b_0 باقیماندهٔ تقسیم $p(x)$ بر $(x - z)$ است. استفاده از (۹.۹.۲) برای محاسبهٔ $p(z)$ و تشکیل چندجمله‌ای خارج‌قسمت $q(x)$ را روش هورنر^۱ نیز می‌خوانند.

اگر z یک ریشهٔ $p(x)$ باشد، آنگاه $b_0 = 0$ و $p(x) = (x - z)q(x)$ برای پیدا کردن ریشه‌های دیگر $p(x)$ ، می‌توانیم جستجوی خود را به ریشه‌های $q(x)$ محدود کنیم. این فرآیند کاهش را تقلیل می‌نامند و آن را باید با احتیاط به کار برد، نکته‌ای که بعداً به آن باز خواهیم گشت.

روش نیوتن. اگر بخواهیم برای یافتن یک ریشهٔ $p(x)$ ، روش نیوتن را به‌کار ببریم، باید بتوانیم هم $p(x)$ و هم $p'(x)$ را در هر نقطهٔ z محاسبه کنیم. با توجه به (۱۲.۹.۲) داریم:

$$p'(x) = (x - z)q'(x) + q(x)$$

$$p'(z) = q(z) \quad (۱۳.۹.۲)$$

از (۱۰.۹.۲) و (۱۳.۹.۲) در توجیه روش نیوتن برای ریشه‌یابی چندجمله‌یها که ذیلاً آمده است استفاده می‌کنیم.

الگوریتم $Polynew(a, n, x_0, \varepsilon, itMaxroot, b, ier)$

۱. توجه: a بردار ضرایب است، $itmax$ بیشترین تعداد بارستهای است که باید محاسبه شود، b بردار ضرایب چندجمله‌ای تقلیل‌یافته است، و ier یک شاخص برای خطاست.

$$itnum := ۱.۲$$

$$b_n := c = a_n, z := x_0.۳$$

$$۴. برای $k := n - ۱, \dots, ۱$ ، $b_k := a_k + zb_{k+۱}$ و $c := b_k + zc$$$

$$۵. $b_0 := a_0 + zb_1$$$

$$۶. اگر $c = 0$ ، آنگاه $ier = ۲$ قرار داده خارج شوید.$$

$$۷. $x_1 = x_0 - b_0/c$$$

$$۸. اگر $\varepsilon \leq |x_1 - x_0|$ ، آنگاه $ier := 0$ و $root := x_1$ قرار داده خارج شوید.$$

$$۹. اگر $itnum = itmax$ ، آنگاه $ier := ۱$ قرار داده و خارج شوید.$$

$$۱۰. در غیر این صورت، $itnum := itnum + ۱$ و $x_1 := x_0$ ، و به مرحلهٔ ۳ بروید.$$

مسائل پایداری. چندجمله‌یهای بسیاری هستند که ریشه‌های آنها نسبت به تغییرات کوچک ضرایب حساسیت زیادی دارند. در بعضی از این مسائل که دارای ریشه‌های چندگانه‌اند، تعجب‌آور نیست که این ریشه‌ها نسبت به تغییرات کوچک ضرایب حساسیت زیادی نشان دهند. ولی چندجمله‌یهای زیادی هستند که فقط ریشه‌های سادهٔ کاملاً از هم جدا دارند و در آنها باز هم ریشه‌ها نسبت به اختلالات کوچک بسیار حساس‌اند. فرمولهایی ذیلاً داده شده‌اند که این حساسیت را توضیح می‌دهند. و مثالهای عددی نیز آورده شده‌اند.

برای این بحث نظری، توابع زیر را معرفی می‌کنیم

$$p(x) = a_0 + a_1x + \dots + a_nx^n \quad a_n \neq 0$$

$$q(x) = b_0 + b_1x + \dots + b_nx^n \quad (۱۴.۹.۲)$$

و اختلال $p(x)$ را با

$$p(x; \varepsilon) = p(x) + \varepsilon q(x) \quad (۱۵.۹.۲)$$

تعریف می‌نماییم. ریشه‌های $p(x; \varepsilon)$ را با $z_1(\varepsilon), \dots, z_n(\varepsilon)$ نشان می‌دهیم که برحسب چندگانگی تکرار می‌شوند، گیریم $z_j = z_j(0)$, $j = 1, \dots, n$ ریشه متناظر $p(x) = p(x, 0)$ باشند. روشن است که صفرهای یک چندجمله‌ای توابع پیوسته‌ای از ضرایب چندجمله‌ای هستند [برای مثال هنریچی (۱۹۷۴، ص ۲۸۱) را ببینید]. در نتیجه، $z_j(\varepsilon)$ تابعی پیوسته از ε است. آنچه ما می‌خواهیم معلوم کنیم این است که وقتی ε نزدیک صفر است، ریشه $z_j(\varepsilon)$ با چه تندی ε تغییر می‌کند.

مثال

$$p(x; \varepsilon) = (x - 1)^3 - \varepsilon \quad p(x) = (x - 1)^3 \quad \varepsilon > 0$$

پس ریشه‌های $p(x)$ عبارت‌اند از $z_1 = z_2 = z_3 = 1$. ریشه‌های $p(x; \varepsilon)$ عبارت‌اند از:

$$z_1(\varepsilon) = 1 + \sqrt[3]{\varepsilon} \quad z_2(\varepsilon) = 1 + w \cdot \sqrt[3]{\varepsilon} \quad z_3(\varepsilon) = 1 + w^2 \cdot \sqrt[3]{\varepsilon}$$

با $w = \frac{1}{\sqrt[3]{-1 + i\sqrt{3}}}$. برای هر سه ریشه $p(x; \varepsilon)$:

$$|z_j(\varepsilon) - 1| = \sqrt[3]{\varepsilon}$$

برای روشن شدن مطلب، گیریم $\varepsilon = 0.001$. در این صورت

$$p(x; \varepsilon) = x^3 - 3x^2 + 3x - 1.001$$

که تغییر نسبتاً کوچکی در $p(x)$ است. ولی برای ریشه‌ها

$$|z_j(\varepsilon) - 1| = 0.1$$

که تغییر نسبتاً بزرگی در ریشه‌های $z_j = 1$ است. ما اکنون برآوردهای کلیتری برای $z_j(\varepsilon) - z_j$ می‌دهیم.

حالت ۱. z_j یک ریشه ساده $p(x)$ است، پس $p'(z_j) \neq 0$. با استفاده از نظریه توابع با متغیر مختلط، می‌دانیم که $z_j(\varepsilon)$ را می‌توان به صورت یک سری توانی نوشت؛

$$z_j(\varepsilon) = z_j + \sum_{l=1}^{\infty} \gamma_l \varepsilon^l \quad (16.9.2)$$

برای برآورد $z_j(\varepsilon) - z_j$ ، یک فرمول برای اولین جمله $\gamma_1 \varepsilon$ سری به دست می‌آوریم. برای شروع، به راحتی می‌توان دید که

$$\gamma_1 = z_j'(\varepsilon)$$

برای محاسبه $z_j'(\varepsilon)$ ، از اتحاد زیر، که برای کلیه مقادیر به اندازه کافی کوچک ε برقرار است، مشتق می‌گیریم

$$p(z_j(\varepsilon)) + \varepsilon q(z_j(\varepsilon)) = 0$$

به دست می‌آوریم

$$p'(z_j(\varepsilon))z_j'(\varepsilon) + q(z_j(\varepsilon)) + \varepsilon q'(z_j(\varepsilon))z_j'(\varepsilon) = 0$$

$$z_j'(\varepsilon) = \frac{-q(z_j(\varepsilon))}{p'(z_j(\varepsilon)) + \varepsilon q'(z_j(\varepsilon))} \quad (17.9.2)$$

با قراردادن $\varepsilon = 0$ ، به دست می‌آوریم

$$\gamma_1 = z_j'(\varepsilon) = -\frac{q(z_j)}{p'(z_j)}$$

به رابطه (16.9.2) برمی‌گردیم

$$z_j(\varepsilon) = z_j - \frac{q(z_j)}{p'(z_j)}\varepsilon + \sum_{l=2}^{\infty} \gamma_l \varepsilon^l$$

$$\left| z_j(\varepsilon) - \left[z_j - \frac{q(z_j)}{p'(z_j)}\varepsilon \right] \right| \leq K\varepsilon^2 \quad |\varepsilon| \leq \varepsilon. \quad (18.9.2)$$

برای بعضی ثابتهای $\varepsilon_0 > 0$ و $K > 0$. در برآورد $z_j(\varepsilon)$ برای مقدار کوچک ε ، از رابطه زیر استفاده می‌کنیم

$$z_j(\varepsilon) \doteq z_j - \frac{q(z_j)}{p'(z_j)} \varepsilon \quad (19.9.2)$$

ضریب ε نشان می‌دهد که تغییر $z_j(\varepsilon)$ نسبت به ε چه اندازه سریع است؛ اگر این نسبت زیاد باشد ریشه z_j را بدوضع گویند.

حالت ۲. m چندگانگی z_j از 1 بزرگتر است. با به‌کار بردن تکنیکهای مربوط به حالت (۱)، برای یک مقدار $\varepsilon_0 > 0$ و $K > 0$ می‌توانیم به‌دست آوریم

$$|z_j(\varepsilon) - [z_j + \gamma_1 \varepsilon^{1/m}]| \leq K |\varepsilon|^{2/m} \quad |\varepsilon| \leq \varepsilon_0. \quad (20.9.2)$$

m مقدار ممکن برای γ_1 وجود دارد که m ریشه مختلط

$$\gamma_1^m = \frac{-m!q(z_j)}{p^{(m)}(z_j)}$$

هستند.

مثال چندجمله‌ای ساده زیر را در نظر می‌گیریم

$$\begin{aligned} p(x) &= (x-1)(x-2)\dots(x-7) \\ &= x^7 - 28x^6 + 322x^5 - 1960x^4 + 6769x^3 - 13132x^2 \\ &\quad + 13068x - 5040 \end{aligned} \quad (21.9.2)$$

برای اختلال، می‌گیریم

$$q(x) = x^6 \quad \varepsilon = -0.002$$

پس برای ریشه $z_j = j$

$$p'(z_j) = \prod_{l \neq j} (j-l) \quad q(z_j) = j^6$$

از (۱۹.۹.۲)، برآورد زیر را داریم

$$z_j(\varepsilon) \doteq j + \frac{0.002j^6(-1)^{j-1}}{(j-1)!(7-j)!} = j + \delta(j) \quad (22.9.2)$$

جدول ۱۶.۲ مقادیر $\delta(j)$ از (۲۲.۹.۲)

j	$\delta(j)$
۱	$۲,۷۸E-۶$
۲	$-۱,۰۷E-۳$
۳	$۳,۰۴E-۲$
۴	$-۲,۲۸E-۱$
۵	$۶,۵۱E-۱$
۶	$-۷,۷۷E-۱$
۷	$۳,۲۷E-۱$

جدول ۱۷.۲ ریشه‌های $p(x; \varepsilon)$ برای (۲۱-۹-۲)

j	$z_j(\varepsilon)$	$z_j(\varepsilon) - z_j(0)$
۱	$۱,۰۰۰۰۰۰۲۸$	$۲,۸۰E-۶$
۲	$۱,۹۹۸۹۳۸۲$	$-۱,۰۶E-۳$
۳	$۳,۰۳۳۱۲۵۳$	$۳,۳۱E-۲$
۴	$۳,۸۱۹۵۶۹۲$	$-۱,۸۰E-۱$
۵	$۵,۴۵۸۶۷۵۸ + ۰,۵۴۰۱۲۵۷۸i$	
۶	$۵,۴۵۸۶۷۵۸ - ۰,۵۴۰۱۲۵۷۸i$	
۷	$۷,۲۳۳۰۱۲۸$	$۲,۳۳E-۱$

مقادیر عددی $\delta(j)$ در جدول ۱۶.۲، داده شده‌اند. خطای نسبی در ضریب x^6 برابر است با $7,1E-5 = \frac{z_j(0.2) - z_j(0)}{28}$ ، ولی خطاهای نسبی در ریشه‌ها بسیار بزرگتر است. در واقع، اندازه بعضی اختلالات $\delta(j)$ ما را به شک می‌اندازد که آیا برآورد خطی (۲۲.۹.۲) معتبر است؟ ریشه‌های واقعی $p(x) + \varepsilon q(x)$ در جدول ۱۷.۲ داده شده‌اند، که به اختلالات پیش‌بینی شده خیلی نزدیک‌اند. بیشترین انحراف در ریشه‌ها، به ازای $j = 5$ و $j = 6$ رخ داده است. این ریشه‌ها مختلط‌اند که با برآورد خطی (۲۲.۹.۲) پیش‌بینی نشده بود. در این دو حالت، ε در خارج همگرایی سری توانی (۱۶.۹.۲) واقع است، زیرا این سری فقط دارای ضرایب حقیقی

$$\gamma_l = \frac{1}{l!} z_j^{(l)}(0)$$

است که از مشتق‌گیری (۱۷.۹.۲) به دست می‌آیند.

یک چندجمله‌ای را که ریشه‌های آن نسبت به تغییرات نسبی کوچک ضرایب ناپایدارند، بد وضع گوئیم. خیلی از این چندجمله‌یها به‌طور طبیعی در کاربردها، پدید می‌آیند. مثال قبلی، دشواری تعیین بدوضع بودن یا نبودن یک چندجمله‌ای را فقط با یک بررسی سطحی نشان می‌دهد.

تقلیل چندجمله‌ای. مسألهٔ دیگری که رخ می‌دهد تقلیل یک چندجمله‌ای به یک چندجمله‌ای از درجهٔ کمتر است، فرایندی که به دنبال (۱۲.۹.۲) تعریف شده است. چون صفرها، دقیقاً پیدا نمی‌شوند، چندجمله‌ای از درجهٔ پایینتر (۱۱.۹.۲) که با خارج کردن آخرین ریشه به‌دست آمده، معمولاً در تمام ضرایب خود دارای خطاست. از مثال اخیر روشن می‌شود، که این امر می‌تواند یک اختلال قابل توجهی در ریشه‌های رده‌هایی از چندجمله‌یها ایجاد نماید. ویلکینسن^۱ (۱۹۶۳) اثرات این تقلیل را تحلیل نموده و استراتژی کلی زیر را توصیه نموده است: (۱) مسأله را با پیدا کردن ریشه‌های با اندازهٔ کوچکتر شروع، و به پیدا کردن ریشه‌های بزرگتر ختم کنید؛ (۲) پس از اینکه تقریب تمام ریشه‌ها را یافتید، با استفاده از مقادیر یافته شده به‌عنوان حدسهای اولیه، و به‌کار بردن چندجمله‌ای اصلی، حل را دوباره تکرار نمایید. بحث کامل را می‌توان در ویلکینسن (۱۹۶۳، ۵۵-۶۵) پیدا کرد.

مثال یافتن ریشه‌های چندجمله‌ای از درجهٔ ۶ لاگر را در نظر می‌گیریم.

$$p(x) = x^6 - 36x^5 + 450x^4 - 2400x^3 + 5400x^2 - 4320x + 720$$

الگوریتم نیوتن بخش اخیر برای یافتن ریشه‌ها به‌کار گرفته شده و با پذیرفتن هر ریشهٔ جدید عمل تقلیل انجام شده است. ریشه‌ها به دو طریق محاسبه شده‌اند: (۱) از بزرگتر به کوچکتر و (۲) از کوچکتر به بزرگتر. محاسبات با حساب با دقت ساده با رایانهٔ IBM-360 انجام پذیرفته و نتایج عددی در جدول ۱۸.۲ داده شده‌اند. مقایسهٔ ستونهایی که با روش‌های (۱) و (۲) مشخص شده‌اند، بروشنی برتری روشی را که ریشه‌ها را در جهت صعودی محاسبه می‌کند نشان می‌دهد. اگر نتایج روش (۱) به‌عنوان حدسهای اولیه برای بارست بعدی با چندجمله‌ای اصلی به‌کار برده شوند ریشه‌های تقریبی که پیدا می‌شوند دقیقتر از جوابهای روش (۲) هستند؛ ستون جدول را که با روش (۳) در بالای ستون مشخص شده، ببینید. این جدول مجدداً اهمیت بارستن با چندجمله‌ای اصلی را برای از بین بردن اثرات فرایند تقلیل نشان می‌دهد. روشهای دیگری برای تقلیل یک چندجمله‌ای وجود دارند که یکی از آنها پیدا کردن ریشهٔ بزرگتر در ابتدای امر است. برای یک بحث کامل پیترز^۲ و ویلکینسن (۱۹۷۱، بخش ۵) را ببینید. یک الگوریتم برای تقلیل مرکب داده شده است که احتیاج به پیدا کردن ریشه‌ها به ترتیبی خاص را

جدول ۱۸.۲ مثال از تقلیل کثیرالجمله

واقعی	روش (۱)	روش (۲)	روش (۳)
۱۵,۹۸۲۸۷	۱۵,۹۸۲۸۷	۱۵,۹۸۲۷۹	۱۵,۹۸۲۸۷
۹,۸۳۷۴۶۷	۹,۸۳۷۴۷۱	۹,۸۳۷۴۶۹	۹,۸۳۷۴۶۷
۵,۷۷۵۱۴۴	۵,۷۷۵۷۶۴	۵,۷۷۵۲۰۷	۵,۷۷۵۱۴۴
۲,۹۹۲۷۳۶	۲,۹۹۱۰۸۰	۲,۹۹۲۷۱۰	۲,۹۹۲۷۳۶
۱,۱۸۸۹۳۲	۱,۱۹۰۹۳۷	۱,۱۸۸۹۳۲	۱,۱۸۸۹۳۲
۰,۲۲۲۸۴۶۶	۰,۲۲۱۹۴۲۹	۰,۲۲۲۸۴۶۶	۰,۲۲۲۸۴۶۶

منتفی می‌سازد. در این مقاله، مولفان، قوانین به‌کار بردن تقلیل ضمنی

$$q(x) = \frac{p(x)}{(x - z_1) \dots (x - z_r)}$$

را برای حذف ریشه‌های z_1, \dots, z_r که قبلاً محاسبه شده‌اند، نیز مورد بحث قرار می‌دهند. ما این بحث را قبلاً، در (۴.۴.۲)، در رابطه با روش مولر آورده‌ایم.

روشهای کلی ریشه‌یابی چندجمله‌یها. الگوریتم‌های ریشه‌یابی زیادی به‌ویژه برای چندجمله‌یها طراحی شده‌اند. تعداد زیادی از آنها مشروحاً در کتابهای دژون^۱ و هنریچی (۱۹۶۹) و هنریچی (۱۹۷۴، فصل ۶) و هاوسهولدر (۱۹۷۰) آمده است. انواع این روشها آنقدر زیادند که نمی‌توانیم در اینجا به شرح آنها بپردازیم.

در یک رده وسیع از این روشهای مهم قضایای جایابی را که در (۳.۹.۲) - (۶.۹.۲) شرح داده شده‌اند به‌کار می‌برند، تا به‌طور بارسستی، ریشه‌ها را به ناحیه‌های مجزا و همیشه کوچکتر، که اغلب دایره‌اند، جدا کنند. معروفترین این روشها، احتمالاً روش لمر-شور^۲ است [هاوسهولدر (۱۹۷۰)، بخش ۲-۷] را ببینید. این روشها بطور خطی همگرا هستند، و به همین دلیل، اغلب آنها را با روشهای همگرای سریعت، مانند روش نیوتن، ترکیب می‌کنند. وقتی که ریشه‌ها به ناحیه‌های مجزا، از هم جدا شدند، روشهای سریعتی برای پیدا کردن ریشه در هر ناحیه به‌کار گرفته می‌شود. برای بحث کلی در مورد این روشهای ریشه‌یابی به هنریچی (۱۹۷۴)، بخش ۶-۱۰) مراجعه کنید.

روشهای دیگری که صورت الگوریتم‌های گسترده - کار بردی پیدا کرده‌اند روش جنکینز و تراوب^۳ و روش لاگر است. برای اولی هاوسهولدر (۱۹۷۰، ص ۱۷۳)، جنکینز و تراوب (۱۹۷۰)، (۱۹۷۲) را ببینید. برای روش لاگر، به هاوسهولدر (۱۹۷۰، بخش ۴-۵) و کاهان^۴ (۱۹۶۷) مراجعه کنید.

یک روش عددی، آسان در کاربرد، براین پایه بنا شده است که بتوان ویژه مقدارهای ماتریس را محاسبه کرد. برای یک چندجمله‌ای داده شده $p(x)$ ممکن است بتوان به سادگی ماتریسی ساخت که $p(x)$ چندجمله‌ای مشخصه آن باشد (مسئله ۲ از فصل ۹ را ببینید). چون نرم افزار عالی برای حل مسئله ویژه مقدار وجود دارد، این نرم افزار را می توان برای پیدا کردن ریشه های چندجمله‌ای $p(x)$ به کار برد.

۱۰.۲ دستگاه معادلات غیرخطی

در این بخش و بخش بعد سروکار ما با حل عددی دستگاه معادلات غیرخطی چند متغیره است. این مسائل کاربردهای گسترده ای دارند و به شکلهای گوناگونی هستند. روشهای متنوعی برای حل چنین دستگاهها وجود دارد، لذا ما فقط خود موضوع را معرفی می کنیم و نظریه های کلی و بعضی روشهای عددی را که به سادگی قابل برنامه نویسی هستند ارائه خواهیم داد. برای بسط کامل تحلیل عددی حل دستگاههای غیرخطی، به چند قضیه از جبر خطی عددی نیاز داریم که تا فصل ۷-۹ آورده نشده اند. برای سادگی بیان و بهتر فهمیدن، این نظریه فقط برای دو معادله ارائه شده است:

$$f_1(x_1, x_2) = 0 \quad f_2(x_1, x_2) = 0 \quad (۱.۱۰.۲)$$

وقتی که مفاهیم اساسی روشن شد، تعمیم مطلب برای n معادله n مجهولی ساده خواهد بود. برای کمک بیشتر صورت و جواب (۱.۱۰.۲) را با نماد برداری در نظر می گیریم:

$$f(\mathbf{x}) = 0 \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad f(\mathbf{x}) = \begin{bmatrix} f_1(x_1, x_2) \\ f_2(x_1, x_2) \end{bmatrix} \quad (۲.۱۰.۲)$$

به حل (۱.۱۰.۲) می توان مانند یک فرایند دو مرحله ای نگاه کرد: (۱) خمهای صفر رویه های $z = f_1(x_1, x_2)$ و $z = f_2(x_1, x_2)$ را در صفحه $x_1 x_2$ پیدا می کنیم، و (۲) نقاط برخورد این خمهای صفر را در صفحه $x_1 x_2$ به دست می آوریم. این دیدگاه در بخش بعد برای تعمیم روش نیوتن برای حل (۱.۱۰.۲) به کار رفته است.

نظریه نقطه ثابت. با تعمیم قسمتی از نظریه بارسستی نقطه ثابت، بخش ۲-۵، مطلب را آغاز می کنیم. فرض می کنیم مسئله ریشه یابی (۱.۱۰.۲) به شکل هم ارز آن به صورت زیر دوباره نوشته شده است

$$x_1 = g_1(x_1, x_2) \quad x_2 = g_2(x_1, x_2) \quad (۳.۱۰.۲)$$

جواب آن را با نماد

$$\alpha = \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix}$$

نشان می‌دهیم. بارستی نقطه ثابت

$$x_{1,n+1} = g_1(x_{1,n}, x_{2,n}) \quad x_{2,n+1} = g_2(x_{1,n}, x_{2,n}) \quad (۴.۱۰.۲)$$

را مطالعه می‌کنیم. با استفاده از نماد برداری می‌نویسیم

$$\mathbf{x}_{n+1} = \mathbf{g}(\mathbf{x}_n) \quad (۵.۱۰.۲)$$

که در آن

$$\mathbf{x}_n = \begin{bmatrix} x_{1,n} \\ x_{2,n} \end{bmatrix} \quad \mathbf{g}(\mathbf{x}) = \begin{bmatrix} g_1(x_1, x_2) \\ g_2(x_1, x_2) \end{bmatrix}$$

مثال حل دستگاه

$$f_1 \equiv 3x_1^2 + 4x_2^2 - 1 = 0 \quad f_2 \equiv x_2^2 - 8x_1^2 - 1 = 0 \quad (۶.۱۰.۲)$$

برای جواب α نزدیک به $(x_1, x_2) = (-0.5, 0.25)$ را در نظر می‌گیریم. این دستگاه را با بارست زیر حل می‌کنیم

$$\begin{bmatrix} x_{1,n+1} \\ x_{2,n+1} \end{bmatrix} = \begin{bmatrix} x_{1,n} \\ x_{2,n} \end{bmatrix} - \begin{bmatrix} 0.016 & -0.17 \\ 0.52 & -0.26 \end{bmatrix} \begin{bmatrix} 3x_{1,n}^2 + 4x_{2,n}^2 - 1 \\ x_{2,n}^2 - 8x_{1,n}^2 - 1 \end{bmatrix} \quad (۷.۱۰.۲)$$

منشأ این صورت جدید (۶.۱۰.۲) بعداً داده خواهد شد. نتایج عددی (۷.۱۰.۲) در جدول ۱۹.۲ داده شده است. واضح است که این بارستها به سرعت همگرا هستند.

برای تحلیل همگرایی (۵.۱۰.۲)، دو معادله (۴.۱۰.۲) را از معادلات متناظر

$$\alpha_1 = g_1(\alpha_1, \alpha_2) \quad \alpha_2 = g_2(\alpha_1, \alpha_2)$$

جدول ۱۹.۲ مثال (۷.۱۰.۲) از بارست نقطه - ثابت

n	$x_{1,n}$	$x_{2,n}$	$f_1(x_{1,n}, x_{2,n})$	$f_2(x_{1,n}, x_{2,n})$
۰	-۰٫۵	۰٫۲۵	۰٫۰	۱٫۵۶E - ۲
۱	-۰٫۴۹۷۳۴۳۷۵۰	۰٫۲۵۴۰۶۲۵۰۰	۲٫۴۳E - ۴	۵٫۴۶E - ۴
۲	-۰٫۴۹۷۲۵۴۷۹۴	۰٫۲۵۴۰۷۷۹۲۲	۹٫۳۵E - ۶	۲٫۱۲E - ۵
۳	-۰٫۴۹۷۲۵۱۳۴۳	۰٫۲۵۴۰۷۸۵۶۶	۳٫۶۴E - ۷	۸٫۲۶E - ۷
۴	-۰٫۴۹۷۲۵۱۲۰۸	۰٫۲۵۴۰۷۸۵۹۲	۱٫۵۰E - ۸	۳٫۳۰E - ۸

که شامل جوابهای درست α هستند، کم می‌کنیم. قضیه مقدار میانی برای توابع دو متغیره (قضیه ۵-۱ با $n = 1$) را برای این تفاضلات به‌کار می‌بریم تا به دست آوریم

$$\alpha_i - x_{i,n+1} = \frac{\partial g_i(\xi_{1,n}^{(i)}, \xi_{2,n}^{(i)})}{\partial x_1} (\alpha_1 - x_{1,n}) + \frac{\partial g_i(\xi_{1,n}^{(i)}, \xi_{2,n}^{(i)})}{\partial x_2} (\alpha_2 - x_{2,n})$$

برای $i = 1, 2$. نقاط $\xi_n^{(i)} = (\xi_{1,n}^{(i)}, \xi_{2,n}^{(i)})$ روی قطعه خطی است که α را به x_n متصل می‌سازد. در شکل ماتریسی، این معادلات خطا چنین خواهند شد

$$\begin{bmatrix} \alpha_1 - x_{1,n+1} \\ \alpha_2 - x_{2,n+1} \end{bmatrix} = \begin{bmatrix} \frac{\partial g_1(\xi_n^{(1)})}{\partial x_1} & \frac{\partial g_1(\xi_n^{(1)})}{\partial x_2} \\ \frac{\partial g_2(\xi_n^{(2)})}{\partial x_1} & \frac{\partial g_2(\xi_n^{(2)})}{\partial x_2} \end{bmatrix} \begin{bmatrix} \alpha_1 - x_{1,n} \\ \alpha_2 - x_{2,n} \end{bmatrix} \quad (۸.۱۰.۲)$$

گیریم G_n معرف ماتریس (۸.۱۰.۲) باشد در این صورت این معادله را می‌توان دوباره چنین نوشت

$$\alpha - x_{n+1} = G_n(\alpha - x_n) \quad (۹.۱۰.۲)$$

در اینجا بجاست که ماتریس ژاکوبی را برای توابع g_1 و g_2 معرفی کنیم:

$$G(x) = \begin{bmatrix} \frac{\partial g_1(x)}{\partial x_1} & \frac{\partial g_1(x)}{\partial x_2} \\ \frac{\partial g_2(x)}{\partial x_1} & \frac{\partial g_2(x)}{\partial x_2} \end{bmatrix} \quad (۱۰.۱۰.۲)$$

در (۹.۱۰.۲)، اگر x_n نزدیک به α باشد، G_n نزدیک به $G(\alpha)$ می‌شود. این امر موجب می‌شود اندازه یا نرم ماتریس $G(\alpha)$ در تحلیل همگرایی (۹.۱۰.۲) نقش اساسی داشته باشد. ماتریس $G(\alpha)$ نقش $g'(\alpha)$ را در نظریه بخش ۵.۲، بازی می‌کند. برای اندازه‌گیری اندازه خطاهای $\alpha - x_n$ و ماتریس‌های G_n و $G(\alpha)$ از نرم‌های برداری و ماتریسی (۱۶.۱.۱) و (۱۹.۱.۱) فصل ۱ استفاده می‌کنیم.

قضیه ۹.۲.۲. D مجموعه‌ای بسته، کراندار و محدب در صفحه باشد. D را محدب گوییم، اگر برای هر دو نقطه در D ، قطعه خطی هم که آنها را به هم وصل می‌کند در D باشد. فرض کنید مؤلفه‌های $g(x)$ در تمام نقاط D پیوسته مشتق‌پذیر باشند. همچنین فرض کنید

$$g(D) \subset D \quad (۱۱.۱۰.۲)$$

۲.

$$\lambda = \max_{x \in D} \|G(x)\|_{\infty} < 1 \quad (۱۲.۱۰.۲)$$

آنگاه

(الف) $x = g(x)$ یک ریشه یکتای $\alpha \in D$ دارد.

(ب) برای هر نقطه اولیه داده شده $x_0 \in D$ ، بارست $(۵.۱۰.۲)$ در D به α می‌گراید.

(ج)

$$\|\alpha - x_{n+1}\|_{\infty} \leq (\|G(\alpha)\|_{\infty} + \varepsilon_n) \|\alpha - x_n\|_{\infty} \quad (۱۳.۱۰.۲)$$

با $\varepsilon_n \rightarrow 0$ هرگاه $n \rightarrow \infty$

برهان (الف) وجود نقطه ثابت α را می‌توان با اثبات اینکه دنباله بارستهای $\{x_n\}$ از $(۵.۱۰.۲)$ در D همگراست، نشان داد. ما اثبات وجود را به عنوان یک مسأله به خواننده واگذار می‌کنیم و فقط یکتایی α را ثابت می‌نماییم.

فرض کنید α و β هر دو، نقطه‌های ثابت $g(x)$ در D باشند. در این صورت

$$\alpha - \beta = g(\alpha) - g(\beta) \quad (۱۴.۱۰.۲)$$

قضیه مقدار میانی را برای مؤلفه نام به کار می‌بریم تا به دست آوریم

$$g_i(\alpha) - g_i(\beta) = \nabla g_i(\xi^{(i)})(\alpha - \beta) \quad i = 1, 2 \quad (۱۵.۱۰.۲)$$

که در آن

$$\nabla g_i(x) = \left[\frac{\partial g_i}{\partial x_1} \quad \frac{\partial g_i}{\partial x_2} \right]$$

و $\xi^{(i)} \in D$ روی پاره خطی که α را به β وصل می‌کند قرار دارد. چون $\lambda < 1$ و $\|G(x)\| \leq \lambda$ ، با توجه به تعریف نرم خواهیم داشت،

$$\left| \frac{\partial g_i(x)}{\partial x_1} \right| + \left| \frac{\partial g_i(x)}{\partial x_2} \right| \leq \lambda < 1, \quad x \in D, \quad i = 1, 2$$

از ترکیب این رابطه با (۱۵.۱۰.۲)،

$$\begin{aligned} |g_i(\alpha) - g_i(\beta)| &\leq \lambda \|\alpha - \beta\|_\infty \\ \|g(\alpha) - g(\beta)\|_\infty &\leq \lambda \|\alpha - \beta\|_\infty \end{aligned} \quad (۱۶.۱۰.۲)$$

این رابطه وقتی با (۱۴.۱۰.۲) ترکیب شود نتیجه می‌دهد؛

$$\|\alpha - \beta\|_\infty \leq \lambda \|\alpha - \beta\|_\infty$$

که تنها به‌ازای $\alpha = \beta$ امکان دارد، و یکتایی α را در D نشان می‌دهد.

(ب) شرط (۱۱.۱۰.۲) تضمین می‌کند که اگر $x_0 \in D$ هر $x_n \in D$ اکنون $x_{n+1} = g(x_n)$ را از $\alpha = g(\alpha)$ کم می‌کنیم و به‌دست می‌آوریم

$$\alpha - x_{n+1} = g(\alpha) - g(x_n)$$

نتیجه (۱۶.۱۰.۲) برای هر دو نقطه‌ای از D صادق است. با به‌کار بردن آن

$$\|\alpha - x_{n+1}\|_\infty \leq \lambda \|\alpha - x_n\|_\infty \quad (۱۷.۱۰.۲)$$

و به‌طور استقرایی

$$\|\alpha - x_n\|_\infty \leq \lambda^n \|\alpha - x_0\|_\infty \quad (۱۸.۱۰.۲)$$

چون $\lambda < 1$ ، رابطهٔ اخیر نشان می‌دهد که وقتی $n \rightarrow \infty$ ، آنگاه $x_n \rightarrow \alpha$. (ج) با توجه به (۹.۱۰.۲) و استفاده از (۲۱.۱.۱)،

$$\|\alpha - x_{n+1}\|_\infty \leq \|G_n\|_\infty \|\alpha - x_n\|_\infty \quad (۱۹.۱۰.۲)$$

وقتی $n \rightarrow \infty$ ، نقاط $\xi_n^{(i)}$ که در محاسبهٔ G_n به‌کار رفته‌اند همگی به α میل خواهند کرد، زیرا روی پاره‌خط واصل x_n و α قرار دارند پس وقتی $n \rightarrow \infty$ ، آنگاه $\|G_n\|_\infty \rightarrow \|G(\alpha)\|_\infty$. نتیجهٔ (۱۳.۱۰.۲) را می‌توان از (۱۹.۱۰.۲) با فرض $\varepsilon_n = \|G_n\|_\infty - \|G(\alpha)\|_\infty$ به‌دست آورد. ■

قضیه قبل، تعمیم قضیه ۶.۲ به توابع دو متغیره است. فرع ذیل تعمیم قضیه ۷.۲ است.

فرع ۱۰.۲ گزیریم α یک نقطه ثابت $g(x)$ باشد و فرض می‌کنیم مؤلفه‌های $g(x)$ در یک همسایگی α پیوسته مشتقپذیر باشند. به علاوه فرض می‌کنیم

$$\|G(\alpha)\|_{\infty} < 1 \quad (20.10.2)$$

در این صورت برای x_n که به اندازه کافی نزدیک به α انتخاب شده باشد، بارست $x_{n+1} = g(x_n)$ به سمت α می‌گراید و نتایج قضیه ۹.۲ در یک ناحیه بسته، کراندار و محدب حول α معتبر است. ■

ما اثبات این فرع را به عنوان یک مسأله واگذار می‌کنیم. بر پایه نتایج فصل ۷، همگرایی خطی x_n به α باز هم درست است اگر قدر مطلق همه ویژه مقدرهای $G(\alpha)$ کوچکتر از یک باشند، که می‌توان نشان داد فرض ضعیفتری از (۲۰.۱۰.۲) است.

مثال مثال قبلی (۷.۱۰.۲) را ادامه می‌دهیم. به سادگی

$$G(\alpha) \doteq \begin{bmatrix} 0.38920 & 0.000401 \\ 0.008529 & -0.006613 \end{bmatrix}$$

حساب می‌شود. بنابراین

$$\|G(\alpha)\|_{\infty} \doteq 0.393$$

پس شرط (۲۰.۱۰.۲) از قضیه برقرار است. از (۱۳.۱۰.۲)، رابطه

$$\frac{\|\alpha - x_{n+1}\|_{\infty}}{\|\alpha - x_n\|_{\infty}} \leq \|G_n\|_{\infty} \doteq 0.393$$

برای تمام مقادیر به اندازه کافی بزرگ n ، تقریباً درست است.

فرض کنید A یک ماتریس ثابت و ناتکین 2×2 باشد. پس می‌توانیم (۱.۱۰.۲) را به شکل

زیر دوباره بنویسیم

$$x = x + Af(x) \equiv g(x) \quad (21.10.2)$$

مثال (۷.۱۰.۲) این روند را نشان می‌دهد. برای آنکه ببینیم A چه شرایطی باید داشته باشد، ماتریس ژاکوبی را تشکیل می‌دهیم. به سادگی

$$G(\mathbf{x}) = I + AF(\mathbf{x})$$

به دست می‌آید. که $F(\mathbf{x})$ ماتریس ژاکوبی f_1 و f_2 است.

$$F(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} \end{bmatrix} \quad (۲۲.۱۰.۲)$$

می‌خواهیم A را به گونه‌ای انتخاب کنیم که (۲۰.۱۰.۲) برقرار شود. و برای همگرایی سریع، می‌خواهیم $\|G(\alpha)\|_\infty \doteq 0$ یا

$$A \doteq -F(\alpha)^{-1}$$

ماتریس در (۷.۱۰.۲) با استفاده از

$$A \doteq -F(\mathbf{x}_n)^{-1}$$

بدین طریق انتخاب شده است. از این رابطه به فکر استفاده از پیوسته نو کردن A ، یعنی $A = -F(\mathbf{x}_n)^{-1}$ ، می‌افتیم. روشی که به دست می‌آید چنین است

$$\mathbf{x}_{n+1} = \mathbf{x}_n - F(\mathbf{x}_n)^{-1} \mathbf{f}(\mathbf{x}_n) \quad n \geq 0 \quad (۲۳.۱۰.۲)$$

این روش را در بخش بعد مطالعه می‌کنیم.

۱۱.۲ روش نیوتن برای دستگاههای غیرخطی

همانند روش نیوتن برای حل یک معادله، برای حل دستگاه معادلات غیر خطی بیش از یک راه برای نگرستن و پیدا کردن روش نیوتن وجود دارد. ما با یک روش تحلیلی آغاز می‌کنیم و سپس یک نگرش هندسی ارائه خواهیم داد.

قضیه تیلر برای توابع دومتغیره را برای هر یک از معادله‌های $f_i(x_1, x_2) = 0$ با بسط $f_i(\alpha)$ ، $i = 1, 2$ ، حول \mathbf{x}_0 به کار می‌بریم:

$$0 = f_i(\alpha) = f_i(\mathbf{x}_0) + (\alpha_1 - x_{1,0}) \frac{\partial f_i(\mathbf{x}_0)}{\partial x_1} + (\alpha_2 - x_{2,0}) \frac{\partial f_i(\mathbf{x}_0)}{\partial x_2} + \frac{1}{2} \left[(\alpha_1 - x_{1,0}) \frac{\partial}{\partial x_1} + (\alpha_2 - x_{2,0}) \frac{\partial}{\partial x_2} \right]^2 f_i(\xi^{(i)}) \quad (۱.۱۱.۲)$$

$\xi^{(i)}$ روی پاره خط و اصل بین x_0 و α واقع است. اگر از جمله‌های درجه دوم صرف نظر کنیم تقریب زیر را به دست می‌آوریم

$$\begin{aligned} 0 &\doteq f_1(x_0) + (\alpha_1 - x_{1,0}) \frac{\partial f_1(x_0)}{\partial x_1} + (\alpha_2 - x_{2,0}) \frac{\partial f_1(x_0)}{\partial x_2} \\ 0 &\doteq f_2(x_0) + (\alpha_1 - x_{1,0}) \frac{\partial f_2(x_0)}{\partial x_1} + (\alpha_2 - x_{2,0}) \frac{\partial f_2(x_0)}{\partial x_2} \end{aligned} \quad (2.11.2)$$

در شکل ماتریسی

$$0 \doteq f(x_0) + F(x_0)(\alpha - x_0) \quad (3.11.2)$$

که $F(x_0)$ ماتریس ژاکوبی f است که در (۲۲.۱۰.۲) داده شده است. اگر رابطه فوق را نسبت به α حل کنیم داریم:

$$\alpha \doteq x_0 - F(x_0)^{-1} f(x_0) \equiv x_1$$

تقریب x_1 بایستی یک بهبود x_0 باشد، به شرطی که x_0 به اندازه کافی به α نزدیک انتخاب شده باشد. این عمل ما را به روشی بارستی می‌رساند که ابتدا در پایان بخش قبلی به دست آمده بود

$$x_{n+1} = x_n - F(x_n)^{-1} f(x_n) \quad n \geq 0 \quad (4.11.2)$$

این روش نیوتن برای حل دستگاه غیرخطی $f(x) = 0$ است.

در کاربرد واقعی، $F(x_n)$ را وارون نمی‌کنیم، به‌ویژه برای دستگاه‌های بیش از دو معادله. به جای آن، یک دستگاه خطی را برای یک جمله اصلاحی x_n حل می‌کنیم:

$$F(x_n) \delta_{n+1} = -f(x_n)$$

$$x_{n+1} = x_n + \delta_{n+1} \quad (5.11.2)$$

این عمل زمان محاسبات را بسیار کوتاهتر می‌کند، زیرا حدود یک سوم عملیات وارون‌سازی $F(x_n)$ کار می‌برد. بخش‌های ۱.۸ و ۲.۸ را برای بحث در حل عددی دستگاه معادلات خطی ملاحظه کنید.

یک راه هندسی برای پیدا کردن روش نیوتن، شبیه به تقریب خط مماسی که برای یک معادله غیرخطی در بخش ۲.۲ به کار رفت، وجود دارد. نمودار معادله

$$z = f_i(x_0) + (x_1 - x_{1,0}) \frac{\partial f_i(x_0)}{\partial x_1} + (x_2 - x_{2,0}) \frac{\partial f_i(x_0)}{\partial x_2} \equiv p_i(x_1, x_2)$$

در فضا صفحه‌ای است مماس بر نمودار $z = f_i(x_1, x_2)$ در نقطه x_0 برای $i = 1, 2$. اگر x_0 نزدیک به α باشد، آنگاه این صفحات مماس باید تقریبات خوبی برای رویه‌های متناظر $z = f_i(x_1, x_2)$ در $x = (x_1, x_2)$ نزدیک به α باشند. در این صورت محل برخورد خم‌های صفر صفحات مماس $z = p_i(x_1, x_2)$ باید تقریب خوبی برای محل برخورد متناظر α از خم‌های صفر رویه‌های اصلی $z = f_i(x_1, x_2)$ باشد. این به نتیجه (۲.۱۱.۲) می‌انجامد. محل برخورد خم‌های صفر $z = p_i(x_1, x_2)$ ، $i = 1, 2$ ، نقطه x_1 است.

مثال دستگاه زیر را در نظر می‌گیریم

$$f_1 \equiv 4x_1^2 + x_2^2 - 4 = 0 \quad f_2 \equiv x_1 + x_2 - \sin(x_1 - x_2) = 0$$

دستگاه فقط دو ریشه دارد، یکی نزدیک به $(1, 0)$ و دیگری قرینه آن نسبت به مرکز نزدیک به $(-1, 0)$. با استفاده از (۴.۱۱.۲) با $x_0 = (1, 0)$ ، نتایج جدول ۲۰.۲ را به دست می‌آوریم.

تحلیل همگرایی. برای تحلیل همگرایی روش نیوتن، (۴.۱۱.۲)، آن را مانند یک روش بارستی نقطه ثابت با

$$g(x) = x - F(x)^{-1}f(x) \quad (۶.۱۱.۲)$$

در نظر می‌گیریم. همچنین فرض می‌کنیم

درمیان $F(\alpha)$ مخالف صفر است

که مشابه این فرض است که α یک ریشه ساده است وقتی با یک معادله تنها نظیر معادله مربوط به قضیه ۱.۲ سروکار داریم. در این صورت می‌توان نشان داد که ژاکوبی $G(x)$ از (۶.۱۱.۲) در $x = \alpha$ صفر است (مسئله ۵۳ را ببینید)؛ در نتیجه، شرط (۲۰.۱۰.۲) برقرار است.

جدول ۲۰.۲ مثال از روش نیوتن

n	$x_{1,n}$	$x_{2,n}$	$f_1(x_n)$	$f_2(x_n)$
۰	۱٫۰	۰٫۰	۰٫۰	۱٫۵۹E - ۱
۱	۱٫۰	-۰٫۱۰۲۹۲۰۷۱۵۴	۱٫۰۶E - ۲	۴٫۵۵E - ۳
۲	۰٫۹۹۸۶۰۸۷۵۹۸	-۰٫۱۰۵۵۳۰۷۲۳۹	۱٫۴۶E - ۵	۶٫۶۳E - ۷
۳	۰٫۹۹۸۶۰۶۹۴۴۱	-۰٫۱۰۵۵۳۰۴۹۲۳	۱٫۳۲E - ۱۱	۱٫۸۷E - ۱۲

در این صورت فرج ۲-۱۰ ایجاب می‌کند که x_n به α همگرا باشد به شرطی که x_0 به اندازه کافی به α نزدیک باشد. به علاوه می‌توان نشان داد که این بارسیتی از درجه دو است. به‌ویژه فرمولهای (۱.۱۱.۲) و (۴.۱۱.۲) را می‌توان ترکیب کرده نتیجه

$$\|\alpha - x_{n+1}\|_\infty \leq B \|\alpha - x_n\|_\infty^2 \quad n \geq 0 \quad (7.11.2)$$

را برای مقدار ثابتی مانند $B > 0$ به‌دست‌آورد.

شکلهای گوناگون روش نیوتن. روش نیوتن در مقایسه با سایر روشهای حل دستگاه معادلات غیرخطی، دارای محاسن و معایبی است. از جمله محاسن آن دارا بودن شکل ساده و انعطاف‌پذیری زیاد برای کاربرد آن در مسائل بسیار متنوع است. اگر نخواهیم زحمت محاسبه مشتقات جزئی در برنامه رایانه‌یی را به خود بدهیم می‌توانیم از تفاضلات تقریبی استفاده کنیم. به‌عنوان مثال، معمولاً از

$$\frac{\partial f_i(x_1, x_2)}{\partial x_1} \approx \frac{f_i(x_1 + \varepsilon, x_2) - f_i(x_1, x_2)}{\varepsilon} \quad (8.11.2)$$

با عدد خیلی کوچک ε استفاده می‌کنیم. برای بحث مفصل در انتخاب ε ، دنیس و اشناابل^۱ (۱۹۸۳، صص ۹۴-۹۹) را ببینید.

اولین عیب روش نیوتن این است که روشهای دیگری هستند که (۱) کار با آنها ارزانتر است، و یا (۲) برای بعضی از رده‌های مسائل کار با آنها آسانتر است. برای یک دستگاه m معادله m مجهولی غیرخطی، هر بارست برای روش نیوتن بطور کلی مستلزم $m^2 + m$ محاسبه تابع است. به علاوه، روش نیوتن برای هر بارست به حل یک دستگاه خطی m معادله m مجهولی نیاز دارد، که به بهای حدود $\frac{1}{4}m^2$ عمل حساب، برای هر دستگاه خطی تمام می‌شود. روشهای دیگری وجود دارند که در همگرایی ریاضی، یک اندازه یا تقریباً یک اندازه سرعت دارند ولی به تعداد کمتری محاسبه تابع و عملیات حساب رایانه‌ای نیاز دارند. این روشها را معمولاً روشهای نیوتنی نما، شبه‌نیوتنی یا نیوتنی اصلاح شده می‌گویند. برای یک توضیح کلی درباره بسیاری از این روشها دنیس و اشناابل (۱۹۸۳) را ببینید. یک اصلاح ساده روش نیوتن تثبیت ماتریس ژاکوبی برای چند، مثلاً k مرحله است؛

$$x_{rk+j+1} = x_{rk+j} - F(x_{rk})^{-1} f(x_{rk+j}) \quad j = 0, 1, \dots, k-1 \quad (9.11.2)$$

برای $r = 0, 1, 2, \dots$. این بدان معنی است که دستگاه خطی

$$F(x_{rk}) \delta_{rk+j+1} = -f(x_{rk+j})$$

$$x_{rk+j+1} = x_{rk+j} + \delta_{rk+j+1}, \quad j = 0, 1, \dots, k-1, \quad (10.11.2)$$

خیلی بهتر از روش اصلی نیوتن (۵.۱۱.۲) حل می‌شود. حل دستگاه خطی از مرتبه m در حدود $\frac{1}{2}m^2$ عمل در اولین حالت، $j = 0$ ، نیاز دارد. ولی در هر یک از حالت‌های بعدی، $j = 1, \dots, k-1$ ، فقط $\frac{1}{2}m^2$ عمل حساب برای حل آن لازم است. بخش ۱.۸ را برای توضیح کامل ببینید. سرعت همگرایی (۹.۱۱.۲) آهسته‌تر از سرعت همگرایی روش اصلی (۴.۱۱.۲) است، ولی زمان محاسبه واقعی روش اصلاح شده غالباً خیلی کمتر از زمان روش اصلی است. برای توضیحات بیشتر این مسأله پوترا و پتاک^۱ (۱۹۸۴، ص ۱۱۹) را ببینید.

یک مشکل دوم با روش نیوتن و بسیاری از روش‌های دیگر، این است که اغلب x_0 باید بطور قابل قبولی به α نزدیک باشد تا همگرایی حاصل شود. اصلاحاتی از روش نیوتن وجود دارند که برای انتخاب‌های بد x_0 ، همگرایی را تسریع می‌کنند. برای مثال، تعریف می‌کنیم

$$x_{n+1} = x_n + s d_n \quad d_n = -F'(x_n)^{-1} f(x_n) \quad (11.11.2)$$

و $s > 0$ را طوری انتخاب می‌کنیم که

$$\|f(x_n + s d_n)\|_2^2 = \sum_{j=1}^m [f_j(x_n + s d_n)]^2 \quad (12.11.2)$$

می‌نیم شود. انتخاب $s = 1$ در (۱۱.۱۱.۲) روش نیوتن را به دست می‌دهد، ولی این مقدار ممکن است بهترین انتخاب نباشد. در بعضی حالات، لازم است s بسیار کوچکتر از یک باشد، حداقل در ابتدا، تا همگرایی تضمین شود. برای بحث مفصل‌تر دنیس و اشناپل (۱۹۸۳، فصل ۶) را ببینید. برای تحلیلی از بعضی برنامه‌های حل دستگاه‌های غیرخطی، هیرت^۲ (۱۹۸۲) را ببینید. نامبرده درباره مشکلات چنین نرم‌افزارهایی نیز بحث نموده است.

۱۲.۲ بهینه‌سازی نامقید

بهبینه‌سازی به پیدا کردن ماکسیم یا مینیم یک تابع پیوسته $f(x_1, \dots, x_n)$ گفته می‌شود. این مسأله‌ای بسیار مهم است، به ویژه در مهندسی صنایع جدید، علم مدیریت و سایر زمینه‌ها. در این بخش از بعضی روش‌ها بحث می‌شود و محاسبه مینیم یا ماکسیم یک تابع $f(x_1, \dots, x_n)$ نشان داده می‌شود. در اینجا هیچ الگوریتم دقیقی داده نشده است زیرا این کار مفصلی را می‌طلبد.

برای عرضهٔ بسیاری از مسائل از نماد بردار استفاده شده است تا به‌طور کلی نتایج برای تعداد m متغیر به‌دست آید. ما فقط بهینه‌سازی غیر مقید را مورد توجه قرار می‌دهیم، که در آن هیچ محدودیتی برای (x_1, \dots, x_m) وجود ندارد. همچنین فقط برای سادگی فرض می‌کنیم $f(x_1, \dots, x_m)$ برای تمام مقادیر (x_1, \dots, x_m) معین است.

چون رفتار تابع $f(x)$ ممکن است خیلی تغییر کند، مسأله باید محدودتر شود. یک نقطهٔ α یک مینیمم موضعی اکید f خوانده می‌شود اگر برای همهٔ مقادیر x نزدیک به α و $x \neq \alpha$ نابرابری $f(x) > f(\alpha)$ برقرار باشد. ما خود را به یافتن مینیمم موضعی اکید $f(x)$ محدود می‌نماییم. معمولاً یک حدس اولیهٔ x_0 از α معلوم است و تابع $f(x)$ دوبار پیوسته مشتق‌پذیر نسبت به متغیرهای x_1, \dots, x_m فرض می‌شود.

بیانی دیگر به‌صورت دستگاه غیرخطی. با فرض مشتق‌پذیری، یک شرط لازم برای اینکه α مینیمم موضعی اکید باشد

$$\frac{\partial f(\alpha)}{\partial x_i} = 0 \quad i = 1, \dots, m \quad (1.12.2)$$

است. پس دستگاه غیرخطی

$$\frac{\partial f(x)}{\partial x_i} = 0 \quad i = 1, \dots, m \quad (2.12.2)$$

را می‌توان حل کرد و بررسی نمود که هر جواب محاسبه شده ماکسیمم یا مینیمم موضعی است یا هیچکدام. برای نمادگذاری، بردار گزادیان را معرفی می‌کنیم

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_m} \end{bmatrix}$$

با استفاده از این بردار، دستگاه (۲.۱۲.۲) به شکل فشرده‌تر زیر نوشته می‌شود

$$\nabla f(x) = 0 \quad (3.12.2)$$

برای حل (۳.۱۲.۲)، می‌توان از روش (۴.۱۱.۲) نیوتن و همچنین سایر روشهای ریشه‌یابی استفاده کرد. با به‌کار بردن دستگاههای غیرخطی روش نیوتن رابطهٔ زیر به‌دست می‌آید

$$x_{n+1} = x_n - H(x_n)^{-1} \nabla f(x_n) \quad n \geq 0 \quad (4.12.2)$$

که در آن $H(x)$ ماتریس هسه‌ای تابع f است،

$$H(x)_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}, \quad 1 \leq i, j \leq m$$

اگر α یک مینیمم موضعی اکید f باشد، با استفاده از قضیه (۱۲.۱۱.۱) تیلر می‌توان نشان داد که $H(\alpha)$ یک ماتریس ناتکین است؛ پس برای x های نزدیک به α ناتکین خواهد بود. برای همگرایی، تحلیل روش نیوتن در بخش قبل را می‌توان به‌کار برد و همگرایی مرتبه دوم x_n به α را مشروط بر آنکه x به اندازه کافی نزدیک به α انتخاب شده باشد، ثابت نمود.

معایب عمده روش بارستی (۴.۱۲.۲) همانهایی هستند که برای روش نیوتن در حل دستگاههای غیرخطی، در آخرین قسمت بخش قبل داده شد. روشهای بهبودسازی کاراتری وجود دارند که در آنها برای پیدا کردن تقریب α ، فقط از $f(x)$ و $\nabla f(x)$ استفاده می‌شود. ممکن است این روشها به تعداد بارستهای بیشتری نیاز داشته باشند ولی معمولاً زمان کلی محاسبات بسیار کمتر از زمان محاسبه در روش نیوتن خواهد بود. به علاوه این روشها در پی یافتن همگرایی، در مجموعه بزرگتری از مقادیر اولیه x هستند.

روشهای کاهش. فرض کنید می‌کوشیم یک تابع $f(x)$ را مینیمم کنیم. بیشتر روشها برای این کار، بر پایه فرایند بارستی دو مرحله‌ای کلی زیر بنا شده‌اند.

مرحله ک ۱: در نقطه x_n ، یک امتداد d_n می‌گیریم به‌گونه‌ای که وقتی x در امتداد d_n از x_n دور می‌شود، $f(x)$ کاهش یابد.

مرحله ک ۲: گیریم $x_{n+1} = x_n + s d_n$ ، که s به گونه‌ای انتخاب شده است که تابع

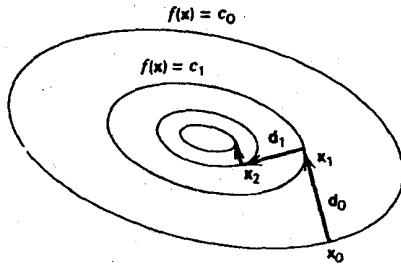
$$\varphi(s) = f(x_n + s d_n) \quad s \geq 0 \quad (5.12.2)$$

را مینیمم سازد. معمولاً s کوچکترین مینیمم نسبی مثبت $\varphi(s)$ انتخاب می‌شود. این گونه روشها، روشهای کاهش خوانده می‌شوند. با هر بارست

$$f(x_{n+1}) < f(x_n)$$

روشهای کاهش در شرایطی کلیتر از شرایط در روش (۴.۱۲.۲) نیوتن، ضامن همگرایی هستند. رویه تراز،

$$C = \{x \mid f(x) = f(x_0)\}$$



شکل ۸.۲. نمایش روش تندترین کاهش

را در نظر می‌گیریم، فقط تکه همبندی از آن، مثلاً C' را که شامل x باشد اختیار می‌کنیم. اگر C' کراندار و α را در درون خود داشته باشد، آنگاه روشهای کاهش در شرایط بسیار کلی، همگرا خواهند بود. این امر در شکل ۸.۲ برای حالت دو متغیره نشان داده شده است. چندین خم تراز $f(x_1, x_2) = c$ برای مجموعه‌ای از مقادیر c که به $f(\alpha)$ نزدیک می‌شوند، نشان داده شده است. بردارهای d_n امتدادهایی هستند که در امتداد آنها $f(x)$ کاهش می‌یابد. راههای زیادی برای انتخاب امتدادهای d_n وجود دارند و شناخته‌ترین آنها در زیر می‌آیند:

۱. روش تندترین کاهش. در اینجا $d_n = -\nabla f(x_n)$. این امتدادی است که در جهت آن $f(x)$ تندترین کاهش را وقتی از x_n دور می‌شود، خواهد داشت. این خط مشی خوبی در نزدیک x_n است ولی غالباً به یک خط مشی ضعیف برای همگرایی سریع به α تبدیل می‌شود.

۲. روشهای شبه نیوتنی. به این روشها می‌توان به عنوان تقریبهای روش (۴.۱۲.۲) نیوتن نگاه کرد. در این روشها از تقریبهای $H(x_n)$ و $H(x_n)^{-1}$ که محاسبه آنها آسان است، استفاده می‌شود و این روشها نیز کاهش می‌دهند. معروفترین مثال، روش دیویدسن - فلچر - پاول و روشهای برویدن^۱ هستند.

۳. روش گرادیان مزدوج. در این روش برای تولید امتدادهای d_n ، از تعمیم مفهوم در یک فضای برداری استفاده می‌شود. که d_n ها به طریق بهینه‌ای به تابع $f(x)$ که باید بهینه شود مربوط می‌شوند. در فصل ۸ روش گرادیان مزدوج برای حل دستگاههای معادله‌های خطی به‌کار گرفته شده است. روشهای دیگری برای مینیم کردن یک تابع وجود دارند، ولی تعداد آنها آن قدر زیاد است که نمی‌توان در اینجا آورد. به عنوان منابع کلی برای مفاهیم بالا، مراجعه کنید به: دنیس اشنابل (۱۹۸۳) و فلچر (۱۹۸۰)، لوتنبرگر^۲ (۱۹۸۴) و گیل^۳ و همکاران (۱۹۸۰). یک روش مهم و کاملاً متفاوت برای مینیم کردن یک تابع روش سارکی است که در نلدر^۵ و مید^۶ (۱۹۶۵) داده

1. Davidson-Fletcher-Powell

2. Broyden

3. Luenberger

4. Gill et al

5. Nelder

6. Mead

شده و بحث آن در گیل و همکاران (۱۹۸۱، ص ۹۴) و وودز^۱ (۱۹۸۵، فصل ۲) داده شده است. در این روش فقط از مقادیر تابع (نه مقادیر مشتق) استفاده می‌شود، و ظاهراً برای توابع نوفه‌ای بسیار مناسب است.

یک طرح مهم برای به وجود آوردن برنامه‌هایی برای حل مسائل بهینه‌سازی و دستگاه‌های غیرخطی در Argonne National Laboratory در جریان است. این بسته برنامه MINPACK خوانده می‌شود و نسخه اول آن در دسترس است [موره و همکاران (۱۹۸۰) و موره و همکاران (۱۹۸۴)] را ببینید. این بسته شامل برنامه‌هایی برای دستگاه‌های غیرخطی و مسائل کمترین مربعات غیرخطی است. نسخه‌های دیگر قرار است شامل برنامه‌هایی برای مسائل بهینه‌سازی مقید و نامقید باشد.

بحث در آثار خواندنی

نوشته‌های زیادی در مورد محاسبه ریشه‌های یک معادله وجود دارد. برای توضیحاتی بیشتر از آنچه که در این کتاب آمده است، به کتابهای هاوزولدر (۱۹۷۰)، آستروفسکی (۱۹۷۳)، تراؤب (۱۹۶۴) مراجعه کنید. روش نیوتن یکی از روشهایی است که بسیار به کار گرفته شده است و گسترش آن مرهون افراد زیادی است. برای یک گزارش تاریخی در سهم نیوتن، رافسون، و کوشی در این روش به گولداستاین (۱۹۷۷، صص ۶۴ و ۲۷۸) مراجعه کنید.

برای برنامه‌های رایانه‌ای، اغلب افراد، هنوز برنامه‌ای را که شخصاً تهیه دیده و مناسب برای کار خودشان است به کار می‌برند. ولی، باید یکی از برنامه‌هایی را که برای منظوره‌های کلی در سالهای اخیر نوشته و توسعه داده شده است و در کتابخانه‌های نرم‌افزاری بازرگانی در دسترس هستند در نظر گرفت. این برنامه‌ها معمولاً دقیق و کارا هستند و استفاده از آنها آسان است. در میان این برنامه‌های کلی آنها که از برنت (۱۹۷۳) و دکر (۱۹۶۹) گرفته شده‌اند، مشهورتر از همه هستند، و توسعه بیشتر آنها هنوز ادامه دارد، مانند له^۲ (۱۹۸۵) در کتابخانه‌های رایانه‌ای IMSL و NAG این برنامه‌ها و برنامه‌های عالی ریشه‌یابی دیگر موجود است.

ریشه‌یابی چندجمله‌یها موضوعی بسیار قدیمی است که به زمانهای یونان باستان می‌رسد. روشها و نوشته‌های زیادی در مورد آنها وجود دارد، و در دو سه دهه اخیر روشهای جدید زیادی برای آنها یافته‌اند. برای آشنایی با این موضوع کتابهای مقدماتی در این زمینه، کتابهای هنریچی^۳ و دوژون^۴ (۱۹۶۹)، هنریچی (۱۹۷۴، فصل ۶)، هاوزولدر (۱۹۷۰)، تراؤب (۱۹۶۴) و کتابنامه‌های آنها را ملاحظه نمائید. مقاله ویلکینسن (۱۹۸۴) بعضی مشکلات عملی حل مسأله ریشه‌یابی چند جمله‌یها را با رایانه نشان می‌دهد. برنامه‌های رایانه‌ای دقیق، کارا، خودکار و قابل اعتماد برای

ریشه‌یابی چندجمله‌یها نوشته شده‌اند. از این جمله‌اند (الف) برنامه‌های جنکینز^۱ (۱۹۷۵) و جنکینز و تراوب (۱۹۷۰)، (۱۹۷۲)، و (ب) برنامه ZERPOL اثر اسمیت (۱۹۶۷) که بر مبنای روش لاگر^۲ نوشته شده است. هاوسهولدر [۱۹۷۰، ص ۱۷۶]، کاهان (۱۹۶۷) را ببینید. این برنامه‌های خودکار، هم از نظر ریاضی و هم از نظر الگوریتمی، آن قدر پیچیده‌اند که نمی‌توان در یک کتاب مقدماتی، چون این کتاب به معرفی آنها پرداخت. مع هذا کار با آنها با ارزش است. بسیاری از افراد نمی‌توانند برنامه‌ای بنویسند که از نظر سرعت و دقت بتواند با آنها رقابت کند. به‌ویژه دقت، بسیار اهمیت دارد، زیرا مسأله ریشه‌یابی چندجمله‌ای، می‌تواند نسبت به خطای گرد کردن، همان گونه که در مثالهای قبلی نشان داده شد بسیار حساس باشد.

مطالعه روشهای عددی برای حل دستگاههای معادله‌های غیرخطی و مسائل بهینه‌سازی، اخیراً یک زمینه بسیار متداول تحقیق گردیده است. برای آشنایی با روشهای عددی حل دستگاههای غیرخطی، بیکر و فیلیس^۳ (۱۹۸۱، فصل ۱)، اورتگا و راینبولت^۴ (۱۹۷۰) و راینبولت (۱۹۷۴) را ببینید. برای تعمیم این روشها به معادلات دیفرانسیل و انتگرالی غیرخطی، بیکر و فیلیس (۱۹۸۱)، کاتوروویچ (۱۹۴۸) [که مقاله‌ای کلاسیک در این زمینه است]، کاتوروویچ و اکیلوف (۱۹۶۴) و رال (۱۹۶۹) را ببینید. برای یک بازنگری روشهای عددی در بهینه‌سازی، دنیس (۱۹۸۴) و پاول (۱۹۸۲) را ببینید. آشنایی کلی در اشنا بل و دنیس (۱۹۸۳)، فلچر (۱۹۸۰) و (۱۹۸۱)، گیل و همکاران (۱۹۸۱)، لوتنبرگر^۵ (۱۹۸۴) داده شده است. به‌عنوان یک مثال از تحقیقات اخیر در نظریه بهینه‌سازی و توسعه نرم‌افزار، باگز و همکاران (۱۹۸۵) را ببینید. برای برنامه‌های کامپیوتری هیبرت (۱۹۸۲) و موره و همکاران (۱۹۸۴) را ببینید.

مراجع

- Baker, C., and C. Phillips, eds. (1981). *The Numerical Solution of Nonlinear Problems*. Clarendon Press, Oxford, England.
- Boggs, P., R. Byrd, and R. Schnabel, eds. (1985). *Numerical Optimization 1984*. Society for Industrial and Applied Mathematics, Philadelphia.
- Brent, R. (1973). *Algorithms for Minimization Without Derivatives*. Prentice-Hall, Englewood Cliffs, N.J.
- Byrne, G., and C. Hall, eds. (1973). *Numerical Solution of Systems of Nonlinear Algebraic Equations*. Academic Press, New York.

1. Jenkins 2. Laguerre 3. Baker and Phillips
4. Ortega and Rheinboldt 5. Luenberger

- Dejon, B., and P. Henrici, eds. (1969). *Constructive Aspects of the Fundamental Theorem of Algebra*. Wiley, New York.
- Dekker, T. (1969). Finding a zero by means of successive linear interpolation. In B. Dejon and P. Henrici (eds.), *Constructive Aspects of the Fundamental Theorem of Algebra*, pp. 37–51. Wiley, New York.
- Dennis, J. (1984). A user's guide to nonlinear optimization algorithms. *Proc. IEEE*, **72**, 1765–1776.
- Dennis, J., and R. Schnabel (1983). *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*. Prentice-Hall, Englewood Cliffs, N.J.
- Fletcher, R. (1980). *Practical Methods of Optimization*, Vol. 1, *Unconstrained Optimization*. Wiley, New York.
- Fletcher, R. (1981). *Practical Methods of Optimization*, Vol. 2, *Constrained Optimization*. Wiley, New York.
- Forsythe, G. (1969). What is a satisfactory quadratic equation solver? In B. Dejon and P. Henrici (eds.), *Constructive Aspects of the Fundamental Theorem of Algebra*, pp. 53–61. Wiley, New York.
- Gill, P., W. Murray, and M. Wright (1981). *Practical Optimization*. Academic Press, New York.
- Goldstine, H. (1977). *A History of Numerical Analysis*. Springer-Verlag, New York.
- Henrici, P. (1974). *Applied and Computational Complex Analysis*, Vol. 1. Wiley, New York.
- Hiebert, K. (1982). An evaluation of mathematical software that solve systems of nonlinear equations. *ACM Trans. Math. Softw.*, **11**, 250–262.
- Householder, A. (1970). *The Numerical Treatment of a Single Nonlinear Equation*. McGraw-Hill, New York.
- Jenkins, M. (1975). Algorithm 493: Zeroes of a real polynomial. *ACM Trans. Math. Softw.*, **1**, 178–189.
- Jenkins, M., and J. Traub (1970). A three state algorithm for real polynomials using quadratic iteration. *SIAM J. Numer. Anal.*, **7**, 545–566.
- Jenkins, M., and J. Traub (1972). Algorithm 419—Zeros of a complex polynomial. *Commun. ACM*, **15**, 97–99.
- Kahan, W. (1967). Laguerre's method and a circle which contains at least one zero of a polynomial. *SIAM J. Numer. Anal.*, **4**, 474–482.
- Kantorovich, L. (1948). Functional analysis and applied mathematics. *Usp. Mat. Nauk*, **3**, 89–185.
- Kantorovich, L., and G. Akilov (1964). *Functional Analysis in Normed Spaces*. Pergamon, London.
- Le, D. (1985). An efficient derivative-free method for solving nonlinear equations. *ACM Trans. Math. Softw.*, **11**, 250–262.

- Luenberger, D. (1984). *Linear and Nonlinear Programming*, 2nd ed. Wiley, New York.
- Moré, J., B. Garbow, and K. Hillstrom (1980). *User Guide for MINPACK-1*. Argonne Nat. Lab. Rep. ANL-80-74.
- Moré, J., and D. Sorenson. Newton's method. In *Studies in Numerical Analysis*, G. Golub (ed.), pp. 29-82. Math. Assoc. America, Washington, D.C.
- Moré, J., D. Sorenson, B. Garbow, and K. Hillstrom (1984). The MINPACK project. In *Sources and Development of Mathematical Software*, Cowell (ed.), pp. 88-111. Prentice-Hall, Englewood Cliffs, N.J.
- Nelder, A., and R. Mead (1965). A simplex method for function minimization. *Comput. J.*, 7, 308-313.
- Ortega, J., and W. Rheinboldt (1970). *Iterative Solution of Nonlinear Equations in Several Variables*. Academic Press, New York.
- Ostrowski, A. (1973). *Solution of Equations in Euclidean and Banach Spaces*. Academic Press, New York.
- Peters, G., and J. Wilkinson (1971). Practical problems arising in the solution of polynomial equations. *J. Inst. Math. Its Appl.* 8, 16-35.
- Potra, F., and V. Ptak (1984). *Nondiscrete Induction and Iterative Processes*. Pitman, Boston.
- Powell, M., ed. (1982). *Nonlinear Optimization 1981*. NATO Conf. Ser. Academic Press, New York.
- Rall, L. (1969). *Computational Solution of Nonlinear Operator Equations*. Wiley, New York.
- Rheinboldt, W. (1974). *Methods for Solving Systems of Nonlinear Equations*. Society for Industrial and Applied Mathematics, Philadelphia.
- Smith, B. (1967). ZERPOL: A zero finding algorithm for polynomials using Laguerre's method. Dept. of Computer Science, Univ. Toronto, Toronto, Ont., Canada.
- Traub, J. (1964). *Iterative Methods for the Solution of Equations*. Prentice-Hall, Englewood Cliffs, N.J.
- Whitley, V. (1968). Certification of algorithm 196: Muller's method for finding roots of an arbitrary function. *Commun. ACM* 11, 12-14.
- Wilkinson, J. (1963). *Rounding Errors in Algebraic Processes*. Prentice-Hall, Englewood Cliffs, N.J.
- Wilkinson, J. (1984). The perfidious polynomial. In *Studies in Numerical Analysis*, G. Golub (ed.). Math. Assoc. America, Washington, D.C.
- Woods, D. (1985). An interactive approach for solving multi-objective optimization problems. Ph.D. dissertation, William Marsh Rice Univ., Houston, Tex.

مسائل

۱. مثالهای مقدماتی برای $f(x) = a - (1/x)$ به حاصلضرب نامتناهی زیر مربوطاند

$$\prod_{j=0}^{\infty} (1 + r^{2^j}) \equiv \lim_{n \rightarrow \infty} [(1+r)(1+r^2)(1+r^4) \dots (1+r^{2^n})]$$

با استفاده از فرمولهای (۲.۰۶) و (۲.۰۹) می‌توانیم مقدار حاصلضرب نامحدود فوق را محاسبه کنیم. این مقدار چیست، و چه شرطی برای r باید برقرار باشد تا حاصلضرب فوق همگرا شود؟ راهنمایی: بگیرد $r = r_0$ و x_n را برحسب x_0 و r_0 بنویسید.

۲. یک برنامه بنویسید که در آن از الگوریتم نیم‌سازی که در بخش ۱.۲ داده شده استفاده شود. برنامه را برای محاسبه ریشه‌های حقیقی معادلات زیر به‌کار برید. تحمل خطا را برابر $\varepsilon = 10^{-5}$ انتخاب کنید.

(الف) $e^x - 3x^2 = 0$ (ب) $x^3 = x^2 + x + 1$

(ج) $e^x = \frac{1}{0.1 + x^2}$ (د) $x = 1 + 0.3 \cos(x)$

۳. با استفاده از برنامه مسئله ۲ (الف) کوچکترین ریشه مثبت $x - \tan(x) = 0$ و (ب) نزدیکترین ریشه این معادله به $x = 100$ را محاسبه کنید.

۴. الگوریتم نیوتن را که در بخش ۲.۲ داده شده است تکمیل کنید و آن را برای معادلات مسئله ۲ به‌کار گیرید.

۵. روش نیوتن را برای محاسبه ریشه‌هایی که در مسئله ۳ خواسته شده به‌کار گیرید. کوشش کنید که تفاوت‌های پیدا کردن ریشه‌ها را در (الف) و (ب) توضیح دهید.

۶. روش نیوتن را برای محاسبه تنها ریشه

$$x + e^{-Bx^2} \cos(x) = 0$$

با $B > 0$ پارامتری به‌کار برید. مقادیر افزایشی متفاوتی را برای B مثلاً $1, 5, 10, 25, 50$ به‌کار گیرید. از بین انتخاب‌های استفاده شده برای x_0 ، $x_0 = 0$ را انتخاب کرده و هر رفتار غیرعادی را توضیح دهید. از جنبه نظری، روش نیوتن برای هر مقدار x_0 و B همگراست. این موضوع را با محاسباتی که برای مقادیر بزرگتر B کرده‌اید، مقایسه کنید.

۷. یک مسئله جالب ریشه‌یابی چند جمله‌ای، در محاسبه قسط‌السنین پیش می‌آید. مبلغ P_1 تومان در ابتدا سالهای ۱، ۲، ۳، ...، N به حساب گذاشته شده است، و بهره آن در هر سال با نرخ مرکب r (مثلاً 5% ، $r = 0.05$) یعنی با سودی برابر 5% درصد حساب شده است. در ابتدای سالهای

$N_1 + 1, N_1 + 2, \dots, N_1 + N_2$ مبلغ P_2 تومان از حساب برداشته می‌شود پس از آخرین برداشت، حساب دقیقاً صفر می‌شود. رابطه بین متغیرها به صورت زیر است:

$$P_1[(1+r)^{N_1} - 1] = P_2[1 - (1+r)^{-N_2}]$$

اگر $N_1 = 30, N_2 = 20, P_1 = 2000, P_2 = 8000$ چقدر است r ؟ یک روش ریشه‌یابی به انتخاب خود به کار برید.

۸. روش نیوتن - فوریه را برای حل معادلات مسائل ۲ و ۶ به کار برید.
۹. روش خط قاطع را برای حل معادلات مسئله ۲ به کار برید.
۱۰. روش خط قاطع را برای حل معادله مسئله ۶ به کار برید.
۱۱. فرمول خطای (۲.۳.۲) را برای روش خط قاطع زیر ثابت کنید.

$$\alpha - c = -(\alpha - b)(\alpha - a) \frac{f[a, b, \alpha]}{f[a, b]}$$

۱۲. روش نیوتن را برای پیدا کردن ریشه مثبت $a > 0$ در نظر می‌گیریم. با فرض $x_0 > 0$ و $x_0 \neq \sqrt{a}$ نتایج زیر را پیدا کنید.

$$x_{n+1} = \frac{1}{2} \left(x_n + \frac{a}{x_n} \right) \quad (\text{الف})$$

$$x_{n+1}^2 - a = \left[\frac{x_n^2 - a}{2x_n} \right]^2 \quad (\text{ب})$$

(ج) بارستهای $\{x_n\}$ یک دنباله نزولی اکید برای $n \geq 1$ است.

راهنمایی: علامت $x_{n+1} - x_n$ را در نظر بگیرید.

$$e_n = \sqrt{a} - x_n \quad e_{n+1} = -e_n^2 / (2x_n) \quad (\text{د})$$

$$\text{Rel}(x_{n+1}) = \frac{-\sqrt{a}}{2x_n} [\text{Rel}(x_n)]^2$$

که در آن $\text{Rel}(x_n)$ خطای نسبی x_n است.

(ه) اگر $x_0 \geq \sqrt{a}$ و $| \text{Rel}(x_0) | \leq 0.1$ ، کران $\text{Rel}(x_2)$ را پیدا کنید.

۱۳. روش نیوتن یک روش متداول در محاسبه ریشه‌های دوم در رایانه است. برای استفاده از روش رایانه‌ای در محاسبه \sqrt{a} یک حدس اولیه x_0 باید انتخاب شود و مناسبتر از همه این است که به جای آزمایش همگرایی از تعداد ثابتی بارست استفاده کنیم. برای قطعیت، فرض می‌کنیم که حساب رایانه‌ای ما دودویی و جزء اعشاری شامل ۴۸ بیت دودویی باشد. می‌نویسیم

$$a = \hat{a} \cdot 2^e \quad \frac{1}{2} \leq \hat{a} < 1$$

این را به سادگی می‌توان به شکل

$$a = b \cdot 2^f \quad \frac{1}{4} \leq b < 1$$

که در آن f یک عدد صحیح زوج است، بدل کرد. پس

$$\sqrt{a} = \sqrt{b} \cdot 2^{f/2} \quad \frac{1}{2} \leq \sqrt{b} < 1$$

و هرگاه \sqrt{b} معلوم باشد، \sqrt{a} به شکل ممیز شناور استاندارد خواهد بود.

این امر مسأله را به محاسبه \sqrt{b} برای $\frac{1}{4} \leq b < 1$ بدل می‌کند. فرمول درونیابی خطی

$$x_n = \frac{1}{3}(2b + 1) \quad \frac{1}{4} \leq b \leq 1$$

را به عنوان یک حدس اولیه برای بارست نیوتن در محاسبه \sqrt{b} به کار برید. کران خطا را $\sqrt{b} - x_n$ بگیریید. تعداد بارستهای لازم را برای برقراری

$$0 \leq x_n - \sqrt{b} \leq 2^{-48}$$

که حد درجه دقت ماشین برای b در یک محاسبه خاص است، تخمین بزنید. [توجه داشته باشید که تأثیر خطای گردکردن نادیده گرفته شده است]. چگونه می‌توان انتخاب x_n را بهبود بخشید؟

۱۴. بارستهای نیوتن را برای حل $x^2 - 1 = 0$ به طور عددی محاسبه کنید و از $x_n = 1.000000$ استفاده کنید. سرعت همگرایی نتیجه را پیدا کنید و توضیح دهید.

۱۵. (الف) روش نیوتن را برای تابع

$$f(x) = \begin{cases} \sqrt{x} & x \geq 0 \\ -\sqrt{-x} & x < 0 \end{cases}$$

با ریشه $x = 0$ به کار برید. رفتار بارستها چگونه است؟ آیا همگرا هستند، اگر جواب مثبت است، با چه نرخ؟

(ب) مانند (الف) عمل کنید ولی برای

$$f(x) = \begin{cases} \sqrt[3]{x^2} & x \geq 0 \\ -\sqrt[3]{x^2} & x < 0 \end{cases}$$

۱۶. یک دنباله $\{x_n\}$ را به α همگرای زیر خطی گویند اگر

$$|\alpha - x_{n+1}| \leq c_n |\alpha - x_n| \quad n \geq 0$$

با $c_n \rightarrow 0$ هرگاه $n \rightarrow \infty$. نشان دهید که در این حالت

$$\lim_{n \rightarrow \infty} \frac{|\alpha - x_n|}{|x_{n+1} - x_n|} = 1$$

بنابراین وقتی $n \rightarrow \infty$ ، اعتبار $|x_{n+1} - x_n| \doteq |\alpha - x_n|$ با n افزایش می‌یابد.

۱۷. در روش نیوتن برای پیدا کردن یک ریشهٔ $f(x) = 0$ ، مانند α ، گاهی برای همگرایی نیاز است که حدس اولیهٔ x_0 خیلی به ریشهٔ α نزدیک باشد. ثابت کنید که برای $\alpha = \pi/2$ ، ریشهٔ از تابع

$$f(x) = \cos(x) + \sin^2(5^\circ x)$$

مسئله چنین است.

یک برآورد تقریبی از کوچکی $|x_0 - \alpha|$ برای همگرایی به α به دست آورید.

راهنمایی: (۶.۲.۲) را ملاحظه کنید.

۱۸. برنامه‌ای بنویسید که روش مولر را انجام دهد. آن را برای ریشه‌یابی در مسائل ۲، ۳ و ۶ به کار گیرید.

۱۹. نشان دهید که $x = 1 + \tan^{-1}(x)$ دارای یک جواب α است. یک بازهٔ $[a, b]$ پیدا کنید

که شامل α باشد به گونه‌ای که برای هر $x \in [a, b]$ بارست

$$x_{n+1} = 1 + \tan^{-1}(x_n) \quad n \geq 0$$

به α همگرا شود. چند بارست اول را محاسبه و نرخ همگرایی را پیدا کنید.

۲۰. مانند مسئلهٔ ۱۹ ولی برای بارست

$$x_{n+1} = 3 - 2 \log(1 + e^{-x_n}) \quad n \geq 0$$

عمل کنید.

۲۱. برای پیدا کردن ریشهٔ $f(x) = 0$ روش بارستی، معادله را مجدداً به شکل زیر می‌نویسیم:

$$x = x + cf(x) \equiv g(x)$$

به ازای ثابت $c \neq 0$. اگر α یک ریشهٔ $f(x)$ و $f'(\alpha) \neq 0$ باشد، چگونه باید انتخاب شود تا

دنبالهٔ $x_{n+1} = g(x_n)$ به α همگرا شود؟

۲۲. معادلهٔ

$$x = d + hf(x)$$

را در نظر می‌گیریم که در آن d یک ثابت داده شده و $f(x)$ یک تابع پیوسته از x است. اگر $h = 0$ ، یک ریشه $\alpha = d$ است. نشان دهید که برای تمام مقادیر h به اندازه کافی کوچک این معادله یک ریشه $\alpha(h)$ دارد. با چه شرطی، اگر به شرطی نیاز باشد، یکتایی ریشه $\alpha(h)$ در بازه‌ای حول d تضمین می‌شود؟

۲۳. بارست $x_{n+1} = 2 - (1+c)x_n + cx_n^2$ برای بعضی از مقادیر c به $\alpha = 1$ همگرا می‌شود [به شرطی که x به اندازه کافی نزدیک به α انتخاب شده باشد]. این مقادیر c را پیدا کنید. برای چه مقداری از c ، همگرایی از درجه دو است؟

۲۴. کدام یک از بارستهای زیر به نقطه ثابت داده شده α همگراست. (به شرطی که x به اندازه کافی به α نزدیک باشد)؟ اگر بارست همگرا باشد مرتبه همگرایی را پیدا کنید، در حالت همگرایی خطی، نرخ همگرایی را پیدا کنید.

$$\alpha = 2 \quad x_{n+1} = -16 + 6x_n + \frac{12}{x_n} \quad (\text{الف})$$

$$\alpha = 3^{1/3} \quad x_{n+1} = \frac{2}{3}x_n + \frac{1}{x_n^2} \quad (\text{ب})$$

$$\alpha = 3 \quad x_{n+1} = \frac{12}{1+x_n} \quad (\text{ج})$$

۲۵. نشان دهید که

$$x_{n+1} = \frac{x_n(x_n^2 + 3a)}{3x_n^2 + a} \quad n \geq 0$$

برای محاسبه \sqrt{a} یک روش مرتبه سه است. با فرض اینکه x به اندازه کافی به α نزدیک انتخاب شده باشد، حد زیر را محاسبه کنید.

$$\lim_{n \rightarrow \infty} \frac{\sqrt{a} - x_{n+1}}{(\sqrt{a} - x_n)^3}$$

۲۶. با استفاده از قضیه ۸.۲، نشان دهید که فرمول (۱۱.۴.۲) یک روش بارستی همگرایی مرتبه ۳ است.

۲۷. یک فرمول بارستی به شکل زیر تعریف می‌کنیم

$$x_{n+1} = z_{n+1} - \frac{f(z_{n+1})}{f'(x_n)}, \quad z_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

نشان دهید که مرتبه همگرایی $\{x_n\}$ به α حداقل ۳ است.

راهنمایی: قضیه ۲-۸ را به‌کار برید و بگیرید

$$g(x) = h(x) - \frac{f(h(x))}{f'(x)}, \quad h(x) = x - \frac{f(x)}{f'(x)}$$

۲۸. یک اصلاح دیگری برای روش نیوتن، مشابه با روش خط قاطع، وجود دارد، که از تقریب متفاوتی برای مشتق $f'(x_n)$ در آن استفاده می‌شود. تعریف می‌کنیم

$$x_{n+1} = x_n - \frac{f(x_n)}{D(x_n)} \quad D(x_n) = \frac{f(x_n + f(x_n)) - f(x_n)}{f(x_n)} \quad n \geq 0$$

این روش تک-نقطه‌ای، روش استیفنسن^۱ نامیده می‌شود. فرض می‌کنیم $f'(\alpha) \neq 0$ ، نشان دهید که این روش از مرتبه دو است.

راهنمایی: بارست را به شکل $x_{n+1} = g(x_n)$ بنویسید. از $f(x) = (x - \alpha)h(x)$ با $h(\alpha) \neq 0$ استفاده کنید و این فرمول را برای $g(x)$ برحسب $h(x)$ محاسبه کنید. سپس، قضیه ۲-۸ را به‌کار برید.

۲۹. جدول زیر مقادیر بارستهای یک بارست همگرای خطی $x_{n+1} = g(x_n)$ را نشان می‌دهد. خواسته‌های زیر را تخمین بزنید: (الف) نرخ همگرایی خطی، (ب) نقطه ثابت α ، و (ج) خطای $x_5 - \alpha$.

n	x_n
۰	۱,۰۹۴۹۲۴۲
۱	۱,۲۰۹۲۷۵۱
۲	۱,۲۸۰۷۹۱۷
۳	۱,۳۲۵۴۹۴۳
۴	۱,۳۵۳۴۳۳۹
۵	۱,۳۷۰۸۹۶۲

۳۰. می‌توان نشان داد که الگوریتم ایتکن که در بخش ۲-۶ داده شد از نظر سرعت همگرایی از درجه دوم است. گیریم که بارست اصلی $x_{n+1} = g(x_n)$ ، $n \geq 0$ باشد. فرمول (۸.۶.۲) را می‌توان به شکل هم‌ارز زیر مجدداً نوشت

$$\alpha \hat{=} \hat{x}_{n-2} = x_{n-2} + \frac{(x_{n-1} - x_{n-2})^2}{(x_{n-1} - x_{n-2}) - (x_n - x_{n-1})} \quad n \geq 2$$

برای آزمایش سرعت همگرایی برون‌بهای ایتکن، دنباله متناظر زیر را در نظر می‌گیریم

$$z_{n+1} = z_n + \frac{[g(z_n) - z_n]^2}{[g(z_n) - z_n] - [g(g(z_n)) - g(z_n)]} \quad n \geq 0$$

مقادیر z_n ، مقادیر پی در پی \hat{x}_n مسأله هستند که از الگوریتم ایتکن به دست آمده‌اند. برای ۱ یا $g'(\alpha) \neq 0$ نشان دهید که z_n با سرعت درجه دو به α می‌گراید. این درست است حتی اگر $|g'(\alpha)| > 1$ و بارست اصلی و اگر باشد. راهنمایی: سعی نکنید که قضیه ۸.۲ را مستقیماً به‌کار ببرید، زیرا خیلی پیچیده می‌شود و به جای آن بنویسید

$$g(x) = (x - \alpha)h(x) \quad h(\alpha) = g'(\alpha) \neq 0$$

این رابطه را به‌کار ببرید تا نشان دهید که برای بعضی توابع $H(x)$ که حول $x = \alpha$ کراندار باشد

$$\alpha - z_{n+1} = H(z_n)(\alpha - z_n)^2$$

۳۱. دنباله زیر را در نظر می‌گیریم

$$x_n = \alpha + \beta\rho^n + \gamma\rho^{2n} \quad n \geq 0, \quad |\rho| < 1$$

با $\beta, \gamma \neq 0$ که با یک نرخ خطی ρ به α همگراست. گیریم \hat{x}_{n-2} برون‌یاب ایتکن باشد:

$$\hat{x}_{n-2} = x_n - \frac{(x_n - x_{n-1})^2}{(x_n - x_{n-1}) - (x_{n-1} - x_{n-2})} \quad n \geq 0$$

نشان دهید که

$$\hat{x}_{n-2} = \alpha + a\rho^{2n} + b\rho^{4n} + c_n\rho^{6n}$$

که در آن c_n وقتی $n \rightarrow \infty$ کراندار است. عباراتی برای a و b بیابید. دنباله $\{\hat{x}_n\}$ با یک نرخ خطی ρ^2 به α همگراست.

۳۲. گیریم $f(x)$ یک ریشه چندگانه α ، مثلاً از مرتبه $m > 1$ داشته باشد. نشان دهید که

$$K(x) = \frac{f(x)}{f'(x)}$$

دارای یک ریشه ساده α است. چرا این امر در دشواری اصلی محاسبه عددی ریشه‌های چندگانه، یعنی بازه بزرگ عدم اطمینان α ، به ما کمکی نمی‌کند؟

۳۳. روش نیوتن را برای محاسبه هر چه دقیقتر ریشه‌های چند جمله‌بیهای زیر به‌کار ببرید. چندگانگی هر ریشه را برآورد کنید و در صورت لزوم، برای بهبود مقادیر محاسبه شده خود، راه دیگری را امتحان کنید.

$$x^4 - 32x^3 + 96x^2 + 4608x - 3456 \quad (\text{الف})$$

$$x^5 + 9x^4 - 162x^3 - 1458x^2 + 6561x + 59049 \quad (\text{ب})$$

۳۴. از برنامهٔ مسئلهٔ ۲ در حل معادلات زیر برای ریشهٔ $\alpha = 1$ استفاده کنید. بازهٔ اولیهٔ $[0, 3]$ را به‌کار برید و در همهٔ حالات $\varepsilon = 10^{-5}$ را به‌عنوان تحمل توقف برنامه در نظر بگیرید. نتایج را با آنچه در بخش ۲-۸ با روش برنت به‌دست‌آمده بود مقایسه نمایید.

$$(x-1)[1+(x-1)^2] = 0 \quad (\text{i})$$

$$x^2 - 1 = 0 \quad (\text{ii})$$

$$-1 + x(3 + x(-3 + x)) = 0 \quad (\text{iii})$$

$$(x-1) \exp(-1/(x-1)^2) = 0 \quad (\text{iv})$$

۳۵. با استفاده از کران بالای (۶.۹.۲) و توصیهٔ (۷.۹.۲)، کران پایین (۶.۹.۲) را ثابت کنید.

۳۶. گیریم $p(x)$ یک چندجمله‌ای از درجهٔ n باشد. گیریم ریشه‌های متمایز آن: $\alpha_1, \dots, \alpha_r$ به ترتیب از درجهٔ چندگانگی m_1, \dots, m_r باشند،
(الف) نشان دهید که

$$\frac{p'(x)}{p(x)} = \sum_{j=1}^r \frac{m_j}{x - \alpha_j}$$

(ب) گیریم c عددی باشد که برای آن $p'(c) \neq 0$. نشان دهید که یک ریشهٔ α دارد که در رابطهٔ زیر صدق می‌کند

$$|\alpha - c| \leq n \left| \frac{p(c)}{p'(c)} \right|$$

۳۷. برای چندجمله‌ای

$$p(x) = a_0 + a_1x + \dots + a_nx^n \quad a_n \neq 0$$

تعریف می‌کنیم

$$R = \frac{|a_0| + |a_1| + \dots + |a_{n-1}|}{|a_n|}$$

نشان دهید که هر ریشه x از $p(x) = 0$ در رابطه

$$|x| \leq \text{Max}\{R, \sqrt[n]{R}\}$$

صدق می‌کند.

۳۸. یک برنامه رایانه‌ای بنویسید که چندجمله‌یهای $p(x)$ زیر را برای مقادیر داده شده x محاسبه کند و نوفه را در مقادیر $p(x)$ به دست آورد. برای هر مقدار x ، $p(x)$ را در حساب با دقت ساده و دقت مضاعف محاسبه نمایید؛ تفاضل آنها را به عنوان نوفه در حساب با دقت ساده، بر اثر خطای گرد کردن در محاسبه $p(x)$ تلقی نمایید. هم فرمول معمولی (۱.۹.۲) و هم قاعده هورنر (۸.۹.۲) را برای محاسبه هر چندجمله‌ای به کار برید. این محاسبات باید نشان دهند که نوفه در دو حالت فوق متفاوت است.

(الف) $p(x) = x^4 - 57x^3 - 047x^2 + 29865x - 261602$ ؛ $-3 \leq x \leq 5$ با طول گام x برابر 0.1 .

(ب) $p(x) = x^4 - 54x^3 + 1056x^2 - 8954x + 27951$ ، $1 \leq x \leq 1.2$ ، با طول گام x برابر 0.01 .

توجه: برای کامپیوترهای مختلف ممکن است مقادیر کوچکتر یا بزرگتر h مناسب باشند. همچنین، برای مقایسه معتبرتر، قبل از استفاده از حساب با دقت مضاعف، ضرایب با دقت ساده را وارد کنید.

۳۹. حساب مختلط و روش نیوتن را برای محاسبه یک ریشه مختلط

$$p(z) = z^4 - 3z^3 + 20z^2 + 44z + 54$$

که نزدیک به $z_0 = 2.5 + 4.5i$ باشد به کار برید.

۴۰. برنامه‌ای برای محاسبه ریشه‌های چندجمله‌یهای زیر، تا اندازه ممکن دقیق، بنویسید.

(الف) $676039x^{12} - 1939938x^{10} + 2078505x^8 - 1021020x^6 + -18018x^4 + 231$
 $225225x^4$

(ب) $x^4 - 4096152422706631x^3 + 3284222335022705x^2$
 $+ 4703847577293368x - 5715767664977294$

۴۱. از یک بسته برنامه ریشه‌یابی برای چندجمله‌یها در محاسبه ریشه‌های چندجمله‌یهای مسائل ۳۸، ۳۹ و ۴۰ استفاده کنید.

۴۲. در مثال $f(x) = (x-1)(x-2)\dots(x-7)$ ، در (۲۱.۹.۲) بخش ۹.۲، اختلالی به اندازه $\varepsilon_i x^i$ برای ضریب x^i ایجاد می‌کنیم که در آن ε_i به گونه‌ای انتخاب شده است که اختلال نسبی در ضریب x^i همان اختلال مثال داده شده در کتاب برای ضریب x^6 است. نظریه خطی‌سازی (۱۹.۹.۲) چه اختلالی را در ریشه‌ها پیشگویی می‌کند؟ تغییر در کدام ضریب بزرگترین تغییرها را در ریشه‌ها ایجاد می‌کند؟

۴۳. چندجمله‌ای $f(x) = x^5 - 300x^4 - 126x + 5005$ دارای یک ریشه $\alpha = 5$ است. اثر تغییر ضریب x^5 را از $1 + \varepsilon$ به 1 روی ریشه α برآورد نمایید.

۴۴. نتیجه پایداری (۱۹.۹.۲) برای ریشه‌های چندجمله‌ای را می‌توان برای توابع دیگر تعمیم داد. گیریم α یک ریشه ساده $f(x) = 0$ باشد و گیریم $f(x)$ و $g(x)$ حول α پیوسته مشتقپذیر باشند. تابع $F_\varepsilon(x) = f(x) + \varepsilon g(x)$ را تعریف می‌کنیم. گیریم $\alpha(\varepsilon)$ یک ریشه $F_\varepsilon(x) = 0$ متناظر با $\alpha = \alpha(0)$ برای مقدار کوچک ε باشد. برای اینکه بینیم چنین $\alpha(\varepsilon)$ ای وجود دارد و برای اثبات اینکه این تابع پیوسته مشتقپذیر است، قضیه تابع ضمنی برای توابع یک متغیره را به‌کار می‌بریم. با استفاده از این قضیه، (۱۹.۹.۲) را برای حالت کلی تعمیم دهید.

۴۵. با استفاده از نتیجه پایداری در مسأله ۴۴، ریشه $\alpha(\varepsilon)$ معادله

$$x - \tan(x) + \varepsilon = 0$$

را برآورد نمایید. دو حالت را صریحاً برای ریشه α از معادله $x - \tan(x) = 0$ در نظر بگیرید؛

$$(1) \alpha \in (0.5\pi, 1.5\pi) \quad (2) \alpha \in (3.15\pi, 3.25\pi)$$

۴۶. دستگاه زیر را در نظر می‌گیریم

$$x = \frac{0.5}{1 + (x + y)^2} \quad y = \frac{0.5}{1 + (x - y)^2}$$

یک ناحیه کراندار D بیابید که برای آن فرضهای قضیه (۹.۲) برقرار باشند.

راهنمایی: علامت مؤلفه‌های ریشه α چه خواهد بود؟ همچنین مقادیر ماکسیمم ممکن برای

x و y در فرمولهای فوق چه خواهد بود؟

۴۷. دستگاه زیر را در نظر می‌گیریم

$$x = 1 + h \cdot \frac{e^{-x^2}}{1 + y^2} \quad y = 0.5 + h \cdot \tan^{-1}(x^2 + y^2)$$

نشان دهید که اگر h به اندازه کافی کوچک انتخاب شده باشد، این دستگاه دارای یک ریشه

یکتای α در ناحیه مستطیلی شکل خواهد داشت. به‌علاوه، نشان دهید که یک بارست ساده به

شکل (۴.۱۰.۲) به این جواب همگرا خواهد بود.

۴۸. فرع ۱۰.۲ را ثابت کنید.

راهنمایی: از پیوستگی مشتقات جزئی مؤلفه‌های $g(x)$ استفاده کنید.

۴۹. ثابت کنید که بارستهای x_n در قضیه ۹. به یک جواب $x = g(x)$ همگراست.

راهنمایی: مجموع نامتناهی

$$x_0 + \sum_{n=0}^{\infty} [x_{n+1} - x_n]$$

را در نظر بگیرید. حاصل جمعهای جزئی آن چنین اند

$$x_0 + \sum_{n=0}^{N-1} [x_{n+1} - x_n] = x_N$$

پس اگر این سری نامتناهی مثلاً به α همگرا شود، آنگاه x_N به α همگرا می‌شود. با نشان دادن رابطه

$$\|x_{n+1} - x_n\|_{\infty} \leq \lambda \|x_n - x_{n-1}\|_{\infty}$$

و استفاده از آن، نشان دهید که سری نامتناهی فوق همگرایی مطلق است. سپس نشان دهید که α یک نقطه ثابت $g(x)$ است،

۵۰. با استفاده از روش نیوتن برای دستگاههای غیرخطی، دستگاه غیرخطی

$$x^2 + y^2 = 4 \quad x^2 - y^2 = 1$$

را حل کنید. به سادگی معلوم می‌شود که جوابهای حقیقی دستگاه، $(\pm\sqrt{2.5}, \pm\sqrt{1.5})$ است. به عنوان یک حدس اولیه، $(x_0, y_0) = (1.6, 1.2)$ را به کار برید. ۵۱. دستگاه

$$x^2 + xy^3 = 9 \quad 3x^2y - y^3 = 4$$

را با استفاده از روش نیوتن برای دستگاههای غیرخطی، حل کنید. هر یک از حدسهای اولیه $(2, -2.5), (-2, 2.5), (-1.2, -2.5), (1.2, 2.5)$ را به کار برید. ببینید کدام ریشه است که این روش به آن همگرا می‌شود، تعداد بارستها و سرعت همگرایی را نیز پیدا کنید.

۵۲. با استفاده از روش نیوتن برای دستگاههای غیرخطی، دستگاههای غیرخطی زیر را برای تمام ریشه‌های آنها حل کنید. از نمودارهای آنها برای تعیین حدسهای اولیه استفاده کنید.

$$x^2 + y^2 - 2x - 2y + 1 = 0 \quad x + y - 2xy = 0$$

(الف)

$$x^2 + 2xy + y^2 - x + y - 4 = 0 \quad (ب)$$

$$5x^2 - 6xy + 5y^2 + 16x - 16y + 12 = 0$$

۵۳. ثابت کنید که ژاکوبی

$$g(x) = x - F(x)^{-1}f(x)$$

در هر ریشه α از معادله $f(x) = 0$ به شرطی که $F(\alpha)$ ناسنگین باشد، برابر صفر است. از ترکیب این رابطه با فرع ۱۰.۲ بخش ۱۰.۲ اثبات همگرایی روش نیوتن حاصل می‌شود.

۵۴. روش نیوتن (۴.۱۲.۲) را برای پیدا کردن مقدار مینیمم تابع

$$f(x) = x_1^4 + x_1x_2 + (1 + x_2)^2$$

به کار برید. حدسهای اولیه مختلفی را به کار برید و رفتار همگرایی را مشاهده کنید.

نظریه درونیابی

مفهوم درونیابی عبارت است از انتخاب یک تابع $p(x)$ از رده‌ای از توابع داده شده به نحوی که نمودار $y = p(x)$ از یک مجموعه متناهی از نقاط داده شده بگذرد. در بیشتر این فصل تابع درونیاب $p(x)$ را به یک چندجمله‌ی محدود می‌کنیم.

نظریه درونیابی چندجمله‌ی شماری کاربرد مهم دارد. در این کتاب، کاربرد نخست آن، تهیه ابزارهای ریاضی برای ایجاد روشهایی است در زمینه‌های نظریه تقریب، انتگرال‌گیری عددی و حل عددی معادلات دیفرانسیل. کاربرد دوم آن، تهیه وسایلی است برای کار با توابعی که به شکل جدول داده شده‌اند. به عنوان مثال، تقریباً همه با درونیابی خطی در یک جدول لگاریتم از جبر دبیرستانی آشنایی داریم. از جنبه محاسباتی شکلهای مناسبی برای درونیابی چندجمله‌ی با داده‌های جدولی، به دست می‌آوریم و خطای حاصل را تحلیل خواهیم نمود. معلوم است که با استفاده گسترده از ماشینهای حساب و رایانه امروزه به استفاده خیلی کمتر از گذشته نزدیک به درونیابی جدولی، نیاز است. ما آن را از این جهت آورده‌ایم که فرمولهای به دست آمده در موارد دیگر هنوز مفیدند و نیز درونیابی جدولی، مثالها و تمرینات مناسبی در اختیار ما می‌گذارد.

این فصل را با آوردن دو موضوع دیگر خاتمه می‌دهیم. این دو موضوع عبارت‌اند از: (۱) توابع درونیاب چندجمله‌ی تکه‌یی، بویژه توابع براز، و (۲) درونیابی با توابع مثلثاتی.

۱.۳ نظریه درونیابی چندجمله‌یی

گیریم x_0, x_1, \dots, x_n اعداد حقیقی یا مختلط متمایز باشند و y_0, y_1, \dots, y_n مقادیر تابع متناظر آنها. اکنون به مطالعه مسأله پیدا کردن یک چندجمله‌یی $p(x)$ می‌پردازیم که داده‌های معلوم

$$p(x_i) = y_i \quad i = 0, 1, \dots, n \quad (1.1.3)$$

را درونیابی می‌کند. آیا چنین چندجمله‌یی وجود دارد؟ اگر جواب مثبت است، درجه آن چیست؟ آیا یکتاست؟ فرمول تولید $p(x)$ از داده‌های فوق چیست؟

چند جمله‌یی کلی درجه m

$$p(x) = a_0 + a_1x + \dots + a_mx^m$$

$m + 1$ پارامتر مستقل a_0, a_1, \dots, a_m دارد. چون (۱.۱.۳)، $n + 1$ شرط برای $p(x)$ پدید می‌آورد، منطقی است که ابتدا حالت $m = n$ را در نظر بگیریم. پس می‌خواهیم a_0, a_1, \dots, a_n را به گونه‌ای بیابیم که داشته باشیم

$$a_0 + a_1x_0 + a_2x_0^2 + \dots + a_nx_0^n = y_0$$

⋮

$$a_0 + a_1x_n + a_2x_n^2 + \dots + a_nx_n^n = y_n \quad (2.1.3)$$

این یک دستگاه $n + 1$ معادله خطی با $n + 1$ مجهول است، و حل آن کاملاً هم‌ارز با حل مسأله درونیابی چندجمله‌یی است. با نمادهای برداری و ماتریسی، این دستگاه چنین نوشته می‌شود:

$$Xa = y$$

با

$$X = [x_i^j] \quad i, j = 0, 1, \dots, n \quad (3.1.3)$$

$$a = [a_0, a_1, \dots, a_n]^T, y = [y_0, \dots, y_n]^T$$

ماتریس X ، ماتریس واندرموند نامیده می‌شود.

قضیه ۱.۳ برای $n + 1$ نقطه متمایز داده شده x_0, \dots, x_n و $n + 1$ عرض داده شده y_0, \dots, y_n یک چندجمله‌یی $p(x)$ از درجه نایزگتر از n هست که مقادیر y_i را در x_i ، $i = 0, 1, \dots, n$ درونیابی می‌کند. این چندجمله‌یی در مجموعه تمام چندجمله‌یهای از درجه حداکثر n ، یکتاست.

برهان سه برهان برای این قضیه مهم داده شده است. هرکدام اطلاعات مورد نیازی را به دست می‌دهد و کاربردهای مهمی در مسائل دیگر درونیابی خواهد داشت.

(i) می‌توان نشان داد که برای ماتریس X در (۳.۱.۳)،

$$\det(X) = \prod_{0 \leq j < i \leq n} (x_i - x_j) \quad (4.1.3)$$

(مسئله ۱ را ببینید). این نشان می‌دهد که $\det(X) \neq 0$ ، زیرا نقاط x_i متمایزند. پس X ناکین است و دستگاه $Xa = y$ یک جواب یکتای a دارد. این، وجود و یکتایی چندجمله‌یی درونیاب از درجه نایزگتر از n را اثبات می‌کند.

(ii) به موجب یک قضیه استانده جبر خطی (قضیه ۲.۷ فصل ۷)، دستگاه $Xa = y$ دارای یک جواب یکتاست اگر و فقط اگر دستگاه همگن $Xb = 0$ فقط جواب نمایان $b = 0$ داشته باشد. بنابراین، فرض می‌کنیم برای یک مقدار b ، دستگاه $Xb = 0$. با استفاده از b تعریف می‌کنیم

$$p(x) = b_0 + b_1x + \dots + b_nx^n$$

با توجه به دستگاه $Xb = 0$ ، داریم

$$p(x_i) = 0 \quad i = 0, 1, \dots, n$$

چندجمله‌یی $p(x)$ ، $n+1$ صفر دارد و درجه $p(x)$ نایزگتر از n است. این ممکن نیست مگر اینکه $p(x) \equiv 0$. پس همه ضرایب $b_i = 0$ ، $i = 0, 1, \dots, n$ ، که اثبات را کامل می‌کند.

(iii) اکنون چندجمله‌یی درونیاب را به طور صریح ارائه می‌دهیم. برای شروع، مسئله درونیابی خاصی را در نظر می‌گیریم که در آن برای یک مقدار i ، $0 \leq i \leq n$ ،

$$y_i = 1 \quad y_j = 0 \quad i \neq j$$

یک چندجمله‌یی می‌خواهیم از درجه نایزگتر از n با n صفر x_j ، $i \neq j$. پس

$$p(x) = c(x - x_0) \dots (x - x_{i-1})(x - x_{i+1}) \dots (x - x_n)$$

برای مقدار ثابتی مانند c . شرط $p(x_i) = 1$ نتیجه می‌دهد که

$$c = [(x_i - x_0) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_i - x_n)]^{-1}$$

این چندجمله‌یی خاص به شکل

$$l_i(x) = \prod_{j \neq i} \left(\frac{x - x_j}{x_i - x_j} \right) \quad i = 0, 1, \dots, n \quad (5.1.3)$$

نوشته می‌شود. برای حل مسئله کلی درونیابی (۱.۱.۳)، می‌نویسیم

$$p(x) = y_0 l_0(x) + y_1 l_1(x) + \dots + y_n l_n(x)$$

با استفاده از ویژگیهای خاص چندجمله‌یهای $l_i(x)$ ، دیده می‌شود که تابع $p(x)$ در رابطه (۱.۱.۳) صدق می‌کند. همچنین، درجه $p(x)$ نابزرگتر از n است، زیرا تمام $l_i(x)$ ها از درجه n هستند. برای اثبات یکتایی، گیریم $q(x)$ چندجمله‌یی دیگری از درجه نابزرگتر از n باشد که در رابطه (۱.۱.۳) صدق می‌کند. تعریف می‌کنیم

$$r(x) = p(x) - q(x)$$

در این صورت درجه $r(x)$ نابزرگتر از n است، و

$$r(x_i) = p(x_i) - q(x_i) = y_i - y_i = 0 \quad i = 0, 1, \dots, n$$

چون $r(x)$ دارای $n + 1$ صفر است، باید داشته باشیم $r(x) \equiv 0$. این ثابت می‌کند که $p(x) \equiv q(x)$ و برهان کامل می‌شود. ■

از ویژگی یکتایی در بسیاری از مطالبی که بعداً خواهد آمد استفاده خواهیم کرد. فرمولهای دیگری برای مسئله درونیابی (۱.۱.۳) به دست خواهیم آورد، و یکتایی می‌گوید که با همه فرمولها همان یک چندجمله‌ای یکتا به دست می‌آید. همچنین، بدون یکتایی، امکان داشت دستگاه خطی (۲.۱.۳) به طور یکتا حل پذیر نباشد، این امر با استفاده از قضایای جبر خطی، وجود داده‌هایی چون بردارهای y را ایجاب می‌کرد که برای آنها هیچ چند جمله‌ای درونیابی با درجه نابزرگتر از n وجود نداشته باشد.

فرمول

$$p_n(x) = \sum_{i=0}^n y_i l_i(x) \quad (6.1.3)$$

فرمول لاگرانژ برای چندجمله‌یی درونیاب نامیده می‌شود.

$$p_1(x) = \frac{x - x_1}{x_0 - x_1} y_0 + \frac{x - x_0}{x_1 - x_0} y_1 = \frac{(x_1 - x)y_0 + (x - x_0)y_1}{x_1 - x_0} \quad \text{مثال}$$

$$p_2(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} y_0 + \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} y_1 + \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} y_2$$

چندجمله‌یی از درجهٔ نایزگتر از ۲ که از سه نقطهٔ $(0, 1)$ ، $(-1, 2)$ و $(1, 3)$ می‌گذرد عبارت است از

$$P_2(x) = \frac{(x+1)(x-1)}{(0+1)(0-1)} \cdot 1 + \frac{(x-0)(x-1)}{(-1-0)(-1-1)} \cdot 2 + \frac{(x-0)(x+1)}{(1-0)(1+1)} \cdot 3 \\ = 1 + \frac{1}{4}x + \frac{3}{4}x^2$$

اگر تابع $f(x)$ داده شده باشد، می‌توانیم با استفاده از چندجمله‌یی درونیاب

$$p_n(x; f) \equiv p_n(x) = \sum_{i=0}^n f(x_i) l_i(x) \quad (7.1.3)$$

تقریبی برای آن بسازیم.

این چندجمله‌ای تابع $f(x)$ را در نقاط x_0, \dots, x_n درونیابی می‌کند. برای مثال، بعداً تابع $f(x) = \log_1 x$ را با درونیابی خطی در نظر می‌گیریم. قضیهٔ اساسی که در تحلیل خطای درونیابی مورد استفاده قرار می‌گیرد قضیهٔ زیر است. نماد $\mathcal{H}\{a, b, c, \dots\}$ معرف کوچکترین بازهٔ شامل تمام اعداد حقیقی a, b, c, \dots است.

قضیهٔ ۲.۳ بگیریم x_0, \dots, x_n اعداد حقیقی متمایز باشند و f یک تابع حقیقی داده شده با $n+1$ مشتق پیوسته در بازهٔ $I_t = \mathcal{H}\{t, x_0, \dots, x_n\}$ باشد که t مقدار حقیقی داده‌شده‌ای است. آنگاه عددی مانند $\xi \in I_t$ با ویژگی زیر وجود دارد:

$$f(t) - \sum_{j=0}^n f(x_j) l_j(t) = \frac{(t - x_0) \dots (t - x_n)}{(n+1)!} f^{(n+1)}(\xi) \quad (8.1.3)$$

برهان توجه کنید که قضیهٔ وقتی t یک نقطهٔ گرهی باشد درست است، زیرا در آن حالت هر دو طرف (۸.۱.۳) برابر صفر می‌شود. فرض کنید t نقطهٔ گرهی نباشد. تعریف می‌کنیم:

$$E(x) = f(x) - p_n(x) \quad p_n(x) = \sum_{j=0}^n f(x_j) l_j(x)$$

$$G(x) = E(x) - \frac{\Psi(x)}{\Psi(t)} E(t) \quad x \in I_t \text{ هر بازای هر } (9.1.3)$$

که در آن

$$\Psi(x) = (x - x_0) \dots (x - x_n)$$

تابع $G(x)$ در بازه I_t ، $n + 1$ بار پیوسته مشتقپذیر است زیرا $E(x)$ و $\Psi(x)$ چنین اند. همچنین

$$G(x_i) = E(x_i) - \frac{\Psi(x_i)}{\Psi(t)} E(t) = 0 \quad i = 0, 1, \dots, n$$

$$G(t) = E(t) - E(t) = 0$$

بنابراین G ، $n + 2$ صفر متمایز در I_t دارد. با استفاده از قضیه مقدار میانی، G' دارای $n + 1$ صفر متمایز است. استدلال استقرایی نشان می‌دهد که به ازای $j = 0, 1, \dots, n + 1$ $G^{(j)}(x)$ دارای $n + 2 - j$ صفر در I_t است. گیریم ξ یک صفر $G^{(n+1)}(x)$ باشد.

$$G^{(n+1)}(\xi) = 0$$

چون

$$E^{(n+1)}(x) = f^{(n+1)}(x)$$

$$\Psi^{(n+1)}(x) = (n + 1)!$$

به دست می‌آوریم

$$G^{(n+1)}(x) = f^{(n+1)}(x) - \frac{(n + 1)!}{\Psi(t)} E(t)$$

اگر به جای x ، ξ بگذاریم و معادله را نسبت به $E(t)$ حل کنیم داریم،

$$E(t) = \frac{\Psi(t)}{(n + 1)!} \cdot f^{(n+1)}(\xi)$$

که همان قضیه مطلوب است.

این طریق پیدا کردن فرمول خطا ممکن است «دشوار» به نظر آید، ولی این یک راه معمولی به دست آوردن بعضی از فرمولهای خطاست.

مثال برای $n = 1$ با استفاده از x به جای t

$$f(x) - \frac{(x_1 - x)f(x_0) + (x - x_0)f(x_1)}{x_1 - x_0} = \frac{(x - x_0)(x - x_1)}{2} f''(\xi_x) \quad (10.1.3)$$

که در آن $\xi_x \in \mathcal{H}\{x_0, x_1, x\}$ زیرنمایه x در ξ_x صریحاً نشان می‌دهد که ξ به x بستگی دارد؛ معمولاً زیرنمایه را برای راحتی نمی‌نویسیم.

اکنون حالت $n = 1$ را برای تکنیک معمول دبیرستانی در درونیابی خطی در یک جدول لگاریتم به کار می‌بریم. گیریم

$$f(x) = \log_{10} x$$

پس $f''(x) = -\log_{10} e/x^2$ ، $\log_{10} e \doteq 0.434$. در جدول معمولاً باید $x_0 < x < x_1$ پس

$$E(x) = \frac{(x - x_0)(x_1 - x)}{2} \cdot \frac{\log_{10} e}{\xi^2} \quad x_0 \leq \xi \leq x_1$$

و از این، کرانهای بالا و پایین به دست می‌آیند

$$\frac{\log_{10} e}{x_1^2} \cdot \frac{(x - x_0)(x_1 - x)}{2} \leq E(x) \leq \frac{\log_{10} e}{x_0^2} \cdot \frac{(x - x_0)(x_1 - x)}{2}$$

این نشان می‌دهد که تابع خطای $E(x)$ بسیار شبیه به یک چندجمله‌یی درجه دوم است، به ویژه اگر فاصله $h = x_1 - x_0$ نسبتاً کوچک باشد. برای یک کران یکنواخت در $[x_0, x_1]$ ، برای $x_0 \geq 1$ همان‌گونه که در جدول لگاریتم معمول است،

$$\text{Max}_{x_0 \leq x \leq x_1} (x_1 - x)(x - x_0) = \frac{h^2}{4}$$

$$|\log_{10} x - p_1(x)| \leq \frac{h^2 \cdot 0.434}{8 x_0^2} = \frac{0.0542 h^2}{x_0^2} \leq 0.0542 h^2 \quad (11.1.3)$$

ملاحظه می‌کنیم که در جدول، برای x نزدیک به 10^3 خطای درونیابی بسیار کوچکتر از x نزدیک به 1 است. همچنین است خطای ماکسیمم نزدیک به نقطه وسط $[x_0, x_1]$.

در یک جدول چهار رقمی، $h = 0.1$

$$|\log_{10} x - p_1(x)| \leq 5.42 \times 10^{-6} \quad 1 \leq x_0 < x_1 \leq 10$$

چون درایه‌های جدول تا چهار رقم داده شده‌اند (مثل $10^3 = \log_{10} 2$)، این نتیجه به اندازه کافی دقیق است. پس اگر این نتیجه به دست آمده تا این اندازه دقیق است، چرا به یک جدول دقیقتر پنج رقمی نیاز داریم؟ برای اینکه اثرات خطاهای گردکردن موجود در درایه‌های جدول را ندیده گرفته‌ایم. برای مثال، در $10^3 = \log_{10} 2$.

$$|\log_{10} 2 - 0.3010| \leq 0.00005$$

این، بر خطای درونیابی وقتی x_0 یا x_1 برابر 2 باشد غلبه دارد.

تحلیل خطای گردکردن در درونیابی خطی. گیریم

$$f(x_0) = f_0 + \epsilon_0 \quad f(x_1) = f_1 + \epsilon_1$$

که f_0 و f_1 درایه‌های جدول و ϵ_0 و ϵ_1 خطاهای گردکردن هستند. فرض می‌کنیم برای یک ϵ معلوم

$$|\epsilon_0|, |\epsilon_1| \leq \epsilon$$

در مورد جدول لگاریتم چهار رقمی، $\epsilon = 0.00005$ می‌خواهیم عبارت

$$\mathcal{E}(x) = f(x) - \frac{(x_1 - x)f_0 + (x - x_0)f_1}{x_1 - x_0} \quad x_0 \leq x \leq x_1 \quad (12.1.3)$$

را کراندار کنیم. با استفاده از $f_i = f(x_i) - \epsilon_i$

$$\begin{aligned} \mathcal{E}(x) &= f(x) - \frac{(x_1 - x)f(x_0) + (x - x_0)f(x_1)}{x_1 - x_0} \\ &\quad + \frac{(x_1 - x)\epsilon_0 + (x - x_0)\epsilon_1}{x_1 - x_0} \\ &\equiv E(x) + R(x) \end{aligned} \quad (13.1.3)$$

$$E(x) = \frac{(x - x_0)(x - x_1)}{2} f''(\xi) \quad \xi \in [x_0, x_1]$$

خطای $\mathcal{E}(x)$ مجموع خطاهای نظری درونیابی $E(x)$ و تابع $R(x)$ است که $R(x)$ به ϵ_0 و ϵ_1 بستگی دارد. چون $R(x)$ یک خط مستقیم است، ماکسیم آن روی $[x_0, x_1]$ در یکی از دو نقطه انتهایی قرار می‌گیرد،

$$\text{Max}_{x_0 \leq x \leq x_1} |R(x)| = \text{Max}\{|\epsilon_0|, |\epsilon_1|\} \leq \epsilon \quad (14.1.3)$$

با $x_1 = x_0 + h$ داریم

$$|\mathcal{E}(x)| \leq \frac{h^2}{8} \text{Max}_{x_0 \leq t \leq x_1} |f''(t)| + \text{Max}\{|\epsilon_0|, |\epsilon_1|\} \quad (15.1.3)$$

مثال در مثال لگاریتم قبلی، که از جدول چهار رقمی استفاده شده است

$$|\mathcal{E}(x)| \leq 5,42 \times 10^{-6} + 5 \times 10^{-5} \doteq 5,5 \times 10^{-5}$$

در یک جدول پنج رقمی، $h = 0.001$ ، $\epsilon = 0.000005$ و

$$|\mathcal{E}(x)| \leq 5.42 \times 10^{-8} + 5 \times 10^{-6} \approx 5.05 \times 10^{-6} \quad x_0 \leq x \leq x_1$$

تنها خطای مهم در استفاده از یک جدول پنج رقمی خطای گردکردن است. در واقع، به نظر مفید می‌آید که جدول پنج رقمی را به جدول شش رقمی تبدیل کنیم بدون اینکه اندازه فاصله h تغییر داده شود. در این صورت یک ماکسیمم خطا برابر $10^{-7} \times 5.42$ برای $\mathcal{E}(x)$ خواهیم داشت بدون اینکه محاسبات افزایش چشمگیری داشته باشند. این دلایل در مورد خطای گردکردن برای درونیابی چند جمله‌ای از درجه بالاتر نیز عمومیت می‌یابد اگرچه نتیجه آن برای $\max |R(x)|$ کمی پیچیده‌تر می‌شود (مسئله ۸ را ببینید).

توجه کنید که در هیچ یک از نتایج این بخش خطاهای گردکردن جدیدی که در محاسبه $p_n(x)$ رخ می‌دهند به حساب نیامده است. این خطاها به کمک قضایایی که در بخش بعد می‌آیند، مینیمم می‌شوند.

۲.۳ تفاضلات منقسم نیوتن

شکل لاگرانژی چندجمله‌ای درونیابی را می‌توان برای درونیابی یک تابع که به شکل جدولی داده شده است به کار برد؛ از جداول موجود در آبراموویتس و استگون^۱ (۱۹۶۴، فصل ۲۵) $l_i(x)$ آسانتر محاسبه می‌شود. ولی شکلهای دیگری وجود دارند که بسیار مناسبترند و در این بخش و بخش بعدی بحث می‌شوند. با شکل لاگرانژی، مشکل می‌توان از یک چندجمله‌ای درونیابی به یک چندجمله‌ای دیگر با یک درجه بیشتر رسید. مقایسه چندجمله‌یهای درونیابی با درجات مختلف، برای انتخاب درجه چندجمله‌ای که باید به کار رود تکنیک مفیدی است. فرمولهایی که در این بخش مطالعه می‌شوند برای نقاط شبکه‌ی نامتساوی الفاصله $\{x_i\}$ هستند. بدین لحاظ این فرمولها، برای درونیابی معکوس در جدول مناسب‌اند، نکته‌ای است که بعداً نشان خواهیم داد. این فرمولها، در بخش ۳.۳ برای نقاط شبکه‌ی متساوی الفاصله تخصیص یافته‌اند.

می‌نویسیم:

$$p_n(x) = p_{n-1} + C(x) \quad C(x) = \text{جمله تصحیحی} \quad (1.2.3)$$

پس، در حالت کلی، $C(x)$ یک چندجمله‌ای از درجه n است، زیرا معمولاً درجه p_{n-1} برابر

$n - 1$ است و درجهٔ p_n مساوی n . همچنین داریم:

$$C(x_i) = p_n(x_i) - p_{n-1}(x_i) = f(x_i) - f(x_i) = 0 \quad i = 0, 1, \dots, n-1$$

بنابراین

$$C(x) = a_n(x - x_0) \dots (x - x_{n-1})$$

چون $p_n(x_n) = f(x_n)$ ، از رابطهٔ (۱.۲.۳) به دست می‌آوریم:

$$a_n = \frac{f(x_n) - p_{n-1}(x_n)}{(x_n - x_0) \dots (x_n - x_{n-1})}$$

به دلایلی که در زیر آمده‌است، این ضریب a_n را تفاضل منقسم مرتبهٔ n ام نیوتن مربوط به f می‌نامند، و آن را با نماد زیر نشان می‌دهند.

$$a_n \equiv f[x_0, x_1, \dots, x_n]$$

پس فرمول درونیابی ما چنین می‌شود

$$p_n(x) = p_{n-1}(x) + (x - x_0) \dots (x - x_{n-1}) f[x_0, x_1, \dots, x_n] \quad (۲.۲.۳)$$

برای آنکه اطلاعات بیشتری در مورد a_n به دست آوریم، به فرمول (۷.۱.۳) لاگرانژ برای $p_n(x)$ برمی‌گردیم. می‌نویسیم

$$\Psi_n(x) = (x - x_0) \dots (x - x_n) \quad (۳.۲.۳)$$

بنابراین

$$\Psi'_n(x_i) = (x_i - x_0) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_i - x_n)$$

و اگر x یک نقطهٔ گرهی نباشد

$$p_n(x) = \sum_{j=0}^n \frac{\Psi_n(x)}{(x - x_j) \Psi'_n(x_j)} \cdot f(x_j) \quad (۴.۲.۳)$$

از آنجا که a_n ضریب x^n در $p_n(x)$ است، فرمول لاگرانژ را برای به دست آوردن ضریب x^n به کار می‌بریم. با نگاه کردن به هر جملهٔ درجهٔ n ام در فرمول (۴.۲.۳)، به دست می‌آوریم

$$f[x_0, x_1, \dots, x_n] = \sum_{j=0}^n \frac{f(x_j)}{\Psi'_n(x_j)} \quad (۵.۲.۳)$$

از این فرمول، یک ویژگی مهم برای تفاضل منقسم به دست می‌آوریم. گیریم (i_0, i_1, \dots, i_n) یک جایگشت $(0, 1, \dots, n)$ باشد. به سادگی به دست می‌آوریم:

$$\sum_{j=0}^n \frac{f(x_j)}{\Psi'_n(x_j)} = \sum_{j=0}^n \frac{f(x_{i_j})}{\Psi'_n(x_{i_j})}$$

زیرا مجموع دوم، تنها یک جابه‌جایی در مجموع اول است. از طرف دیگر برای هر جایگشت (i_0, \dots, i_n) ، مانند $(0, 1, \dots, n)$

$$f[x_0, x_1, \dots, x_n] = f[x_{i_0}, x_{i_1}, \dots, x_{i_n}] \quad (6.2.3)$$

یک فرمول مفید دیگر برای محاسبه $f[x_0, x_1, \dots, x_n]$ فرمول زیر است:

$$f[x_0, x_1, \dots, x_n] = \frac{f[x_1, \dots, x_n] - f[x_0, \dots, x_{n-1}]}{x_n - x_0} \quad (7.2.3)$$

که نامگذاری تفاضل منقسم را نیز توجیه می‌کند. این نتیجه را می‌توان از فرمول (5.2.3) یا فرمول دیگری که ذیلاً برای $p_n(x)$ بیان می‌شود، ثابت کرد:

$$p_n(x) = \frac{(x_n - x)p_{n-1}^{(0, n-1)}(x) + (x - x_0)p_{n-1}^{(1, n)}(x)}{x_n - x_0} \quad (8.2.3)$$

که در آن $p_{n-1}^{(0, n-1)}(x)$ چندجمله‌ای از درجه نایزگتر از $n-1$ است که $f(x)$ را در $\{x_0, \dots, x_{n-1}\}$ درونیایی می‌کند و $p_{n-1}^{(1, n)}(x)$ چندجمله‌ای است که $f(x)$ را در $\{x_1, \dots, x_n\}$ درونیایی می‌کند. اثبات فرمولهای (7.2.3) و (8.2.3) در مسأله ۱۳ آمده است.

با بازگشت به فرمول (2.2.3)، فرمولهای زیر را خواهیم داشت:

$$p_0(x) = f(x_0)$$

$$p_1(x) = f(x_0) + (x - x_0)f[x_0, x_1]$$

⋮

$$p_n(x) = f(x_0) + (x - x_0)f[x_0, x_1] + (x - x_0)(x - x_1)f[x_0, x_1, x_2]$$

$$+ \dots + (x - x_0) \dots (x - x_{n-1})f[x_0, x_1, \dots, x_n] \quad (9.2.3)$$

جدول ۱.۳ قالب برای ساختن تفاضلات منقسم $f(x)$

x_i	$f(x_i)$	$f[x_i, x_{i+1}]$	$f[x_i, x_{i+1}, x_{i+2}] \dots$
x_0	f_0	$f[x_0, x_1]$	
x_1	f_1	$f[x_1, x_2]$	$f[x_0, x_1, x_2] \dots$
x_2	f_2	$f[x_2, x_3]$	$f[x_1, x_2, x_3] \dots$
x_3	f_3	$f[x_3, x_4]$	$f[x_2, x_3, x_4] \dots$
x_4	f_4	$f[x_4, x_5]$	$f[x_3, x_4, x_5] \dots$
x_5	f_5		
\vdots	\vdots	\vdots	\vdots

این فرمول را فرمول تفاضل منقسم نیوتن برای چندجمله‌ای درونیاب می‌نامند. این فرمول برای محاسبات خیلی بهتر از فرمول لاگرانژ است (اگرچه گونه‌های دیگری از فرمول لاگرانژ وجود دارند که از این فرمول نیوتن کاراترند).

برای ساختن تفاضلات منقسم، قالبی را که در جدول ۱.۳ داده شده به کار می‌گیریم. صورت هرکسر تفاضل، از تفاضل دو درایه مجاور در ستونهای سمت چپ ستونی که دارید می‌سازید به دست می‌آید.

مثال یک جدول تفاضل منقسم برای $f(x) = \sqrt{x}$ می‌سازیم که در جدول ۲.۳ نشان داده

جدول ۲.۳ مثالی برای ساختن تفاضلات منقسم

x_i	$f(x_i)$	$f[x_i, x_{i+1}]$	$D^1 f[x_i]$	$D^2 f[x_i]$	$D^3 f[x_i]$
۲٫۰	۱٫۴۱۴۲۱۴	۰٫۳۴۹۲۴			
۲٫۱	۱٫۴۴۹۱۳۸	۰٫۳۴۱۰۲	-۰٫۰۴۱۱۰	۰٫۰۰۹۱۶۷	
۲٫۲	۱٫۴۸۳۲۴۰	۰٫۳۳۳۳۵	-۰٫۰۳۸۳۵	۰٫۰۰۸۳۳۳	-۰٫۰۰۲۰۸۴
۲٫۳	۱٫۵۱۶۵۷۵	۰٫۳۲۶۱۸	-۰٫۰۳۵۸۵		
۲٫۴	۱٫۵۴۹۱۹۳				

شده است. از نماد $D^r f(x_i) = f[x_i, x_{i+1}, \dots, x_{i+r}]$ استفاده کرده‌ایم. توجه داشته باشید که درایه‌های جدول در ستون $f(x_i)$ در رقم هفتم خطای گردکردن دارند و این امر بر دقت تفاضلات منقسم حاصل اثر می‌گذارد. بحثی دربارهٔ نتایج خطای گردکردن در محاسبهٔ $p_n(x)$ برای هر دو فرمول لاگرانژ و تفاضلات منقسم نیوتن، در پاؤل (۱۹۸۱، ص ۵۱) داده شده است.

یک الگوریتم ساده برای ساختن تفاضلات

$$f(x_0), f[x_0, x_1], f[x_0, x_1, x_2], \dots, f[x_0, x_1, \dots, x_n]$$

که در محاسبهٔ فرمول (۹.۲.۳) نیوتن مورد نیازند، می‌توان ارائه کرد.

الگوریتم $Divdif(d, x, n)$

۱. توجه: d و x بردارهایی با درایه‌های به ترتیب، $f(x_i)$ و x_i و $i = 0, 1, \dots, n$ هستند.

هنگام خروج، d_i شامل $f[x_0, \dots, x_i]$ است.

۲. مرحلهٔ ۴ را برای $i = 1, 2, \dots, n$ عمل کنید.

۳. مرحلهٔ ۴ را برای $j = n, n-1, \dots, i$ عمل کنید.

$$d_j := (d_j - d_{j-1}) / (x_j - x_{j-i}) \quad ۴.$$

۵. از الگوریتم خارج شوید.

برای محاسبهٔ فرمول نیوتن چندجمله‌ای درونیاب (۹.۲.۳)، ضرب تودرتوی چندجمله‌ای (۸.۹.۲) از فصل ۲ را با تغییری ساده می‌آوریم.

الگوریتم $Interp(d, x, n, t, p)$

۱. توجه: در ورودی، d و x به ترتیب بردارهایی شامل $f[x_0, \dots, x_i]$ و x_i و $i = 0, 1, \dots, n$ هستند.

هنگام خروج، p شامل مقدار $p_n(t)$ چندجمله‌ای درجهٔ n است که f را در x درونیابی می‌کند.

$$p := d_n \quad ۲.$$

۳. از مرحلهٔ ۴ برای $i = n-1, n-2, \dots, 0$ عمل کنید.

$$p := d_i + (t - x_i)p \quad ۴.$$

۵. از الگوریتم خارج شوید.

جدول ۳.۳ مثالی در استفاده از فرمول نیوتن

x	$p_1(x)$	$p_2(x)$	$p_3(x)$	$p_4(x)$	\sqrt{x}
۲٫۰۵	۱٫۴۳۱۶۷۶	۱٫۴۳۱۷۷۹	۱٫۴۳۱۷۸۲	۱٫۴۳۱۷۸۲	۱٫۴۳۱۷۸۲
۲٫۱۵	۱٫۴۶۶۶۰۰	۱٫۴۶۶۶۲۹۲	۱٫۴۶۶۶۲۸۸	۱٫۴۶۶۶۲۸۸	۱٫۴۶۶۶۲۸۸
۲٫۴۵	۱٫۵۷۱۳۷۲	۱٫۵۶۴۸۹۹	۱٫۵۶۵۲۶۰	۱٫۵۶۵۲۴۷	۱٫۵۶۵۲۴۸

مثال برای $f(x) = \sqrt{x}$ ، در جدول ۳.۳، مقادیر $p_n(x)$ برای مقادیر مختلف x و n داده شده‌اند. $p_4(x)$ چندجمله‌ای بالاترین درجه، از مقادیر تابع در نقاط شبکه از $x_0 = ۲٫۰$ تا $x_4 = ۲٫۴$ استفاده می‌کند. تفاضلات منقسم لازم در مثال آخر، در جدول ۲.۳ داده شده‌اند.

وقتی مقدار x خارج از $\mathcal{H}\{x_0, x_1, \dots, x_n\}$ می‌افتد، غالباً می‌گوییم که $p_n(x)$ تابع $f(x)$ را برون‌یابی می‌کند. در مثال اخیر، به بی‌دقتی زیاد در مقدار برون‌یابی $p_4(۲٫۴۵)$ در مقایسه با مقادیر $p_3(۲٫۱۵)$ ، $p_3(۲٫۰۵)$ توجه نمایید. ولی در این کتاب، واژه درونیابی همواره، امکان بیرون افتادن x از بازه $\mathcal{H}\{x_0, x_1, \dots, x_n\}$ را نیز شامل می‌شود.

اغلب مقدار تابع $f(x)$ را داریم، و می‌خواهیم x مقدار متناظر آن را به دست آوریم. این را درونیابی وارون می‌نامند. معمولاً استفاده‌کنندگان از جداول لگاریتم برای محاسبه آنتی لگاریتم یک عدد x ، آن را به کار می‌برند. برای محاسبه x ، آن را مانند یک متغیر وابسته و $y = f(x)$ را یک متغیر مستقل می‌گیریم. اگر جدولی از مقادیر (x_i, y_i) ، $i = 0, 1, \dots, n$ ، داده شده باشد، یک چندجمله‌ای $p_n(y)$ می‌سازیم که مقادیر x_i را در y_i ، $i = 0, 1, \dots, n$ ، درونیابی کند. در واقع، تابع معکوس $g(y) \equiv f^{-1}(y)$ را درونیابی می‌کنیم؛ فرمول خطای قضیه ۲.۳ با $x = f^{-1}(y)$ به شکل زیر در می‌آید:

$$x - p_n(y) = \frac{(y - y_0) \cdots (y - y_n)}{(n + 1)!} g^{(n+1)}(\xi) \quad (۱۰.۲.۳)$$

به ازای مقداری از $\xi \in \mathcal{H}\{y_0, y_1, \dots, y_n\}$. چنانچه لازم باشد، مشتقات $g(y)$ را می‌توان با مشتق‌گیری از فرمول مرکب $g(f(x)) = x$ محاسبه کرد. برای مثال:

$$g'(y) = \frac{1}{f'(x)} \quad \text{برای } y = f(x)$$

مثال مقادیر تابع بسل $J(x)$ را که از آبراموویتس و استگون، ۱۹۶۴، فصل ۹ استخراج شده در جدول ۴.۳ در نظر می‌گیریم. مقدار x را که به ازای آن $J(x) = ۰٫۱$ ، محاسبه می‌کنیم. جدول ۵.۳، مقادیر $p_n(y)$ را برای $n = 0, 1, \dots, 6$ با $x_0 = ۲٫۰$ به دست می‌دهد.

جدول ۴.۳. مقادیر تابع بسل $J_0(x)$

x	$J_0(x)$
۲٫۰	۰٫۲۲۳۸۹۰۷۷۹۱
۲٫۱	۰٫۱۶۶۶۰۶۹۸۰۳
۲٫۲	۰٫۱۱۰۳۶۲۲۶۶۹
۲٫۳	۰٫۰۵۵۵۳۹۷۸۴۴
۲٫۴	۰٫۰۰۲۵۰۷۶۸۳۲
۲٫۵	-۰٫۰۴۸۳۸۳۷۷۶۴
۲٫۶	-۰٫۰۹۶۸۰۴۹۵۴۴
۲٫۷	-۰٫۱۴۲۴۴۹۳۷۰۰
۲٫۸	-۰٫۱۸۵۰۳۶۰۳۳۴
۲٫۹	-۰٫۲۲۴۳۱۱۵۴۵۸

چند جمله‌ای $p_n(y)$ که تابع وارون $J_0(x)$ را درونیایی می‌کند، $g(y)$ می‌نامیم. جواب $x = ۲٫۲۱۸۶۸۳۸$ تا هشت رقم معنی‌دار درست است.

یک فرمول خطای درونیایی با استفاده از تفاضلات منقسم، گیریم t یک عدد حقیقی متمایز از نقاط گرهی x_0, \dots, x_n باشد. چند جمله‌ای $p_{n+1}(x)$ که $f(x)$ را در نقاط x_0, \dots, x_n و t درونیایی می‌کند می‌سازیم.

$$\begin{aligned}
 p_{n+1}(x) &= f(x_0) + (x - x_0)f[x_0, x_1] + \dots \\
 &\quad + (x - x_0) \dots (x - x_{n-1})f[x_0, \dots, x_n] \\
 &\quad + (x - x_0) \dots (x - x_n)f[x_0, x_1, \dots, x_n, t] \\
 &= p_n(x) + (x - x_0) \dots (x - x_n)f[x_0, \dots, x_n, t]
 \end{aligned}$$

جدول ۵.۳. مثال از درونیایی وارون

n	$p_n(y)$	$p_n(y) - p_{n-1}(y)$	$g[y_0, \dots, y_n]$
۰	۲٫۰		۲٫۰
۱	۲٫۲۱۶۲۷۵۴۲۵	$۲٫۱۶E - ۱$	-۱٫۷۴۵۶۹۴۲۸۲
۲	۲٫۲۱۸۶۱۹۶۰۸	$۲٫۳۴E - ۳$	۰٫۲۸۴۰۷۴۸۴۰۵
۳	۲٫۲۱۸۶۸۶۲۵۲	$۶٫۶۶E - ۵$	-۰٫۷۷۹۳۷۱۱۸۱۲
۴	۲٫۲۱۸۶۸۳۳۴۴	$-۲٫۹۱E - ۶$	۰٫۷۶۴۸۹۸۶۷۰۴
۵	۲٫۲۱۸۶۸۳۹۶۴	$۶٫۲۰E - ۷$	-۱٫۶۷۲۳۵۷۲۶۴
۶	۲٫۲۱۸۶۸۳۷۷۳	$-۱٫۹۱E - ۷$	۳٫۴۷۷۳۳۳۱۲۶

چون $p_{n+1}(t) = f(t)$ ، می‌گیریم $x = t$ و به دست می‌آوریم:

$$f(t) - p_n(t) = (t - x_0) \dots (t - x_n) f[x_0, \dots, x_n, t] \quad (۱۱.۲.۳)$$

که در آن $p_n(t)$ به سمت چپ معادله برده شده است. این رابطه فرمول دیگری را برای خطای فرمول خطای قبلی (۸.۱.۳) مقایسه و $\Psi_n(t)$ ، ضریب مشترک را حذف کنیم، برای مقداری مانند

$$\xi \in \mathcal{H}\{x_0, x_1, \dots, x_n, t\}$$

$$f[x_0, x_1, \dots, x_n, t] = \frac{f^{(n+1)}(\xi)}{(n+1)!}$$

برای این که نتیجه فوق را نسبت به شناسه‌های آن متقارن سازیم، معمولاً می‌گیریم $t = x_{n+1}$ ، $n = m - 1$ و فرمول زیر را برای یک مقدار $\xi \in \mathcal{H}\{x_0, \dots, x_m\}$ به دست می‌آوریم:

$$f[x_0, x_1, \dots, x_m] = \frac{f^{(m)}(\xi)}{m!} \quad (۱۲.۲.۳)$$

با این نتیجه، فرمول (۹.۲.۳) نیوتن، مانند سری بریده شده تیلر $f(x)$ حول x_0 به نظر می‌آید، به شرطی که $x_n - x_0$ خیلی بزرگ نباشد.

مثال از جدول تفاضلات منقسم ۲.۳ برای تابع $f(x) = \sqrt{x}$ ، داریم:

$$f[2_0, 2_1, \dots, 2_4] = -0.002084$$

چون $f^{(4)}(x) = -15/(16x^3\sqrt{x})$ ، می‌توان نشان داد که

$$\frac{f^{(4)}(2_3 10^3)}{4!} \doteq -0.002084$$

پس در این حالت در (۱۲.۲.۳)، داریم $\xi \doteq 2_3 1$.

مثال اگر فرمول (۱۲.۲.۳) برای برآورد مشتقات $g(y)$ ، تابع وارون $J(x)$ ، که در یکی از مثالهای قبل داده شد، استفاده شود، مشتقات $g^{(n)}(y)$ سریعاً با n افزایش می‌یابد. برای مثال،

$$g[y_0, \dots, y_6] \doteq 3_4 8$$

$$g^{(6)}(\xi) \doteq (6!)(3_4 8) \doteq 2500$$

به ازای یک ξ در $[0, 2239, 968, 0]$. برآوردهای مشابهی را می‌توان برای مشتقات دیگر انجام داد.

برای تعمیم تعریف تفاضل منقسم نیوتن، در حالتی که بعضی یا تمام نقاط گره برهم منطبق می‌شوند، فرمول دیگری را برای تفاضل منقسم معرفی می‌کنیم.

قضیه ۳.۳ ارمیت - جنوکی^۱ گیریم x_0, \dots, x_n متمایز باشند و $f(x)$ در بازه $\mathcal{H}\{x_0, x_1, \dots, x_n\}$ n بار پیوسته مشتقپذیر باشد. آنگاه

$$f[x_0, x_1, \dots, x_n] = \int_{\tau_n} \dots \int f^{(n)}(t_0 x_0 + \dots + t_n x_n) dt_1 \dots dt_n \quad (13.2.3)$$

که در آن

$$\tau_n = \{(t_1, t_2, \dots, t_n) \mid \forall t_i \geq 0, \sum_1^n t_i \leq 1\} \quad (14.2.3)$$

$$t_0 = 1 - \sum_1^n t_i$$

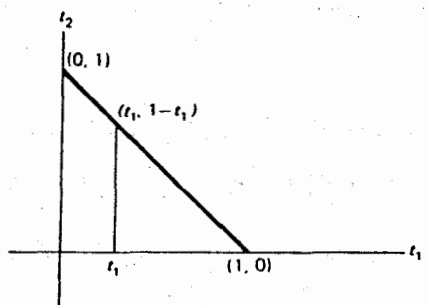
توجه داشته باشید که $t_0 \geq 0$ و $\sum_1^n t_i = 1$.

برهان نشان می‌دهیم که (۱۳.۲.۳) به ازای $n = 1, 2$ صحیح است و از این دو حالت، برهان استقرایی در حالت کلی روشن می‌شود.

۱. $n = 1$. آنگاه $\tau_1 = [0, 1]$.

$$\begin{aligned} \int_0^1 f'(t_0 x_0 + t_1 x_1) dt_1 &= \int_0^1 f'(x_0 + t_1(x_1 - x_0)) dt_1 \\ &= \frac{1}{x_1 - x_0} f(x_0 + t_1(x_1 - x_0)) \Big|_{t_1=0}^{t_1=1} \\ &= \frac{f(x_1) - f(x_0)}{x_1 - x_0} = f[x_0, x_1] \end{aligned}$$

۲. $n = 2$. آنگاه τ_2 مثلثی است با رأسهای $(0, 0)$, $(0, 1)$ و $(1, 0)$ که در شکل ۱.۳ نشان داده شده است.

شکل ۱.۳ ناحیه τ_T

$$\begin{aligned}
 & \int_{\tau_n} \int f''(t_1 x_0 + t_2 x_1 + t_3 x_2) dt_1 dt_2 \\
 &= \int_0^1 \int_0^{1-t_1} f''(x_0 + t_1(x_1 - x_0) + t_2(x_2 - x_0)) dt_2 dt_1 \\
 &= \int_0^1 \frac{1}{x_2 - x_0} \left[f'(x_0 + t_1(x_1 - x_0) + t_2(x_2 - x_0)) \right]_{t_2=0}^{t_2=1-t_1} dt_1 \\
 &= \frac{1}{x_2 - x_0} \left[\int_0^1 f'(x_2 + t_1(x_1 - x_2)) dt_1 \right. \\
 &\quad \left. - \int_0^1 f'(x_0 + t_1(x_1 - x_0)) dt_1 \right] \\
 &= \frac{1}{x_2 - x_0} \{ f[x_1, x_2] - f[x_0, x_1] \} = f[x_0, x_1, x_2]
 \end{aligned}$$

با روندی مشابه برای حالت کلی عمل می‌کنیم. با یک بار انتگرالگیری مسأله را به یک بُعد کمتر کاهش می‌دهیم. سپس به فرضهای استقراء استناد نموده (۷.۲.۳) را به کار می‌بریم تا برهان کامل شود. ■

اکنون می‌توانیم (۱۳.۲.۳)، را در $f[x_0, x_1, \dots, x_n]$ به کار ببریم. می‌بینیم که اگر $f(x)$ در $\mathcal{H}\{x_0, \dots, x_n\}$ n بار پیوسته مشتقپذیر باشد، آنگاه $f[x_0, \dots, x_n]$ تابعی است پیوسته از n متغیر x_0, \dots, x_n ، خواه این متغیرها متمایز باشند یا نباشند. به عنوان مثال، اگر فرض کنیم که تمام نقاط به x_0 بدل شوند، آنگاه برای تفاضل منقسم مرتبه n ,

$$\begin{aligned}
 f[x_0, \dots, x_0] &= \int_{\tau_n} \dots \int f^{(n)}(x_0) dt_1 \dots dt_n \\
 &= f^{(n)}(x_0) \cdot \text{Vol}(\tau_n)
 \end{aligned}$$

که در آن $\text{Vol}(\tau_n)$ حجم τ_n است. با توجه به مسأله ۱۵ داریم $\text{Vol}(\tau_n) = 1/n!$ و بنابراین

$$f[x_0, \dots, x_n] = \frac{f^{(n)}(x_0)}{n!} \quad (۱۵.۲.۳)$$

این رابطه را می‌توانستیم مستقیماً از رابطه (۱۲.۲.۳) پیشینی کنیم. ولی اگر فقط بعضی از نقاط گرهی یکی می‌شدند، باید برای توجیه وجود $f[x_0, \dots, x_n]$ از رابطه (۱۳.۲.۳) استفاده می‌کردیم. در کاربردهای انتگرالگیری عددی، باید بدانیم که آیا

$$\frac{d}{dx} f[x_0, \dots, x_n, x] \quad (۱۶.۲.۳)$$

وجود دارد یا نه. اگر f ، $n+2$ بار پیوسته مشتقپذیر باشد، آنگاه می‌توانیم قضیه ۳.۳ را به‌کار ببریم. با استفاده از قضایای مشتقگیری و انتگرالگیری نسبت به یک پارامتر انتگرالده می‌توانیم وجود (۱۶.۲.۳) را نتیجه بگیریم. از راه مستقیم‌تر

$$\begin{aligned} \frac{d}{dx} f[x_0, \dots, x_n, x] &= \lim_{h \rightarrow 0} \frac{f[x_0, \dots, x_n, x+h] - f[x_0, \dots, x_n, x]}{h} \\ &= \lim_{h \rightarrow 0} \frac{f[x_0, \dots, x_n, x+h] - f[x, x_0, \dots, x_n]}{h} \\ &= \lim_{h \rightarrow 0} f[x, x_0, \dots, x_n, x+h] = f[x, x_0, x_1, \dots, x_n, x] \end{aligned}$$

$$\frac{d}{dx} f[x_0, x_1, \dots, x_n, x] = f[x_0, x_1, \dots, x_n, x, x] \quad (۱۷.۲.۳)$$

وجود و پیوستگی طرف راست، با استفاده از (۱۳.۲.۳) تضمین می‌شود.

نظریه‌ای پر بار شامل درونبایی چندجمله‌ای و تفاضلات منقسم وجود دارد، ولی ما در اینجا موضوع را با یک قضیه سر راست پایانی خاتمه می‌دهیم. اگر $f(x)$ یک چندجمله‌ای از درجه m باشد، آنگاه

$$f[x_0, \dots, x_n, x] = \begin{cases} \text{چندجمله‌ای از درجه } m-n-1 & n < m-1 \\ a_m & n = m-1 \\ 0 & n > m-1 \end{cases} \quad (۱۸.۲.۳)$$

که در آن، (جملات از درجه پایینتر) $f(x) = a_m x^m + \dots$ برای برهان قضیه به مسأله ۱۴ مراجعه کنید.

۳.۳ تفاضلات متناهی و جدول جهتدار فرمولهای درونیابی

در این بخش، تفاضلات پیشرو و پسرو را معرفی می‌کنیم، همراه با فرمولهای درونیابی که از آنها استفاده می‌شود. این تفاضلات را کلاً تفاضلات متناهی می‌نامند و آنها را در ساختن فرمولهای درونیابی برای جداولی که طولها، یعنی $\{x_i\}$ ها متساوی الفاصله‌اند، به‌کار می‌برند. این فرمولهای درونیابی در حل عددی معادلات دیفرانسیل معمولی و جزئی نیز کاربرد دارند. بعلاوه، از تفاضلات متناهی برای مسأله تعیین بیشترین درجه چندجمله‌ای درونیابی، که بتوان با اطمینان به‌کار برد، بر پایه دقتی که درایه‌های جدول دارند، استفاده می‌شود. و تفاضلات متناهی را برای کشف نوفه در داده‌ها، وقتی نوفه نسبت به خطای گردکردن یا خطای عدم قطعیت اندازه‌گیری فیزیکی زیاد باشد، می‌توان به‌کار برد. این موضوع در بخش ۴.۳ بحث شده است.

برای یک $h > 0$ مفروض، تعریف می‌کنیم

$$\Delta_h f(z) = f(z+h) - f(z)$$

معمولاً h از قراین مطلب معلوم می‌شود، و می‌نویسیم

$$\Delta f(z) = f(z+h) - f(z) \quad (1.3.3)$$

این مقدار را تفاضل پیشرو f در z ، و Δ را عملگر تفاضل پیشرو می‌نامند. ما همیشه با نقاط گرهی متساوی الفاصله $x_i = x_0 + ih$ ، $i = 0, 1, \dots, n$ کار می‌کنیم. پس می‌نویسیم

$$\Delta f(x_i) = f(x_{i+1}) - f(x_i)$$

یا به گونه خلاصه‌تر

$$\Delta f_i = f_{i+1} - f_i \quad f(x_i) = f_i \quad (2.3.3)$$

برای $r \geq 0$ ، تعریف می‌کنیم

$$\Delta^{r+1} f(z) = \Delta^r f(z+h) - \Delta^r f(z) \quad (3.3.3)$$

با $\Delta^0 f(z) = f(z)$ ، $\Delta^r f(z)$ تفاضل پیشرو مرتبه r ام در z است. تفاضلات پیشرو به‌سادگی محاسبه می‌شوند، و مثالهایی بعداً در ارتباط با فرمول درونیابی داده خواهد شد.

ابتدا، با به‌کار بردن فرمول تفاضلات منقسم نیوتن، نتایجی برای عملگر تفاضل پیشرو به‌دست می‌آوریم.

لم ۱ برای $k \geq 0$

$$f[x_0, x_1, \dots, x_k] = \frac{1}{k!h^k} \Delta^k f. \quad (۴.۳.۳)$$

برهان برای $k = 0$ ، درستی نتیجه بدیهی است. برای $k = 1$

$$f[x_0, x_1] = \frac{f_1 - f_0}{x_1 - x_0} = \frac{1}{h} \Delta f.$$

که (۴.۳.۳) را نشان می‌دهد. گیریم حکم (۴.۳.۳) برای تمام تفاضلات پیشرو از مرتبه $k \leq r$ درست است. در این صورت برای $k = r + 1$ ، با استفاده از (۷.۲.۳)

$$f[x_0, x_1, \dots, x_{r+1}] = \frac{f[x_1, \dots, x_{r+1}] - f[x_0, \dots, x_r]}{x_{r+1} - x_0}$$

با به‌کارگرفتن فرض استقراء، این برابر است با

$$\frac{1}{(r+1)h} \left[\frac{1}{r!h^r} \Delta^r f_1 - \frac{1}{r!h^r} \Delta^r f_0 \right] = \frac{1}{(r+1)!h^{r+1}} \Delta^{r+1} f.$$

اکنون فرمول (۹.۲.۳) درونیایی نیوتن را به صورت فرمولی که از تفاضلات پیشرو به جای تفاضلات منقسم، استفاده می‌کند درمی‌آوریم. برای یک مقدار داده شده x که به‌ازای آن چندجمله‌ای درونیاب را محاسبه خواهیم کرد. برای نشان دادن وضعیت x نسبت به x_0 تعریف می‌کنیم

$$\mu = \frac{x - x_0}{h}$$

برای مثال، $\mu = ۱.۶$ یعنی x در $\frac{6}{10}$ فاصله x_1 تا x_2 قرار دارد. به فرمولی برای

$$(x - x_0) \dots (x - x_k)$$

برحسب متغیر μ ، نیاز داریم:

$$x - x_j = x_0 + \mu h - (x_0 + jh) = (\mu - j)h$$

$$(x - x_0) \dots (x - x_k) = \mu(\mu - 1) \dots (\mu - k) h^{k+1} \quad (۵.۳.۳)$$

از ترکیب (۴.۳.۳) و (۵.۳.۳) با فرمول درونیایی تفاضل منقسم (۹.۲.۳)، به‌دست می‌آوریم:

$$p_n(x) = f_0 + \mu h \cdot \frac{\Delta f_0}{h} + \mu(\mu - 1) h^2 \cdot \frac{\Delta^2 f_0}{2!h^2} \\ + \dots + \mu(\mu - 1) \dots (\mu - n + 1) h^n \cdot \frac{\Delta^n f_0}{n!h^n}$$

ضرایب دوجمله‌ای را تعریف می‌کنیم،

$$\binom{\mu}{k} = \frac{\mu(\mu-1)\dots(\mu-k+1)}{k!} \quad k > 0 \quad (۶.۳.۳)$$

و $\binom{\mu}{0} = 1$ آنگاه

$$p_n(x) = \sum_{j=0}^n \binom{\mu}{j} \Delta^j f_0 \quad \mu = \frac{x-x_0}{h} \quad (۷.۳.۳)$$

این شکل تفاضل پیشرو نیوتن برای چندجمله‌ای درونیاب است.

مثال برای $n = 1$

$$p_1(x) = f_0 + \mu \Delta f_0 \quad (۸.۳.۳)$$

این فرمولی است که اغلب در جدول هنگام درونیابی خطی از آن استفاده می‌شود.

برای $n = 2$

$$p_2(x) = f_0 + \mu \Delta f_0 + \frac{\mu(\mu-1)}{2} \Delta^2 f_0 \quad (۹.۳.۳)$$

که یک شکل ساده محاسبه‌پذیر از چندجمله‌ای درونیاب درجه دوم است.

تفاضلات پیشرو به روندی مشابه آنچه که در تفاضلات منقسم گفته شده ساخته می‌شوند، ولی در اینجا عملهای تقسیم وجود ندارند (جدول ۶.۳ را ببینید).

مثال تفاضلات پیشرو برای $f(x) = \sqrt{x}$ در جدول ۷.۳ داده شده‌اند. مقادیر چندجمله‌ای درونیاب همان مقادیری هستند که با استفاده از فرمول تفاضلات منقسم نیوتن پیدا شده‌اند، ولی فرمول تفاضلات پیشرو (۷.۳.۳) برای محاسبه به مراتب آسانتر است.

مثال $p_n(x)$ را با استفاده از جدول ۷.۳ برای $n = 1, 2, 3, 4$ با $x = 2.15$ محاسبه کنید. توجه کنید که $\sqrt{2.15} = 1.4662878$ و $\mu = 1.5$.

$$p_1(x) = 1.414214 + 1.5(0.034924) = 1.414214 + 0.52386 = 1.4666$$

$$p_2(x) = p_1(x) + \frac{(1.5)(0.5)}{2} (-0.000822) = 1.4666 - 0.00030825$$

$$= 1.466292$$

جدول ۶.۳ قالب ساختن تفاضلات پیشرو

x_i	f_i	Δf_i	$\Delta^2 f_i$	$\Delta^3 f_i$...
x_0	f_0	Δf_0			
x_1	f_1	Δf_1	$\Delta^2 f_0$		
x_2	f_2	Δf_2	$\Delta^2 f_1$	$\Delta^3 f_0$	
x_3	f_3	Δf_3	$\Delta^2 f_2$	$\Delta^3 f_1$	\vdots
x_4	f_4	Δf_4	$\Delta^2 f_3$	$\Delta^3 f_2$	
x_5	f_5				
\vdots	\vdots				

$$p_2(x) = p_2(x) + \frac{(1.5)(0.5)(-0.5)}{6} (0.0000055) = 1.466292 - 0.0000034 = 1.466288$$

$$p_2(x) = p_2(x) + \frac{(1.5)(0.5)(-0.5)(-1.5)}{24} (-0.000005) = 1.466288 + 0.0000012 = 1.466288$$

جمله‌های تصحیحی به‌سادگی محاسبه می‌شوند، و با مشاهده اندازه آنها، یک تصوّر دقیق‌کلی در این که چه وقت درجه n به‌قدر کافی بزرگ است به‌دست می‌آید. توجه کنید که هفت رقم دقت در جدول مقادیر \sqrt{x} حداکثر به یک رقم دقت در تفاضل پیشرو $\Delta^4 f_0$

جدول ۷.۳ جدول تفاضلات پیشرو برای $f(x) = \sqrt{x}$

x_i	f_i	Δf_i	$\Delta^2 f_i$	$\Delta^3 f_i$	$\Delta^4 f_i$
۲.۰	۱.۴۱۴۲۱۴				
		۰.۳۴۹۲۴			
۲.۱	۱.۴۴۹۱۳۸		-۰.۰۰۰۸۲۲		
		۰.۳۴۱۰۲		۰.۰۰۰۰۵۵	
۲.۲	۱.۴۸۳۲۴۰		-۰.۰۰۰۷۶۷		-۰.۰۰۰۰۰۵
		۰.۳۳۳۳۵		۰.۰۰۰۰۵۰	
۲.۳	۱.۵۱۶۵۷۵		-۰.۰۰۰۷۱۷		
		۰.۳۲۶۱۸			
۲.۴	۱.۵۴۹۱۹۳				

می‌انجامد. تفاضلات پیشرو از مرتبه بیشتر از ۳ تقریباً همه نتیجه خطاهای گردکردن در تفریق درایه‌های جدول است؛ در نتیجه درونیابی در این جدول باید به چندجمله‌بیهای از درجه کمتر از چهار محدود شود. این موضوع با یک توجیه نظری بیشتری در بخش بعد خواهد آمد.

شکلهای دیگری از تفاضلات و فرمولهای درونیابی متناظر با آنها وجود دارد.

تفاضل پسرو را چنین تعریف می‌کنیم:

$$\begin{aligned}\nabla f(z) &= f(z) - f(z - h) \\ \nabla^{r+1} f(z) &= \nabla^r f(z) - \nabla^r f(z - h) \quad r \geq 1\end{aligned}\quad (10.3.3)$$

نتایج کاملاً مشابهی با تفاضلات پیشرو می‌توان به دست آورد. و ما فرمول درونیابی تفاضل پسرو نیوتن را به دست می‌آوریم.

$$p_n(x) = f_0 + \binom{-\nu}{1} \nabla f_0 + \binom{-\nu+1}{2} \nabla^2 f_0 + \dots + \binom{-\nu+n-1}{n} \nabla^n f_0 \quad (11.3.3)$$

در این فرمول، مانند قبل، گره‌های درونیابی عبارت‌اند از $x_0, x_{-1}, x_{-2}, \dots, x_{-n}$ با $x_{-j} = x_0 - jh$ عدد ν با

$$\nu = \frac{x_0 - x}{h}$$

داده می‌شود، و بدین معناست که وقتی از این فرمول استفاده می‌نماییم x معمولاً کوچکتر از x_0 است. نمودار تفاضلات پسرو را می‌توان به طریق مشابه با نمودار تفاضلات پیشرو ساخت. فرمول تفاضل پسرو در فصل ۶، برای مطالعه خانواده فرمولهای آدامز^۱ (به پاس احترام به جان کوچ آدامز ستاره‌شناس انگلیسی قرن نوزدهم) در حل عددی معادلات دیفرانسیل به‌کار برده شده است.

فرمولهای تفاضلی دیگری و فرمولهای درونیابی متناظر با آنها را می‌توان پیدا کرد. چون کاربرد آنها بسیار کمتر از آنهاست که بیان شد، خوانندگان را به هیلدبرانت^۲ ارجاع می‌دهیم.

۴.۳ خطاها در داده‌ها و تفاضلات پیشرو

برای پیدا کردن نوفه در داده‌های فیزیکی، مادامی که نوفه نسبت به حدود خطای تجربی بزرگ است می‌توانیم از یک جدول تفاضل پیشرو، استفاده کنیم. با چند لم مقدماتی شروع می‌کنیم.

لم ۲ $\Delta^r f(x_i) = h^r f^{(r)}(\xi_i)$ برای مقداری از $x_i \leq \xi_i \leq x_{i+r}$.

برهان با استفاده از لم ۱ و (۱۲.۲.۳).

$$\Delta^r f_i = h^r r! f[x_i, \dots, x_{i+r}] = h^r r! \frac{f^{(r)}(\xi_i)}{r!} = h^r f^{(r)}(\xi_i)$$

لم ۳ برای هر دو تابع دلخواه f و g ، و هر دو مقدار ثابت α و β ،

$$\Delta^r(\alpha f(x) + \beta g(x)) = \alpha \Delta^r f(x) + \beta \Delta^r g(x) \quad r \geq 0$$

برهان حکم برای $r = 0$ یا $r = 1$ بدیهی است. فرض کنید حکم برای هر $r \leq n$ درست است، و آن را برای $r = n + 1$ ثابت کنید؛ با به کار بردن تعریف (۳.۳.۳) و فرض استقرا:

$$\begin{aligned} \Delta^{n+1}[\alpha f(x) + \beta g(x)] &= \Delta^n[\alpha f(x+h) + \beta g(x+h)] \\ &\quad - \Delta^n[\alpha f(x) + \beta g(x)] \\ &= \alpha \Delta^n f(x+h) + \beta \Delta^n g(x+h) \\ &\quad - \alpha \Delta^n f(x) - \beta \Delta^n g(x) \end{aligned}$$

سپس با جابه‌جا کردن جملات، داریم:

$$\begin{aligned} &\alpha[\Delta^n f(x+h) - \Delta^n f(x)] + \beta[\Delta^n g(x+h) - \Delta^n g(x)] \\ &= \alpha \Delta^{n+1} f(x) + \beta \Delta^{n+1} g(x) \end{aligned}$$

لم ۲ می‌گوید که اگر مشتقات $f(x)$ کراندار باشند، یا اگر در مقایسه با h^{-n} به سرعت افزایش نیابند، آنگاه تفاضلات پیشرو $\Delta^n f(x)$ ، با افزایش n ، بایستی کوچکتر شوند. ما بعد از خطاهای گردکردن و خطاهای دیگری با اندازه بزرگتر از خطاهای گردکردن را در نظر می‌گیریم. می‌نویسیم

$$f(x_i) = \tilde{f}_i + e(x_i) \quad i = 0, 1, 2, \dots \quad (1.4.3)$$

که \tilde{f}_i مقداری است که در جدول برای ساختن جدول تفاضل پیشرو به کار می‌بریم. پس:

$$\begin{aligned}\Delta^r \tilde{f}_i &= \Delta^r f(x_i) - \Delta^r e(x_i) \\ &= h^r f^{(r)}(\xi_i) - \Delta^r e(x_i)\end{aligned}\quad (2.4.3)$$

همان گونه که در جدول تفاضل پیشرو برای $f(x) = \sqrt{x}$ نشان داده شد، اولین عبارت با افزایش r کوچکتر می‌شود.

برای درک بهتر از رفتار $\Delta^r e(x_i)$ ، حالت ساده زیر را در نظر می‌گیریم که در آن

$$e(x_i) = \begin{cases} 0 & i \neq k \\ \varepsilon & i = k \end{cases}\quad (3.4.3)$$

تفاضل پیشرو این تابع در جدول ۸.۳ داده شده است. می‌توان ثابت کرد که ستون $\Delta^r e(x_i)$ به شکل زیر است:

$$0, \dots, 0, \varepsilon, -\binom{r}{1}\varepsilon, \binom{r}{2}\varepsilon, -\binom{r}{3}\varepsilon, \dots, (-1)^{r+1}\varepsilon, 0, \dots\quad (4.4.3)$$

بنابراین اثر خطای گردکردن، وقتی تفاضلات درجات بالاتر شکل می‌گیرند، پخش می‌شود و افزایش می‌یابد.

درباره خطاهای گردکردن که یک تابع کلی خطا را مشخص می‌کنند، مانند تابع خطا در (۱.۴.۳)، به اثرات آنها می‌توان به دیده مجموع توابع (۳.۴.۳) نگاه کرد. چون مقادیر $e_i(x)$ به طور کلی از نظر اندازه و علامت تغییر می‌یابند، اثرات آنها ظاهراً به تصادف روی هم می‌افتند. ولی اندازه تفاضلات آنها همچنان افزایش می‌یابد، و در نتیجه، تفاضلات مرتبه بالای مقادیر جدول \tilde{f}_i سرانجام بی‌فایده می‌شوند. هنگامی که اندازه تفاضلات $\Delta^r \tilde{f}_i$ با ازدیاد r شروع به افزایش می‌کند، به احتمال زیاد، خطای گردکردن بر آنها غلبه می‌کند و این تفاضلات نباید به کار روند. لذا از یک فرمول درونیابی با درجه کمتر از r باید استفاده شود.

کشف نوفه در داده‌ها. همین تحلیل را می‌توان برای کشف و تصحیح خطاهای منفردی که نسبت به خطاهای گردکردن بزرگ هستند، به کار برد. چون (۲.۴.۳) می‌گوید که اثر خطاها سرانجام غالب می‌شود، ما دنبال الگویی مانند (۴.۴.۳) می‌گردیم. روش کلی در جدول ۹.۳ نشان داده شده است. از (۲.۴.۳)

$$\Delta^r e(x_i) = \Delta^r f(x_i) - \Delta^r \tilde{f}_i$$

جدول ۸.۳ تفاضلات پیشرو برای تابع خطای $e(x)$

x_i	$e(x_i)$	$\Delta e(x_i)$	$\Delta^2 e(x_i)$	$\Delta^3 e(x_i)$
\vdots	\vdots	\circ	\vdots	\circ
x_{k-2}	\circ	\vdots	\circ	\vdots
		\circ		ε
x_{k-1}	\circ		ε	
		ε		-3ε
x_k	ε		-2ε	
		$-\varepsilon$		3ε
x_{k+1}	\circ		ε	
		\circ		$-\varepsilon$
x_{k+2}	\circ		\circ	\circ
\vdots	\vdots	\vdots	\vdots	\vdots

جدول ۹.۳ مثالی از کشف یک خطای منفرد در داده‌ها

\tilde{f}_i	$\Delta \tilde{f}_i$	$\Delta^2 \tilde{f}_i$	$\Delta^3 \tilde{f}_i$	خطاهای تخمین	تخمین $\Delta^3 f(x_i)$
$\circ.۱۰۳۹۶$					
$\circ.۱۲۰۹۶$	$\circ.۰۱۷۰۰$				
		$-\circ.۰۰۰۰۱۴$			
$\circ.۱۳۷۸۲$	$\circ.۰۱۶۸۶$	$-\circ.۰۰۰۰۱۷$	$-\circ.۰۰۰۰۰۳$	\circ	$-\circ.۰۰۰۰۰۳$
			$-\circ.۰۰۰۰۰۲$	\circ	$-\circ.۰۰۰۰۰۲$
$\circ.۱۵۴۵۱$	$\circ.۰۱۶۶۹$	$-\circ.۰۰۰۰۱۹$	$\circ.۰۰۰۰۰۶$	ε	$-\circ.۰۰۰۰۰۲$
$\circ.۱۷۱۰۱$	$\circ.۰۱۶۵۰$	$-\circ.۰۰۰۰۱۳$			
			$-\circ.۰۰۰۰۲۵$	-3ε	$-\circ.۰۰۰۰۰۲$
$\circ.۱۸۷۳۸$	$\circ.۰۱۶۳۷$	$-\circ.۰۰۰۰۳۸$			
			$\circ.۰۰۰۰۲۱$	3ε	$-\circ.۰۰۰۰۰۲$
$\circ.۲۰۳۳۷$	$\circ.۰۱۵۹۹$	$-\circ.۰۰۰۰۱۷$			
			$-\circ.۰۰۰۰۱۰$	$-\varepsilon$	$-\circ.۰۰۰۰۰۲$
$\circ.۲۱۹۱۹$	$\circ.۰۱۵۸۲$	$-\circ.۰۰۰۰۲۷$			
$\circ.۲۳۴۷۴$	$\circ.۰۱۵۵۵$				

با استفاده از $r = 3$ و انتخاب یکی از درایه‌های خطا به‌طور دلخواه، مثلاً اولین، داریم

$$\varepsilon = -\circ.۰۰۰۰۰۲ - (\circ.۰۰۰۰۰۶) = -\circ.۰۰۰۰۰۸$$

این خطا را امتحان می‌کنیم تا ببینیم چگونه ستون $\Delta^3 \tilde{f}_i$ را تغییر می‌دهد. (جدول ۱۰.۳ را ببینید).

جدول ۱۰.۳ تصحیح یک خطای داده‌ها

$\Delta^2 \tilde{f}_i$	$\Delta^2 e(x_i)$	$\Delta^2 f(x_i)$
-۰.۰۰۰۰۰۲	۰.۰	-۰.۰۰۰۰۰۲
۰.۰۰۰۰۰۶	-۰.۰۰۰۰۰۸	-۰.۰۰۰۰۰۲
-۰.۰۰۰۰۲۵	۰.۰۰۰۰۲۴	-۰.۰۰۰۰۰۱
۰.۰۰۰۰۲۱	-۰.۰۰۰۰۲۴	-۰.۰۰۰۰۰۳
-۰.۰۰۰۰۱۰	۰.۰۰۰۰۰۸	-۰.۰۰۰۰۰۲

با یک انتخاب دیگر ϵ ، مثلاً $\epsilon = -۰.۰۰۰۰۰۷$ ، بهبودی حاصل نمی‌شود، اگرچه نتایج ممکن است به همان خوبی بماند.

به عقب برمی‌گردیم، درایه $\tilde{f}_i = ۰.۱۸۷۳۸$ بایستی چنین باشد:

$$f(x_i) = \tilde{f}_i + e(x_i) = ۰.۱۸۷۳۸ + (-۰.۰۰۰۰۰۸) = ۰.۱۸۷۳۰$$

در جدولی که دویا سه خطای منفرد وجود داشته باشد، تفاضلات مراتب بالاتر آنها ممکن است تداخل کنند و کشف خطاها را مشکلتر سازند (مسأله ۲۲ را ببینید).

۵.۳ نتایج دیگری درباره خطای درونیابی

فرمول خطا را دوباره در نظر می‌گیریم.

$$f(x) - p_n(x) = \frac{(x - x_0) \dots (x - x_n)}{(n+1)!} f^{(n+1)}(\xi_x) \quad \xi_x \in \mathcal{H}\{x_0, \dots, x_n, x\} \quad (۱.۵.۳)$$

فرض می‌کنیم که $f(x)$ در یک بازه I ، که شامل $\mathcal{H}\{x_0, x_1, \dots, x_n, x\}$ است، به ازای جمیع مقادیر x مورد نظر، $n+1$ بار پیوسته مشتق‌پذیر است. چون ξ_x مجهول است، باید برای محاسبه (۱.۵.۳)، به جای $f^{(n+1)}(\xi_x)$ قرار دهیم

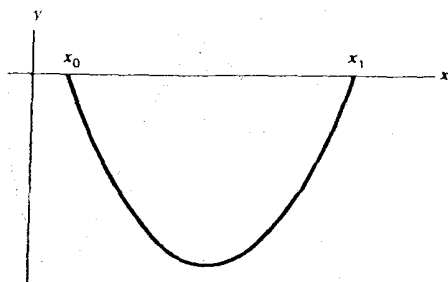
$$c_{n+1} = \text{Max}_{t \in I} |f^{(n+1)}(t)| \quad (۲.۵.۳)$$

توجه خود را به کراندارنمودن چندجمله‌ای

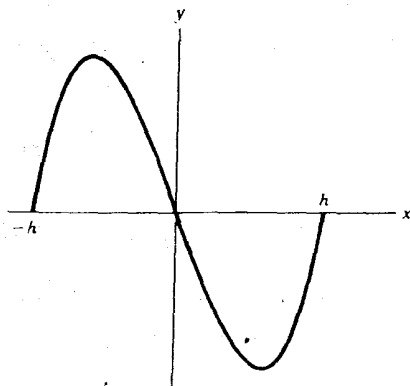
$$\Psi_n(x) = (x - x_0) \dots (x - x_n) \quad (۳.۵.۳)$$

متمرکز می‌کنیم. از (۱.۵.۳)، به دست می‌آوریم

$$\text{Max}_{x \in I} |f(x) - p_n(x)| \leq \frac{c_{n+1}}{(n+1)!} \text{Max}_{x \in I} |\Psi_n(x)| \quad (۴.۵.۳)$$



شکل ۲.۳ $y = \Psi_1(x)$



شکل ۳.۳ $y = \Psi_2(x)$

ما فقط حالتی را که گره‌ها متساوی‌فاصله هستند در نظر می‌گیریم: $x_j = x_0 + jh$, $j = 0, 1, \dots, n$. ابتدا حالت‌های ویژه‌ای از مقادیر n را در نظر می‌گیریم و سپس برای حالت کلی n اظهار نظر می‌کنیم.

حالت ۱. $n = 1$. $\Psi_1(x) = (x - x_0)(x - x_1)$. به‌سادگی دیده می‌شود

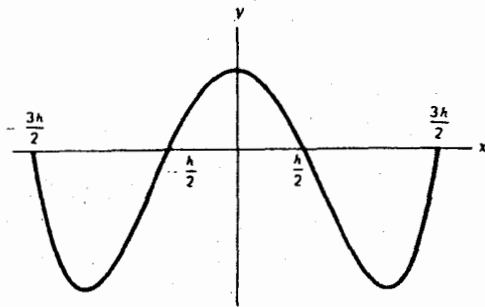
$$\text{Max}_{x_0 \leq x \leq x_1} |\Psi_1(x)| = \frac{h^2}{4}$$

برای توضیح به شکل ۲.۳ نگاه کنید.

حالت ۲. $n = 2$. برای کراندار نمودن $\Psi_2(x)$ در $[x_0, x_2]$ ، آن را بر محور x ها انتقال می‌دهیم تا چند جمله‌ای هم‌ارز زیر را پیدا کنیم

$$\hat{\Psi}_2(x) = (x + h)x(x - h)$$

که نمودار آن در شکل ۳.۳ نشان داده شده است. شکل و اندازه آن درست مانند چند جمله اصلی



شکل ۴.۳ $y = \hat{\Psi}_r(x)$

$\Psi_r(x)$ است، ولی از نظر تحلیلی ساده‌تر می‌توان آن را کراندار کرد. با استفاده از $\hat{\Psi}(x)$ به سادگی به دست می‌آوریم

$$\text{Max}_{x_1 - h/2 \leq x \leq x_1 + h/2} |\Psi_r(x)| = 0.375h^2$$

$$\text{Max}_{x_0 \leq x \leq x_2} |\Psi_r(x)| = \frac{2\sqrt{3}}{9}h^2 \doteq 0.385h^2 \quad (5.5.3)$$

بنابراین قرار گرفتن x در بازه $[x_0, x_2]$ نزدیک x_1 مهم نیست، اگرچه برای درونیابیهای درجه بالاتر، این یک تفاوتی به بار خواهد آورد. از ترکیب (۵.۵.۳) و (۴.۵.۳) داریم

$$\text{Max}_{x_0 \leq x \leq x_2} |f(x) - p_r(x)| \leq \frac{\sqrt{3}}{27}h^2 \cdot \text{Max}_{x_0 \leq t \leq x_2} |f^{(r)}(t)| \quad (6.5.3)$$

$$\sqrt{3}/27 \doteq 0.064$$

حالت ۳. $n = 3$. مانند گذشته، چندجمله‌ای را برای مقارن کردن گره‌ها نسبت به مرکز انتقال

می‌دهیم، به دست می‌آوریم:

$$\hat{\Psi}_r(x) = (x^2 - \frac{9}{4}h^2)(x^2 - \frac{1}{4}h^2)$$

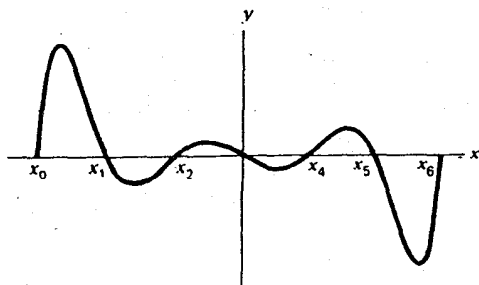
نمودار $\hat{\Psi}_r(x)$ در شکل ۴.۳ نشان داده شده است. با این تغییر

$$\text{Max}_{x_1 \leq x \leq x_2} |\Psi_r(x)| = \frac{9}{16}h^4 \doteq 0.56h^4$$

$$\text{Max}_{x_0 \leq x \leq x_2} |\Psi_r(x)| = h^4$$

بنابراین در درونیابی $f(x)$ در نقطه x ، گره‌ها باید به گونه‌ای انتخاب شوند که $x_1 < x < x_2$. در این صورت از (۱.۵.۳)،

$$\text{Max}_{x_1 \leq x \leq x_2} |f(x) - p_r(x)| \leq \frac{3h^4}{128} \cdot \text{Max}_{x_0 \leq t \leq x_2} |f^{(r)}(t)| \quad (7.5.3)$$



شکل ۵.۳ $y = \Psi_\epsilon(x)$

حالت ۴. برای حالت کلی $n > 3$ ، رفتاری که هم‌اکنون برای $n = 3$ نشان داده شده، تشدید می‌شود. برای مثال، نمودار $\Psi_\epsilon(x)$ را که در شکل ۵.۳ نشان داده شده است در نظر می‌گیریم. مانند قبل، می‌توانیم نشان دهیم

$$\text{Max}_{x_1 \leq x \leq x_2} |\Psi_\epsilon(x)| \doteq 12,36h^4$$

$$\text{Max}_{x_4 \leq x \leq x_6} |\Psi_\epsilon(x)| \doteq 95,8h^4$$

برای آنکه خطای درونیایی به حداقل برسد، گره‌های درونیایی باید به گونه‌ای انتخاب شوند که نقطهٔ درونیایی x در حد ممکن به نقطهٔ وسط $[x_0, x_n]$ نزدیک باشد. وقتی n افزایش می‌یابد، برای درونیایی در گره‌های متساوی‌الفاصله، تأکید زیادی می‌شود که نقاط طوری انتخاب شوند که x نزدیک وسط $[x_0, x_n]$ قرار گیرد.

مثال درونیایی درجهٔ پنج $J_5(x)$ را در $x = 2,45$ در نظر می‌گیریم، با مقادیر $J_5(x)$ که در جدول ۴.۳، بخش ۲.۳ داده شده‌اند. اول $x_0 = 2,4$ و $x_5 = 2,9$ را به‌کار می‌بریم. به‌دست می‌آوریم:

$$P_5(2,45) = -0,232267384 \quad \text{خطا} = -4,9 \times 10^{-9}$$

ثانیاً $x_0 = 2,2$ و $x_5 = 2,7$ را به‌کار می‌بریم. آنگاه

$$P_5(2,45) = -0,232267423 \quad \text{خطا} = 1,0 \times 10^{-9}$$

در اثر نزدیکی وضعیت x به وسط $[x_0, x_5]$ ، خطا پنج بار کوچکتر است.

در جداول آبراموویتس و استگون^۱ (۱۹۶۴) تعداد زیادی ارقام معنی‌دار داده شده‌اند، و h

فاصله شبکه به همان نسبت کوچک نیست. در نتیجه، درونیابی با درجه بالا باید به کار گرفته شود. اگرچه این کار بیشتری برای استفاده کننده از جدول ایجاد می نماید، ولی موجب می شود که جدول در شکل فشرده تر جای بسیار کمتر بگیرد و می توان جداول بیشتری برای توابع بیشتری در یک حجم جای داد. وقتی که با استفاده از این جداول، درونیابی با درجه بالا انجام می گیرد، گره ها باید طوری انتخاب شوند که در صورت امکان x نزدیک به $(x_0 + x_n)/2$ قرار گیرد.

مسئله تقریب. وقتی از رایانه استفاده می کنیم، معمولاً ترجیح می دهیم که یک تقریب تحلیلی تابع را ذخیره کنیم نه یک جدول مقادیری که با آن درونیابی می کنیم. تقریبیابی یک تابع داده شده $f(x)$ در یک بازه داده شده $[a, b]$ را با استفاده از چند جمله ایهای درونیابی در نظر می گیریم. به ویژه، حالتی را در نظر می گیریم که چند جمله ای $p_n(x)$ با درونیابی $f(x)$ در زیرتقسیمهای متساوی الفاصله $[a, b]$ تولید شده است.

برای هر $n \geq 1$ ، تعریف می کنیم $h = (b - a)/n$ ، $x_j = a + jh$ ، $j = 0, 1, \dots, n$.
گیریم $p_n(x)$ چند جمله ای درونیاب $f(x)$ در x_0, \dots, x_n است. سپس می پرسیم که آیا

$$\max_{a \leq x \leq b} |f(x) - p_n(x)| \quad (۸.۵.۳)$$

با $n \rightarrow \infty$ به صفر میل می کند یا نه؟ جواب این است که، لزوماً نه. برای توابع بسیاری، مثلاً e^x در $[0, 1]$ ، خطا در (۸.۵.۳) وقتی $n \rightarrow \infty$ به صفر همگراست (مسئله ۲۴ را ببینید). ولی توابع دیگری وجود دارند، که کاملاً خوشرفتارند و همگرا نیستند.

معروفترین مثال عدم همگرایی مثالی است که کارل رونگه^۱ آورده است. گیریم

$$f(x) = \frac{1}{1+x^2} \quad -5 \leq x \leq 5 \quad (۹.۵.۳)$$

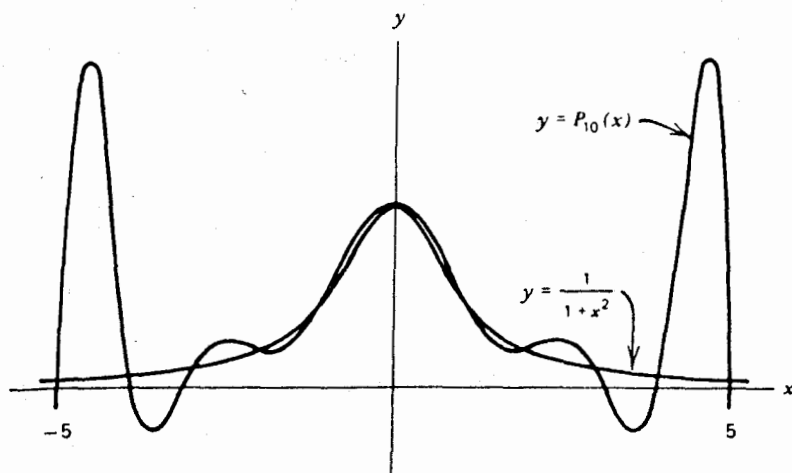
در کتاب آیزکسون و کِلر^۲ (۱۹۶۶، صص ۲۷۵، ۲۷۹) نشان داده شده است که برای هر $3.64 < |x| < 5$ ،

$$\sup_{n \geq k} |f(x) - p_n(x)| = \infty, \quad k \geq 0 \quad (۱۰.۵.۳)$$

بنابراین برای چنین x هرگاه $n \rightarrow \infty$ به $f(x)$ همگرا نمی شود. در ابتدا این امر مغایر با احساس ما به نظر می رسد، ولی اساس این امر رفتار چند جمله ای $y = \Psi_n(x)$ نزدیک به دو سر بازه $[a, b] = [x_0, x_n]$ است. این موضوع در نمودار $p_{10}(x)$ و $f(x)$ در شکل ۶.۳

1. Carl Runge

2. Isaacson and Keller



شکل ۶.۳ درونیابی $1/(1+x^2)$

به خوبی نمایان است. اگرچه درونیابی با نقاط متساوی الفاصله ممکن است دنباله‌ای همگرا از چندجمله‌یهای درونیابی ایجاد ننماید، مجموعه‌های مناسبتری از نقاط شبکه‌ای $\{x_j\}$ وجود دارند که برای تمام توابع پیوسته مشتق‌پذیر تقریبات خوبی به بار می‌آورند. این شبکه در بخش ۷.۴ از فصل ۴ دربارهٔ نظریهٔ تقریب، توضیح داده شده است.

۶.۳ درونیابی ارمیت

در بسیاری از کاربردها، مناسب است که چندجمله‌ای $p(x)$ چنان در نظر گرفته شود که تابع $f(x)$ را درونیابی نماید، و بعلاوه $p'(x)$ چندجمله‌ای $f'(x)$ ، مشتق تابع، را نیز درونیابی کند. در این کتاب، کاربرد ابتدایی این موضوع، به‌عنوان یک ابزار ریاضی در مبحث انتگرالگیری عددی گاوسی در فصل ۵، آمده است. این‌گونه درونیابی برای مطالعهٔ روشهای عددی در حل بعضی مسائل معادلات دیفرانسیل نیز مناسب‌اند.

موضوع را با در نظر گرفتن یک قضیهٔ وجود برای مسئلهٔ اساسی درونیابی شروع می‌کنیم

$$p(x_i) = y_i \quad p'(x_i) = y'_i \quad i = 1, \dots, n \quad (۱.۶.۳)$$

که در آن x_1, \dots, x_n گره‌های متمایز (حقیقی یا مختلط) و $y_1, \dots, y_n, y'_1, \dots, y'_n$ داده‌ها هستند (نمادگذاری از $n+1$ نقطهٔ $\{x_0, \dots, x_n\}$ به n نقطهٔ $\{x_1, \dots, x_n\}$ تغییر یافته است تا مانند کاربرد در فصل ۵ باشد). $2n$ شرط برای (۱.۶.۳) گذاشته شده‌است. پس دنبال یک چندجمله‌ای $p(x)$ حداکثر از درجهٔ $2n-1$ هستیم.

برای بحث درباره وجود و یکتایی $p(x)$ ، سومین برهان قضیه ۱.۳ را تعمیم می‌دهیم. برای هماهنگی با نمادگذاری بخش ۱.۳، می‌گیریم

$$\begin{aligned}\Psi_n(x) &= (x - x_1) \dots (x - x_n) \\ l_i(x) &= \frac{(x - x_1) \dots (x - x_{i-1})(x - x_{i+1}) \dots (x - x_n)}{(x_i - x_1) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_i - x_n)} \\ &= \frac{\Psi_n(x)}{(x - x_i)\Psi'_n(x_i)} \\ \tilde{h}_i(x) &= (x - x_i)[l_i(x)]^2 \\ h_i(x) &= [1 - 2l'_i(x_i)(x - x_i)][l_i(x)]^2\end{aligned}\quad (۲.۶.۳)$$

در این صورت به ازای $i, j = 1, \dots, n$

$$\begin{aligned}h'_i(x_j) &= \tilde{h}_i(x_j) = 0 \quad 1 \leq i, j \leq n \\ h_i(x_j) &= \tilde{h}'_i(x_j) = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}\end{aligned}\quad (۳.۶.۳)$$

چندجمله‌ای درونیاب برای (۱.۶.۳) با رابطه زیر داده می‌شود

$$H_n(x) = \sum_{i=1}^n y_i h_i(x) + \sum_{i=1}^n y'_i \tilde{h}_i(x) \quad (۴.۶.۳)$$

برای نشان دادن یکتایی $H_n(x)$ ، فرض می‌کنیم چندجمله‌ای دیگری با درجه نابزرگتر از $2n-1$ باشد که در (۱.۶.۳) صدق می‌کند. تابع $R = H - G$ را تعریف می‌کنیم. لذا با توجه به (۱.۶.۳)

$$R(x_i) = R'(x_i) = 0 \quad i = 1, 2, \dots, n$$

که R یک چندجمله‌ای از درجه نابزرگتر از $2n-1$ است با n ریشه دوگانه x_1, x_2, \dots, x_n . این مطلب فقط وقتی درست است که به ازای یک چندجمله‌ای $q(x)$ داشته باشیم

$$R(x) = q(x)(x - x_1)^2 \dots (x - x_n)^2$$

اگر $q(x) \not\equiv 0$ ، آنگاه درجه $R(n) \leq 2n$ که یک تناقض است. پس باید $R(x) \equiv 0$.

برای پیدا کردن یک شکل محاسبه پذیرتر از (۴.۶.۳) و یک عبارت خطا، ابتدا چندجمله‌ای درونیاب $f(x)$ را در نقاط z_{2n}, \dots, z_2, z_1 در نظر می‌گیریم که به شکل تفاضل منقسم نیوتن نوشته شده است:

$$p_{2n-1}(x) = f(z_1) + (x - z_1)f[z_1, z_2] + \dots + (x - z_1) \dots (x - z_{2n-1})f[z_1, \dots, z_{2n}] \quad (۵.۶.۳)$$

برای خطا

$$f(x) - p_{2n-1}(x) = (x - z_1) \dots (x - z_{2n})f[z_1, \dots, z_{2n}, x] \quad (۶.۶.۳)$$

در فرمول (۵.۶.۳) می‌توانیم گره‌ها را منطبق برهم فرض کنیم و باز فرمول برقرار خواهد بود. به‌ویژه فرض کنید

$$z_1, z_2 \rightarrow x_1, \quad z_3, z_4 \rightarrow x_2, \dots, z_{2n-1}, z_{2n} \rightarrow x_n$$

تا به دست آوریم

$$p_{2n-1}(x) = f(x_1) + (x - x_1)f[x_1, x_1] + (x - x_1)^2 f[x_1, x_1, x_2] + \dots + (x - x_1)^2 \dots (x - x_{n-1})^2 (x - x_n)f[x_1, x_1, \dots, x_n, x_n] \quad (۷.۶.۳)$$

این یک چندجمله‌ای است با درجهٔ نابزرگتر از $2n - 1$. برای خطای آن، از (۶.۶.۳) وقتی که $z_1, z_2 \rightarrow x_1, \dots, z_{2n-1}, z_{2n} \rightarrow x_n$ حد می‌گیریم. چون تفاضل منقسم با فرض آنکه f به اندازه کافی مشتق‌پذیر باشد پیوسته است

$$f(x) - p_{2n-1}(x) = (x - x_1)^2 \dots (x - x_n)^2 f[x_1, x_1, \dots, x_n, x_n, x] \quad (۸.۶.۳)$$

ادعا می‌کنیم - $p_{2n-1}(x) = H_n(x)$. برای اثبات، گیریم $f(x)$ تابعی $2n + 1$ بار پیوسته مشتق‌پذیر باشد، در این صورت ملاحظه می‌کنیم که:

$$f(x_i) - p_{2n-1}(x_i) = 0 \quad i = 1, 2, \dots, n$$

همچنین

$$f'(x) - p'_{2n-1}(x) = (x - x_1)^2 \dots (x - x_n)^2 \frac{d}{dx} f[x_1, x_1, \dots, x_n, x_n, x] + 2f[x_1, x_1, \dots, x_n, x_n, x] \sum_{i=1}^n \left[(x - x_i) \prod_{\substack{j=1 \\ j \neq i}}^n (x - x_j)^2 \right]$$

$$f'(x_i) - p'_{2n-1}(x_i) = 0 \quad i = 1, \dots, n$$

بنابراین درجه p_{2n-1} نابزرگتر از $2n-1$ است و p_{2n-1} و p'_{2n-1} در رابطه (۱۰.۶.۳) در داده‌های $y_i = f(x_i)$ و $y'_i = f'(x_i)$ صدق می‌کنند. از یکتایی چندجمله‌ای درونیاب ارمیت داریم $p_{2n-1} = H_n$. بنابراین (۷.۶.۳) یک فرمول تفاضل منقسم برای محاسبه $H_n(x)$ است و (۸.۶.۳) یک فرمول خطا به صورت

$$f(x) - H_n(x) = [\Psi_n(x)]^r f[x_1, x_1, \dots, x_n, x_n, x] \quad (9.6.3)$$

به دست می‌دهد:

با استفاده از رابطه (۱۳.۲.۳) برای تعمیم (۱۲.۲.۳)، به دست می‌آوریم:

$$f(x) - H_n(x) = [\Psi_n(x)]^r \frac{f^{(2n)}(\xi_n)}{(2n)!} \quad \xi_n \in \mathcal{H}\{x_1, \dots, x_n, x\} \quad (10.6.3)$$

مثال صورت درونیابی ارمیت که در بیشترین جاها مورد استفاده قرار می‌گیرد شاید چندجمله‌ای درجه ۳ ارمیت باشد، که مسأله زیر را حل می‌کند:

$$\begin{aligned} p(a) &= f(a) & p'(a) &= f'(a) \\ p(b) &= f(b) & p'(b) &= f'(b) \end{aligned} \quad (11.6.3)$$

فرمول (۴.۶.۳) چنین می‌شود:

$$\begin{aligned} H_3(x) &= \left[1 + 2 \frac{x-a}{b-a} \right] \left[\frac{b-x}{b-a} \right]^2 \cdot f(a) + \left[1 + 2 \frac{b-x}{b-a} \right] \left[\frac{x-a}{b-a} \right]^2 \cdot f(b) \\ &+ \frac{(x-a)(b-x)^2}{(b-a)^2} f'(a) - \frac{(x-a)^2(b-x)}{(b-a)^2} f'(b) \end{aligned} \quad (12.6.3)$$

فرمول تفاضل منقسم (۷.۶.۳) چنین می‌شود:

$$\begin{aligned} H_3(x) &= f(a) + (x-a)f'(a) + (x-a)^2 f[a, a, b] \\ &+ (x-a)^2(x-b)f[a, a, b, b] \end{aligned} \quad (13.6.3)$$

$$f[a, a, b] = \frac{f[a, b] - f'(a)}{b - a}$$

$$f[a, a, b, b] = \frac{f'(b) - 2f[a, b] + f'(a)}{(b - a)^2}$$

فرمول (۱۳.۶.۳) را می‌توان با الگوریتم ضرب تودرتو، که مشابه با *Interp* است که در بخش ۲.۳ آمده است، محاسبه کرد.

فرمول خطا برای (۱۲.۶.۳) یا (۱۳.۶.۳) چنین است:

$$f(x) - H_r(x) = (x - a)^r(x - b)^r f[a, a, b, b, x] \quad (14.6.3)$$

$$= \frac{(x - a)^r(x - b)^r}{r!} \cdot f^{(r)}(\xi_x) \quad \xi_x \in \mathcal{H}\{a, b, x\}$$

$$\text{Max}_{a \leq x \leq b} |f(x) - H_r(x)| \leq \frac{(b - a)^r}{r!} \cdot \text{Max}_{a \leq t \leq b} |f^{(r)}(t)| \quad (15.6.3)$$

کاربرد دیگری از چند جمله‌ای درجه سوم ارمیت در بخش بعد خواهد آمد و یک مثال عددی در جدول ۱۲.۳ آن بخش داده شده است.

مسئله درونیابی کلی ارمیت. مسئله ساده (۱.۶.۳) ارمیت را به مسئله زیر تعمیم می‌دهیم: یک چندجمله‌ای $p(x)$ پیدا کنید که در روابط زیر صدق کند.

$$p^{(i)}(x_1) = y_1^{(i)} \quad i = 0, 1, \dots, \alpha_1 - 1$$

$$\vdots$$

$$p^{(i)}(x_n) = y_n^{(i)} \quad i = 0, 1, \dots, \alpha_n - 1 \quad (16.6.3)$$

اعداد $y_j^{(i)}$ داده‌های معلوم‌اند و تعداد شرایط برای $p(x)$ در x_j برابر α_j است، $j = 1, \dots, n$.
تعریف می‌کنیم

$$N = \alpha_1 + \dots + \alpha_n$$

آنگاه در میان چندجمله‌بیهای با درجه نایزگتر از $N - 1$ چندجمله‌ای یکتایی وجود دارد که در (۱۶.۶.۳) صدق می‌کند. برهان به‌عنوان مسئله ۲۵، واگذار شده است. همه نتایج قبلی، مانند

(۴.۶.۳) و (۹.۶.۳) را می‌توان تعمیم داد. به‌عنوان یک حالت خاص جالب، فرض می‌کنیم $\alpha_1 = N$ و $n = 1$. این بدان معناست که $p(x)$ باید در

$$p^{(i)}(x_1) = f^{(i)}(x_1) \quad i = 0, 1, \dots, N-1$$

صدق کند. برای نمادگذاری مناسبتر، به جای $y_1^{(i)}, f^{(i)}(x_1)$ گذارده‌ایم. در این صورت با استفاده از (۱۵.۲.۳)، شکل تفاضل منقسم نیوتن برای چندجمله‌ای درونیاب ارمیت چنین می‌شود

$$\begin{aligned} p(x) &= f(x_1)(x - x_1)f[x_1, x_1] + (x - x_1)^2 f[x_1, x_1, x_1] + \dots \\ &\quad + (x - x_1)^{N-1} f[x_1, \dots, x_1] \\ &= f(x_1) + (x - x_1)f'(x_1) + \frac{(x - x_1)^2}{2!} f''(x_1) + \dots \\ &\quad + \frac{(x - x_1)^{N-1}}{(N-1)!} f^{(N-1)}(x_1) \end{aligned}$$

که چندجمله‌ای تیلر f حول x_1 نیز هست.

۷.۳ درونیابی چندجمله‌یی تکه‌یی

از اوایل دهه ۱۹۶۰، موضوع توابع چندجمله‌یی تکه‌یی، به‌ویژه توابع برازا، عمومیت روزافزون یافته است. این توابع چندجمله‌یی به انحاء گوناگون در نظریه تقریب، گرافیک رایانه‌یی، برازش داده‌ها، انتگرالگیری و مشتقگیری عددی و حل عددی معادلات انتگرالی، دیفرانسیل، و مشتقات جزئی به‌کار رفته‌اند. ما فقط از دیدگاه نظریه درونیابی به توابع چندجمله‌ای تکه‌ای نگاه می‌کنیم، ولی کاربردهای زیادی در چند زمینه دارند که درونیابی به‌طور جنبی در آن زمینه‌ها ظاهر می‌شود. برای یک تابع چندجمله‌ای تکه‌ای $p(x)$ ، یک شبکه از نقاط گرهی وابسته وجود دارد:

$$-\infty < x_0 < x_1 < \dots < x_n < \infty \quad (۱.۷.۳)$$

نقاط x_j اغلب گره یا نقطه شکستگی نامیده می‌شوند. تابع $p(x)$ یک چندجمله‌ای در هر یک از زیربازه‌های:

$$(-\infty, x_0], [x_0, x_1], \dots, [x_n, \infty) \quad (۲.۷.۳)$$

است، اگرچه زیربازه‌های $(-\infty, x_0]$ و $[x_n, \infty)$ اغلب در آنها نیستند. $p(x)$ را یک چندجمله‌ای تکه‌ای مرتبه r گویند اگر درجه $p(x)$ در هر یک از زیربازه‌های (۲.۷.۳) کوچکتر از r باشد. هیچ قید پیوستگی برای $p(x)$ و مشتقات آن وجود ندارد، اگرچه معمولاً $p(x)$ پیوسته است. در این بخش، خود را بیشتر به توابع چندجمله‌ای تکه‌ای درجه ۳ (مرتبه چهار) محدود می‌کنیم. این عمومیترین مورد در کاربردهاست، و معین بودن مرتبه، بیان مطلب را آسان می‌نماید.

یک راه رده‌بندی مسائل درونیابی چندجمله‌ای تکه‌ای تقسیم آنها به موضعی یا کلی است. برای یک مسأله موضعی، چندجمله‌ای $p(x)$ در هر زیربازه $[x_{i-1}, x_i]$ به وسیله داده‌های درونیابی، در نقاط گرهی داخل و مجاور $[x_{i-1}, x_i]$ ، کاملاً معین می‌شود. ولی برای یک مسأله کلی انتخاب $p(x)$ در هر $[x_{i-1}, x_i]$ به همه داده‌های درونیابی بستگی دارد. مطالعه مسائل کلی تا اندازه‌ای پیچیده‌تر است؛ بهترین مثالها، توابع براز هستند که بعداً تعریف می‌شوند.

مسائل درونیابی موضعی. فرض کنید که می‌خواهیم $f(x)$ را در بازه $[a, b]$ تقریب بزنیم. با انتخاب یک شبکه

$$a = x_0 < x_1 < \dots < x_n = b \quad (3.7.3)$$

که اغلب متساوی‌الفاصله است شروع می‌کنیم. اولین حالت تابع درونیابی چندجمله‌ای تکه‌ای ما بر پایه استفاده از درونیابی چندجمله‌ای معمولی در هر زیربازه $[x_{i-1}, x_i]$ استوار است. گیریم چهارگره درونیابی در هر زیربازه $[x_{i-1}, x_i]$ داده شده است،

$$x_{i-1} \leq z_{i,1} < z_{i,2} < z_{i,3} < z_{i,4} \leq x_i \quad i = 1, \dots, n \quad (4.7.3)$$

$p(x)$ را یک چندجمله‌ای از درجه نایزگتر از ۳ در $x_{i-1} < x < x_i$ تعریف می‌کنیم که $f(x)$ را در $z_{i,1}, \dots, z_{i,4}$ درونیابی کند. اگر

$$z_{i,1} = x_{i-1} \quad z_{i,4} = x_i \quad i = 1, 2, \dots, n \quad (5.7.3)$$

آنگاه $p(x)$ در $[a, b]$ پیوسته است. در گره‌های (۴.۷.۳)، این تابع درونیابی را، تابع چندجمله‌ای تکه‌ای لاگرانژ می‌نامیم و آن را با $L_n(x)$ نشان می‌دهیم.

با توجه به فرمول خطای (۸.۱.۳) در درونیابی چندجمله‌ای، فرمول خطا در $L_n(x)$ در رابطه زیر صدق می‌کند

$$f(x) - L_n(x) = \frac{(x - z_{i,1})(x - z_{i,2})(x - z_{i,3})(x - z_{i,4})}{4!} f^{(4)}(\xi_i) \quad (6.7.3)$$

به ازای $x_{i-1} \leq x \leq x_i$ با $i = 1, \dots, n$ ، $x_{i-1} < \xi_i < x_i$. حالت خاص تساوی فاصله‌ها را در نظر بگیرید. گیریم

$$\delta_i = \frac{(x_i - x_{i-1})}{3} \quad z_{i,j} = x_{i-1} + (j-1)\delta_i \quad j = 1, 2, 3, 4 \quad (7.7.3)$$

در این حالت (7.5.3) و (6.7.3) به ازای $i = 1, 2, \dots, n$ ، نتیجه می‌دهند

$$|f(x) - L_n(x)| \leq \frac{\delta_i^4}{24} \cdot \text{Max}_{x_{i-1} \leq t \leq x_i} |f^{(4)}(t)| \quad x_{i-1} < x < x_i \quad (8.7.3)$$

از این رابطه می‌توان دید، برای آنکه خطا در $a \leq x \leq b$ در یک سطح باقی بماند، فاصله‌گذاری δ_i باید براساس اندازه مشتق $f^{(4)}(t)$ در نقاط $[x_{i-1}, x_i]$ انتخاب شود. بنابراین اگر تابع $f(x)$ رفتار متغیری در بازه $[a, b]$ داشته باشد، تابع چندجمله‌ای تکه‌ای $L_n(x)$ را می‌توان با تنظیم مناسب شبکه (3.7.3)، به نحوی انتخاب کرد که همان رفتار را دنبال کند. درونیابی چندجمله‌ای معمولی در $[a, b]$ ، این انعطاف‌پذیری را ندارد، و این یک دلیل به کار بردن درونیابی چندجمله‌ای تکه‌ای است. برای مواردی که نقاط شبکه متساوی‌فاصله (3.7.3) را به کار می‌بریم، نتیجه (8.7.3) همگرایی را تضمین می‌نماید، حال آنکه در درونیابی چندجمله‌ای معمولی ممکن است، مانند مثال (9.5.3) رونگ همگرایی برقرار نباشد.

مثال برای $f(x) = e^x$ در $[0, 1]$ ، فرض کنید می‌خواهیم فاصله‌گذاریها برابر باشند و خطای ماکسیمم از 10^{-6} کوچکتر باشد. با استفاده از $\delta_i = \delta$ در (8.7.3)، می‌خواهیم داشته باشیم:

$$\frac{\delta^4}{24} e \leq 10^{-6}$$

$$\delta \leq 0.055, \quad n = \frac{1}{\delta} \geq 12.6$$

شاید به لحاظ ماهیت پایداری (8.7.3)، استفاده از $n = 6$ کافی باشد. در اینجا شش زیربازه درونیابی درجه سوم داریم.

از نظر نیاز به ذخیره‌سازی درونیابی چندجمله‌ای تکه‌ای لاگرانژ، با فرض (7.7.3)، باید چهار تکه اطلاعات در هر زیربازه $[x_{i-1}, x_i]$ نگهداری کنیم. این امر در کل ذخیره‌سازی به $4n$ جا در حافظه برای $L_n(x)$ و همچنین $n - 1$ جای اضافی برای نقاط شکستگی (3.7.3) نیاز دارد. انتخاب اینکه چگونه اطلاعات نگهداری شوند بستگی به چگونگی محاسبه $L_n(x)$ دارد. اگر

مشقات $L_n(x)$ خواسته شده باشد، مناسبترین شکل ذخیره‌سازی $L_n(x)$ ، شکل چندجمله‌ای تیلر در هر زیربازه $[x_{j-1}, x_j]$ است:

$$L_n(x) = a_j + b_j(x - x_{j-1}) + c_j(x - x_{j-1})^2 + d_j(x - x_{j-1})^3 \quad (9.7.3)$$

$$x_{j-1} \leq x \leq x_j$$

این محاسبه بهتر است به شکل تودرتو انجام گیرد. ضرایب $\{a_j, b_j, c_j, d_j\}$ از شکل‌های استانده چندجمله‌ای درونیابی درجه سوم به‌سادگی به‌دست می‌آیند.

یک شکل دوم تابع درونیابی چندجمله‌ای تکه‌ای که به‌طور گسترده‌ای از آن استفاده می‌شود، بر پایه چندجمله‌ای درجه ۳ ارمیت $(12.6.3) - (15.6.3)$ است. در هر زیربازه $[x_{i-1}, x_i]$ ، گیریم $Q_n(x)$ ، چندجمله‌ای درجه ۳ ارمیت باشد که $f(x)$ را در x_{i-1} و x_i درونیابی می‌کند. تابع $Q_n(x)$ یک چندجمله‌ای درجه سوم تکه‌ای بر شبکه $(3.7.3)$ است، و چون در شرایط درونیابی، هم $Q_n(x)$ و هم $Q'_n(x)$ در $[a, b]$ پیوسته هستند، پس معمولاً هموارتر از $L_n(x)$ است. برای خطای، $Q_n(x)$ در $[x_{i-1}, x_i]$ ، $(15.6.3)$ را به‌کار می‌بریم تا به‌دست آوریم

$$|f(x) - Q_n(x)| \leq \frac{h_i^4}{384} \cdot \text{Max}_{x_{i-1} \leq t \leq x_i} |f^{(4)}(t)| \quad x_{i-1} \leq x \leq x_i \quad (10.7.3)$$

با $1 \leq i \leq n$ ، وقتی این رابطه با رابطه $(8.7.3)$ برای خطا در $L_n(x)$ مقایسه شود، ممکن است چنین به‌نظر آید که درونیابی درجه ۳ تکه‌ای ارمیت بهتر است. ولی این موضوعی گمراه‌کننده است. برای روشن‌تر شدن موضوع، گیریم شبکه $(3.7.3)$ متساوی‌الفاصله با $x_i - x_{i-1} = h$ به ازای تمام i ها باشد. گیریم $L_n(x)$ بر پایه $(3.7.3)$ درست شده باشد و در $(8.7.3)$ داشته باشیم $\delta = h/3$. توجه داشته باشید که $Q_n(x)$ بر پایه $2n + 2$ تکه از داده‌ها درباره $f(x)$ ، یعنی $\{f(x_i), f'(x_i) \mid i = 0, 1, \dots, n\}$ درست شده است. و $L_n(x)$ بر پایه $3n + 1$ تکه از داده‌ها درباره $f(x)$ درست شده است. برای مقایسه خطای $L_n(x)$ با خطای $Q_{n_2}(x)$ با فرض $n_2 = 1.5n_1$ این دو را معادل می‌گیریم. در این صورت کرانه‌های خطا که از $(8.7.3)$ و $(10.7.3)$ به‌دست می‌آیند دقیقاً یکی می‌شوند.

چون تفاوتی در خطا وجود ندارد، شکل تابع چندجمله‌ای تکه‌ای بستگی به کاربردی دارد که روش برای آن بکار گرفته می‌شود. در کاربردهای انتگرال‌گیری عددی، تابع تکه‌ای لاگرائز مناسبترین است؛ همچنین این تابع در حل بعضی معادلات انتگرال‌تکین، توسط روش‌های انتگرال‌گیری حاصلضرب بخش ۶.۵ از فصل پنج به‌کار برده می‌شود. تابع تکه‌ای ارمیت برای حل بعضی مسائل

معادلات دیفرانسیل مفید است. برای مثال، این یک تابع رایج در روش عنصر متناهی، برای حل مسائل مقدار مرزی معادلات دیفرانسیل مرتبه دوم است؛ استرنگ و فیکس^۱ (۱۹۷۳، فصل ۱) را ببینید. مثالهای عددی برای مقایسه $L_n(x)$ و $Q_n(x)$ در جدولهای ۱۱.۳ و ۱۲.۳، پس از معرفی توابع برازا داده شده‌اند.

توابع برازا. مانند قبل، شبکه

$$a = x_0 < x_1 < \dots < x_n = b$$

را در نظر می‌گیریم. $s(x)$ را تابع برازا از مرتبه $m \geq 1$ می‌نامیم اگر در دو ویژگی زیر صدق کند.

(۱- $s(x)$ یک چندجمله‌ای از درجه کوچکتر از m در هر زیربازه $[x_{i-1}, x_i]$ باشد.

(۲- $s^{(r)}(x)$ در $[a, b]$ برای $0 \leq r \leq m - 2$ پیوسته باشد.

مشق تابع برآزی مرتبه m یک تابع برازا از مرتبه $m - 1$ است و به همین گونه است برای انتگرالهای آن. اگر پیوستگی در (۱- $s(x)$ تا $s^{(m-1)}$ ادامه یابد، آنگاه می‌توان ثابت کرد که $s(x)$ یک چندجمله‌ای از درجه نابزرگتر از $m - 1$ در $[a, b]$ است (مسئله ۳۳ را ببینید). توابع برازای درجه ۳ (از مرتبه $m = 4$) به دلایل گوناگون رایجترین توابع برازا هستند. اینها توابع همواری برای برازش داده‌ها هستند، و هنگامی که در درونیابی به‌کار برده می‌شوند، رفتار نوسانی را که مشخصه درونیابی چندجمله‌ای با درجه بالاست ندارند. یک دلیل دیگر برای شکل‌های خاص درونیابی برای درجه ۳ در (۲۷.۲.۳) و در مسئله ۳۸ داده شده است.

برای مسئله درونیابی، می‌خواهیم یک تابع برازای $s(x)$ از درجه ۳ پیدا کنیم که داشته باشیم.

$$s(x_i) = y_i \quad i = 0, 1, \dots, n \quad (11.7.3)$$

بررسی را شروع می‌کنیم با اینکه ببینیم پس از قیدهای (۱۱.۷.۳) چند درجه آزادی برای انتخاب $s(x)$ باقی می‌ماند. تکنیکی که به‌کار برده می‌شود مستقیماً به یک وسیله عملی برای محاسبه $s(x)$ نمی‌انجامد، ولی بینش بیشتری به‌دست می‌دهد. می‌نویسیم

$$s(x) = a_i + b_i x + c_i x^2 + d_i x^3 \quad x_{i-1} \leq x \leq x_i \quad i = 1, \dots, n \quad (12.7.3)$$

در اینجا $4n$ ضریب مجهول $\{a_i, b_i, c_i, d_i\}$ وجود دارند. قیدهای مربوط به $s(x)$ عبارت‌اند از (۱۱.۷.۳) و محدودیت‌های پیوستگی (۲- $s(x)$).

$$s^{(j)}(x_i + 0) = s^{(j)}(x_i - 0) \quad i = 1, \dots, n - 1, \quad j = 0, 1, 2 \quad (13.7.3)$$

این قیود بر روی هم تعداد

$$n + 1 + 3(n - 1) = 4n - 2$$

شرط ایجاد می‌نمایند. پس از $4n$ مجهول حداقل دو درجه آزادی در انتخاب ضرایب (۱۲.۷.۳) وجود دارد. باید انتظار داشته باشیم که با تحمیل شرایط اضافی به $s(x)$ بتوانیم یک تابع برآزای درونیاب یکتای $s(x)$ به دست آوریم.

اکنون روشی برای ساختن $s(x)$ ارائه می‌نماییم. ابتدا نماد زیر را معرفی می‌کنیم

$$M_i = s''(x_i) \quad i = 0, 1, \dots, n \quad (14.7.3)$$

چون $s(x)$ در $[x_i, x_{i+1}]$ از درجه ۳ است، $s''(x)$ خطی است و بنابراین

$$s''(x) = \frac{(x_{i+1} - x)M_i + (x - x_i)M_{i+1}}{h_i} \quad i = 0, 1, \dots, n - 1 \quad (15.7.3)$$

که در آن $h_i = x_{i+1} - x_i$. با این فرمول، $s''(x)$ در $[x_0, x_n]$ پیوسته است. دوبار انتگرال می‌گیریم تا به دست آوریم

$$s(x) = \frac{(x_{i+1} - x)^2 M_i + (x - x_i)^2 M_{i+1}}{6h_i} + C(x_{i+1} - x) + D(x - x_i)$$

که C و D دلخواه‌اند. شرط درونیاب (۱۱.۷.۳) ایجاب می‌کند:

$$C = \frac{y_i}{h_i} - \frac{h_i M_i}{6} \quad D = \frac{y_{i+1}}{h_i} - \frac{h_i M_{i+1}}{6}$$

$$s(x) = \frac{(x_{i+1} - x)^2 M_i + (x - x_i)^2 M_{i+1}}{6h_i} + \frac{(x_{i+1} - x)y_i + (x - x_i)y_{i+1}}{h_i}$$

$$- \frac{h_i}{6} [(x_{i+1} - x)M_i + (x - x_i)M_{i+1}]$$

$$x_i \leq x \leq x_{i+1}, \quad 0 \leq i \leq n - 1 \quad (16.7.3)$$

این فرمول پیوستگی $s(x)$ را در $[a, b]$ ، و نیز شرط درونیاب (۱۱.۷.۳) را نتیجه می‌دهد. برای تعیین مقادیر ثابت M_n, \dots, M_0 فرض می‌کنیم $s'(x)$ در x_1, \dots, x_{n-1} پیوسته است

$$\lim_{x \searrow x_i} s'(x) = \lim_{x \nearrow x_i} s'(x) \quad i = 1, \dots, n - 1 \quad (17.7.3)$$

در $[x_i, x_{i+1}]$ داریم

$$s'(x) = \frac{-(x_{i+1} - x)^2 M_i + (x - x_i)^2 M_{i+1}}{2h_i} + \frac{y_{i+1} - y_i}{h_i} - \frac{(M_{i+1} - M_i)h_i}{6} \quad (۱۸.۷.۳)$$

و در $[x_{i-1}, x_i]$ داریم

$$s'(x) = \frac{-(x_i - x)^2 M_{i-1} + (x - x_{i-1})^2 M_i}{2h_{i-1}} + \frac{y_i - y_{i-1}}{h_{i-1}} - \frac{(M_i - M_{i-1})h_{i-1}}{6}$$

با استفاده از (۱۷.۷.۳) و اندکی دستکاری، به ازای $i = 1, \dots, n-1$ خواهیم داشت:

$$\frac{h_{i-1}}{6} M_{i-1} + \frac{h_i + h_{i-1}}{3} M_i + \frac{h_i}{6} M_{i+1} = \frac{y_{i+1} - y_i}{h_i} - \frac{y_i - y_{i-1}}{h_{i-1}} \quad (۱۹.۷.۳)$$

این روابط $n-1$ معادله برای $n+1$ مجهول M_0, \dots, M_n به دست می دهند. معمولاً دو شرط انتهایی در x_0 و x_n تعیین می کنیم تا دو درجه آزادی موجود در (۱۹.۷.۳) را از میان برداریم. حالت ۱. شرایط مشتق در نقاط انتهایی. می خواهیم $s(x)$ در شرایط

$$s'(x_0) = y'_0, \quad s'(x_n) = y'_n \quad (۲۰.۷.۳)$$

که y'_0, y'_n ثابتهای داده شده اند، صدق کنند. با استفاده از این شرایط همراه با (۱۸.۷.۳)، برای $i = 0$ و $i = n-1$ معادلات اضافی زیر را به دست می آوریم

$$\frac{h_0}{3} M_0 + \frac{h_0}{6} M_1 = \frac{y_1 - y_0}{h_0} - y'_0$$

$$\frac{h_{n-1}}{6} M_{n-1} + \frac{h_{n-1}}{3} M_n = y'_n - \frac{y_n - y_{n-1}}{h_{n-1}}$$

از ترکیب با (۱۹.۷.۳) یک دستگاه معادلات خطی خواهیم داشت به صورت:

$$AM = D \quad (۲۱.۷.۳)$$

با

$$D^T = \left[\frac{y_1 - y_0}{h_0} - y'_0, \frac{y_2 - y_1}{h_1} - \frac{y_1 - y_0}{h_0}, \dots, \frac{y_n - y_{n-1}}{h_{n-1}} - \frac{y_{n-1} - y_{n-2}}{h_{n-2}}, y'_n - \frac{y_n - y_{n-1}}{h_{n-1}} \right]$$

$$M^T = [M_0, M_1, \dots, M_n]$$

$$A = \begin{bmatrix} \frac{h_0}{3} & \frac{h_0}{6} & 0 & 0 & \dots & 0 \\ \frac{h_0}{6} & \frac{h_0 + h_1}{3} & \frac{h_1}{6} & & & \\ 0 & \frac{h_1}{6} & \frac{h_1 + h_2}{3} & \frac{h_2}{6} & & \\ \vdots & & & \ddots & & \\ 0 & & & & \frac{h_{n-2}}{6} & \frac{h_{n-2} + h_{n-1}}{3} & \frac{h_{n-1}}{6} \\ & & & & \frac{h_{n-1}}{6} & \frac{h_{n-1}}{3} & 0 \end{bmatrix} \quad (22.7.3)$$

این ماتریس، متقارن، معین مثبت و غالب قطری است، و دستگاه خطی $AM = D$ به‌طور یکتا حل پذیر است. این دستگاه به‌سادگی و به‌سرعت با تقریباً $8n$ عمل حساب حل می‌شود (مطلب مربوط به دستگاه‌های سه قطری در بخش ۳.۸ را ببینید).

$s(x)$ ، تابع برازای درجه ۳ که به‌دست آوردیم، گاهی درونیاب برازای درجه ۳ کامل خوانده می‌شود، که آن را با $s_{c,n}(x)$ نشان می‌دهیم. تحلیل خطای آن نیاز به بحث دامنه‌داری دارد، ما فقط قضایای دور (۱۹۷۸، صص ۶۸-۶۹)، را نقل می‌کنیم.

قضیه ۴.۳ گیریم $f(x)$ در $a \leq x \leq b$ چهاربار پیوسته مشتق‌پذیر باشد و دنباله افزایشی زیر داده شده باشد

$$\tau_n : a = x_0^{(n)} < x_1^{(n)} < \dots < x_n^{(n)} = b$$

و تعریف می‌کنیم

$$|\tau_n| = \max_{1 \leq i \leq n} (x_i^{(n)} - x_{i-1}^{(n)})$$

گیریم $s_{c,n}(x)$ درونیاب برازای درجه ۳ کامل $f(x)$ در افراز τ_n باشد:

$$s_{c,n}(x_i^{(n)}) = f(x_i^{(n)}) \quad i = 0, 1, \dots, n$$

$$s'_{c,n}(a) = f'(a) \quad s'_{c,n}(b) = f'(b)$$

در این صورت به ازای مقادیر ثابت مناسب c_j ،

$$\max_{a \leq x \leq b} |f^{(j)}(x) - s_{c,n}^{(j)}(x)| \leq c_j |\tau_n|^{r-j} \cdot \max_{a \leq x \leq b} |f^{(r)}(x)| \quad (23.7.3)$$

برای $j = 0, 1, 2$ با فرض اضافی

$$\sup_n \left\{ \frac{|\tau_n|}{\min_{1 \leq i \leq n} (x_i^{(n)} - x_{i-1}^{(n)})} \right\} < \infty$$

نتیجه (۲۳.۷.۳) برای $j = 3$ نیز برقرار است. ثابتهای قابل قبول عبارت‌اند از:

$$c_0 = \frac{5}{384} \quad c_1 = \frac{1}{24} \quad c_2 = \frac{3}{8} \quad (24.7.3)$$

برهان برهانهای بسیاری از این قضایا همراه با قضایای دیگری مربوط به $s_{c,n}(x)$ را می‌توان در دور (۱۹۷۸)، صص ۶۹-۶۸ یافت. ■

اگر در رابطه (۲۳.۷.۳)، بگیریم $j = 0$ ، خواهیم دید که برای یک شبکه یکنواخت τ_n ، نرخ همگرایی متناسب با $1/n^2$ است. این را درونیابی درجه سوم تکه‌ای لاگرانژ و ارمیت هم داشتیم، ولی ضریب ثابت c_0 تقریباً سه برابر کوچکتر است. پس درونیاب برازی درجه ۳ کامل باید تا اندازه‌ای تقریب برتری باشد، همان‌گونه که نتایج جداول ۱۱.۳ تا ۱۴.۳ این موضوع را نشان می‌دهند.

دلیل دیگر برای استفاده از $s_c(x)$ ، ویژگی بهینگی زیر است. گیریم $g(x)$ تابع دلخواهی دوبار پیوسته مشتق‌پذیر در $[a, b]$ باشد، و به علاوه در شرایط درونیابیهای (۱۱.۷.۳) و (۲۰.۷.۳) صدق کند. بنابراین

$$\int_a^b |s''_c(x)|^2 dx \leq \int_a^b |g''(x)|^2 dx \quad (25.7.3)$$

که تساوی فقط در حالت $g(x) = s_c(x)$ برقرار است. بنابراین $s_c(x)$ ، در بین تمام توابعی که در شرایط درونیابی (۱۱.۷.۳) و (۲۰.۷.۳) صدق می‌کنند «کمترین نوسان» را دارد. برای اثبات این قضیه، گیریم $k(x) = s_c(x) - g(x)$ ، و می‌نویسیم

$$\begin{aligned} \int_a^b |g''(x)|^2 dx &= \int_a^b |s''_c(x) - k''(x)|^2 dx \\ &= \int_a^b |s''_c(x)|^2 dx - 2 \int_a^b s''_c(x)k''(x) dx \\ &\quad + \int_a^b |k''(x)|^2 dx \end{aligned}$$

با انتگرالگیری جزء به جزء، و استفاده از شرایط درونیابی و ویژگیهای $s_c(x)$ می‌توانیم نشان دهیم

$$\int_a^b s''_c(x)k''(x)dx = 0 \quad (26.7.3)$$

و بنابراین

$$\int_a^b |g''(x)|^2 dx = \int_a^b |s''_c(x)|^2 dx + \int_a^b |s''_c(x) - g''(x)|^2 dx \quad (27.7.3)$$

این رابطه (۲۵.۷.۳) را ثابت می‌کند. تساوی در (۲۵.۷.۳) فقط وقتی اتفاق می‌افتد که در $[a, b]$ ، $s''_c(x) - g''(x) \equiv 0$ یا هم‌ارز با آن، $s_c(x) - g(x)$ خطی باشد. در این صورت شرایط درونیاب ایجاب می‌نماید $s_c(x) - g(x) \equiv 0$. ما بحث بیشتر در این موضوع را به مسئله ۳۸ احاله می‌کنیم.

حالت ۲. شرط «نبود یک گره». وقتی مشتقات $f'(a)$ و $f'(b)$ در دست نباشند، به شرایط انتهایی دیگری برای $s(x)$ نیاز داریم تا دستگاه معادله‌های (۱۹.۷.۳) کامل شود. این امر با شرط پیوسته بودن $s^{(3)}(x)$ در نقاط x_1 و x_{n-1} انجام می‌شود. این شرط هم‌ارز با این است که خواسته شود $s(x)$ یک تابع برازای درجه ۳ با نقاط گره‌ی $\{x_0, x_1, x_2, \dots, x_{n-2}, x_n\}$ باشد، حال آنکه هنوز هم به درونیابی در نقاط گره‌ی $\{x_0, x_1, x_2, \dots, x_{n-1}, x_n\}$ نیاز داریم. این کار دستگاه (۱۹.۷.۳) را به $n - 3$ معادله بدل می‌کند و درونیابی در نقاط x_1, x_{n-1} دو معادله جدید به دست می‌دهند (پیدا کردن آنها را به مسئله ۳۴ احاله می‌کنیم). مجدداً یک دستگاه خطی سه قطری $AM = D$ به دست می‌آوریم، اگرچه ماتریس A بعضی ویژگیهای خوب (۲۲.۷.۳) را نخواهد داشت. تابع برازای به دست آمده در اینجا را با $s_{nk}(x)$ نشان می‌دهند که زیرنمایه آن به معنای شرط «نبود یک گره» است. یک تحلیل همگرایی مشابه آنچه در قضیه ۴.۳ آمد، می‌توان برای s_{nk} نوشت. برای این بحث دبور (۱۹۷۸، ص ۲۱۱)، (۱۹۸۵) را ببینید.

روشهای دیگری برای معرفی شرایط نقطه انتهایی، وقتی $f'(a)$ و $f'(b)$ مجهول باشند، وجود دارد. بحث درباره بعضی از اینها را می‌توان در دبور (۱۹۷۸، ص ۵۶) یافت. ولی، طرح ارائه شده در اینجا ساده‌ترین شکل در کاربرد است و به‌طور گسترده‌ای مورد استفاده قرار می‌گیرد. در حالات خاص بعضی شرایط نقاط انتهایی ساده‌تری وجود دارند که می‌توان از آنها استفاده نمود، یکی از اینها را در مسئله ۳۸ مطرح کرده‌ایم. ولی در حالت کلی، شرایط نقاط انتهایی به شکلی که قبلاً آمد، لازم است، تا نرخ همگرایی را که در قضیه ۴.۳ داده شد تأمین کند.

مثالهای عددی. گیریم $f(x) = \tan^{-1}x$ و $0 \leq x \leq 5$. در جدول ۱۱.۳ خطاهای

$$E_i = \max_{0 \leq x \leq 5} |f^{(i)}(x) - L_n^{(i)}(x)| \quad i = 0, 1, 2, 3 \quad (28.7.3)$$

جدول ۱۱.۳ درونیابی درجه سوم تکه‌ای لاگرانژ: $L_n(x)$

n	E_0	نسبت	E_1	نسبت	E_2	نسبت	E_3	نسبت
۲	$۱٫۲۰E - ۲$		$۱٫۲۲E - ۱$		$۷٫۸۱E - ۱$		۲٫۳۲	
		۳٫۳		۲٫۱		۱٫۵		۱٫۲
۴	$۳٫۶۲E - ۳$		$۵٫۸۳E - ۲$		$۵٫۲۴E - ۱$		۱٫۹۵	
		۱۱٫۴		۶٫۱		۳٫۲		۱٫۶
۸	$۳٫۱۸E - ۴$		$۹٫۷۵E - ۳$		$۱٫۶۴E - ۱$		۱٫۱۹	
		۱۶٫۹		۸٫۱		۳٫۹		۱٫۷
۱۶	$۱٫۸۸E - ۵$		$۱٫۱۱E - ۳$		$۴٫۲۱E - ۲$		۰٫۶۸۲	
		۱۴٫۵		۷٫۳		۳٫۷		۱٫۹
۳۲	$۱٫۳۰E - ۶$		$۱٫۶۱E - ۴$		$۱٫۱۴E - ۲$		۰٫۳۵۹	

جدول ۱۲.۳ درونیابی درجه سوم تکه‌ای ارمیت: $Q_n(x)$

n	E_0	نسبت	E_1	نسبت	E_2	نسبت	E_3	نسبت
۳	$۲٫۶۴E - ۲$		$۵٫۱۸E - ۲$		$۴٫۹۲E - ۱$		۲٫۰۶	
		۵٫۶		۳٫۰		۲٫۳		۱٫۵
۶	$۴٫۷۳E - ۳$		$۱٫۷۴E - ۲$		$۲٫۱۴E - ۱$		۱٫۳۳	
		۱۶٫۰		۸٫۰		۳٫۶		۱٫۵
۱۲	$۲٫۹۵E - ۴$		$۲٫۱۷E - ۳$		$۵٫۹۱E - ۲$		۰٫۸۹۱	
		۱۳٫۱		۶٫۷		۳٫۶		۱٫۹
۲۴	$۲٫۲۶E - ۵$		$۳٫۲۵E - ۴$		$۱٫۶۶E - ۲$		۰٫۴۷۵	
		۱۶٫۰		۸٫۰		۴٫۰		۲٫۰
۴۸	$۱٫۴۱E - ۶$		$۴٫۰۶E - ۵$		$۴٫۱۸E - ۳$		۰٫۲۴۱	

داده شده‌اند که در آنها $L_n(x)$ تابع درجه ۳ تکه‌ای لاگرانژ است، که $f(x)$ را در نقاط $x_j = a + jh$ همگرایی نوشته شده است، نرخ کاهش خطا را وقتی n دو برابر می‌شود، مشخص می‌نمایند. توجه کنید که نرخ همگرایی $L_n^{(2)}$ به $f^{(2)}$ ، با h^{4-2} ، $i = 0, 1, 2, 3$ ، متناسب است. این را می‌توان دقیقاً اثبات کرد و یک راهنمایی برای اثبات آن در مسأله ۳۲ داده شده است.

در جدول ۱۲.۳، خطاهای مشابه برای تابع درجه ۳ تکه‌ای ارمیت $Q_n(x)$ که $f(x)$ را درونیابی می‌کند داده شده است.

جدول ۱۳.۳ درونیابی برازی درجهٔ سوم کامل: $s_{c,n}(x)$

n	E_0	نسبت	E_1	نسبت	E_2	نسبت	E_3	نسبت
۶	$7.09E-3$		$2.45E-2$		$1.40E-1$		$1.06E^0$	
		۲۱٫۹		۱۰٫۷		۴٫۸		۲٫۶
۱۲	$3.24E-4$		$2.28E-3$		$2.90E-2$		$4.09E-1$	
		۱۰٫۶		۵٫۶		۲٫۹		۱٫۶
۲۴	$3.06E-5$		$4.09E-4$		$9.84E-3$		$2.53E-1$	
		۲۰٫۷		۹٫۷		۴٫۶		۲٫۱
۴۸	$1.48E-6$		$4.22E-5$		$2.13E-3$		$1.22E-1$	
		۱۶٫۴		۸٫۱		۴٫۰		۲٫۰
۹۶	$9.04E-8$		$5.19E-6$		$5.30E-4$		$6.09E-2$	

توجه کنید که در اینجا هم خطاها با

$$\max_{a \leq x \leq b} |f^{(i)}(x) - Q_n^{(i)}(x)| \leq ch^{r-i} \quad i = 0, 1, 2, 3$$

سازگار است، که باز می‌توان آن را برای بعضی از مقادیر $c > 0$ ثابت کرد.

همان‌گونه که قبلاً به دنبال (۱۰.۷.۳) بیان شد، توابع $L_n(x)$ و $Q_m(x)$ ، $m = 1.5n$ ، دارای دقتی قابل مقایسه، در تقریب $f(x)$ هستند، و جداول ۱۱.۳ و ۱۲.۳ این امر را تأیید می‌کنند. برعکس، $Q'_m(x)$ در مقایسه با $L'_n(x)$ تقریبی دقیق‌تر از $f'(x)$ است. در مسألهٔ ۳۲ توضیحی برای آن داده شده است.

در جدول ۱۳.۳ نتایج استفاده از درونیاب برازی درجهٔ سوم کامل را داده‌ایم. برای مقایسهٔ $L_n(x)$ با $Q_m(x)$ برای مقادیر قابل مقایسهٔ داده‌های $f(x)$ ، همان تعداد نقاط درونیابی متساوی‌فاصله را که در $L_n(x)$ به‌کار رفته است به‌کار می‌بریم.

مثال یک مثال آموزندهٔ دیگر انتخاب $f(x) = x^4$ ، $0 \leq x \leq 1$ ، است. در همهٔ فرمولهای درونیابی قبلی، $f^{(4)}(x) = 24$ به صورت یک ضریب در فرمولهای خطا وجود دارد. چون $f^{(4)}(x) = 24$ ثابت است، خطا در هر سه شکل درونیابی در رابطهٔ زیر صدق می‌کند

$$\max_{0 \leq x \leq 1} |x^4 - f_n(x)| = c_j h^{4-j} \quad j = 0, 1, 2, 3 \quad (29.7.3)$$

ثابتهای c_j با شکل درونیابی که به‌کار رفته است تغییر می‌یابند. در محاسبات فعلی، خطاها درست مانند (۲۹.۷.۳) رفتار می‌کنند، و بنابراین وسیلهٔ دیگری برای مقایسهٔ روشها به‌دست می‌دهند. نتایج جدول ۱۴.۳ را فقط برای دقیقترین حالات آورده‌ایم.

جدول ۱۴.۳ مقایسه سه شکل درونیابی درجه سوم تکه‌یی

روش	E_0	E_1	E_2	E_3
لاگرانژ $n = 32$	$1.10E - 8$	$6.788E - 6$	$2.932E - 3$	0.375
ارمیت $n = 48$	$1.18E - 8$	$1.74E - 6$	$8.68E - 4$	0.250
برازا $n = 96$	$7.36E - 10$	$2.12E - 7$	$1.09E - 4$	0.625

این مثالها نشان می‌دهند که درونیابی برازای درجه سوم در بعضی حالات به طور چشمگیری دقیقتر است. این مثالها همچنین نشان می‌دهند که تمام این روشها، از نظر دقت، احتمالاً مناسب‌اند، و همه آنها با یک نرخ همگرا هستند. بنابراین، تصمیم اینکه کدام یک از روشها به کار گرفته شود، به عوامل دیگری که معمولاً در زمینه کاربرد پیش می‌آیند بستگی دارد. ثابت شده است که توابع برازا برای مسائل برازش داده‌ها، و برازش منحنی بسیار مفیدند، و توابع لاگرانژ و ارمیت، به ترتیب، برای تقریب تحلیلی در حل معادلات انتگرالی و دیفرانسیل. تمام این شکلهای تقریب چند جمله‌ای تکه‌یی برای تمام این کاربردها سودمندند و شخص باید براساس نیازهایی که در حل مسأله مورد نظر دارد، شکل تقریب را انتخاب کند.

B- برازاها. یک راه نمایش توابع برازای درجه سوم در $(12.7.3)$ ، $(13.7.3)$ داده شده است که در آن یک چند جمله‌ای درجه سه در هر زیربازه داده شده است. این روش برای مسائل درونیابی، مانند آنچه در $(16.7.3)$ داده شده، کافی است، ولی برای بیشتر کاربردها، روشهای بهتری برای نمایش توابع برازای درجه سه وجود دارد. مانند گذشته، توجه ما به توابع برازای درجه سه با نقاط گرهی $\{x_0, x_1, \dots, x_n\}$ است.

تعریف می‌کنیم

$$x_+^r = \begin{cases} 0 & x < 0 \\ x^r & x \geq 0 \end{cases} \quad (30.7.3)$$

این یک تابع برازا از مرتبه $r+1$ است و فقط یک گره $x=0$ دارد. این تابع را می‌توان برای به دست آوردن یک نمایش دوم توابع برازا به کار برد. گیریم $s(x)$ یک تابع برازا از مرتبه m با گره‌های

$\{x_0, \dots, x_n\}$ باشد. آنگاه برای $x_0 \leq x \leq x_n$.

$$s(x) = p_{m-1}(x) + \sum_{j=1}^{n-1} \beta_j (x - x_j)_+^{m-1} \quad (۳۱.۷.۳)$$

که $p_{m-1}(x)$ یک چند جمله‌ای یکتا انتخابی از درجهٔ نابزرگتر از $m-1$ است و $\beta_1, \dots, \beta_{n-1}$ ضرایب تعیین شدهٔ یکتا هستند. اثبات این فرمول به مسألهٔ ۳۷ محوّل شده است. اگر این نمایش برای حل مسائل دیگر به‌کار برده شود، چندین اشکال وجود دارند. جذبتین اشکال این است که این نمایش به طرحهای عددی بدوضع می‌انجامد. بدین دلیل، نمایش دیگری از $s(x)$ ارائه می‌دهیم، که ویژگیهای عددی خیلی بهتری دارد. برای ساده‌کردن موضوع، ما فقط برازاهای درجهٔ سه را در نظر می‌گیریم. با افزایش گره‌های $\{x_0, \dots, x_n\}$ شروع می‌کنیم. نقاط گرهی دیگری به شکل دلخواه انتخاب می‌کنیم:

$$x_{-2} < x_{-1} < x_0 < x_1 < x_2 < x_3 < \dots < x_n < x_{n+1} < x_{n+2} < x_{n+3} \quad (۳۲.۷.۳)$$

برای $i = -3, -2, \dots, n-1$ یک تفاضل منقسم مرتبهٔ چهار

$$B_i(x) = (x_{i+4} - x_i) f_x[x_i, x_{i+1}, x_{i+2}, x_{i+3}, x_{i+4}] \quad (۳۳.۷.۳)$$

را برای تابع

$$f_x(t) = (t - x)_+^3 \quad (۳۴.۷.۳)$$

تعریف می‌کنیم. تابع $B_i(x)$ یک B -برازا نامیده می‌شود، که مخفف تابع برازای پایه است. به‌عنوان شکل دیگر (۳۳.۷.۳)، فرمول (۵.۲.۳) تفاضلات منقسم را به‌کار می‌بریم و به‌دست می‌آوریم

$$B_i(x) = (x_{i+4} - x_i) \sum_{j=i}^{i+3} \frac{(x_j - x)_+^3}{\Psi'_i(x_j)}$$

$$\Psi_i(x) = (x - x_i)(x - x_{i+1})(x - x_{i+2})(x - x_{i+3})(x - x_{i+4}) \quad (۳۵.۷.۳)$$

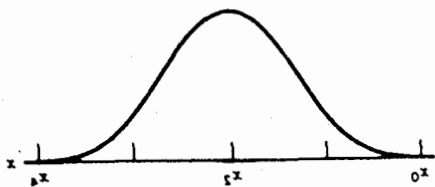
این نشان می‌دهد که $B_i(x)$ ، یک برازای درجهٔ سه با گره‌های x_i, \dots, x_{i+4} است. یک نمودار B -برازا در شکل ۷.۳ داده شده است.

بعضی از ویژگیهای مهم B -برازا را ذیل خلاصه می‌کنیم.

قضیهٔ ۵.۳ - برازای درجهٔ سه در روابط زیر صدق می‌کند:

(الف)

$$B_i(x) = 0 \quad \text{خارج از } x_i < x < x_{i+4} \quad (۳۶.۷.۳)$$

شکل ۳.۷.۳ - برزای $B_i(x)$

(ب)

$$0 \leq B_i(x) \leq 1 \quad \text{برای همهٔ مقادیر } x. \quad (۳۷.۷.۳)$$

(ج)

$$\sum_{i=-2}^{n-1} B_i(x) = 1 \quad x_0 \leq x \leq x_n \quad (۳۸.۷.۳)$$

(د)

$$\int_{x_i}^{x_{i+2}} B_i(x) dx = \frac{(x_{i+2} - x_i)}{4} \quad (۳۹.۷.۳)$$

(ه) اگر $s(x)$ یک تابع برزای درجهٔ سه باگره‌های $\{x_0, \dots, x_n\}$ باشد، آنگاه برای $x_0 \leq x \leq x_n$

$$s(x) = \sum_{i=-2}^{n-1} \alpha_i B_i(x) \quad (۴۰.۷.۳)$$

که انتخاب $\alpha_{n-1}, \dots, \alpha_{-2}$ یکتاست.

برهان (الف) برای $x \leq x_i$ تابع $f_x(t)$ یک چندجمله‌ای درجهٔ سه در بازهٔ $x_i \leq t \leq x_{i+2}$ است. بنابراین تفاضل منقسم مرتبهٔ چهار آن صفر است. به ازای $x \geq x_{i+2}$ تابع $f_x(t)$ در بازهٔ

$$B_i(x) = 0 \quad \text{متحد با صفر، و بنابراین}$$

(ب) دیور (۱۹۷۸، ص ۱۳۱) را ببینید.

(ج) با استفاده از رابطهٔ بازگشتی برای تفاضلات منقسم

$$B_i(x) = f_x[x_{i+1}, x_{i+2}, x_{i+3}, x_{i+4}] - f_x[x_i, x_{i+1}, x_{i+2}, x_{i+3}] \quad (۴۱.۷.۳)$$

سپس فرض می‌کنیم $x_k \leq x \leq x_{k+1}$. در این صورت $B_k(x), \dots, B_{k-2}(x), B_{k-3}(x)$ تنها B -برازاهایی هستند که ممکن است در x ناصفر باشند. با استفاده از (۴۱.۷.۳)

$$\begin{aligned} \sum_{i=-3}^{n-1} b_i(x) &= \sum_{i=k-3}^k B_i(x) \\ &= \sum_{i=k-3}^k (f_x[x_{i+1}, x_{i+2}, x_{i+3}, x_{i+4}] - f_x[x_i, x_{i+1}, x_{i+2}, x_{i+3}]) \\ &= f_x[x_{k+1}, x_{k+2}, x_{k+3}, x_{k+4}] - f_x[x_{k-3}, x_{k-2}, x_{k-1}, x_k] \\ &= 1 - 0 = 1 \end{aligned}$$

در مرحله آخر از این استفاده می‌کنیم که اولاً: $f_x(t)$ در $[x_{k+1}, x_{k+4}]$ از درجه ۳ است، پس تفاضل منقسم براساس (۱۸.۲.۳)، برابر ۱ است؛ و ثانیاً: $f_x(t) \equiv 0$ در $[x_{k-3}, x_k]$.
(د) دیور (۱۹۷۸، ص ۱۵۱) را ببینید.

(ه) مفهوم B - برازاها را شونبرگ^۱ پایه‌گذاری کرده، و فرمول (۴۰.۷.۳) از اوست. برای اثبات، دیور (۱۹۷۸، ص ۱۱۳) را ببینید.

به دلیل (۳۶.۷.۳)، مجموع (۴۰.۷.۳) حداکثر چهار جمله ناصفر دارد. برای $x_k \leq x < x_{k+1}$

$$s(x) = \sum_{i=k-3}^k \alpha_i B_i(x) \quad (۴۲.۷.۳)$$

به علاوه با استفاده از (۳۷.۷.۳) و (۳۸.۷.۳)،

$$\min\{\alpha_{k-3}, \alpha_{k-2}, \alpha_{k-1}, \alpha_k\} \leq s(x) \leq \max\{\alpha_{k-3}, \alpha_{k-2}, \alpha_{k-1}, \alpha_k\}$$

که نشان می‌دهد مقدار $s(x)$ با ضرایب B - برازاهای نزدیک به x ، کراندار شده است. از این جهت، (۴۰.۷.۳) یک نمایش موضعی $s(x)$ ، برای هر $x \in [x_0, x_n]$ است.

یک برداشت کلیتر از B - برازا و همچنین بعضی ویژگیهایی که اینجا حذف شده، در (فصلهای ۹-۱۱، دیور، ۱۹۷۸) داده شده است. در آنجا برنامه‌هایی برای محاسبه B - برازاها نیز داده شده است.

یک تعمیم مهم برازاها وقتی به وجود می‌آید که گرهما بتوانند یکی شوند. به ویژه، گیریم بعضی از گرهما در (۳۳.۷.۳) یکی شده‌اند. در این حالت، مادام که $x_i < x_{i+4}$ تابع $B_i(x)$ یک چندجمله‌ای درجه سه تکه‌ای است. اگر دوگره یکی شوند، تعداد مشتقات پیوسته در نقطه گرهی چندگانه، از دو به یک کاهش می‌یابد. اگر سه نقطه گرهی یکی شوند، معنی آن این خواهد بود که $B_i(x)$ فقط پیوسته خواهد بود. با این کار، (۴۰.۷.۳) نمایشی می‌شود برای همه

چندجمله‌بیهیای درجه سه تکه‌ای. در این صورت، تمام توابع چندجمله‌ای تکه‌ای توابع برازا هستند و بالعکس. این امر در دپور (۱۹۷۸) کاملاً تحقیق شده است.

۸.۳ درونیابی مثلثاتی

توابع دوره‌ای یک رده بسیار مهم از توابع، هستند. تابع $f(t)$ را دوره‌ای با دوره تناوب τ خوانند اگر

$$f(t + \tau) = f(t) \quad -\infty < t < \infty$$

و این رابطه برای هیچ مقدار مثبت کوچکتر از τ درست نباشد. معروفترین توابع دوره‌ای، توابع مثلثاتی هستند. توابع دوره‌ای درکاربردها بسیار فراوانند. این امر توجه ما را به درونیابی مناسبی برای داده‌هایی که از این توابع به دست آمده‌اند، توجیه می‌کند. علاوه بر آن، ما از این موضوع برای معرفی تبدیل سریع فوریه، که برای حل بسیاری از مسائل که شامل داده‌های توابع دوره‌ای هستند، استفاده می‌کنیم. با انتخاب یک مقیاس مناسب برای متغیر مستقل همیشه می‌توان دوره نوسان τ را 2π گرفت:

$$f(t + 2\pi) = f(t) \quad -\infty < t < \infty \quad (۱.۸.۳)$$

این تابعهای $f(t)$ را با استفاده از چندجمله‌بیهیای مثلثاتی تقریب می‌زنیم

$$p_n(t) \approx a_0 + \sum_{j=1}^n a_j \cos(jt) + b_j \sin(jt) \quad (۲.۸.۳)$$

اگر $a_n \neq 0$ یا $b_n \neq 0$ ، آنگاه این تابع $p_n(t)$ را چندجمله‌ای مثلثاتی از درجه n خوانند. می‌توان با استفاده از فرمولهای جمع مثلثاتی، نشان داد که یک صورت هم‌ارز آن به شکل زیر است:

$$p_n(t) = \alpha_0 + \sum_{j=1}^n \alpha_j [\cos(t)]^j + \beta_j [\sin(t)]^j \quad (۳.۸.۳)$$

که تا اندازه‌ای، استفاده از کلمه چندجمله‌ای را برای چنین تابعی توجیه می‌کند. چندجمله‌ای $p_n(t)$ دارای دوره تناوب 2π یا کسر صحیحی از آن است.

در مطالعه مسائل درونیابی با جوابی به صورت $p_n(t)$ ، باید $2n + 1$ شرط درونیابی برقرار کنیم، زیرا $p_n(t)$ شامل $2n + 1$ ضریب α_j و β_j است. چون تابع $f(t)$ و چندجمله‌ای $p_n(t)$ دوره‌ای هستند گره‌های درونیابی را، در بازه $0 \leq t < 2\pi$ (یا هم‌ارز آن)، $-\pi \leq t < \pi$ یا $0 < t \leq 2\pi$ می‌گیریم. بنابراین گره‌های درونیابی را به شکل زیر فرض می‌کنیم،

$$0 \leq t_0 < t_1 < \dots < t_{2n} < 2\pi \quad (۴.۸.۳)$$

و می‌خواهیم که $p_n(t)$ به گونه‌ای انتخاب شود که در روابط زیر صدق نماید

$$p_n(t_i) = f(t_i) \quad i = 0, 1, \dots, 2n \quad (5.8.3)$$

بعداً نشان می‌دهیم که این مسأله دارای جوابی یکتاست.

این مسأله درونیایی یک جواب صریح دارد که با فرمول لاگرانژ (۶.۱.۳) در درونیایی چندجمله‌ای، قابل مقایسه است؛ این موضوع در مسأله ۴۱ نشان داده شده است. به جای آنکه از این راه آغاز کنیم، ابتدا (۴.۸.۳) و (۵.۸.۳) را به مسأله هم‌ارز آن که شامل چندجمله‌یها و توابعی از یک متغیر مختلط است، برمی‌گردانیم. این صورت جدید، بهترین قالب ریاضی برای درونیایی چندجمله‌یهای مثلثاتی است.
با استفاده از فرمول اویلر

$$e^{i\theta} = \cos \theta + i \sin \theta \quad i = \sqrt{-1} \quad (6.8.3)$$

به دست می‌آوریم

$$\cos \theta = \frac{e^{i\theta} + e^{-i\theta}}{2} \quad \sin \theta = \frac{e^{i\theta} - e^{-i\theta}}{2i} \quad (7.8.3)$$

با استفاده از این روابط در (۲.۸.۳)، خواهیم داشت

$$p_n(t) = \sum_{j=-n}^n c_j e^{ijt} \quad (8.8.3)$$

ضرایب با روابط زیر به هم بستگی دارند

$$c_0 = a_0 \quad c_{-j} = \frac{1}{i^j} (a_j - ib_j) \quad c_j = \frac{1}{i^j} (a_j + ib_j) \quad 1 \leq j \leq n$$

اگر $\{c_j\}$ ها داده شده باشند، ضرایب $\{a_j, b_j\}$ به سادگی از حل معادلات اخیر به دست می‌آیند. اگر بگیریم $z = e^{it}$ ، می‌توانیم (۸.۸.۳) را به صورت یک تابع مختلط بنویسیم

$$P_n(z) = \sum_{j=-n}^n c_j z^j \quad (9.8.3)$$

تابع $z^n P_n(z)$ یک چندجمله‌ای از درجه نایزگتر از $2n$ است.

برای بیان مسأله درونیایی (۴.۸.۳) و (۵.۸.۳) به صورتی دیگر، گیریم $z_j = e^{it_j} = z$ ، $0, \dots, 2n$ با قید (۴.۸.۳)، اعداد z_j نقاط متمایزی روی دایره واحد $|z| = 1$ در صفحه

مختلط هستند. مسأله درونیابی چنین است:

$$P_n(z_j) = f(t_j) \quad j = 0, \dots, 2n \quad (10.8.3)$$

برای دیدن یکتایی جواب این مسأله، توجه نمایید که با تعریف $Q(z) = z^{2n}P_n(z)$ این معادله هم‌ارز است با

$$Q(z_j) = z_j^{2n}f(t_j) \quad j = 0, \dots, 2n$$

که یک مسأله درونیابی چندجمله‌ای با $2n+1$ نقطه گرهی متمایز z_0, \dots, z_{2n} است. قضیه ۱.۳ نشان می‌دهد که یک جواب یکتا وجود دارد. همچنین فرمول (۶.۱.۳) لاگرانژ به $Q(z)$ و از آنجا به $P(z)$ تعمیم می‌یابد.

دلایل چندی، هم نظری و هم عملی، برای برگرداندن درونیابی مثلثاتی به شکل متغیر مختلط وجود دارند. مهمترین آنها از نظر ما این است که درونیابی و تقریب با چندجمله‌بیهای مثلثاتی کاملاً به موضوع توابع مشتق‌پذیر متغیر مختلط مربوط است، و قسمت بیشتر نظریه از این چشم‌انداز بهتر فهمیده می‌شود. ما این نظریه را بسط نمی‌دهیم، یک بحث کامل آن را می‌توان در هنریچی (۱۹۸۶، فصل ۱۳) و زیگموند (۱۹۵۹، فصل ۱۰) یافت.

درونیابی متساوی‌الفاصله. جالبترین شکل درونیابی در کاربرد، آن است که نقاط شبکه t_j متساوی الفاصله انتخاب شوند. به عبارت دقیقتر، تعریف می‌کنیم

$$t_j = j \frac{2\pi}{2n+1} \quad j = 0, \pm 1, \pm 2, \dots \quad (11.8.3)$$

نقاط t_0, \dots, t_{2n} در (۴.۸.۳) صدق می‌کنند، و نقاط $z_j = e^{it_j}$ ، $j = 0, \dots, 2n$ ، نقاط متساوی‌الفاصله روی دایره واحد $|z| = 1$ هستند. همچنین توجه داریم که نقاط z_j وقتی j به اندازه $2n+1$ افزایش یابد، تکرار می‌شوند.

ما اکنون صورت دیگری از شکل لاگرانژی $p_n(t)$ را، هنگامی که گرههای $\{t_j\}$ در (۱۱.۸.۳) صدق می‌کنند، ارائه می‌دهیم. با لم زیر شروع می‌کنیم.

لم ۴ به ازای همه اعداد صحیح k ،

$$\sum_{j=0}^{2n} e^{ikt_j} = \begin{cases} 2n+1 & e^{it_k} = 1 \\ 0 & e^{it_k} \neq 1 \end{cases} \quad (12.8.3)$$

شرط $e^{it_k} = 1$ هم‌ارز است با k مضرب صحیحی از $2n+1$ است.

برهان گیریم $z = e^{itk}$ ، با استفاده از (۱۱.۸.۳)، $e^{ikt_j} = e^{ijtk}$ ، و مجموع در (۱۲.۸.۳) چنین می‌شود

$$S = \sum_{j=0}^{2n} z^j$$

اگر $z = 1$ ، آنگاه $S = 2n + 1$. اگر $z \neq 1$ ، آنگاه از فرمول سری هندسی (۸.۱.۱) داریم

$$S = \frac{z^{2n+1} - 1}{z - 1}$$

با استفاده از (۱۱.۸.۳)، $e^{2\pi ki} = 1$ ؛ بنابراین $S = 0$.

شرط درونیایی (۱۰.۸.۳) را می‌توان به شکل زیر نوشت

$$\sum_{k=-n}^n c_k e^{ikt_j} = f(t_j) \quad j = 0, 1, \dots, 2n \quad (13.8.3)$$

برای پیدا کردن ضرایب c_k ، از لم ۴ استفاده می‌کنیم. معادله z ام را در e^{-ilt_j} ضرب، سپس نسبت به z جمع می‌کنیم و l را مقید می‌کنیم که در $-n \leq l \leq n$ صدق کند. این عمل نتیجه می‌دهد

$$\sum_{j=0}^{2n} \sum_{k=-n}^n c_k e^{i(k-l)t_j} = \sum_{j=0}^{2n} e^{-ilt_j} f(t_j) \quad (14.8.3)$$

ترتیب جمع را وارونه و سپس از لم ۴ استفاده می‌کنیم تا به دست آوریم

$$\sum_{j=0}^{2n} e^{i(k-l)t_j} = \begin{cases} 0 & k \neq l \\ 2n + 1 & k = l \end{cases}$$

با گذاردن این مقدار در (۱۴.۸.۳)، خواهیم داشت

$$c_l = \frac{1}{2n + 1} \sum_{j=0}^{2n} e^{-ilt_j} f(t_j) \quad l = -n, \dots, n \quad (15.8.3)$$

ضرایب $\{c_{-n}, \dots, c_n\}$ تبدیل متناهی فوریه داده‌های $\{f(t_0), \dots, f(t_{2n})\}$ خوانده می‌شوند. این ضرایب یک فرمول صریح برای چندجمله‌ای درونیاب مثلثاتی $p_n(t)$ (۸.۸.۳) به دست می‌دهند. فرمول (۱۵.۸.۳) به ضرایب فوریه $f(t)$ مربوط است:

$$\gamma_l = \frac{1}{2\pi} \int_0^{2\pi} e^{-ilt} f(t) dt \quad -\infty < l < \infty \quad (16.8.3)$$

اگر از قاعده انتگرالگیری عددی ذوزنقه‌ای [بخش ۱.۵ را ببینید] برای انتگرالهای فوق استفاده شود، با استفاده از $1 + 2n$ تقسیم جزئی $[2\pi, 0]$ ، (۱۵.۸.۳) نتیجه می‌شود، مشروط بر آنکه $f(t)$ در $[0, 2\pi]$ متناوب باشد. ما اکنون در همگرایی $p_n(t)$ به $f(t)$ بحث می‌کنیم.

قضیه ۶.۳ بگیریم $f(t)$ یک تابع پیوسته، دوره‌ی و 2π ضربی صحیح از دوره تناوب آن باشد. تعریف می‌کنیم

$$\rho_n(f) = \inf_{\deg(q) \leq n} [\text{Max}_{0 \leq t \leq 2\pi} |f(t) - q(t)|] \quad (17.8.3)$$

که در آن $q(t)$ یک چندجمله‌ی مثلثاتی است. آنگاه تابع درونیاب $p_n(t)$ از (۸.۸.۳) و (۱۵.۸.۳) در رابطه زیر صدق می‌کند

$$\text{Max}_{0 \leq t \leq 2\pi} |f(t) - p_n(t)| \leq c[\ln(n+2)]\rho_n(f) \quad n \geq 0 \quad (18.8.3)$$

ثابت c مستقل از f و n است.

برهان چون اثبات نسبتاً پیچیده است به زیگموند (۱۹۵۹، ص ۱۹ فصل ۱۰) مراجعه کنید. ■

کمیت $\rho_n(f)$ خطای مینیماکس نامیده می‌شود (فصل ۴ را ببینید)، و از راههای گوناگون می‌توان آن را برآورد کرد. مهمترین قید $\rho_n(f)$ شاید همان قیدی است که جکسن^۱ نهاده است. فرض می‌کنیم $f(t)$ در بازه $[0, 2\pi]$ ، k بار پیوسته مشتقپذیر باشد، $k \geq 0$ ، و علاوه فرض می‌کنیم، به ازای مقداری از $1 \geq \alpha \geq 0$ ، $f^{(k)}(t)$ در شرط

$$|f^{(k)}(t_1) - f^{(k)}(t_2)| \leq c_f |t_1 - t_2|^\alpha \quad 0 \leq t_1, t_2 \leq 2\pi$$

صدق می‌کند. (این را شرط هولدر^۲ گویند). پس

$$\rho_n(f) \leq \frac{c_k(f)}{n^{k+\alpha}} \quad n \geq 1 \quad (19.8.3)$$

که $c_k(f)$ مستقل از n است. برای اثبات میناردوس^۳ (۱۹۶۷، ص ۵۵) را ببینید. یک شق دیگر فرمول خطای (۱۸.۸.۳) را هنریچی (۱۹۸۶، فرج ۱۳۶۶) داده است که از ضرایب سری فوریه (۱۶.۸.۳) برای $f(t)$ استفاده کرده است.

جدول ۱۵.۳ خطا در درونیابی چندجمله‌ای مثلثاتی

n	E_n	n	E_n
۱	$۵,۳۹E - ۱$	۶	$۶,۴۶E - ۶$
۲	$۹,۳۱E - ۲$	۷	$۴,۰۱E - ۷$
۳	$۱,۱۰E - ۲$	۸	$۲,۲۲E - ۸$
۴	$۱,۱۱E - ۳$	۹	$۱,۱۰E - ۹$
۵	$۹,۱۱E - ۵$	۱۰	$۵,۰۰E - ۱۱$

مثال تقریب زدن $f(t) = e^{\sin(t)}$ را با استفاده از تابع درونیاب $p_n(t)$ در نظر می‌گیریم. خطای ماکسیمم

$$E_n = \text{Max}_{0 \leq t \leq 2\pi} |f(t) - p_n(t)|$$

به ازای مقادیر مختلف n در جدول ۱۵.۳ داده شده است. همگرایی سریع است.

تبدیل سریع فوریه. تقریب $f(t)$ با $p_n(t)$ در مثال قبل، برای مقادیر کوچک n بسیار دقیق بود. اما، محاسبه تبدیل متناهی فوریه (۱۵.۸.۳) در کاربردهای دیگر، اغلب نیاز به مقادیر بزرگ n دارد. اکنون روشی را معرفی می‌کنیم که در کاهش هزینه محاسبه $\{c_t\}$ ، وقتی n بزرگ باشد، بسیار سودمند است. به جای استفاده از فرمول (۱۵.۸.۳)، فرمول هم‌ارز آن، یعنی

$$d_k = \frac{1}{m} \sum_{j=0}^{m-1} w_m^{jk} \cdot f_j \quad w_m = e^{-2\pi i/m} \quad k = 0, 1, \dots, m-1 \quad (20.8.3)$$

را با داده‌های مفروض $\{f_0, \dots, f_{m-1}\}$ در نظر می‌گیریم. این فرمول را تبدیل متناهی فوریه از مرتبه m می‌نامند. در فرمول (۱۵.۸.۳)، می‌گیریم $m = 2n + 1$. k مجاز است هر عدد صحیحی باشد، با توجه به

$$d_{k+m} = d_k \quad -\infty < k < \infty \quad (21.8.3)$$

بنابراین کافی است d_0, \dots, d_{m-1} یا هر m ضریب پیایی دیگر d_k را محاسبه کنیم. برای مقایسه فرمول (۲۰.۸.۳) با شکل دیگری که در زیر خواهد آمد، هزینه محاسبه d_0, \dots, d_{m-1} را، با استفاده از (۲۰.۸.۳) حساب می‌کنیم. برای محاسبه d_k ، گیریم $z_k = w_m^k$ پس

$$d_k = \frac{1}{m} \sum_{j=0}^{m-1} f_j z_k^j \quad (22.8.3)$$

با استفاده از ضرب تودرتو، در اینجا به $m - 1$ ضرب و $m - 1$ جمع نیاز داریم. تقسیم بر m را ندیده می‌گیریم زیرا اغلب عوامل دیگری به کار برده می‌شوند. محاسبهٔ z_k فقط به یک ضرب نیاز دارد زیرا $z_k = w_m z_{k-1}$ ، $k \geq 2$. هزینهٔ کلی محاسبهٔ d_{m-1}, \dots, d_0 برابر m^2 ضرب و $m(m - 1)$ جمع است.

برای معرفی منظور اصلی از تبدیل سریع فوریه، گیریم $m = pq$ که p و q اعدادی صحیح و مثبت بزرگتر از یک هستند. تعریف (۲۰.۸.۳) را به شکل هم‌ارز زیر می‌نویسیم

$$d_k = \frac{1}{p} \sum_{l=0}^{p-1} \frac{1}{q} \sum_{g=0}^{q-1} w_m^{k(l+pg)} f_{l+pg}$$

از $w_m^p = \exp(-2\pi i/q) = w_q$ استفاده می‌کنیم. پس

$$d_k = \frac{1}{p} \sum_{l=0}^{p-1} w_m^{kl} \left[\frac{1}{q} \sum_{g=0}^{q-1} w_q^{kg} f_{l+pg} \right] \quad k = 0, 1, \dots, m - 1$$

می‌نویسیم

$$e_k^{(l)} = \frac{1}{q} \sum_{g=0}^{q-1} w_q^{kg} f_{l+pg} \quad 0 \leq l \leq p - 1 \quad (23.8.3)$$

$$d_k = \frac{1}{p} \sum_{l=0}^{p-1} w_m^{kl} e_k^{(l)} \quad 0 \leq k \leq m - 1 \quad (24.8.3)$$

وقتی $\{e_k^{(l)}\}$ معلوم شد، هر مقدار d_k ، با استفاده از قاعدهٔ ضرب تودرتوی (۲۲.۸.۳)، به p ضرب نیاز خواهد داشت. اگر فرض کنیم همهٔ $e_k^{(l)}$ ‌ها قبلاً حساب شده‌اند، به دست آوردن (۲۴.۸.۳) به mp ضرب نیاز خواهد داشت. تعدادی در همین حدود هم عمل جمع وجود خواهد داشت. اکنون نظر خود را به محاسبهٔ $e_k^{(l)}$ معطوف می‌داریم. زیرنمایهٔ k از 0 تا $m - 1$ تغییر می‌کند، ولی به محاسبهٔ همهٔ آنها نیاز نیست. توجه کنید که

$$e_{k+q}^{(l)} = \frac{1}{q} \sum_{g=0}^{q-1} w_q^{(k+q)g} f_{l+pg} = e_k^{(l)}$$

زیرا $w_q^q = 1$. بنابراین فقط لازم است $e_k^{(l)}$ برای $k = 0, 1, \dots, q - 1$ محاسبه شود و سپس این مقادیر خود تکرار می‌شوند. برای هر l ، یک تبدیل متناهی فوریه به ازای داده‌های $\{f_l, f_{l+p}, \dots, f_{l+p(q-1)}\}$ ، $0 \leq l \leq p - 1$ ، است. بنابراین محاسبهٔ $\{e_k^{(l)}\}$ به محاسبهٔ p تبدیل متناهی فوریه از مرتبهٔ q (یعنی برای داده‌های به طول q) تبدیل می‌شود.

تبدیل سریع فوریه استفاده مکرر از این نحوه عمل است برای تبدیل محاسبه به تبدیلهای متناهی فوریه با مراتب پایین و پایینتر. برای توضیح بیشتر، فرض کنید r عدد صحیح و $m = 2^r$. در مرحله اول، گیریم $p = 2$ و $q = 2^{r-1}$. پس محاسبه (۲۴.۸.۳) به $2m$ ضرب به اضافه هزینه محاسبه دو تبدیل متناهی فوریه از مرتبه $q = 2^{r-1}$ نیاز دارد. برای هر یک از اینها، فرایند را به طور بازگشتی تکرار می‌کنیم. در این فرآیند، r ردیف وجود دارد که سرانجام به محاسبه تبدیل متناهی فوریه از مرتبه یک می‌انجامد. تعداد کل ضربها با رابطه زیر داده می‌شود

$$2m + 2 \left[2 \left(\frac{m}{2} \right) \right] + 4 \left[2 \left(\frac{m}{4} \right) \right] + \dots + 2^r \left[2 \left(\frac{m}{2^r} \right) \right]$$

که برابر است با

$$2rm = 2m \cdot \log_2 m$$

بنابراین تعداد عملیات متناسب است با $m \cdot \log_2 m$ ، در مقابل m^2 عمل در الگوریتم ضرب تودرتوی (۲۲.۸.۳). وقتی m بزرگ، مثلاً 2^{10} یا بیشتر باشد، در زمان صرفه‌جویی هنگفتی می‌شود. برای حالت خاص $m = 2^r$ ، یک حساب دقیقتر نشان می‌دهد که در واقع فقط به $m \log_2 m$ ضرب نیاز است، و یک روند دیگری وجود دارد که فقط به نصف این تعداد نیاز دارد. برای سایر مقادیر m همان‌گونه که قبلاً ذکر شد، تعمیمهایی از روند قبلی وجود دارند که باز هم نیاز به عملیاتی متناسب با $m \log_2 m$ دارند، ولی حالت $m = 2^r$ به بزرگترین صرفه‌جویی می‌انجامد. برای بحث بیشتر در این موضوع، خواننده را به هنریچی (۱۹۸۶، فصل ۱۳) ارجاع می‌دهیم. در آنجا همچنین بحثی دارد در مورد پایداری الگوریتم وقتی خطای گردکردن به حساب آمده باشد. استفاده از تبدیل سریع فوریه انقلابی در بسیاری از موضوعات به وجود آورده است و محاسباتی را که قبلاً عملی نبوده‌اند ممکن ساخته است.

بحث در آثار خواندنی

همان‌گونه که در مقدمه گفته شد، نظریه درونیابی پایه‌ای برای توسعه روشها، در انتگرالگیری عددی و مشتقگیری، نظریه تقریب، و حل عددی معادلات دیفرانسیل است. درباره هر یک از این موضوعها در فصلهای آتی بحث شده است و نوشتارهای مربوط به آنها در جای خود آورده شده است. مطالب بیشتری از نظریه درونیابی در دبور (۱۹۷۸)، دیویس (۱۹۶۳)، هنریچی (۱۹۸۲)، فصل ۵ و (۷) و هیلدبرانت (۱۹۵۶) داده شده است. برای گزارشهای تاریخی بسیاری از موضوعات این فصل به گولدشتاین (۱۹۷۷) مراجعه نمایید.

عرضه کامپیوترهای رقمی، انقلابی در آنالیز عددی، از جمله نظریه درونیابی، ایجاد نمود. قبل

از استفاده از کامپیوترهای رقمی، محاسبات دستی الزامی بود. بدین معنی که از روشهایی عددی استفاده می‌شد که نیاز کمتری به محاسبه داشت. این روشها، غالباً پیچیده‌تر از روشهایی بودند که امروزه در کامپیوتر استفاده می‌شوند، و از ویژگیهای ریاضی منحصر به هر مسأله، بهره‌برداری خاص می‌نمایند. در این روشها از جداول استفاده بسیاری می‌کردند تا از تکرار محاسباتی که توسط دیگران انجام شده، جلوگیری شود؛ فرمولهای درونیابی بر پایه تفاضلات متناهی به‌طور وسیعی به‌کار گرفته می‌شدند. یک مبحث گسترده‌ای به نام حساب تفاضلات متناهی ابداع شده بود که در حل مسائل در زمینه‌های گوناگون آنالیز عددی و ریاضیات کاربردی به‌کار می‌رفت. برای معرفی کلی این روش در آنالیز عددی هیلد برانت (۱۹۵۹) و مراجعی را که در آن داده شده ملاحظه کنید.

استفاده از کامپیوترهای رقمی نیاز به زمینه‌های دیگر را برای نظریه درونیابی تغییر داده و به‌طور وسیعی نیاز به فرمولهای درونیابی بر پایه تفاضلات متناهی را تقلیل داده است. ولی هنوز یک جای مهم هم برای محاسبات دستی و هم برای استفاده از جداول ریاضی وجود دارد، بویژه برای توابع پیچیده‌تر در فیزیک ریاضی. هر کس که با آنالیز عددی کار می‌کند باید یک کتاب مقدماتی از این جداول مثل جدول معروف CRC را داشته باشد. دفتر ملی جدولهای استاندارد آبراموویچ و استگان (۱۹۶۴) مرجع بسیار عالی برای توابع غیرمقدماتی است. دسترسی به ماشین حسابهای دستی پیشرفته و ریزکامپیوترها، محاسبات دستی (یا شخصی) جدیدی به‌وجود آورده است.

نظریه تقریب چندجمله‌ای تکه‌ای از دهه ۱۹۶۰ به بعد بسیار متداول شده، و در بسیاری از زمینه‌ها از آن استفاده می‌شود. برای مثال، استرنگ و فیکس (۱۹۷۳، فصل ۱) را برای کاربرد در حل مسائل مقادیر مرزی در معادلات دیفرانسیل معمولی ببینید، و پاولیدیس^۱ (۱۹۸۲، فصل ۱۰.۱۲) را برای کاربرد در گرافیک رایانه‌ای ملاحظه کنید. بیشترین استفاده از توابع چندجمله‌ای تکه‌ای در حول توابع برازا دور می‌زند. آغاز نظریه توابع برازا را عمدتاً به شونبرک^۲ در مقاله‌های ۱۹۴۶ او نسبت می‌دهند و کارهای او برای کمک به توسعه این موضوع برجسته بوده است (مثلاً شونبرک (۱۹۷۳) را ببینید). اکنون نوشته‌های مفصلی درباره توابع برازا وجود دارد که شامل کارهای افراد و گروههاست. برای یک مطالعه کلی، آلبرگ و همکاران^۳ (۱۹۶۷) و دوبر (۱۹۷۸) و شومیکر^۴ (۱۹۸۱) را ببینید. بعضی از معروفترین نرم‌افزارهای رایانه‌ای، برای استفاده از توابع برازا، بر پایه برنامه‌های دوبر (۱۹۷۸) تهیه شده‌اند. نسخه‌هایی از آنها در کتابخانه‌های آنالیز عددی IMSL و NAG وجود دارند.

تبدیلات متناهی فوریه، درونیابی مثلثاتی و مباحث وابسته به آنها، نسبتاً قدیمی‌اند؛ برای مثال گولدشتاین (۱۹۷۷، ص ۲۳۸) را برای کارهای گاوس در زمینه درونیابی مثلثاتی، ببینید. از آنجا که سری فوریه و تبدیلات فوریه ابزارهای مهمی در ریاضیات کاربردی هستند، جای تعجب نیست

1. Pavlidis

2. Schoenberg

3. Ahlberg et al

4. Schumaker

که توجه بسیار زیادی به تقریبهای خاص آنها وجود دارد. به دنبال مقاله مشهور کولی^۱ و تیوکی^۲ (۱۹۶۵) درباره تبدیل سریع فوریه، کاربرد تبدیلات سریع فوریه و موضوعات وابسته به آنها فوق العاده افزایش یافته است. به عنوان مثال از این راه روشهایی بسیار سریع برای حل معادلات دیفرانسیل با مشتقات جزئی لاپلاس روی ناحیه‌های مستطیلی پیدا شده‌اند، که ما در فصل ۸ بیشتر درباره آنها بحث خواهیم کرد. برای یک توضیح کلاسیک از درونیابی مقدماتی زیگموند (۱۹۵۹، فصل ۱۰)، و برای یک بازنگری جدید در تمام زمینه آنالیز منتهای فوریه هنریچی (۱۹۷۶، فصل ۱۳) را ببینید. نظریه درونیابی چندجمله‌ای چند متغیره، زمینه‌ای است که به سرعت توسعه می‌یابد، به لحاظ کمبود جا، از این کتاب حذف شده است. در روش عنصر منتهای برای حل معادلات دیفرانسیل با مشتقات جزئی استفاده وسیعی از نظریه درونیابی چند متغیره می‌شود و معرفیهای خوبی از این نظریه در کتابهایی در زمینه روش عنصر منتهای یافت می‌شوند. برای مثال جین^۳ (۱۹۸۴)، لاپیدوس^۴ و پیندر^۵ (۱۹۸۲)، میچل^۶ و ویت^۷ (۱۹۷۷) و استرنگ^۸ و فیکس^۹ (۱۹۷۳) را ببینید. اخیراً کار در گرافیک رایانه‌یی موجب پیشرفتهای جدیدی شده است [بارن هیل^{۱۰} (۱۹۷۷) و پاولیدیس (۱۹۸۲، فصل ۱۳) را ببینید].

مراجع

- Abramowitz, M., and I. Stegun, eds. (1964). *Handbook of Mathematical Functions*. National Bureau of Standards, Washington, D.C. (Available now from Dover, New York.)
- Ahlberg, J., E. Nilson, and J. Walsh (1967). *The Theory of Splines and Their Applications*. Academic Press, New York.
- Barnhill, R. (1977). Representation and approximation of surfaces. In *Mathematical Software III*, J. Rice, ed., pp. 69-120. Academic Press, New York.
- De Boor, C. (1978). *A Practical Guide to Splines*. Springer-Verlag, New York.
- De Boor, C. (1985). Convergence of cubic spline interpolation with the not-a-knot condition. Math. Res. Ctr. Tech. Rep. 1876, Madison, Wis.
- CRC Standard Mathematical Tables*. Chem. Rubber Publ. Co., Cleveland. (Published yearly).
- Cooley, J., and J. Tukey (1965). An algorithm for the machine calculation of

1. Cooley	2. Tukey	3. Jain	4. Lapidus
5. Pinder	6. Mitchel	7. Wait	8. Strang
9. Fix	10. Barnhill		

- complex Fourier series. *Math. Comput.*, 19, 297–301.
- Davis, P. (1963). *Interpolation and Approximation*. Ginn (Blaisdell), Boston.
- Goldstine, H. (1977). *A History of Numerical Analysis from the 16th through the 19th Century*. Springer-Verlag, New York.
- Henrici, P. (1982). *Essentials of Numerical Analysis*. Wiley, New York.
- Henrici, P. (1986). *Applied and Computational Complex Analysis*, Vol. 3. Wiley, New York.
- Hildebrand, F. (1956). *Introduction to Numerical Analysis*. McGraw-Hill, New York.
- Isaacson, E., and H. Keller (1966). *Analysis of Numerical Methods*. Wiley, New York.
- Jain, M. (1984). *Numerical Solution of Differential Equations*, 2nd ed., Wiley, New Delhi.
- Lapidus, L., and G. Pinder (1982). *Numerical Solution of Partial Differential Equations in Science and Engineering*. Wiley, New York.
- Lorentz, G., K. Jetter, and S. Riemenschneider (1983). *Birkhoff Interpolation*. In *Encyclopedia of Mathematics and Its Applications*, Vol. 19. Addison-Wesley, Reading, Mass.
- Meinardus, G. (1967). *Approximation of Functions: Theory and Numerical Methods* (transl. L. Schumaker). Springer-Verlag, New York.
- Mitchell, A., and R. Wait (1977). *The Finite Element Method in Partial Differential Equations*. Wiley, London.
- Pavlidis, T. (1982). *Algorithms for Graphics and Image Processing*. Computer Science Press, Rockville, Md.
- Powell, M. (1981). *Approximation Theory and Methods*. Cambridge Univ. Press, Cambridge, England.
- Schoenberg, I. (1946). Contributions to the approximation of equidistant data by analytic functions. *Quart. Appl. Math.*, 4(Part A), 45–99; (Part B), 112–141.
- Schoenberg, I. (1973). *Cardinal Spline Interpolation*. Society for Industrial and Applied Mathematics, Philadelphia.
- Schumaker, L. (1981). *Spline Functions: Basic Theory*. Wiley, New York.
- Stoer, J., and R. Bulirsch (1980). *Introduction to Numerical Analysis*. Springer-Verlag, New York.
- Strang, G., and G. Fix (1973). *An Analysis of the Finite Element Method*. Prentice-Hall, Englewood Cliffs, N.J.
- Zygmund, A. (1959). *Trigonometric Series*, Vols. 1 and 2. Cambridge Univ. Press, Cambridge.

مسائل

۱. ماتریس واندرموند X را که در (۳.۱.۳) داده شده است یادآوری کرده تعریف می‌کنیم.

$$V_n(x) = \det \begin{bmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n-1} & x_{n-1}^2 & \dots & x_{n-1}^n \\ 1 & x & x^2 & \dots & x^n \end{bmatrix}$$

(الف) نشان دهید $V_n(x)$ یک چندجمله‌ای درجه n است، و ریشه‌های آن x_0, \dots, x_{n-1} هستند. فرمول زیر را به دست آورید

$$V_n(x) = (x - x_0) \dots (x - x_{n-1}) V_{n-1}(x_{n-1})$$

راهنمایی: آخرین سطر $V_n(x)$ را توسط درمیان‌های جزئی بسط داده نشان دهید که $V_n(x)$ یک چندجمله‌ای درجه n است و ضریب جمله x^n را پیدا کنید.
(ب) نشان دهید

$$\det(X) \equiv V_n(x_n) = \prod_{0 \leq j < i \leq n} (x_i - x_j)$$

۲. برای توابع پایه $l_{j,n}(x)$ که در (۵.۱.۳) داده شده است، ثابت کنید برای هر $n \geq t$

$$\sum_{j=0}^n l_{j,n}(x) = 1 \quad \text{برای تمام مقادیر } x$$

۳. توابع لاگرانژ $l_0(x), \dots, l_n(x)$ را که در (۵.۱.۳) تعریف شده و سپس به حالت کمی متفاوت (۴.۲.۳) با استفاده از رابطه

$$\Psi_n(x) = (x - x_0) \dots (x - x_n)$$

بازنویس شده است به یاد می‌آوریم. بگیریم $w_j = [\Psi'_n(x_j)]^{-1}$. نشان دهید که چندجمله‌ای $p_n(x)$ ، درونیاب $f(x)$ را می‌توان به شکل زیر نوشت

$$p_n(x) = \frac{\sum_{j=0}^n [w_j f(x_j)] / (x - x_j)}{\sum_{j=0}^n w_j / (x - x_j)}$$

به شرطی که x یک نقطه گرهی نباشد. این عبارت را نمایش گرانیگاهی $p_n(x)$ گویند، که به صورت یک مجموع وزن دار از مقادیر $\{f(x_0), \dots, f(x_n)\}$ داده شده است. برای بحث درباره استفاده از این نمایش، هنریچی (۱۹۸۲، ص ۲۳۷) را ببینید.

۴. درونیابی خطی در جدول مقادیر e^x ، $0 \leq x \leq 2$ با $h = 0.1$ را در نظر می‌گیریم. بگیریم که مقادیر جدول با دقت ۵ رقم بامعنی مانند جدول CRC داده شده باشند. کران خطای درونیابی خطی را، شامل قسمتی که به خطای گردکردن در درایه‌های جدول مربوط می‌شود، به دست آورید.

۵. درونیابی خطی در یک جدول $\cos x$ را که x بر حسب درجه، $0^\circ \leq x \leq 90^\circ$ ، با طول گام (درجه) $h = 1' = \frac{1}{60}$ داده شده است در نظر می‌گیریم. با فرض اینکه درایه‌های جدول با خطای گردکردن تا ۵ رقم بامعنی داده شده باشند، کران خطای کلی درونیابی را پیدا کنید.

۶. فرض کنید که می‌خواهید یک جدول مقادیر $\sin x$ ، $0 \leq x \leq \pi/2$ ، با طول گام h بسازید. فرض کنید که در این جدول درونیابی خطی باید به‌کار گرفته شود و خطای کلی، شامل اثرات خطای گردکردن درایه‌های جدول، باید حداکثر برابر 10^{-6} باشد. h چقدر باید باشد (آن را در یک اندازه مناسب برای کاربرد واقعی انتخاب کنید) و درایه‌های جدول تا چند رقم معنی‌دار باید داده شوند؟

۷. مسأله ۶ را با e^x و $0 \leq x \leq 1$ حل کنید.

۸. مطالب مربوط به اثر خطای گردکردن در درایه‌های جدول بخش ۱.۳ را به درونیابی درجه دوم تعمیم دهید. بگیریم $\varepsilon_i = f(x_i) - \tilde{f}_i$ ، $i = 0, 1, 2$ ، و $\{\varepsilon_0, \varepsilon_1, \varepsilon_2\}$ و $\text{Max}\{|\varepsilon_0|, |\varepsilon_1|, |\varepsilon_2|\}$ نایزرگتر از ε باشد. نشان دهید که نتیجه این خطاهای گردکردن در درونیابی درجه دوم و فرض $x_0 \leq x \leq x_2$ و $x_2 - x_1 = x_1 - x_0 = h$ با 1.25ε کراندار می‌شود.

۹. مسأله ۶ را حل کنید، ولی درونیابی درجه دوم و نتیجه مسأله ۸ را به‌کار برید.

۱۰. فرض کنید می‌خواهید جدولی از مقادیر $f(x) = \log_{10} x$ ، $1 \leq x \leq 10$ ، بسازید و فرض کنید درونیابی درجه دوم باید به‌کار گرفته شود و خطای کلی درونیابی، شامل خطای گردکردن در درایه‌های جدول، باید کوچکتر از 10^{-6} باشد. یک گام فاصله‌گذاری h مناسب انتخاب کنید و تعداد ارقامی را که درایه‌ها باید داشته باشند معلوم کنید. آیا بهتر نیست که فاصله h با تغییر x در $[1, 10]$ تغییر کند؟ اگر جواب مثبت است، یک افزایش مناسب $[1, 10]$ متناظر با مقادیر h را پیشنهاد نمائید. نتیجه مسأله ۸ را برای اثر خطای گردکردن به‌کار برید.

۱۱. بگیریم x_0, \dots, x_n نقاط حقیقی متمایز باشند، و مسأله درونیابی زیر را در نظر می‌گیریم. یک تابع

$$P_n(x) = \sum_{j=0}^n c_j e^{jx}$$

انتخاب می‌کنیم به طوری که

$$P_n(x_i) = y_i \quad i = 0, 1, \dots, n$$

که $\{y_i\}$ ها داده‌های معلوم‌اند. نشان دهید که یک انتخاب یکتا برای مقادیر c_0, \dots, c_n وجود دارد. راهنمایی: مسأله را می‌توان به مسأله معمولی درونیابی چندجمله‌ای تبدیل کرد.

۱۲. پیدا کردن یک تابع گویای $p(x) = (a + bx)/(1 + cx)$ را در نظر بگیرید که در رابطه

$$P(x_i) = y_i \quad i = 1, 2, 3$$

صدق می‌کند و x_1, x_2, x_3 نقاط متمایز هستند. آیا چنین تابع $p(x)$ وجود دارد؟ یا شرایط اضافی دیگری لازم است تا وجود و یکتایی $p(x)$ را تضمین نماید؟ برای یک نظریه کلی درونیابی گویا اشتور و بولیرش^۱ (۱۹۸۰، ص ۵۸) را ببینید.

۱۳. (الف) شکل بازگشتی فرمول درونیابی را که در (۸.۲.۳) داده شده است ثابت کنید.

(ب) فرمول بازگشتی (۷.۲.۳) را برای تفاضلات منقسم ثابت کنید.

۱۴. رابطه‌های (۱۸.۲.۳) مربوط به تفاضل منقسم متغیر یک چندجمله‌ای را ثابت کنید.

۱۵. ثابت کنید که $\text{Vol}(\tau_n) = 1/n!$ که τ_n سادگی تعریف شده در R^n است که در (۱۴.۲.۳) از قضیه ۳.۳ در بخش ۲.۳، تعریف شده است.

راهنمایی: رابطه (۱۳.۲.۳) و قضایای دیگر در تفاضلات منقسم را همراه با انتخاب خاص $f(x)$ به کار برید.

۱۶. گیریم $p_2(x)$ چندجمله‌ای درجه دومی باشد که $f(x)$ را در نقاط متساوی‌فاصله x_0, x_1, x_2 درونیابی می‌کند. فرمولهای خطاهای $f'(x_i) - p_2'(x_i)$ و $x_1 = x_0 + h$ و $x_2 = x_0 + 2h$ را پیدا کنید. فرض کنید $f(x)$ سه بار پیوسته مشتقپذیر باشد، کرانه‌های قابل محاسبه برای این خطاها را به دست آورید.

راهنمایی: فرمول خطای (۱۱.۲.۳) را به کار برید.

۱۷. زیر روال رایانه‌ای را برای پیاده کردن الگوریتم‌های Divdif و Interp که در بخش ۲.۳ داده شده بنویسید و سپس یک برنامه اجراکننده اصلی تهیه کنید که از آن در ساختن جدول درونیابی استفاده شود. جدولی از آبراموویتس و اشتگون (۱۹۶۴) انتخاب و با در نظر گرفتن چندین درجه پایایی n برای چندجمله‌ای درونیابی، برنامه را آزمایش کنید.

۱۸. یک مسأله درونیابی معکوس را با استفاده از جدول $J_0(x)$ که در بخش ۲.۳ داده شده

انجام دهید. مقداری برای x پیدا کنید که $J_0(x) = 0$ ، یعنی یک برآورد دقیق برای ریشه پیدا کنید. دقت خود را برآورد کرده و آنرا با جواب واقعی $x = ۲٫۴۰۴۸۲۵۵۵۷۷$ مقایسه کنید.

۱۹. مشابه لم ۱ در بخش ۳.۳ را برای تفاضلات پسر و به دست آورید. از این لم و صورت منقسم چندجمله‌ای درونیاب (۹.۲.۳) نیوتن برای پیدا کردن فرمول درونیابی تفاضل پسر (۱۱.۳.۳) استفاده کنید.

۲۰. جدول زیر برای مقادیر

$$j_0(x) = \sqrt{\frac{\pi}{2x}} \cdot J_{1/2}(x)$$

را که از آبراموویتس و اشتگون (۱۹۶۴، فصل ۱۰) استخراج شده در نظر می‌گیریم.

x	$j_0(x)$	x	$j_0(x)$
۰٫۰	۱٫۰۰۰۰۰۰	۰٫۷	۰٫۹۲۰۳۱
۰٫۱	۰٫۹۹۸۳۳	۰٫۸	۰٫۸۹۶۷۰
۰٫۲	۰٫۹۹۳۳۵	۰٫۹	۰٫۸۷۰۳۶
۰٫۳	۰٫۹۸۵۰۷	۱٫۰	۰٫۸۴۱۴۷
۰٫۴	۰٫۹۷۳۵۵	۱٫۱	۰٫۸۱۰۱۹
۰٫۵	۰٫۹۵۸۸۵	۱٫۲	۰٫۷۷۶۷۰
۰٫۶	۰٫۹۴۱۰۷	۱٫۳	۰٫۷۴۱۲۰

بر پایه خطاهای گردکردن در درایه‌های جدول، بیشترین درجه درونیابی چندجمله‌ای که در این جدول استفاده شده چه باید باشد؟

راهنمایی: از جدول تفاضلات پیشرو برای پیدا کردن نتایج خطای گردکردن استفاده کنید.

۲۱. داده‌های زیر از یک چندجمله‌ای با درجه نایب‌تر از ۵ استخراج شده است. درجه چندجمله‌ای چه بوده است؟

x	-۲	-۱	۰	۱	۲	۳
$P(x)$	-۵	۱	۱	۱	۷	۲۵

۲۲. داده‌های زیر نوفه‌ای دارند که نسبت به خطای گردکردن بزرگ است. این نوفه را پیدا کنید و داده‌ها را به طریق مناسب تغییر دهید. فقط مقادیر تابع داده شده‌اند، زیرا نقاط گرهی برای محاسبه

جدول تفاضلات پیشرو لازم نیستند.

۳۰۴۳۱۹	۴۱۹۳۲۷	۵۴۵۸۱۱	۶۸۳۱۰۰
۳۲۶۳۱۳	۴۴۳۶۵۵	۵۷۲۴۳۳	۷۱۱۷۰۹
۳۴۸۸۱۲	۴۶۸۵۲۹	۵۹۹۴۷۵	۷۴۰۷۵۶
۳۷۱۸۰۶	۴۹۳۸۵۲	۶۲۶۹۰۹	۷۷۰۱۸۸
۳۹۵۲۸۵	۵۱۹۶۱۵	۶۵۴۷۹۰	۸۰۰۰۰۰

۲۳. برای $f(x) = 1/(1+x^2)$ ، $-5 \leq x \leq 5$ ، با استفاده از $n+1$ نقطه گرهی متساوی‌الفاصله در $[-5, 5]$ ، $P_n(x)$ را بسازید. $P_n(x)$ را در تعداد زیادی نقطه محاسبه و نمودار آن یا نمودار خطا را در $[-5, 5]$ ، مانند شکل ۶.۳، رسم کنید.

۲۴. تابع e^x و تقریب آن توسط یک چندجمله‌ای درونیاب روی $[0, b]$ را در نظر می‌گیریم. برای $n \geq 1$ بگیریم n $x_j = jh$ ، $h = b/n$ و $j = 0, 1, \dots, n$ ، و بگیریم $p_n(x)$ چندجمله‌ای درجه n باشد، که e^x را در نقاط گرهی x_0, \dots, x_n درونیابی کند. ثابت کنید که وقتی $n \rightarrow \infty$

$$\max_{0 \leq x \leq b} |e^x - p_n(x)| \rightarrow 0$$

راهنمایی: نشان دهید که $|\Psi_n(x)| \leq n!h^{n+1}$ ، $0 \leq x \leq b$ ؛ به هر زیربازه $[x_{j-1}, x_j]$ جداگانه نگاه کنید.

۲۵. ثابت کنید که مسأله کلی (۱۶.۶.۳) ارمیت، در بین چندجمله‌یهای از درجه نایزگتر از $N-1$ ، یک جواب یکتای $p(x)$ دارد.

راهنمایی: نشان دهید که مسأله همگن برای دستگاه خطی متناظر فقط دارای جواب صفر است.

۲۶. مسأله ارمیت زیر را در نظر می‌گیریم

$$p^{(j)}(x_i) = y_i^{(j)} \quad i = 1, 2 \quad j = 0, 1, 2$$

که در آن $p(x)$ یک چندجمله‌ای از درجه نایزگتر از ۵ است.

(الف) برای $p(x)$ فرمولی از نوع لاگرانژ به دست آورید که تعمیمی از درونیابی درجه سوم (۱۲.۶.۳) ارمیت باشد.

راهنمایی: برای توابع پایه که در $l(x_2) = l'(x_2) = l''(x_2) = 0$ صدق می‌کنند، از $l(x) = (x-x_2)^2 g(x)$ ، با $g(x)$ از درجه نایزگتر از ۲ استفاده کنید. $g(x)$ را پیدا کنید.

(ب) یک فرمول تفاضلات منقسم نیوتن پیدا کنید که تعمیم فرمول (۱۳.۶.۳) باشد.

(ج) یک فرمول خطا پیدا کنید که تعمیم فرمول (۱۴.۶.۳) باشد.

۲۷. گیریم $p(x)$ یک چندجمله‌ای باشد که مسألهٔ درونیابی ارمیت زیر را حل کند.

$$p^{(j)}(a) = f^{(j)}(a) \quad p^{(j)}(b) = f^{(j)}(b) \quad j = 0, 1, \dots, n-1$$

وجود آن با استدلال مسألهٔ ۲۵ تضمین می‌شود. فرض کنید که $f(x)$ دارای $2n$ مشتق پیوسته

روی $[a, b]$ باشد، نشان دهید که برای $a \leq x \leq b$

$$f(x) - p_n(x) = \frac{(x-a)^n(x-b)^n}{(2n)!} f^{(2n)}(\xi_x)$$

$$a \leq \xi_x \leq b$$

راهنمایی: استدلال در قضیهٔ ۲.۳ را تعمیم دهید.

۲۸. (الف) یک چندجمله‌ای $p(x)$ از درجهٔ نایزگتر از ۲ پیدا کنید که در

$$p(x_0) = y_0 \quad p'(x_0) = y'_0 \quad p'(x_1) = y'_1$$

صدق کند. فرمولی به شکل زیر به دست آورید

$$p(x) = y_0 l_0(x) + y'_0 l_1(x) + y'_1 l_2(x)$$

(ب) فرمولی برای مسألهٔ درونیابی چندجمله‌ای زیر بیابید.

گیریم $i = 0, 1, 2$, $x_i = x_0 + ih$ یک چندجمله‌ای $p(x)$ از درجهٔ نایزگتر از ۴ بیابید که برای آن

$$p(x_i) = y_i \quad i = 0, 1, 2$$

$$p'(x_0) = y'_0 \quad p'(x_2) = y'_2$$

که مقادیر y داده شده‌اند.

۲۹. مسألهٔ پیدا کردن یک چندجمله‌ای درجهٔ دو $p(x)$ که برای آن

$$p(x_0) = y_0 \quad p'(x_1) = y'_1 \quad p(x_2) = y_2$$

را در نظر بگیرید به طوری که $x_0 \neq x_2$ و $\{y_0, y'_1, y'_2\}$ داده‌های معلوم باشند. با فرض

حقیقی بودن گره‌های x_0, x_1, x_2 چه شرایطی باید برقرار باشد تا چنین $p(x)$ ‌ای موجود و یکتا

باشد؟ این مسأله، مسألهٔ ۲۸ (الف) و مسألهٔ بعدی مثالهایی از مسائل درونیابی ارمیت-برکاف^۱

هستند [لورنتس و همکاران (۱۹۸۳) را ببینید].

۳۰. (الف) نشان دهید که یک چندجمله‌ای درجه سه یکتای $p(x)$ وجود دارد که برای آن

$$\begin{aligned} p(x_0) &= f(x_0) & p(x_1) &= f(x_1) \\ p'(x_1) &= f'(x_1) & p''(x_1) &= f''(x_1) \end{aligned}$$

که $f(x)$ یک تابع داده شده است و $x_0 \neq x_1$. فرمولی برای $p(x)$ به دست آورید.
 (ب) گیریم $x_0 = -1, x_1 = 1$ ، با فرض اینکه $f(x)$ روی $[-1, 1]$ چهاربار پیوسته مشتقپذیر باشد، نشان دهید که به ازای $-1 \leq x \leq 1$

$$f(x) - p(x) = \frac{x^4 - 1}{4!} f^{(4)}(\xi_x)$$

برای مقداری $\xi_x \in [-1, 1]$.

راهنمایی: از برهان قضیه ۲.۳ پیروی کنید.

۳۱. برای تابع $f(x) = \sin x, 0 \leq x \leq \pi/2$ ، تابع درجه سه تک‌ای $Q_n(x)$ ارمیت و تابع درجه سه تک‌ای $L_m(x)$ لاگرانژ را برای $m = 3, 6, 12$ پیدا کنید. خطاهای ماکسیم $f(x) - L_m(x), f'(x) - L'_m(x), f(x) - Q_n(x), f'(x) - Q'_n(x)$ را محاسبه کنید. این را می‌توان با دقت قابل قبول در محاسبه خطاها در $8n$ نقطه متساوی‌الفاصله در $[0, \pi/2]$ به دست آورد.

۳۲. (الف) گیریم $p_2(x)$ یک چندجمله‌ای درجه سه باشد که $f(x)$ را در نقاط متساوی‌الفاصله $x_j = x_0 + jh, j = 0, 1, 2, 3$ ، درونیایی می‌کند. با فرض اینکه $f(x)$ به اندازه کافی مشتقپذیر است، کران خطای استفاده از $p'_2(x)$ را به عنوان تقریبی برای $f'(x)$ در $x_0 \leq x \leq x_3$ پیدا کنید.

(ب) گیریم $H_2(x)$ چندجمله‌ای درجه سه ارمیت باشد که $f(x)$ و $f'(x)$ را در x_0 و $x_1 = x_0 + h$ درونیایی می‌کند. کران خطای $f'(x) - H'_2(x)$ را برای $x_0 \leq x \leq x_1$ پیدا کنید.
 (ج) توابع چندجمله‌ای تک‌ای $L_m(x)$ و $Q_n(x)$ بخش ۷.۳ را در نظر بگیرید، $m = \frac{2n}{3}$ ، و کران خطاهای $f'(x) - L'_m(x)$ و $f'(x) - Q'_n(x)$ را پیدا کنید. آن را برای حالت خاص $f(x) = x^4$ به کار برید و جوابهای خود را با نتایج عددی که در جدول ۱۴.۳ داده شده‌اند مقایسه کنید.

۳۳. گیریم $s(x)$ یک تابع برآزا از مرتبه m باشد. گیریم b یک گره و $s(x)$ یک چندجمله‌ای از درجه نابزرگتر از $m - 1$ روی $[a, b]$ و $[b, c]$ باشد. نشان دهید که اگر $s^{(m-1)}(x)$ در $x = b$ پیوسته باشد، آنگاه $s(x)$ برای $a \leq x \leq c$ یک چندجمله‌ای از درجه نابزرگتر از $m - 1$ است.

۳۴. شرایطی را برای برازای درونیاب درجه سهی $s(x)$ پیدا کنید که از شرایط «نه یک گره»ی نقاط انتهایی نتیجه می‌شوند. به بحث حالت (۲) پس از (۲۷.۷.۳) رجوع کنید.

۳۵. پیدا کردن تابع درونیاب برازای درجه ۳ برای داده‌های

x	۰	۱	۲	۲٫۵	۳	۴
y	۱٫۴	۰٫۶	۱٫۰	۰٫۶۵	۰٫۶	۱٫۰

را در نظر می‌گیریم. شرط «نه یک گره» را برای به دست آوردن شرایط مرزی مکمل (۱۹.۷.۳) به دست آورید. تابع نتیجه $s(x)$ را رسم کنید. آن را با درونیابی خطی تکه‌ای، که نقاط پشت سرهم (x_i, y_i) را با پاره خطها به هم متصل می‌کند، مقایسه کنید.

۳۶. برنامه‌ای بنویسید که نرخ همگرایی درونیابی برازای درجه ۳ (مانند جدول ۱۳.۳) را با شرایط مرزی گوناگون، بررسی کند. بسیاری از مراکز رایانه‌یی بسته‌هایی برای تولید این درونیابها دارند که شرایط مرزی را مشخص می‌کنند. اگر این بسته را ندارید، یک دستگاه خطی حلال با (۱۹.۷.۳) و شرایط مرزی اضافی به‌کار برید. شرایط مرزی زیر را بررسی کنید؛ (الف) مشتقات مانند (۲۰.۷.۳) داده شده‌اند؛ (ب) شرط «نه یک گره» که در مسأله ۳۸ داده شده است و (ج) شرایط برازای طبیعی $M_0 = M_n = 0$ در مسأله ۳۸. این برنامه را برای مطالعه همگرایی $s(x)$ به $f(x)$ در حالات زیر به‌کار برید: (۱) $f(x) = e^x$ در $[0, 1]$ ، (۲) $f(x) = \sin(x)$ در $[0, \pi/2]$ ، و (۳) $f(x) = x\sqrt{x}$ در $[0, 1]$. به رفتار خطا نزدیک نقاط انتهایی توجه کنید و آنها را مقایسه نمایید.

۳۷. (الف) گیریم $q(x)$ یک برازای درجه ۳ با تنها یک گره $x = a$ باشد. به علاوه فرض کنید که برای $x \leq a$ ، $q(x) \equiv 0$. نشان دهید که برای مقداری از c داریم $q(x) = c(x-a)_+^2$.

راهنمایی: برای $x \geq a$ چند جمله‌ای درجه سهی $s(x)$ را به شکل زیر بنویسید

$$s(x) \equiv c_0 + c_1(x-a) + c_2(x-a)^2 + c_3(x-a)^3$$

سپس مفروضات درباره $s(x)$ را به‌کار برید.

(ب) با استفاده از قسمت (الف)، رابطه (۳۱.۷.۳) را برای توابع برازای درجه سوم ($m = 4$) ثابت کنید.

۳۸. محاسبه یک برازای درونیاب درجه سه را براساس (۱۹.۷.۳) با شرایط اضافی زیر در نظر می‌گیریم

$$s''(x_0) = M_0 = 0 \quad s''(x_n) = M_n = 0$$

این مسأله دارای جواب یکتای $\hat{s}(x)$ است. نشان دهید

$$\int_{x_0}^{x_n} [\hat{s}''(x)]^2 dx \leq \int_{x_0}^{x_n} [g''(x)]^2 dx$$

که در آن $g(x)$ یک تابع دوبار پیوسته مشتقپذیر است که در شرایط درونیاب $g(x_i) = y_i$ $i = 0, 1, \dots, n$ صدق می‌کند.

راهنمایی: نشان دهید که (۲۷.۷.۳) برای $\hat{s}(x)$ معتبر است. [این برازای درونیاب را برازای درونیاب درجه سه طبیعی گویند. این یک درونیاب هموار برای داده‌هاست، ولی معمولاً نزدیک نقاط انتهایی به کندی همگراست. برای اطلاعات بیشتر درباره این مطلب به دور (۱۹۷۸، ص ۵۵) مراجعه کنید.]
 ۳۹. برای تعریف یک B -برازای از مرتبه m که در بازه $[x_i, x_{i+m}]$ ناصفر است، تعریف می‌کنیم

$$B_i^{(m)}(x) = (x_{i+m} - x) f_x[x_i, x_{i+1}, \dots, x_{i+m}]$$

با $f_x(t) = (t - x)_+^{m-1}$. یک رابطه بازگشتی برای $B_i^{(m)}(x)$ بر حسب B -برازاهای مرتبه $m - 1$ پیدا کنید.

۴۰. با استفاده از تعریف (۳۳.۷.۳) رفتار B -برازاهای درجه سه را وقتی گرهها مجاز باشند بر هم منطبق شوند، تحقیق کنید. نشان دهید که اگر دوتا از گرهها در $\{x_i, x_{i+1}, \dots, x_{i+2}\}$ بر هم منطبق شوند، $B_i(x)$ فقط یک مشتق پیوسته در گره منطبق شده دارد. همچنین، اگر سه گره بر هم منطبق شوند، نشان دهید که $B_i(x)$ در آن نقطه پیوسته است ولی مشتقپذیر نیست.

۴۱. گیریم $0 \leq t_0 < t_1 < \dots < t_{2n} < 2\pi$ را در نظر می‌گیریم. تعریف می‌کنیم (۵.۸.۳)

$$l_j(t) = \prod_{\substack{k=0 \\ k \neq j}}^{2n} \frac{\sin \frac{1}{2}(t - t_k)}{\sin \frac{1}{2}(t_j - t_k)}$$

برای $j = 0, 1, \dots, 2n$. به سادگی معلوم می‌شود که $\delta_{ij} = l_j(t_i)$ $0 \leq i, j \leq 2n$. نشان دهید $l_j(t)$ یک چندجمله‌ای مثلثاتی از درجه نایزگتر از n است. در این صورت جواب (۵.۸.۳) با رابطه زیر داده می‌شود

$$P_n(t) = \sum_{j=0}^{2n} f(t_j) l_j(t)$$

راهنمایی: استقرا بر n را همراه با اتحادهای مثلثاتی معمولی به‌کار برید.

۴۲. (الف) فرمولهای زیر را ثابت کنید:

$$\sum_{j=1}^{m-1} \sin\left(\frac{2\pi jk}{m}\right) = 0, \quad k \text{ برای همه اعداد صحیح } k, m \geq 2$$

$$\sum_{j=0}^{m-1} \cos\left(\frac{2\pi jk}{m}\right) = \begin{cases} m & k \text{ مضربی از } m \text{ باشد} \\ 0 & k \text{ مضربی از } m \text{ نباشد} \end{cases}$$

(ب) با استفاده از فرمولهای فوق فرمولهایی برای عبارات زیر بیابید

$$\sum_{j=0}^{m-1} \cos\left(\frac{2\pi jk}{m}\right) \cos\left(\frac{2\pi jl}{m}\right) \quad \sum_{j=1}^{m-1} \sin\left(\frac{2\pi jk}{m}\right) \sin\left(\frac{2\pi jl}{m}\right)$$

$$\sum_{j=0}^{m-1} \cos\left(\frac{2\pi jk}{m}\right) \sin\left(\frac{2\pi jl}{m}\right)$$

به ازای $k \leq 0$ و $l \leq m-1$. به فرمولهای به دست آمده روابط متعامد گسسته گویند و نظیر روابط متعامد انتگرالی برای $\{\cos(kx) \text{ و } \sin(kx)\}$ هستند.

۴۳. تبدیل متناهی فوریه از مرتبه m را برای دنباله‌های زیر پیدا کنید

(الف) $0 \leq k \leq m-1 \quad x_k = 1$

(ب) $0 \leq k \leq m-1 \quad x_k = (-1)^k$ زوج m

(ج) $0 \leq k \leq m-1 \quad x_k = k$

تقریب توابع

برای تعیین مقدار اغلب توابع ریاضی ابتدا باید تقریبهای محاسبه‌پذیری برای آنها تهیه کنیم. توابع در کاربردها از راههای گوناگونی تعریف می‌شوند، انتگرال‌ها و سری‌های نامتناهی معمولیترین فرمول‌هایی هستند که در این تعریفها به‌کار می‌روند. این گونه تعریفها برای مشخص‌کردن ویژگیهای تابع سودمندند، ولی معمولاً راه کارآمدی برای تعیین مقدار تابع نیستند. در این فصل ما استفاده از چندجمله‌یها را در تقریب یک تابع داده‌شده بررسی می‌کنیم. ابزارهای مختلف تهیه تقریبهای چندجمله‌یی را توضیح می‌دهیم و آنها را برحسب دقت نسبی‌شان باهم مقایسه می‌کنیم.

برای محاسبه یک تابع $f(x)$ در رایانه، معمولاً داشتن یک تقریب تحلیلی $f(x)$ ، از لحاظ مکان و زمان، کاراتر از ذخیره‌کردن یک جدول و استفاده از درونیابی است. همچنین بهتر است، در بین چندجمله‌یهای که دقت مطلوب در تقریب $f(x)$ را به‌دست می‌دهند، پایینترین درجه ممکن به‌کار برده شود. بخشهای زیر برخی از روشهای تولید تقریب را به‌دست می‌دهند، و معمولاً تقریبهای بهتر آنهایی هستند که تهیه آنها دشوارتر است. مقدار زمان و کوششی که برای تهیه یک تقریب به‌کار می‌رود باید با مقدار استفاده از تقریب نسبت مستقیم داشته باشد. اگر فقط چند مرتبه به‌کار رود اغلب یک سری بریده شده تیلر کافی است. ولی اگر تقریب را باید میلیونها بار افراد زیادی به‌کار برند، آنگاه دقت زیادی در تهیه آن تقریب باید به عمل آید.

غیر از چند جمله‌یها شکل‌های دیگری برای تقریب وجود دارند. توابع گویا خارج قسمتهای چند جمله‌یها هستند، و معمولاً تا حدی صورتهای تقریب کاراتری هستند. ولی از آنجا که چند جمله‌یها شکل مناسب و کارای تقریب را در اختیار می‌گذارند، و به دلیل اینکه نظریه تقریب توابع گویا پیچیده‌تر از نظریه تقریب چند جمله‌یهاست، ما برای بررسی فقط چند جمله‌یها را انتخاب کرده‌ایم. قضایای این فصل را می‌توان در پیدا کردن تقریبهای چند جمله‌یهای تکه‌یی، که تا حدی مشابه توابع درونیاب چند جمله‌یهای تکه‌یی بخش ۷.۳ فصل قبل هستند، به کار برد.

۱.۴ قضیه وایرستراس و قضیه تیلر

برای توجیه استفاده از چند جمله‌یها در تقریب توابع پیوسته، قضیه زیر را معرفی می‌کنیم.

قضیه ۱.۴ (وایرستراس) گیریم $f(x)$ برای $a \leq x \leq b$ پیوسته باشد و $\epsilon > 0$. در این صورت یک چند جمله‌ی $p(x)$ وجود دارد که برای آن

$$|f(x) - p(x)| \leq \epsilon \quad a \leq x \leq b$$

برهان برهانهای زیادی برای این قضیه و تعمیمهای آن وجود دارد. چون این امر برای مطالعه آنالیز عددی، نکته مهمی نیست، ما فقط به یک برهان ساختنی اشاره می‌کنیم. برای برهانهای دیگر، فصل ۶، دیویس (۱۹۶۳) را ببینید.

برای سادگی فرض کنید، $[a, b] = [0, 1]$: با یک تعویض متغیر مناسب، همیشه می‌توان در صورت لزوم بازه را به این حالت تبدیل کرد. تعریف می‌کنیم

$$p_n(x) = \sum_{k=0}^n \binom{n}{k} f\left(\frac{k}{n}\right) x^k (1-x)^{n-k} \quad 0 \leq x \leq 1$$

گیریم $f(x)$ بر $[0, 1]$ کراندار باشد. پس در هر نقطه x که f پیوسته باشد

$$\lim_{n \rightarrow \infty} p_n(x) = f(x)$$

اگر $f(x)$ در هر نقطه x از $[0, 1]$ پیوسته باشد، همگرایی p_n به f بر $[0, 1]$ یکنواخت است، یعنی

$$\max_{0 \leq x \leq 1} |f(x) - p_n(x)| \rightarrow 0 \quad \text{وقتی } n \rightarrow \infty \quad (1.1.4)$$

این مطلب برای پیدا کردن یک چند جمله‌ی که در حکم قضیه صدق کند، راه روشنی به دست می‌دهد. اثبات این نتایج را همراه با ویژگیهای دیگری از چند جمله‌یهای تقریب $p_n(x)$ که

چندجمله‌یی برنشتاین نامیده می‌شوند، می‌توان در صفحات ۱۰۸ تا ۱۱۸ دیویس (۱۹۶۳)، یافت. این چندجمله‌یها از لحاظ رفتار کیفی شباهت خیلی زیادی به تابع $f(x)$ دارند. برای مثال، اگر $f(x)$ ، r بار پیوسته مشتقپذیر بر $[0, 1]$ باشد، آنگاه

$$\max_{0 \leq x \leq 1} |f^{(r)}(x) - p_n^{(r)}(x)| \rightarrow 0 \quad \text{وقتی} \quad n \rightarrow \infty$$

گرچه یک چنین ویژگی کلی تقریب در حد خودش بالارزش است ولی در این حالت، همگرایی (۱.۱.۴) عموماً خیلی کند است. برای مثال، اگر $f(x) = x^2$ ، آنگاه

$$\lim_{n \rightarrow \infty} n[p_n(x) - f(x)] = x(1-x)$$

و بنابراین برای مقادیر بزرگ n

$$p_n(x) - x^2 \doteq \frac{1}{n}x(1-x)$$

■ حتی برای تقریب حالت نمایانی چون $f(x) = x^2$ ، خطا به سرعت کاهش نمی‌یابد.

قضیه تیلر قضیه تیلر قبلاً، در قضیه ۴.۱، بخش ۱.۱ از فصل ۱ ارائه گردید. این اولین ابزار مهم برای تقریب زدن یک تابع است، و اغلب به عنوان یک تقریب اولیه در محاسبه تقریبهای کاراتر به‌کار می‌رود. برای اینکه ببینیم چرا تقریب تیلر آنقدرها کارا نیست، مثال زیر را در نظر می‌گیریم.

مثال خطای تقریب e^x را با استفاده از چندجمله‌یی $p_r(x)$ درجه سوم تیلر بر بازه $[-1, 1]$ ، با بسط در حول $x = 0$ پیدا کنید.

$$p_3(x) = 1 + x + \frac{1}{2}x^2 + \frac{1}{6}x^3$$

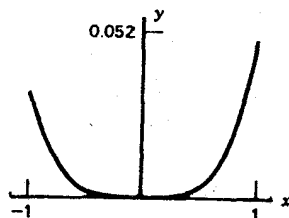
$$e^x - p_3(x) = R_3(x) = \frac{1}{24}x^4 e^\xi \quad (2.1.4)$$

با ξ بین 0 و x .

برای اینکه خطا را با دقت آزمایش کنیم، آن را از بالا و پایین محدود می‌کنیم:

$$\frac{1}{24}x^4 \leq e^x - p_3(x) \leq \frac{e}{24}x^4 \quad 0 \leq x \leq 1$$

$$\frac{e^{-1}}{24}x^4 \leq e^x - p_3(x) \leq \frac{1}{24}x^4 \quad -1 \leq x \leq 0$$

شکل ۱.۴ منحنی خطا برای $p_7(x) \doteq e^x$.

خطا با افزایش $|x|$ ، افزایش می‌یابد، و با محاسبه مستقیم خواهیم داشت:

$$\text{Max}_{-1 \leq x \leq 1} |e^x - P_7| \doteq 0.516 \quad (3.1.4)$$

این خطا در بازه $[-1, 1]$ به‌طور هموار توزیع نشده است (شکل ۱.۴ را ببینید). خطا در نزدیکی مبدأ بسیار کوچکتر از خطا نزدیک به نقاط انتهایی -1 و $+1$ است. این توزیع نامساوی خطا، که نوعاً معرف باقیمانده تیلر است، بدین معناست که معمولاً چندجمله‌بیهای تقریب بسیار بهتری از همان درجه یافت می‌شوند. مثالهای بیشتری در بخش بعد داده شده است.

فضای تابعی $C[a, b]$ مجموعه $C[a, b]$ ی توابع پیوسته و حقیقی مقدار بر بازه $[a, b]$ در بخش ۱.۱ از فصل ۱ معرفی شده است. در این مجموعه معمولاً نرم زیر را به‌کار می‌بریم

$$\|f\|_\infty = \text{Max}_{a \leq x \leq b} |f(x)| \quad f \in C[a, b] \quad (4.1.4)$$

این نرم به نامهای مختلف نرم ماکسیمم، نرم چیشیف، نرم بینهایت، و نرم یکتواخت نیز نامیده شده است. این اندازه طبیعی برای استفاده در نظریه تقریب است، زیرا ما می‌خواهیم

$$\|f - p\|_\infty = \max_{a \leq x \leq b} |f(x) - p(x)| \quad (5.1.4)$$

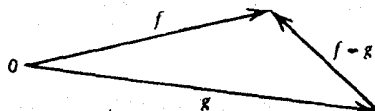
را برای چندجمله‌بیهای گوناگون $p(x)$ ، معین و مقایسه کنیم. نرم دیگری برای $C[a, b]$ در بخش ۳.۴ معرفی شده است، که آن نیز در اندازه‌گیری بزرگی $f(x) - p(x)$ مفید است.

همان‌گونه که در بخش ۱.۱ اشاره کردیم، نرم ماکسیمم در ویژگیهای مخصوص نرم که ذیلاً می‌آید صدق می‌کند:

$$\|f\| = 0 \quad \text{اگر و فقط اگر} \quad f \equiv 0 \quad (6.1.4)$$

$$\|\alpha f\| = |\alpha| \|f\| \quad \alpha \text{ و تمام اسکالرهایی} \quad f \in C[a, b] \quad (7.1.4)$$

$$\|f + g\| \leq \|f\| + \|g\| \quad f, g \in C[a, b] \quad \text{برای تمام مقادیر} \quad (8.1.4)$$



شکل ۲.۴ نمایش تعریف فاصله $D(f, g)$.

اثبات این ویژگیها، برای (۴.۱.۴) کاملاً سراسر است. این ویژگیها نشان می‌دهند که نرم باید تعمیم قدرمطلق یک عدد تلقی شود.

اغلب سودمند است که $C[a, b]$ را یک فضای برداری تلقی کنیم، اگرچه هیچ استفاده مهمی از این ایده نخواهیم کرد. در اینجا بردارها، توابع $f(x)$ ، $a \leq x \leq b$ هستند. فاصله بردار f از بردار g را چنین تعریف می‌کنیم:

$$D(f, g) = \|f - g\| \quad (۸.۱.۴)$$

که با احساس ما از مفهوم فاصله در فضاهای برداری ساده‌تر مطابقت می‌نماید. این مطلب در شکل ۲.۴ نشان داده شده است. با استفاده از نابرابری (۸.۱.۴)،

$$\begin{aligned} \|f - g\| &= \|(f - h) + (h - g)\| \leq \|f - h\| + \|h - g\| \\ D(f, g) &\leq D(f, h) + D(h, g) \end{aligned} \quad (۱۰.۱.۴)$$

این نابرابری به علت تعبیر روشن آن در اندازه‌گیری طولهای اضلاع مثلثی با رئوس f و g و h ، نابرابری مثلثی خوانده می‌شود. فرمول هم‌ارز (۸.۱.۴) را نیز نابرابری مثلثی می‌گویند. یک قضیه مفید دیگر، عکس نابرابری مثلثی است

$$\| \|f\| - \|g\| \| \leq \|f - g\| \quad (۱۱.۱.۴)$$

برای اثبات آن از (۸.۱.۴) استفاده می‌کنیم تا به دست آوریم

$$\|f\| \leq \|f - g\| + \|g\|$$

$$\|f\| - \|g\| \leq \|f - g\|$$

همچنین داریم

$$\|g\| - \|f\| \leq \|g - f\| = \|f - g\|$$

که تساوی اخیر از (۷.۱.۴) با $\alpha = -1$ به دست می‌آید. از ترکیب این دو نابرابری (۱۱.۱.۴) اثبات می‌شود.

یک معرفی کاملتر از فضاهای برداری (فقط برای ابعاد متناهی) و نرم‌های برداری در فصل ۷ داده شده است. و یک هندسه دیگری برای $C[a, b]$ ، همراه با معرفی یک نرم دیگر غیر از (۴.۱.۴)، در بخشهای ۳.۴ و ۴.۴ مطرح شده است. برای مواردی که می‌خواهیم درباره توابعی که چندین مشتق پیوسته دارند، صحبت کنیم، فضای تابعی $C^r[a, b]$ را معرفی می‌نماییم، که متشکل از توابع $f(x)$ با r مشتق پیوسته بر $[a, b]$ است. این فضای برداری فایده جداگانه‌ای دارد ولی ما آن را فقط به عنوان یک وسیله ساده‌کردن نمادی تلقی می‌کنیم.

۲.۴ مسأله تقریب مینیماکس

گیریم $f(x)$ بر $[a, b]$ پیوسته باشد. برای مقایسه چندجمله‌بیهای $p(x)$ ، تقریبهای $f(x)$ که از روشهای گوناگونی به دست آمده‌اند، طبیعی است بپرسیم بهترین دقتی که ممکن است از به‌کار بردن چندجمله‌یی از هر درجه $n \geq 0$ به دست آورد چیست؟ بنابراین به معرفی خطای مینیماکس می‌رسیم

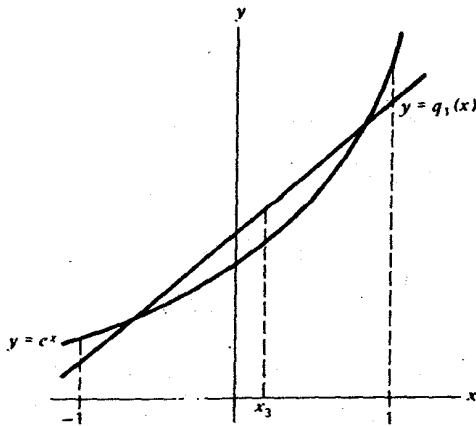
$$\rho_n(f) = \inf_{\deg(q) \leq n} \|f - q\|_\infty \quad (1.2.4)$$

یک چندجمله‌یی $q(x)$ از درجه نایزگتر از n وجود ندارد که $f(x)$ را با خطای ماکسیممی کمتر از $\rho_n(f)$ تقریب زند. با معرفی $\rho_n(f)$ ، می‌خواهیم ببینیم آیا یک چندجمله‌یی چون $q_n^*(x)$ وجود دارد به طوری که

$$\rho_n(f) = \|f - q_n^*\|_\infty \quad (2.2.4)$$

و اگر چنین است، آیا این چندجمله‌یی یکتاست؟ ویژگیهای آن چیست و چگونه می‌توان آن را ساخت؟ تقریب $q_n^*(x)$ را تقریب مینیماکس برای $f(x)$ بر $[a, b]$ خوانند و نظریه آن در بخش ۶.۴ مطرح شده است.

مثال چندجمله‌یی تقریب مینیماکس $q_1^*(x)$ را برای e^x در $-1 \leq x \leq 1$ محاسبه کنید. گیریم $q_1^*(x) = a_0 + a_1x$. برای پیدا کردن a_0 و a_1 ، باید از یک بینش هندسی استفاده کنیم. نمودار $y = e^x$ را با نمودار یک تقریب ممکن $y = q_1(x)$ مانند شکل ۳.۴ در نظر می‌گیریم.



شکل ۳.۴ تقریب مینیمکس خطی برای e^x

فرض می‌کنیم

$$\epsilon(x) = e^x - [a_0 + a_1x] \quad (۳.۲.۴)$$

روشن است که $q_1^*(x)$ و e^x باید در دو نقطه بازه $[-1, 1]$ ، مثلاً در $-1 < x_1 < x_2 < 1$ ، با هم برابر باشند. وگرنه می‌توانیم با جابه‌جا کردن مناسب نمودار $y = q_1^*$ این تقریب را بهبود بخشیم. همچنین

$$\rho_1 = \text{Max}_{-1 \leq x \leq 1} |\epsilon(x)|$$

و $\epsilon(x_1) = \epsilon(x_2) = 0$. با استدلال دیگری که بر پایه جابه‌جا کردن نمودار $y = q_1^*(x)$ استوار است، می‌توانیم نتیجه بگیریم که ماکسیمم خطای ρ_1 دقیقاً در سه نقطه حاصل می‌شود.

$$\epsilon(-1) = \rho_1 \quad \epsilon(1) = \rho_1 \quad \epsilon(x_3) = -\rho_1 \quad (۴.۲.۴)$$

که $x_1 < x_2 < x_3$. چون $\epsilon(x)$ یک مینیمم نسبی در x_3 دارد، داریم $\epsilon'(x_3) = 0$. از ترکیب این چهار معادله، داریم

$$\begin{aligned} e^{-1} - [a_0 - a_1] &= \rho_1 & e - [a_0 + a_1] &= \rho_1 \\ e^{x_3} - [a_0 + a_1x_3] &= -\rho_1 & e^{x_3} - a_1 &= 0 \end{aligned} \quad (۵.۲.۴)$$

این معادلات دارای جوابهای زیرند

$$a_1 = \frac{e - e^{-1}}{2} \doteq 1,1752 \quad x_r = \ln_e(a_1) \doteq 0,1614$$

$$\rho_1 = \frac{1}{2}e^{-1} + \frac{x_r}{2}(e - e^{-1}) \doteq 0,2788$$

$$a_0 = \rho_1 + (1 - x_r)a_1 \doteq 1,2643$$

بنابراین

$$q_1^*(x) = 1,2643 + 1,1752x \quad (۶.۲.۴)$$

$$\rho_1 \doteq 0,2788 \text{ و}$$

با استفاده از آنچه که الگوریتم ریمز خوانده می‌شود، می‌توانیم $q_1^*(x)$ را برای e^x در $[-1, 1]$

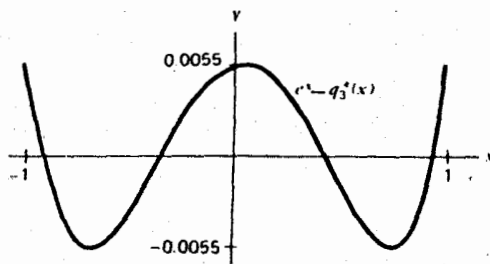
بسازیم:

$$q_3^*(x) = 0,994579 + 0,995668x + 0,542973x^2 + 0,179533x^3 \quad (۷.۲.۴)$$

نمودار خطای آن در شکل ۴.۴ داده شده است، و برخلاف تقریبهای تیلر (شکل ۱.۴ را ببینید) خطا در سراسر بازه تقریب به‌طور هموار توزیع شده است.

مثال خود را با بیان خطاهای تقریب مینیمکس $q_n^*(x)$ و تقریبهای تیلر برای $f(x) = e^x$ در $-1 \leq x \leq 1$ ، خاتمه می‌دهیم. خطاهای ماکسیم برای مقادیر گوناگون n در جدول ۱.۴ داده شده است.

دقت $q_n^*(x)$ به‌طور چشمگیری از دقت $p_n(x)$ بیشتر است و تفاوت با افزایش n زیاد می‌شود. باید توجه داشت که e^x تابعی است که سری تیلر آن در مقایسه با $\log x$ و $\tan^{-1} x$ به سرعت همگرا می‌شود. درباره توابع اخیر، مینیمکس در مقایسه با سری تیلر، بهتر به‌نظر می‌رسد.



شکل ۴.۴ خطا برای تقریب مینیمکس درجه ۳ برای e^x

جدول ۱.۴ خطای تیلر مینیمکس مقدار برای e^x

n	$\ f - p_n\ _\infty$	$\ f - q_n^*\ _\infty$
۱	$۷,۱۸E - ۱$	$۲,۷۹E - ۱$
۲	$۲,۱۸E - ۱$	$۴,۵۰E - ۲$
۳	$۵,۱۶E - ۲$	$۵,۵۲E - ۳$
۴	$۹,۹۵E - ۳$	$۵,۴۷E - ۴$
۵	$۱,۶۲E - ۳$	$۴,۵۲E - ۵$
۶	$۲,۲۶E - ۴$	$۳,۲۱E - ۶$
۷	$۲,۷۹E - ۵$	$۲,۰۰E - ۷$
۸	$۳,۰۶E - ۶$	$۱,۱۱E - ۸$
۹	$۳,۰۱E - ۷$	$۵,۵۲E - ۱۰$

۳.۴ مسأله تقریب کمترین مربعات

به علت دشواری محاسبه تقریب مینیمکس، معمولاً، یک تقریب بینابینی در نظر می‌گیریم که تقریب کمترین مربعات خوانده می‌شود. نماد

$$\|g\|_2 = \sqrt{\int_a^b |g(x)|^2 dx} \quad g \in C[a, b] \quad (۱.۳.۴)$$

را معرفی می‌کنیم. این یک نرم تابعی است که، مانند نرم ماکسیمم (۴.۱.۴)، در ویژگیهای (۶.۱.۴) - (۸.۱.۴) صدق می‌کند. این تعمیم نرم معمولی اقلیدسی برای R^n است که در (۱۷.۱.۱) تعریف شده بود. در بخش بعد، هنگامی که تعریف بالا را تعمیم خواهیم داد، به اثبات نابرابری مثلثی (۸.۱.۴) برمی‌گردیم.

برای تابع مفروض $f \in C[a, b]$ و $n \geq 0$ ، تعریف می‌کنیم

$$M_n(f) = \inf_{\deg(r) \leq n} \|f - r\|_2 \quad (۲.۳.۴)$$

مانند قبل، می‌خواهیم بدانیم آیا یک چندجمله‌یی r_n^* وجود دارد که عبارت

$$M_n(f) = \|f - r_n^*\|_2 \quad (۳.۳.۴)$$

را مینیمم سازد؟ آیا این چندجمله‌یی یکتاست؟ آیا می‌توانیم آن را حساب کنیم؟

برای توجیه بیشتر (۲.۳.۴)، محاسبه خطای متوسط در تقریب $f(x)$ توسط $r(x)$ را در نظر می‌گیریم. برای عدد صحیحی مانند $m \geq 1$ ، نقاط گرهی x_j را چنین تعریف می‌کنیم

$$x_j = a + \left(j - \frac{1}{2}\right) \left(\frac{b-a}{m}\right) \quad j = 1, 2, \dots, m$$

اینها نقاط میانی m زیر بازه متساوی‌الفاصله $[a, b]$ هستند. پس یک خطای متوسط تقریب زد $f(x)$ توسط $r(x)$ بر $[a, b]$ عبارت است از

$$\begin{aligned} E &= \lim_{m \rightarrow \infty} \left\{ \frac{1}{m} \sum_{j=1}^m [f(x_j) - r(x_j)]^2 \right\}^{1/2} \\ &= \lim_{m \rightarrow \infty} \left\{ \frac{1}{b-a} \sum_{j=1}^m [f(x_j) - r(x_j)]^2 \left(\frac{b-a}{m}\right) \right\}^{1/2} \\ &= \frac{1}{\sqrt{b-a}} \sqrt{\int_a^b |f(x) - r(x)|^2 dx} \\ E &= \frac{\|f - r\|_2}{\sqrt{b-a}} \end{aligned} \quad (4.3.4)$$

بنابراین، تقریب کمترین مربعات باید خطای متوسط کوچکی بر $[a, b]$ داشته باشد. کمیت E جذر میانگین مربع خطاها در تقریب $f(x)$ توسط $r(x)$ نامیده می‌شود.

مثال گیریم $f(x) = e^x$ ، $-1 \leq x \leq 1$ ، و گیریم $r_1(x) = b_0 + b_1 x$. مینیمم عبارت

$$\|f - r_1\|_2^2 = \int_{-1}^1 [e^x - b_0 - b_1 x]^2 dx \equiv F(b_0, b_1) \quad (5.3.4)$$

را پیدا کنید. اگر عبارت زیر علامت انتگرال را بسط دهیم و انتگرال را به انتگرال‌های کوچکتر بشکنیم، $F(b_0, b_1)$ یک چندجمله‌یی درجه دوم از دو متغیر b_0 و b_1 خواهد شد،

$$F = \int_{-1}^1 \{e^{2x} + b_0^2 + b_1^2 x^2 - 2b_0 e^x - 2b_1 x e^x + 2b_0 b_1 x\} dx$$

برای پیدا کردن مینیمم، شرط لازم در یک نقطه مینیمم را می‌نویسیم

$$\frac{\partial F}{\partial b_0} = 0 \quad \frac{\partial F}{\partial b_1} = 0$$

به جای مشتق‌گیری از انتگرال اخیر از عبارت زیر علامت انتگرال (۵.۳.۴) مشتق می‌گیریم؛

$$\circ = \frac{\partial F}{\partial b_0} = \int_{-1}^1 \frac{\partial}{\partial b_0} [e^x - b_0 - b_1 x]^2 dx = 2 \int_{-1}^1 [e^x - b_0 - b_1 x](-1) dx$$

$$\circ = \frac{\partial F}{\partial b_1} = 2 \int_{-1}^1 [e^x - b_0 - b_1 x](-x) dx$$

بنابراین

$$b_0 = \frac{1}{2} \int_{-1}^1 e^x dx = \sinh(1) \doteq 1,1752$$

$$b_1 = \frac{3}{2} \int_{-1}^1 x e^x dx = 3e^{-1} \doteq 1,1036$$

$$r_1^*(x) = 1,1752 + 1,1036x \quad (۶.۳.۴)$$

با بررسی مستقیم

$$\|e^x - r_1^*\|_{\infty} \doteq 0,44$$

این یک مقدار بین تقریبهای q_1^* و $p_1(x)$ است که قبلاً به دست آمده بودند. معمولاً تقریب کمترین مربعات یک تقریب یکنواخت نسبتاً خوب، برتر از تقریبهای سری تیلر است.

به عنوان یک مثال دیگر، تقریب کمترین مربعات درجهٔ سوم برای e^x بر $[-1, 1]$ چنین است

$$r_3^*(x) = 0,996294 + 0,997955x + 0,536722x^2 + 0,176139x^3 \quad (۷.۳.۴)$$

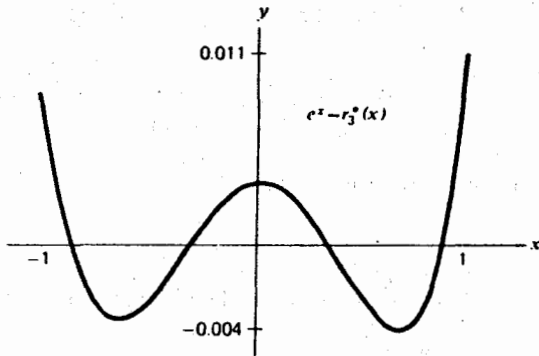
و

$$\|e^x - r_3^*\|_{\infty} = 0,112$$

نمودار خطا در شکل ۵.۴ داده شده است. توجه کنید که خطا در $[-1, 1]$ به اندازهٔ مینیماکس $q_3^*(x)$ در شکل ۴.۴ به طور هموار توزیع نشده است.

مسألهٔ کلی کمترین مربعات مسألهٔ کمترین مربعات (۲.۳.۴) را با دخالت دادن خطاهای میانگین وزن در تقریب $f(x)$ توسط یک چندجمله‌یی $r(x)$ تعمیم می‌دهیم. نظریهٔ کلی وجود، یکتایی، و ساختن تقریبهای کمترین مربعات در بخش ۵.۴ داده شده است. در آنجا از نظریهٔ چندجمله‌یهای متعامد که در بخش بعد خواهد آمد، استفاده شده است.

گیریم $w(x)$ یک تابع وزن نامنفی بر بازهٔ (a, b) باشد، که ممکن است نامتناهی باشد، و برای آن ویژگیهای زیر را فرض می‌کنیم:



شکل ۵.۴ خطا در تقریب کمترین مربعات درجه سوم برای e^x

۱.

$$\int_a^b |x|^n w(x) dx \quad (۸.۳.۴)$$

انتگرالپذیر و برای تمام مقادیر $n \geq 0$ متناهی است؛

۲. فرض می‌کنیم برای یک تابع پیوسته و نامنفی $g(x)$ داشته باشیم؛

$$\int_a^b w(x)g(x)dx = 0 \quad (۹.۳.۴)$$

در این صورت، در (a, b) ، $g(x) \equiv 0$.

مثال توابع وزنی که در این کتاب از آنها بیشتر استفاده شده به شکل زیرند؛

$$w(x) \equiv 1 \quad a \leq x \leq b$$

$$w(x) = \frac{1}{\sqrt{1-x^2}} \quad -1 \leq x \leq 1$$

$$w(x) = e^{-x} \quad 0 \leq x < \infty$$

$$w(x) = e^{-x^2} \quad -\infty < x < \infty$$

اکنون می‌توان برای یک بازه متناهی $[a, b]$ ، مسأله کلی کمترین مربعات را بیان کرد. $g \in C[a, b]$ داده شده است، آیا در میان تمام چندجمله‌بیهای $r(x)$ از درجه نایزگتر از n ، یک چندجمله‌یی

$r_n^*(x)$ وجود دارد که

$$\int_a^b w(x)[f(x) - r(x)]^2 dx \quad (۱۰.۳.۴)$$

را مینیمم سازد؟

تابع $w(x)$ به خطاهای نقاط مختلف بازه $[a, b]$ اهمیت‌های متفاوتی می‌دهد. این امر در مطالعه تقریب‌های نزدیک مینیماکس مفید خواهد بود. برای محاسبه (۱۰.۳.۴) برای یک چندجمله‌ی اختیاری $r(x)$ از درجه نایزگتر از n تعریف می‌کنیم

$$F(a_0, a_1, \dots, a_n) = \int_a^b w(x) \left[f(x) - \sum_{j=0}^n a_j x^j \right]^2 dx \quad (11.3.4)$$

ما می‌خواهیم وقتی که ضرایب $\{a_j\}$ همه اعداد حقیقی را اختیار می‌کنند، F را مینیمم کنیم. یک شرط لازم برای اینکه F در نقطه (a_0, \dots, a_n) مینیمم باشد این است که

$$\frac{\partial F}{\partial a_i} = 0 \quad i = 0, 1, \dots, n \quad (12.3.4)$$

با مشتق‌گیری از انتگرال (۱۱.۳.۴) و استفاده از (۱۲.۳.۴)، دستگاه خطی زیر را به دست می‌آوریم

$$\sum_{j=0}^n a_j \int_a^b w(x) x^{i+j} dx = \int_a^b w(x) f(x) x^i dx \quad i = 0, 1, \dots, n \quad (13.3.4)$$

برای آنکه ببینیم چرا این جواب مسأله کمترین مربعات رضایتبخش نیست، حالت خاص $w(x) \equiv 1$ ، $[a, b] = [0, 1]$ را در نظر می‌گیریم. در این صورت دستگاه خطی چنین می‌شود

$$\sum_{j=0}^n \frac{a_j}{i+j+1} = \int_0^1 f(x) x^i dx \quad i = 0, 1, \dots, n \quad (14.3.4)$$

ماتریس ضرایب، ماتریس هیلبرت از مرتبه $n+1$ است که در (۹.۶.۱) معرفی شده است. جواب دستگاه خطی (۱۴.۳.۴) نسبت به تغییرات کوچک ضرایب یا ثابت‌های طرف راست، بسنار حساس است. پس این راه خوبی برای رسیدن به مسأله کمترین مربعات نیست. در حساب با دقت معمولی در رایانه IBM 3033، حالات $n \geq 4$ اصلاً رضایتبخش نخواهد بود.

۴.۴ چندجمله‌یهای متعامد

همان‌گونه که از نمودارهای x^n بر $[0, 1]$ ، $n \geq 0$ ، مشاهده می‌شود، این یک جمله‌یهای رفتاری خیلی نزدیک به رفتار خطی - وابسته دارند و خیلی شبیه یکدیگرند، و این امر علت ناپایداری دستگاه خطی (۱۴.۳.۴) است. برای اجتناب از این مسأله، پایه دیگری برای چندجمله‌یها در نظر می‌گیریم که بر چندجمله‌یهای متعامد در یک فضای تابعی که در زیر می‌آوریم مبتنی هستند. این نتایج در تقریب توابع، و در بسیاری از کارها در ریاضیات کاربردی، از اهمیت اساسی برخوردارند.

گیریم $w(x)$ همان $w(x)$ در (۸.۳.۴) و (۹.۳.۴) باشد، و حاصلضرب داخلی دو تابع پیوسته f و g را چنین تعریف می‌کنیم

$$(f, g) = \int_a^b w(x)f(x)g(x)dx \quad f, g \in C[a, b] \quad (۱.۴.۴)$$

در این صورت ویژگیهای ساده زیر به آسانی ثابت می‌شوند.

$$۱. (\alpha f, g) = (f, \alpha g) = \alpha(f, g) \quad \text{به ازای همه اسکالرهایی } \alpha$$

$$۲. (f_1 + f_2, g) = (f_1, g) + (f_2, g)$$

$$(f, g_1 + g_2) = (f, g_1) + (f, g_2)$$

$$۳. (f, g) = (g, f)$$

۴. به ازای هر $f \in C[a, b]$ ، $(f, f) \geq 0$ و $(f, f) = 0$ ، اگر و تنها اگر، برای $a \leq x \leq b$ ، $f(x) = 0$.

نرم دو یا نرم اقلیدسی را با رابطه زیر تعریف می‌کنیم

$$\|f\|_2 = \sqrt{\int_a^b w(x)[f(x)]^2 dx} = \sqrt{(f, f)} \quad (۲.۴.۴)$$

این تعریف، در ویژگیهای نرم (۶.۱.۴) - (۸.۱.۴) صدق می‌کند. ولی اثبات نابرابری مثلثی (۸.۱.۴) دیگر بدیهی نیست، و به نابرابری معروف زیر بستگی دارد.

لم (نابرابری کوشی - شوارتس) برای $f, g \in C[a, b]$

$$|(f, g)| \leq \|f\|_2 \|g\|_2 \quad (۳.۴.۴)$$

برهان اگر $g = 0$ ، درستی قضیه بدیهی است. حال $g \neq 0$ را در نظر می‌گیریم. برای هر عدد حقیقی α

$$0 \leq (f + \alpha g, f + \alpha g) = (f, f) + 2\alpha(f, g) + \alpha^2(g, g)$$

چند جمله‌ی طرف راست حداکثر یک ریشه حقیقی دارد، و بنابراین می‌تواند معادله نمی‌تواند مثبت باشد.

$$4|(f, g)|^2 - 4(f, f)(g, g) \leq 0$$

و از این رابطه (۳.۴.۴) نتیجه می‌شود. توجه داشته باشید که در (۳.۴.۴) تساوی تنها وقتی برقرار خواهد بود که مبین صفر باشد. ولی این امر ایجاب می‌کند که یک α^* وجود داشته باشد که به ازای آن چندجمله‌بیهی صفر شود. پس

$$(f + \alpha^*g, f + \alpha^*g) = 0$$

و بنابراین $f = -\alpha^*g$. بنابراین تساوی (۴.۴.۳) برقرار است اگر و تنها اگر (۱): یا f مضربی از g باشد، (۲): یا f یا g متحد با صفر باشد. ■

برای اثبات نابرابری مثلثی (۸.۱.۴)،

$$\begin{aligned} \|f + g\|_2^2 &= (f + g, f + g) = (f, f) + 2(f, g) + (g, g) \\ &\leq \|f\|_2^2 + 2\|f\|_2\|g\|_2 + \|g\|_2^2 = (\|f\|_2 + \|g\|_2)^2 \end{aligned}$$

از طرفین نامساوی جذر می‌گیریم، به دست می‌آید

$$\|f + g\|_2 \leq \|f\|_2 + \|g\|_2 \quad (۴.۴.۴)$$

می‌خواهیم برای چندجمله‌بیهای پایه‌ای غیر از پایه‌ی معمول تک‌جمله‌بیهای $\{1, x, \dots, x^n\}$ بیابیم. یک پایه‌ی متعامد می‌سازیم که تعمیم پایه‌ی متعامد در فضای R^n است (بخش ۱.۷ را ببینید). گوییم f و g متعامدند اگر

$$(f, g) = 0 \quad (۵.۴.۴)$$

قضیه‌ی زیر یک قضیه‌ی وجودی سازنده برای چندجمله‌بیهای متعامد است.

قضیه‌ی ۲.۴ (گرام-اشمیت)^۱ یک دنباله از چندجمله‌بیهای $\{\varphi_n(x) \mid n \geq 0\}$ با درجه‌ی $n = \varphi(n)$ ، به ازای همه‌ی مقادیر n ، وجود دارد به طوری که به ازای همه‌ی مقادیر m و n نا کوچکتر از صفر و $n \neq m$

$$(\varphi_n, \varphi_m) = 0 \quad (۶.۴.۴)$$

به علاوه ما می‌توانیم دنباله را با ویژگیهای اضافی زیر بسازیم

(۱) $(\varphi_n, \varphi_n) = 1$ به ازای هر n ; (۲) ضریب x^n در $\varphi_n(x)$ مثبت باشد. با این ویژگیهای اضافی، دنباله $\{\varphi_n\}$ یکتاست.

برهان یک روش سازنده و بازگشتی برای پیدا کردن عناصر این دنباله ارائه می‌دهیم (این روند، روند گرام-شمیت نامیده می‌شود). گیریم

$$\varphi_0(x) = c$$

c یک ثابت است. c را طوری می‌گیریم که $\|\varphi_0\|_2 = 1$ و $c > 0$. در این صورت

$$(\varphi_0, \varphi_0) = c^2 \int_a^b w(x) dx = 1$$

$$c = \left[\int_a^b w(x) dx \right]^{-1/2}$$

برای ساختن $\varphi_1(x)$ بارابطة زیر شروع می‌کنیم

$$\psi_1(x) = x + a_{1,0} \varphi_0(x)$$

در این صورت

$$(\psi_1, \varphi_0) = 0 \Rightarrow 0 = (x, \varphi_0) + a_{1,0} (\varphi_0, \varphi_0)$$

$$a_{1,0} = -(x, \varphi_0) = \frac{-\int_a^b x w(x) dx}{\left[\int_a^b w(x) dx \right]^{1/2}}$$

تعریف می‌کنیم

$$\varphi_1(x) = \frac{\psi_1(x)}{\|\psi_1\|_2}$$

و توجه داریم که

$$\|\varphi_1\|_2 = 1 \quad (\varphi_1, \varphi_0) = 0$$

و ضریب x مثبت است.

برای ساختن $\varphi_n(x)$ ابتدا تعریف می‌کنیم

$$\psi_n(x) = x^n + a_{n,n-1} \varphi_{n-1}(x) + \dots + a_{n,0} \varphi_0(x) \quad (7.4.4)$$

و ثابتها را طوری انتخاب کنیم که ψ_n بر φ_j ، به ازای $j = 0, \dots, n-1$ عمود شود. در این صورت از $(\psi_n, \varphi_j) = 0$ نتیجه می‌شود

$$a_{n,j} = -(x^n, \varphi_j) \quad j = 0, 1, \dots, n-1 \quad (8.4.4)$$

$\varphi_n(x)$ مورد نظر عبارت است از

$$\varphi_n(x) = \frac{\psi_n(x)}{\|\psi_n\|_2} \quad (9.4.4)$$

و بقیه را با استقرا ادامه می‌دهیم.

مثال برای حالت خاص $w(x) \equiv 1$ و $[a, b] = [-1, 1]$ داریم

$$\varphi_0(x) = \sqrt{\frac{1}{2}} \quad \varphi_1(x) = \sqrt{\frac{3}{2}}x \quad \varphi_2(x) = \sqrt{\frac{5}{2}}(3x^2 - 1)$$

و چندجمله‌بیهای بعدی را می‌توان با فرایند بالا ساخت.

نوشتارهای بسیار زیادی در باره این‌گونه چندجمله‌بیها از جمله فرمول‌های گوناگونی برای این چندجمله‌بیها موجود است. ما الزاماً یک نگاه اجمالی به آنها می‌اندازیم. چندجمله‌بیها معمولاً به شکلی داده شده‌اند که برای آنها $\|\varphi_n\|_2 \neq 1$.

حالتهای خاص

حالت ۱. چندجمله‌بیهای لژاندر. گیریم $w(x) \equiv 1$ در $[-1, 1]$. تعریف می‌کنیم

$$P_n(x) = \frac{(-1)^n}{2^n n!} \cdot \frac{d^n}{dx^n} [(1-x^2)^n] \quad n \geq 1 \quad (10.4.4)$$

با $P_0(x) \equiv 1$. این چندجمله‌بیها در $[-1, 1]$ متعامدند، درجه $P_n(x)$ برابر n است و به ازای همه مقادیر n ، $P_n(1) = 1$. همچنین

$$(P_n, P_n) = \frac{2}{2n+1}$$

$$\varphi_n(x) = \sqrt{\frac{2n+1}{2}} P_n(x) \quad (11.4.4)$$

حالت ۲. چندجمله‌بیهای چیشف. گیریم $w(x) = 1/\sqrt{1-x^2}$ ، $-1 \leq x \leq 1$. در این صورت

$$T_n(x) = \cos(n \cos^{-1} x) \quad n \geq 0 \quad (12.4.4)$$

یک خانواده متعامد از چندجمله‌بیهای با درجه (T_n) برابر n است. برای اینکه بینیم $T_n(x)$ یک چندجمله‌بیهی است، می‌گیریم $\theta = \cos^{-1} x$ ، $0 \leq \theta \leq \pi$. در این صورت

$$T_{n \pm 1}(x) = \cos(n \pm 1)\theta = \cos(n\theta) \cos \theta \mp \sin(n\theta) \sin \theta$$

$$T_{n+1}(x) + T_{n-1}(x) = 2 \cos(n\theta) \cos \theta = 2T_n(x)$$

$$T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x) \quad n \geq 1 \quad (13.4.4)$$

همچنین با محاسبه مستقیم در (۱۲.۴.۴)،

$$T_0(x) \equiv 1 \quad T_1(x) = x$$

با استفاده از (۱۳.۴.۴)

$$T_2(x) = 2x^2 - 1 \quad T_3(x) = 2x(2x^2 - 1) - x = 4x^3 - 3x$$

این چند جمله‌بیهای در شرایط: $T_n(1) = 1$ ، $n \geq 1$ و شرایط زیر نیز صدق می‌کنند

$$(T_n, T_m) = \begin{cases} 0 & n \neq m \\ \pi & n = m = 0 \\ \frac{\pi}{2} & n = m > 0 \end{cases} \quad (14.4.4)$$

چندجمله‌بیهای چبیشف در نظریه تقریب بسیار مهم‌اند، و در بسیاری از زمینه‌های دیگر ریاضیات کاربردی نیز پدید می‌آیند. برای یک بحث کاملتری در مورد آنها، ریولین^۱ (۱۹۷۴) و فاکس و پارکر^۲ (۱۹۶۸) را ببینید. ویژگیهای دیگر چندجمله‌بیهای چبیشف را در بخش بعد ارائه خواهیم داد. حالت ۳. چندجمله‌بیهای لاگر. بگیریم $w(x) = e^{-x}$ ، $[a, b] = [0, \infty)$. در این صورت

$$L_n(x) = \frac{1}{n!e^{-x}} \cdot \frac{d^n}{dx^n} \{x^n e^{-x}\} \quad n \geq 0 \quad (15.4.4)$$

به ازای جمیع مقادیر n ، $\|L_n\|_2 = 1$ ، و $\{L_n\}$ ها در $[0, \infty)$ با تابع وزن e^{-x} متعامدند. یک خانواده از توابع را متعامد گوییم اگر هر عضو خانواده بر تک تک اعضای دیگر خانواده عمود باشد. خانواده را خانواده یکا متعامد خوانیم اگر متعامد باشد و طول هر عضو خانواده یک باشد، یعنی، $\|f\|_2 = 1$. برای مثالهای دیگر از چندجمله‌بیهای متعامد، آبراموویس و اشتگون (۱۹۶۴، فصل ۲۲)، (دیویس (۱۹۶۳)، پیوست)، سگو^۳ (۱۹۶۸) را ببینید.

بعضی ویژگیهای چندجمله‌بیهای متعامد این قضایا در این فصل و برای فصلهای بعد سودمندند.

قضیه ۳.۴ گیریم $\{\varphi_n(x) \mid n \geq 0\}$ یک خانواده از چندجمله‌بیهای متعامد بر (a, b) با تابع وزن $w(x)$ باشد. با چنین خانواده‌ای، همیشه به طور ضمنی فرض می‌کنیم که درجه φ_n برابر n و $n \geq 0$. اگر $f(x)$ یک چندجمله‌بی از درجه m باشد، آنگاه

$$f(x) = \sum_{n=0}^m \frac{(f, \varphi_n)}{(\varphi_n, \varphi_n)} \varphi_n(x) \quad (۱۶.۴.۴)$$

برهان ابتدا نشان می‌دهیم که هر چندجمله‌بی را می‌توان به صورت ترکیبی خطی از چندجمله‌بیهای متعامد از درجه نایبتر نوشت. چون درجه $\varphi_0 = 0$ ، داریم $\varphi_0(x) = c$ ، یک ثابت، و بنابراین

$$1 \equiv \frac{1}{c} \varphi_0(x)$$

چون درجه $(\varphi_1) = 1$ ، از فرایند گرام-شمیت داریم،

$$\varphi_1(x) = c_{1,1}x + c_{1,0}\varphi_0(x) \quad c_{1,1} \neq 0$$

$$x = \frac{1}{c_{1,1}}[\varphi_1(x) - c_{1,0}\varphi_0(x)]$$

با استقرا در فرایند گرام-شمیت

$$\varphi_r(x) = c_{r,r}x^r + c_{r,r-1}\varphi_{r-1}(x) + \dots + c_{r,0}\varphi_0(x) \quad c_{r,r} \neq 0$$

و

$$x^r = \frac{1}{c_{r,r}}[\varphi_r(x) - c_{r,r-1}\varphi_{r-1}(x) - \dots - c_{r,0}\varphi_0(x)]$$

پس هر یک جمله‌بی را می‌توان به صورت ترکیبی از چندجمله‌بیهای متعامد درجه نایبتر بیان کرد. از این امر به سادگی نتیجه می‌شود که یک چندجمله‌بی دلخواه $f(x)$ از درجه m را می‌توان به ازای انتخاب مقادیری مانند b_m, \dots, b_0 به شکل زیر نوشت.

$$f(x) = b_m\varphi_m(x) + \dots + b_0\varphi_0(x)$$

برای محاسبه هر b_i ، طرفین را در $w(x)$ و $\varphi_i(x)$ ضرب نموده روی (a, b) انتگرال می‌گیریم. لذا

$$(f, \varphi_i) = \sum_{j=0}^m b_j(\varphi_j, \varphi_i) = b_i(\varphi_i, \varphi_i)$$

$$b_i = \frac{(f, \varphi_i)}{(\varphi_i, \varphi_i)}$$

که (۱۶.۴.۴) و قضیه را اثبات می‌کند.

$$\int_a^b w(x)B(x)\varphi_n(x)dx \neq 0$$

زیرا روشن است که $B(x)\varphi_n(x) \neq 0$. ولی چون m ، درجه B از n کوچکتر است، از فرع قضیه ۳.۴ نتیجه می‌شود که

$$\int_a^b w(x)B(x)\varphi_n(x)dx = (B, \varphi_n) = 0$$

و این یک تناقض است، پس باید داشته باشیم $m = n$. پس حکم قضیه حاصل می‌شود؛ زیرا $\varphi_n(x)$ حداکثر n ریشه دارد و فرضهای مربوط به x_1, \dots, x_n ایجاب می‌کنند که همه آنها ساده باشند، یعنی $\varphi'_n(x_i) \neq 0$.

مانند قبل فرض می‌کنیم $\{\varphi_n(x) \mid n \geq 0\}$ خانواده‌ای متعامد بر (a, b) با تابع وزن $w(x) \geq 0$ باشد. A_n و B_n را با

$$\varphi_n(x) = A_n x^n + B_n x^{n-1} + \dots \quad (18.4.4)$$

تعریف می‌کنیم. همچنین می‌نویسیم

$$\varphi_n(x) = A_n(x - x_{n,1})(x - x_{n,2}) \dots (x - x_{n,n}) \quad (19.4.4)$$

گیریم

$$a_n = \frac{A_{n+1}}{A_n} \quad \gamma_n = (\varphi_n, \varphi_n) > 0 \quad (20.4.4)$$

قضیه ۵.۴ (رابطه بازگشتی سه‌گانه) گیریم $\{\varphi_n\}$ یک خانواده از چندجمله‌بیهای متعامد بر (a, b) با تابع وزن $w(x) \geq 0$ باشد. در این صورت به ازای $n \geq 1$

$$\varphi_{n+1}(x) = (a_n x + b_n)\varphi_n(x) - c_n \varphi_{n-1}(x) \quad (21.4.4)$$

که در آن

$$b_n = a_n \cdot \left[\frac{B_{n+1}}{A_{n+1}} - \frac{B_n}{A_n} \right] \quad c_n = \frac{A_{n+1}A_{n-1}}{A_n^2} \cdot \frac{\gamma_n}{\gamma_{n-1}} \quad (22.4.4)$$

برهان ابتدا توجه نماید که رابطه بازگشتی سه‌گانه (۱۳.۴.۴) برای چندجمله‌بیهای چیشف یک مثال از (۲۱.۴.۴) است. برای به دست آوردن (۲۱.۴.۴)، اثبات را با در نظر گرفتن چندجمله‌بی

زیر شروع می‌کنیم:

$$\begin{aligned} G(x) &= \varphi_{n+1}(x) - a_n x \varphi_n(x) \\ &= [A_{n+1} x^{n+1} + B_{n+1} x^n + \dots] \\ &\quad - \frac{A_{n+1}}{A_n} x [A_n x^n + B_n x^{n-1} + \dots] \\ &= \left[B_{n+1} - \frac{A_{n+1} B_n}{A_n} \right] x^n + \dots \end{aligned}$$

و درجهٔ G از n نابزرگتر است. به موجب قضیهٔ ۴.۳ برای مجموعهٔ مناسبی از d_n, \dots, d_0 می‌توانیم بنویسیم

$$G(x) = d_n \varphi_n(x) + \dots + d_0 \varphi_0(x)$$

d_i را حساب می‌کنیم

$$d_i = \frac{(G, \varphi_i)}{(\varphi_i, \varphi_i)} = \frac{1}{\gamma_i} [(\varphi_{n+1}, \varphi_i) - a_n (x \varphi_n, \varphi_i)] \quad (23.4.4)$$

به ازای $i \leq n$ داریم $(\varphi_{n+1}, \varphi_i) = 0$ و به ازای $i \leq n-2$ داریم

$$(x \varphi_n, \varphi_i) = \int_a^b w(x) \varphi_n(x) x \varphi_i(x) dx = 0$$

زیرا در این صورت درجهٔ $1 - n \leq (x \varphi_i(x))$. از ترکیب این نتایج

$$d_i = 0 \quad 0 \leq i \leq n-2$$

و بنابراین

$$G(x) = d_n \varphi_n(x) + d_{n-1} \varphi_{n-1}(x)$$

$$\varphi_{n+1}(x) = (a_n x + d_n) \varphi_n(x) + d_{n-1} \varphi_{n-1}(x) \quad (24.4.4)$$

این وجود یک رابطهٔ بازگشتی سه‌گانه را نشان می‌دهد و بقیهٔ کار دستکاری فرمولها برای به‌دست آوردن d_n و d_{n-1} است که مانند b_n و c_n در (۲۲.۴.۴) داده شده‌اند. این ثابتها را در اینجا به‌دست نیاورده‌ایم ولی مقادیر آنها برای بعضی کاربردها مهم‌اند (مسألهٔ ۱۸ را ببینید). ■

مثال ۱. برای چندجمله‌بیهای لاگر

$$L_{n+1}(x) = \frac{1}{n+1} [2n+1-x] L_n(x) - \frac{n}{n+1} L_{n-1}(x) \quad (25.4.4)$$

۲. برای چندجمله‌یهای لژاندر

$$P_{n+1}(x) = \frac{2n+1}{n+1}xP_n(x) - \frac{n}{n+1}P_{n-1}(x) \quad (26.4.4)$$

قضیه ۶.۴ (اتحاد کریستوفل - داربو) برای یک خانواده از چندجمله‌یهای متعامد $\{\varphi_n\}$ با تابع وزن $w(x) \geq 0$

$$\sum_{k=0}^n \frac{\varphi_k(x)\varphi_k(y)}{\gamma_k} = \frac{\varphi_{n+1}(x)\varphi_n(y) - \varphi_n(x)\varphi_{n+1}(y)}{a_n\gamma_n(x-y)}, x \neq y \quad (27.4.4)$$

برهان برهان براساس عملیات روی رابطه بازگشتی سه‌گانه صورت می‌گیرد سگو (۱۹۶۷، ص ۴۳) را ببینید.

۵.۴ مسأله تقریب کمترین مربعات (ادامه)

اکنون به مسأله کلی کمترین مربعات، مینیمم‌سازی (۱۰.۳.۴) در بین همه چندجمله‌یهای از درجه نایزگتر از n برمی‌گردیم. فرض می‌کنیم $\{\varphi_k(x) \mid k \geq 0\}$ یک خانواده از چندجمله‌یهای متعامد با تابع وزن $w(x) \geq 0$ باشد، یعنی

$$(\varphi_n, \varphi_m) = \delta_{n,m} = \begin{cases} 1 & n = m \\ 0 & n \neq m \end{cases}$$

پس یک چندجمله‌ی دلخواه $f(x)$ از درجه نایزگتر از n را می‌توان چنین نوشت

$$r(x) = b_0\varphi_0(x) + \dots + b_n\varphi_n(x) \quad (1.5.4)$$

برای تابع داده شده $f \in C[a, b]$

$$\|f - r\|_2^2 = \int_a^b w(x) \left[f(x) - \sum_{j=0}^n b_j\varphi_j(x) \right]^2 dx \equiv G(b_0, \dots, b_n) \quad (2.5.4)$$

مسأله کمترین مربعات را با مینیمم‌سازی G حل می‌کنیم. مانند قبل، می‌توانیم قرار دهیم

$$\frac{\partial G}{\partial b_i} = 0 \quad i = 0, 1, \dots, n$$

ولی برای به دست آوردن یک نتیجه کاملتر، به راه دیگری می پردازیم. برای هر انتخاب b_0, \dots, b_n ,

$$\begin{aligned} \circ \leq G(b_0, \dots, b_n) &= \left(f - \sum_{j=0}^n b_j \varphi_j, f - \sum_{i=0}^n b_i \varphi_i \right) \\ &= (f, f) - 2 \sum_{j=0}^n b_j (f, \varphi_j) + \sum_i \sum_j b_i b_j (\varphi_i, \varphi_j) \\ &= \|f\|_V^2 - 2 \sum_{j=0}^n b_j (f, \varphi_j) + \sum_{j=0}^n b_j^2 \\ &= \|f\|_V^2 - \sum_{j=0}^n (f, \varphi_j)^2 + \sum_{j=0}^n \left[(f, \varphi_j) - b_j \right]^2 \quad (۳.۵.۴) \end{aligned}$$

که می توان درستی آن را با بسط آخرین جمله بررسی کرد. بنابراین G مینیمم است اگر و تنها اگر

$$b_j = (f, \varphi_j) \quad j = 0, 1, \dots, n$$

پس تقریب کمترین مربعات وجود دارد، یکتاست و با رابطه زیر داده می شود

$$r_n^*(x) = \sum_{j=0}^n (f, \varphi_j) \varphi_j(x)$$

علاوه بر این، از (۲.۵.۴) و (۳.۵.۴)

$$\begin{aligned} \|f - r_n^*\|_V &= \left[\|f\|_V^2 - \sum_{j=0}^n (f, \varphi_j)^2 \right]^{1/2} \\ &= \sqrt{\|f\|_V^2 - \|r_n^*\|_V^2} \quad (۴.۵.۴) \end{aligned}$$

$$\|f\|_V^2 = \|r_n^*\|_V^2 + \|f - r_n^*\|_V^2 \quad (۵.۵.۴)$$

به عنوان یک توضیح مفید برای به دست آوردن $r_{n+1}^*(x)$ از رابطه زیر استفاده کنید

$$r_{n+1}^*(x) = r_n^*(x) + (f, \varphi_{n+1}) \varphi_{n+1}(x) \quad (۶.۵.۴)$$

قضیه ۷.۴ فرض می کنیم $[a, b]$ متناهی باشد،

$$\lim_{n \rightarrow \infty} \|f - r_n^*\|_V = 0 \quad (۷.۵.۴)$$

برهان از تعریف r_n^* به عنوان یک چندجمله‌ای مینیم‌ساز برای $\|f - r_n\|_2$ داریم،

$$\|f - r_n^*\|_2 \geq \|f - r_{n-1}^*\|_2 \geq \dots \geq \|f - r_m^*\|_2 \geq \dots \quad (۸.۵.۴)$$

گیریم $\varepsilon > 0$ دلخواه باشد. پس طبق قضیهٔ وایرستراس، یک چندجمله‌ای $Q(x)$ از درجهٔ m وجود دارد که برای آن

$$\text{Max}_{a \leq x \leq b} |f(x) - Q(x)| \leq \frac{\varepsilon}{c} \quad c = \sqrt{\int_a^b w(x) dx}$$

بنابه تعریف $r_m^*(x)$

$$\begin{aligned} \|f - r_m^*\|_2 &\leq \|f - Q\|_2 = \left[\int_a^b w(x) [f(x) - Q(x)]^2 dx \right]^{1/2} \\ &\leq \left[\int_a^b w(x) \frac{\varepsilon^2}{c^2} dx \right]^{1/2} = \varepsilon \end{aligned}$$

از ترکیب این رابطه با (۸.۵.۴)، به ازای جمیع مقادیر $n > m$

$$\|f - r_n^*\|_2 \leq \varepsilon$$

چون ε دلخواه بود، (۷.۵.۴) اثبات می‌شود. ■

با استفاده از (۵.۵.۴) و یک محاسبهٔ ساده $\|r_n^*\|_2$ ، نامساوی بسط را خواهیم داشت:

$$\|r_n^*\|_2^2 = \sum_{j=0}^n (f, \varphi_j)^2 \leq \|f\|_2^2 \quad (۹.۵.۴)$$

و با استفاده از (۷.۵.۴) در (۵.۵.۴) تساوی پارسوال را به دست می‌آوریم:

$$\|f\|_2 = \left[\sum_{j=0}^{\infty} (f, \varphi_j)^2 \right]^{1/2} \quad (۱۰.۵.۴)$$

قضیهٔ ۷.۴ نمی‌گوید که $\|f - r_n^*\|_{\infty} \rightarrow 0$ ولی اگر فرضهای مشتق‌پذیری بیشتری برای $f(x)$ در نظر گرفته شود، قضایای مربوط به همگرایی یکنواخت r_n^* به f را می‌توان ثابت کرد. یک مثال بعداً داده شده است.

بسطهای چندجمله‌یی لژاندر برای حل مسألهٔ کمترین مربعات بر یک بازهٔ متناهی $[a, b]$ با تابع وزن $w(x) \equiv 1$ می‌توانیم آن را به یک مسأله بر $[-1, 1]$ برگردانیم. تبدیل متغیر

$$x = \frac{b+a+(b-a)t}{2} \quad (11.5.4)$$

بازهٔ $-1 \leq t \leq 1$ را به بازهٔ $a \leq x \leq b$ بدل می‌کند. برای یک تابع داده شدهٔ $f \in C[a, b]$ ، تعریف می‌کنیم

$$F(t) = f\left(\frac{b+a+(b-a)t}{2}\right) \quad -1 \leq t \leq 1 \quad (12.5.4)$$

پس

$$\int_a^b [f(x) - r_n(x)]^2 dx = \left(\frac{b-a}{2}\right) \int_{-1}^1 [F(t) - R_n(t)]^2 dt$$

که $R_n(t)$ از $r_n(x)$ با استفاده از (11.5.4) به دست می‌آید. تغییر متغیر (11.5.4) یک تناظر یک به یک بین چندجمله‌یهای از درجهٔ m بر $[a, b]$ و از درجهٔ m بر $[-1, 1]$ ، برای هر $m > 0$ ، به دست می‌دهد. بنابراین مینیمم کردن $\|f - r_n\|_2$ بر $[a, b]$ هم‌ارز است با مینیمم کردن $\|F - R_n\|_2$ بر $[-1, 1]$. لذا ما توجه خود را به مسألهٔ کمترین مربعات بر $[-1, 1]$ محدود می‌کنیم.

برای تابع مفروض $f \in [-1, 1]$ ، خانوادهٔ یکا متعامد مذکور در قضیهٔ ۲.۴ عبارت است از

$$\varphi_0(x) \equiv \frac{1}{\sqrt{2}}$$

$$\varphi_n(x) = \sqrt{\frac{2n+1}{2}} \cdot \frac{(-1)^n}{2^n n!} \cdot \frac{d^n}{dx^n} [(1-x^2)^n] \quad n \geq 1 \quad (13.5.4)$$

تقریب کمترین مربعات چنین است

$$r_n^*(x) = \sum_{j=0}^n (f, \varphi_j) \varphi_j(x) \quad (f, \varphi_j) = \int_{-1}^1 f(x) \varphi_j(x) dx \quad (14.5.4)$$

که جواب مسألهٔ کلی کمترین مربعات است که در بخش ۳.۴ مطرح شده است. ضرایب (f, φ_j) را ضرایب لژاندر می‌خوانند.

مثال ضرایب بسط (f, φ_j) در (14.5.4) در جدول ۲.۴ برای $f(x) = e^x$ بر $[-1, 1]$ داده شده‌اند. تقریب $r_3^*(x)$ که به شکل چندجمله‌یی استاندارد نوشته شده بود، قبلاً در (۷.۳.۴) داده

جدول ۲.۴ ضرایب بسط لژاندر برای e^x

j	(f, φ_j)	j	(f, φ_j)
۰	۱٫۶۶۱۹۸۵	۳	۰٫۰۳۷۶۶۰
۱	۰٫۹۰۱۱۱۷	۴	۰٫۰۰۴۶۹۸
۲	۰٫۲۲۶۳۰۲	۵	۰٫۰۰۰۴۶۹

شده است. خطای متوسط E در $r_n^*(x)$ را از ترکیب روابط $(۴.۳.۴)$ ، $(۴.۵.۴)$ ، و از جدول ضرایب به دست می آوریم

$$E = \frac{1}{\sqrt{2}} \| e^x - r_n^*(x) \|_2 \doteq ۰٫۰۰۳۴$$

بسطهای چندجمله‌یی چیشف در اینجا $w(x) = 1/\sqrt{1-x^2}$ تابع وزن است و

$$\varphi_0(x) = \frac{1}{\sqrt{\pi}} \quad \varphi_n(x) = \sqrt{\frac{2}{\pi}} T_n(x) \quad n \geq 1 \quad (۱۵.۵.۴)$$

جواب کمترین مربعات چنین است

$$C_n(x) = \sum_{j=0}^n (f, \varphi_j) \varphi_j(x) \quad (f, \varphi_j) = \int_{-1}^1 \frac{f(x) \varphi_j(x)}{\sqrt{1-x^2}} dx \quad (۱۶.۵.۴)$$

با استفاده از تعریف $\varphi_n(x)$ بر حسب $T_n(x)$ ،

$$C_n(x) = \sum_{j=0}^n c_j T_j(x) \quad c_j = \frac{2}{\pi} \int_{-1}^1 \frac{f(x) T_j(x)}{\sqrt{1-x^2}} dx \quad (۱۷.۵.۴)$$

علامت پریم برنماد مجموعیایی بدین معناست که قبل از مجموعیایی، اولین جمله باید نصف شود. بسط چیشف ارتباط نزدیکی به بسط کسینوس فوریه دارد. با استفاده از $x = \cos \theta$ ، $0 \leq \theta \leq \pi$

$$C_n(\cos \theta) = \sum_{j=0}^n c_j \cos(j\theta) \quad (۱۸.۵.۴)$$

$$c_j = \frac{2}{\pi} \int_0^\pi \cos(j\theta) f(\cos \theta) d\theta \quad (۱۹.۵.۴)$$

بنابراین $C_n(\cos \theta)$ برش براز $n + 1$ جمله اول بسط کسینوس فوریه

$$f(\cos \theta) = \sum_{j=0}^{\infty} c_j \cos(j\theta)$$

است. اگر بسط کسینوس فوریه $f(\cos \theta)$ بر $[\pi, 0]$ معلوم باشد، آنگاه با قراردادن $x = \cos^{-1} \theta$ بسط چیشف $f(x)$ به دست می آید.

به دلایلی که بعداً ارائه خواهد شد، تقریب کمترین مربعات چیشف مفیدتر از تقریب کمترین مربعات لژاندر است. بدین دلیل ما قضیه همگرایی برای (۱۷.۵.۴) را با تفصیل بیشتری بیان خواهیم کرد.

قضیه ۸.۴ گیریم $f(x)$ دارای r مشتق پیوسته بر $[-1, 1]$ باشد، با $r \geq 1$. در این صورت به ازای یک ثابت B که به f و r بستگی دارد، برای $C_n(x)$ ، تقریب کمترین مربعات چیشف که در (۱۷.۵.۴) تعریف شده است داریم

$$\|f - C_n\|_{\infty} \leq \frac{B \ln n}{n^r} \quad n \geq 2 \quad (20.5.4)$$

بنابراین با $n \rightarrow \infty$ ، $C_n(x)$ به طور یکنواخت به $f(x)$ میل می کند به شرطی که $f(x)$ پیوسته مشتق پذیر باشد.

برهان برای اثبات، ریولین (۱۹۷۴، قضیه ۳.۳ صفحه ۱۳۴) و میناردوس^۱ (۱۹۷۶، قضیه ۴۵ صفحه ۵۷) را ترکیب نماید.

مثال به عنوان مثال برای بسط چیشف دوباره تقریبهای $f(x) = e^x$ را در نظر می گیریم. برای e^x ضرایب c_j از (۱۷.۵.۴) برابرند با

$$\begin{aligned} c_j &= \frac{2}{\pi} \int_{-1}^1 \frac{e^x T_j(x)}{\sqrt{1-x^2}} dx \\ &= \frac{2}{\pi} \int_0^{\pi} e^{\cos \theta} \cdot \cos(j\theta) d\theta \end{aligned} \quad (21.5.4)$$

فرمول اخیر برای انتگرالگیری عددی بهتر است، زیرا تابع انتگرالده نقطه تکین ندارد. قاعده میانگاهی (نقطه میانی) یا ذوزنقه‌یی به علت دوره‌یی بودن تابع زیر علامت انتگرال روش بسیار عالی خواهد بود (فرع ۱، قضیه ۵.۵، در بخش ۴.۵ را ببینید). با استفاده از انتگرالگیری عددی مقادیر جدول ۳.۴ را به دست می آوریم. با به کار بردن (۱۷.۵.۴) و فرمول $T_j(x)$ خواهیم داشت

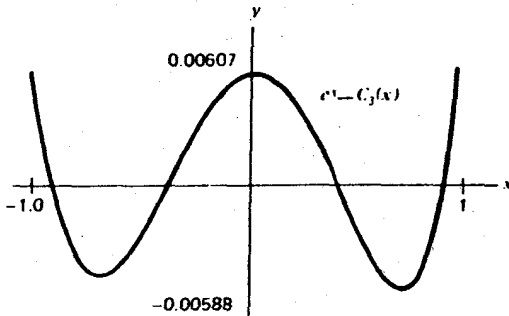
$$C_1(x) = 1.266 + 1.130x$$

$$C_2(x) = 0.994571 + 0.997308x + 0.542991x^2 + 0.177347x^3$$

$$\|e^x - C_1(x)\|_{\infty} = 0.32 \quad \|e^x - C_2(x)\|_{\infty} = 0.00607 \quad (22.5.4)$$

جدول ۳.۴ ضرایب بسط چیشف برای e^x

j	C_j
۰	۲٫۵۳۲۱۳۱۷۶
۱	۱٫۱۳۰۳۱۸۲۱
۲	۰٫۲۷۱۴۹۵۳۴
۳	۰٫۰۴۴۳۳۶۸۵
۴	۰٫۰۰۵۴۷۴۲۴
۵	۰٫۰۰۰۵۴۲۹۳



شکل ۶.۴ خطا در تقریب کمترین مربعات درجه سوم چیشف برای e^x

نمودار $e^x - C_5(x)$ در شکل ۶.۴ داده شده است، و خیلی شبیه به خطای مینیمکس در شکل ۴.۴ است. خطاهای ماکسیمم در این تقریبهای کمترین مربعات چیشف خیلی نزدیک به خطاهای مینیمکس اند و برای اغلب منظوره‌های عملی کفایت می‌کنند.

چندجمله‌یی $C_n(x)$ به شکل (۱۷.۵.۴) را می‌توان به صورت معمول آن، یعنی بر پایه یک جمله‌بندیهای x^j ، درآورد و همان‌گونه که قبلاً در مثال (۲۲.۵.۴) صورت گرفت، آن را محاسبه کرد؛ ولی $C_n(x)$ به شکل (۱۷.۵.۴) سریعتر و دقیقتر محاسبه می‌شود. برای چندجمله‌بندیهای $T_n(x)$ چیشف از رابطه بازگشتی سه‌گانه (۱۳.۴.۴) استفاده می‌کنیم. الگوریتم زیر منسوب به کِلنشاوا است و ما آن را از صفحه ۱۲۵ کتاب ریولین (۱۹۷۴) اخذ کرده‌ایم.

الگوریتم چبوال (مقدار $\text{chebeval}(x, n, a)$)

۱. توجه: با این الگوریتم مقدار $\sum_{j=0}^n a_j T_j(x)$ محاسبه می‌شود

$$z = 2x \quad b_{n+1} = b_{n+2} := 0 \quad ۲.$$

۳. تا مرحله ۵ عمل را برای $j = n, n-1, \dots, 0$ انجام دهید.

$$b_j = zb_{j+1} - b_{j+2} + a_j \quad ۴.$$

۵. z_j بعدی

$$۶. \quad := (b_0 - b_2) / 2 \quad \text{مقدار}$$

این الگوریتم تقریباً به همان اندازه الگوریتم ضرب تودرتو (۸.۹.۲) از بخش ۲.۹ کارایی دارد. مقایسه مشروحتر را به مسأله ۲۵ احاله می‌کنیم. الگوریتمهای مشابهی برای سایر بسطهای چندجمله‌یهای متعامد وجود دارند، که باز هم از رابطه بازگشتی سه‌گانه متناظر استفاده می‌نمایند. برای بررسی خطاهای گردکردن در روند چپوال به (صفحه ۱۲۷) ریولین (۱۹۷۴) و صفحه ۵۷ فکس و پارکر (۱۹۶۸) مراجعه کنید.

۶.۴ تقریبهای مینیماکس

برای یک تقریب یکنواخت خوب تابع داده شده $f(x)$ ، منطقی به نظر می‌آید که خطا در بازه تقریب یکنواخت توزیع شده باشد. به علاوه مثالهای قبلی با $f(x) = e^x$ این نکته را روشن می‌سازند و نشان می‌دهند که خطای ماکسیمم از لحاظ علامت نوسان خواهد کرد. در جدول ۴.۴ آماري از شکلهای گوناگون تقریب e^x در $[-1, 1]$ خلاصه شده است که شامل بعضی از روشهای بخش ۷.۴ نیز می‌باشد. برای نشان دادن اهمیت توزیع یکنواخت خطا با تابع خطایی که علامت آن نوسان می‌نماید، دو قضیه بیان می‌کنیم. اولین قضیه برای برآورد $\rho_n(x)$ ، خطای مینیماکس، بدون اینکه مجبور باشیم تقریب مینیماکس $q_n^*(x)$ را پیدا کنیم مفید است.

جدول ۴.۴ مقایسه تقریبهای گوناگون خطی و درجه سوم برای e^x

خطای ماکسیمم		روش تقریب
درجه سوم	خطی	
۰.۰۵۱۶	۰.۷۱۸	$\rho_n(x)$ چندجمله‌ی تیلر
۰.۱۱۲	۰.۴۳۹	$r_n^*(x)$ کمترین مربعات لژاندر
۰.۰۰۶۰۷	۰.۳۲۲	$C_n(x)$ ، کمترین مربعات چیشیف
۰.۰۰۶۶۶	۰.۳۷۲	$I_n(x)$ ، فرمول درونیابی با نقاط گرهی چیشیف
۰.۰۰۵۵۸	۰.۲۸۶	$F_n(x)$ ، فرمول نوسانی اجباری چیشیف
۰.۰۰۵۵۳	۰.۲۷۹	$q_n^*(x)$ مینیماکس

قضیه ۹.۴ (دولاوله پوسن)^۱ گیریم $f \in C[a, b]$ و $n \geq 0$. فرض کنید یک چندجمله‌یی $Q(x)$ از درجه نایزگتر از n داریم که در رابطه زیر صدق می‌کند

$$f(x_j) - Q(x_j) = (-1)^j e_j \quad j = 0, 1, \dots, n+1 \quad (۱.۶.۴)$$

که تمام e_j ها مخالف صفر و دارای یک علامت‌اند و

$$a \leq x_0 < x_1 < \dots < x_{n+1} \leq b$$

در این صورت

$$\text{Min}_{0 \leq j \leq n+1} |e_j| \leq \rho_n(f) \equiv \|f - q_n^*\|_\infty \leq \|f - Q\|_\infty \quad (۲.۶.۴)$$

برهان کران بالای (۲.۶.۴) از تعریف $\rho_n(f)$ به دست می‌آید. برای اثبات کران پایین، خلاف آن را فرض می‌کنیم و به تناقض می‌رسیم.
فرض کنید

$$\rho_n(f) < \text{Min}_{0 \leq j \leq n+1} |e_j| \quad (۳.۶.۴)$$

لذا طبق تعریف $\rho_n(f)$ ، یک چندجمله‌یی $\rho(x)$ با درجه نایزگتر از n وجود دارد که برای آن

$$\rho_n(f) \leq \|f - P\|_\infty < \text{Min} |e_j| \quad (۴.۶.۴)$$

یک چندجمله‌یی با درجه نایزگتر از n به شکل زیر تعریف می‌کنیم

$$R(x) = Q(x) - P(x)$$

برای سادگی، همه e_j ها را مثبت می‌گیریم؛ استدلال مشابهی برای همه e های منفی نیز برقرار است. برای هر یک از z_j ها، $R(x_j)$ را محاسبه کرده و علامت $R(x_j)$ را در نظر می‌گیریم. ابتدا با استفاده از (۴.۶.۴) داریم؛

$$\begin{aligned} R(x_0) &= Q(x_0) - P(x_0) = [f(x_0) - P(x_0)] - [f(x_0) - Q(x_0)] \\ &= [f(x_0) - P(x_0)] - e_0 < 0 \end{aligned}$$

سپس

$$R(x_1) = Q(x_1) - P(x_1) = [f(x_1) - P(x_1)] + e_1 > 0$$

جدول ۵.۴ ماکسیم‌های نسبی $|e^x - C_2(x)|$

x	$e^x - C_2(x)$
-1.0	0.00497
-0.6919	-0.00511
0.310	0.00547
0.7229	-0.00588
1.0	0.00607

به استقرأ، علامت $R(x_j)$ برابر است با $(-1)^{j+1}$ ، $j = 0, 1, \dots, n+1$. پس $R(x)$ ، $(n+2)$ بار تغییر علامت می‌دهد و در نتیجه $R(x)$ دارای $n+1$ ریشه است. چون درجه $R(x)$ نابزرگتر از n است، این ممکن نخواهد بود مگر آنکه $R(x) \equiv 0$. پس $P \equiv Q$ که با (۱.۶.۴) و (۴.۶.۴) متناقض است. ■

مثال تقریب کمترین مربعات درجه سوم چیبیشف برای e^x بر $[-1, 1]$ را که در (۲۲.۵.۴) با $C_n(x)$ نشان داده شده بود یادآور می‌شویم. ماکسیم خطای آن در بازه $[-1, 1]$ در جدول ۵.۴ داده شده است. این خطاها در مفروضات قضیه ۹.۴ صدق می‌نمایند، و بنابراین

$$0.00497 \leq \rho_2(f) \leq 0.00607$$

مطابق این جدول می‌توانیم نتیجه‌گیری کنیم که $C_2(x)$ به بهترین تقریب ممکن خیلی نزدیک بوده است. توجه داشته باشیم که با محاسبه دقیق و استفاده از (۷.۲.۴)، خواهیم داشت $\rho_2(f) = 0.00553$.

قضیه ۱۰.۴ (قضیه هم‌نوسانی چیبیشف)

گیریم $f \in C[a, b]$ و $n \geq 0$. آنگاه یک چندجمله‌یی یکتای $q_n^*(x)$ از درجه نابزرگتر از n وجود دارد که برای آن

$$\rho_n(f) = \|f - q_n^*\|_\infty$$

این چندجمله‌یی توسط ویژگی زیر به طور یکتا توصیف می‌شود: حداقل $n+2$ نقطه

$$a \leq x_0 < x_1 < \dots < x_{n+1} \leq b$$

وجود دارند که برای آنها

$$f(x_j) - q_n^*(x_j) = \sigma(-1)^j \rho_n(f) \quad j = 0, 1, \dots, n+1 \quad (۵.۶.۴)$$

σ برابر ± 1 است و فقط بستگی به f و n دارد.

برهان برهان کاملاً فنی است و به صورتی پیچیده از راه برهان خلف ثابت می شود. به همین دلیل در اینجا حذف شده است. برای یک بحث کامل (فصل ۷ دیویس^۱ (۱۹۶۳) را ببینید). ■

مثال تقریب مینیمکس درجه سوم $q_3^*(x)$ برای e^x که در (۷.۲.۴) داده شده بود، در احکام این قضیه صدق می نماید، چنانچه می توان از نمودار خطا در شکل ۴.۴ بخش ۴.۲ ملاحظه کرد.

با توجه به این قضیه می توانیم ببینیم که سری تیلر همیشه یک تقریب یکنواخت ضعیف است. خطای سری تیلر

$$f(x) - p_n(x) = \frac{(x - x_0)^{n+1}}{(n+1)!} f^{n+1}(\xi_x) \quad (۶.۶.۴)$$

در بازه تقریب به طور یکنواخت تغییر نمی کند و علامت آن هم خیلی نوسان نمی کند. برای اینکه بهتر نشان دهیم که $q_n^*(x)$ چگونه، با افزایش n ، $f(x)$ را به خوبی تقریب می زند قضیه زیر از د. جکسن^۲ را می آوریم.

قضیه ۱۱.۴ (جکسن) گیریم $f(x)$ به ازای مقداری از $k \geq 0$ دارای k مشتق پیوسته باشد. به علاوه فرض می کنیم که $f^{(k)}(x)$ ، به ازای مقداری مانند $M > 0$ و $0 < \alpha \leq 1$ ، در رابطه زیر صدق می کند

$$\sup_{a \leq x, y \leq b} |f^{(k)}(x) - f^{(k)}(y)| \leq M |x - y|^\alpha \quad (۷.۶.۴)$$

[می گوئیم $f^{(k)}(x)$ با نمای α در شرط هولدر صدق می کند]. در این صورت ثابتی چون d_k ، مستقل از f و n موجود است که برای آن

$$\rho_n(f) \leq \frac{M d_k}{n^{k+\alpha}} \quad n \geq 1 \quad (۸.۶.۴)$$

برهان میناردوس (۱۹۶۷، قضیه ۴۵ ص ۵۷) را ببینید. توجه نمایید که اگر بخواهیم (۷.۶.۴) و فرض k بار پیوسته مشتق پذیری تابع $f(x)$ را نادیده بگیریم، آنگاه فقط $k - 1$ را به جای k در

جدول ۶.۴ مقایسه M_n و ρ_n برای $f(x) = e^x$

n	۲	۳	۴	۵	۶	۷
$M_n(f)$	$۱٫۱۳E-۱$	$۱٫۴۲E-۲$	$۱٫۴۲E-۳$	$۱٫۱۸E-۴$	$۸٫۴۳E-۶$	$۵٫۲۷E-۷$
$\rho_n(f)$	$۴٫۵۰E-۲$	$۵٫۵۳E-۳$	$۵٫۴۷E-۴$	$۴٫۵۲E-۵$	$۳٫۲۱E-۶$	$۲٫۰۰E-۷$

قضیه به کار می‌بریم با $\alpha = 1$ و $\|f^{(k)}\|_\infty = M$. در این صورت حکم قضیه چنین می‌شود

$$\rho_n(f) \leq \frac{d_{k-1}}{n^k} \|f^{(k)}\|_\infty \quad (۹.۶.۴)$$

همچنین باید توجه کرد که اگر f بینهایت بار مشتقپذیر باشد، آنگاه $q_n^*(x)$ ، سریعتر از هر توانی از $\frac{1}{n^k}$ ، $k \geq 1$ ، در بازه $[a, b]$ به طور یکنواخت به $f(x)$ می‌گراید. ■

با توجه به قضیه ۱۲.۴ در بخش بعد، می‌توانیم قضیه زیر را ثابت کنیم. اگر $f(x)$ بر $[a, b]$ ، $n+1$ بار پیوسته مشتقپذیر باشد، آنگاه

$$\rho_n(f) \leq \frac{[(b-a)/2]^{n+1}}{(n+1)! 2^n} \|f^{(n+1)}\|_\infty \equiv M_n(f) \quad (۱۰.۶.۴)$$

اقامه برهان به صورت مسأله ۳۸ به خواننده واگذار شده است. توابع بینهایت مشتقپذیری وجود دارند که برای آنها $M_n(f) \rightarrow \infty$. مع ذلک، برای بیشتر توابعی که زیاد به کار می‌روند، کران (۱۰.۶.۴) برآورد نسبتاً دقیقی از اندازه $\rho_n(f)$ ، به نظر می‌رسد. این برآورد برای $f(x) = e^x$ بر بازه $[-1, 1]$ در جدول ۶.۴ نشان داده شده است. برای برآوردها و کرانهای دیگر $\rho_n(f)$ ، میناردوس (۱۹۶۷، بخش ۲.۶) را ببینید.

۷.۴ تقریبهای نزدیک مینیماکس

با توجه به قضیه همونسانی چبیشف، می‌توانیم روشهایی به دست آوریم که اغلب برآورد خوبی برای تقریب مینیماکس بدهند. با تقریب کمترین مربعات $C_n(x)$ از (۱۷.۵.۴) شروع می‌کنیم. این تقریب اغلب برآورد خوبی برای $q_n^*(x)$ است و ویژگیهای $C_n(x)$ انگیزه دیگر تقریبهای نزدیک مینیماکس ما هستند.

با توجه به (۱۷.۵.۴) داریم،

$$C_n(x) = \sum_{j=0}^n c_j T_j(x) \quad c_j = \frac{2}{\pi} \int_{-1}^1 \frac{f(x) T_j(x) dx}{\sqrt{1-x^2}} \quad (۱۷.۷.۴)$$

که پریم در نماد مجموعیابی بدین معناست که اولین جمله ($j = 0$) قبل از جمع کردن سری باید نصف شود. اگر $f \in C[-1, 1]$ ، آنگاه با استفاده از (۷.۵.۴)

$$f(x) = \sum_{j=0}^{\infty} c_j T_j(x) \quad (2.7.4)$$

با برقراری همگرایی به این معنی که

$$\lim_{n \rightarrow \infty} \int_{-1}^1 \frac{1}{\sqrt{1-x^2}} \left[f(x) - \sum_{j=0}^n c_j T_j(x) \right]^2 dx = 0$$

برای همگرایی یکنواخت، قضیه نسبتاً قوی زیر را داریم

$$\rho_n(f) \leq \|f - C_n\|_{\infty} \leq \left(4 + \frac{4}{\pi^2} \ln(n) \right) \rho_n(f) \quad (3.7.4)$$

برای اثبات، ریولین (۱۹۷۴ ص ۱۳۴) را ببینید. از ترکیب این رابطه با قضیه (۹.۶.۴) جکسن، کران همگرایی (۲۰.۵.۴) قبلی از قضیه ۸.۴ به دست می آید.

اگر $f \in C^r[a, b]$ ، می توان ثابت کرد که مقدار ثابتی چون c وابسته فقط به f و r وجود دارد که برای آن

$$|c_j| \leq \frac{c}{j^r} \quad j \geq 1 \quad (4.7.4)$$

اثبات این فرمول، با در نظر گرفتن c_j ها به عنوان ضرایب فوریه $f(\cos \theta)$ و سپس استفاده از قضایای سری فوریه راجع به نرخ کاهش این ضرایب انجام می گیرد. بنابراین وقتی r بزرگ می شود، ضرایب c_j سریعتر کاهش می یابند.

برای بسط منقطع $C_n(x)$ ، اگر $c_{j+1} \neq 0$ ، و ضرایب c_j به سرعت به صفر بگرایند، داریم

$$f(x) - C_n(x) = \sum_{n+1}^{\infty} c_j T_j(x) \doteq c_{n+1} T_{n+1}(x) \quad (5.7.4)$$

از تعریف T_{n+1} داریم

$$|T_{n+1}(x)| \leq 1 \quad -1 \leq x \leq 1 \quad (6.7.4)$$

همچنین برای $n+2$ نقطه

$$x_j = \cos\left(\frac{j\pi}{n+1}\right) \quad j = 0, 1, \dots, n+1 \quad (7.7.4)$$

داریم

$$T_{n+1}(x_j) = (-1)^j \quad (۸.۷.۴)$$

کران (۶.۷.۴) دقیقاً به $n + ۲$ نقطه، یعنی ماکسیمم تعداد ممکن می‌رسد. با استفاده از این رابطه در (۵.۷.۴)، جمله $c_{n+1}T_{n+1}$ دقیقاً $n + ۲$ ماکسیمم و مینیمم نسبی دارد که همگی اندازه مساوی دارند. بنابراین، از قضیه هموسانی چیشف می‌توانیم انتظار داشته باشیم که $C_n(x)$ خیلی نزدیک به تقریب مینیمکس $q_n^*(x)$ باشد.

مثال مثال (۲۲.۵.۴) برای $f(x) = e^x$ نزدیک به آخر بخش ۵.۴ را به یاد آورید. در آنجا، ضرایب c_j خیلی سریع کاهش می‌یافتند، و

$$\|e^x - C_r\|_{\infty} = ۰.۰۰۰۶۰۷ \quad c_r T_r(x) = ۰.۰۰۵۴۷ T_r(x)$$

$$\|e^x - q_r^*(x)\|_{\infty} = ۰.۰۰۵۵۳$$

توضیحات پاراگراف اخیر با این مثال روشن می‌شود.

مثال می‌توان نشان داد که

$$\tan^{-1} x = ۲ \left[\alpha T_1(x) - \frac{\alpha^3}{3} T_3(x) + \frac{\alpha^5}{5} T_5(x) - \dots \right] \quad (۹.۷.۴)$$

برای $-1 \leq x \leq 1$ با $\alpha = \sqrt{2} - 1 \approx ۰.۴۱۴$ ، همگرایی یکنواخت است. پس

$$C_{r_{n+1}}(x) = ۲ \left[\alpha T_1(x) - \frac{\alpha^3}{3} T_3(x) + \dots + \frac{(-1)^n}{2n+1} \alpha^{2n+1} T_{2n+1}(x) \right] \quad (۱۰.۷.۴)$$

برای خطا،

$$E_{r_{n+1}}(x) \equiv \tan^{-1} x - C_{r_{n+1}}(x) = \frac{۲}{2n+3} (-1)^{n+1} \alpha^{2n+2} T_{2n+2}(x) + ۲ \sum_{j=n+2}^{\infty} \frac{(-1)^j \alpha^{2j+1}}{2j+1} T_{2j+1}(x)$$

کرانه‌های بالا و پایین این عبارت را پیدا می‌کنیم تا خطای $C_{r_{n+1}}(x)$ و خطای مینیمکس $\rho_{r_{n+1}}(f)$ را تخمین بزنیم.

با گرفتن کرانه‌های بالا

$$\left| ۲ \sum_{n+2}^{\infty} \frac{(-1)^j \alpha^{2j+1}}{2j+1} T_{2j+1}(x) \right| \leq ۲ \sum_{n+2}^{\infty} \frac{\alpha^{2j+1}}{2j+1} < \frac{۲}{2n+5} \cdot \frac{\alpha^{2n+5}}{1+\alpha^2}$$

بنابراین کرانهای بالا و پایین $E_{r_{n+1}}$ را، با توجه به اینکه $\alpha^2 / (1 - \alpha^2) \doteq 0.207$ ، به دست می آوریم،

$$\begin{aligned} E_{r_{n+1}}(x) &\leq \frac{2}{2n+3} (-1)^{n+1} \alpha^{2n+2} T_{r_{n+2}}(x) + \frac{2\alpha^{2n+5}}{(2n+5)(1-\alpha^2)} \\ &\leq \frac{2\alpha^{2n+2}}{2n+3} \left[(-1)^{n+1} T_{r_{n+2}}(x) + 0.207 \right] \end{aligned}$$

همچنین

$$E_{r_{n+1}}(x) \geq \frac{\alpha^{2n+2}}{2n+3} \left[(-1)^{n+1} T_{r_{n+2}}(x) - 0.207 \right]$$

بنابراین

$$E_{r_{n+1}}(x) = \frac{\alpha^{2n+2}}{2n+3} (-1)^{n+1} T_{r_{n+2}}(x) \quad (11.7.4)$$

که دارای $2n+4$ ماکسیمم و مینیمم نسبی با علامتهای یک در میان مثبت و منفی است. با توجه به قضیه ۹.۴ و نامساویهای قبلی

$$\frac{2(0.793)\alpha^{2n+2}}{2n+3} \leq \rho_{r_{n+1}}(f) = \rho_{r_{n+2}}(f) \leq \frac{2(1.207)\alpha^{2n+2}}{2n+3} \quad (12.7.4)$$

برای محاسبه عملی ضرایب c_j در (۱.۷.۴) از فرمول

$$c_j = \frac{2}{\pi} \int_0^\pi \cos(j\theta) f(\cos \theta) d\theta \quad (13.7.4)$$

و قاعده انتگرالگیری عددی میانگهی یا دوزنقه‌یی استفاده می‌کنیم. این مطالب قبلاً در (۲۱.۵.۴) نشان داده شده است.

درونیایی در صفرهای چبیشف اگر خطا در تقریب مینیمکس تقریباً برابر $c_{n+1} T_{n+1}(x)$ باشد، همان‌گونه که در (۵.۷.۴) نشان داده شد، خطا در ریشه‌های $T_{n+1}(x)$ در $[-1, 1]$ باید نزدیک صفر باشد. این ریشه‌ها نقاط زیرند

$$x_j = \cos \left[\frac{2j+1}{2n+2} \pi \right] \quad j = 0, 1, \dots, n \quad (14.7.4)$$

گیریم $I_n(x)$ یک چندجمله‌یی با درجه نابزرگتر از n باشد که $f(x)$ را در این گره‌های $\{x_j\}$ درونیایی نماید. چون خطای این درونیاب در x_j ها، صفرهای چبیشف $T_{n+1}(x)$ برابر صفر است، پیوستگی چندجمله‌یی درونیایی نسبت به مقادیر تابع که چندجمله‌یی را تعریف می‌کنند، گویای آن است که $I_n(x)$ باید تقریباً برابر $C_n(x)$ و لذا همچنین تقریباً برابر $q_n^*(x)$ باشد. به عبارت دقیقتر،

شکل لاگرانژی $I_n(x)$ را نوشته و با عملیاتی روی آن به دست می‌آوریم:

$$\begin{aligned} I_n(x) &= \sum_{j=0}^n f(x_j) l_j(x) \\ &= \sum_{j=0}^n C_n(x_j) l_j(x) + \sum_{j=0}^n [f(x_j) - C_n(x_j)] l_j(x) \\ &\doteq C_n(x) \end{aligned}$$

$$I_n(x) \doteq C_n(x) \quad (۱۵.۷.۴)$$

زیرا $f(x_j) - C_n(x_j) \doteq c_{n+1} T_{n+1}(x_j) = 0$ عبارت $I_n(x)$ را می‌توان با الگوریتم‌های *Divdif* و *Interp* بخش ۲.۳ حساب کرد. مواظب باشید: اگر $C_{n+1} = 0$ ، خطای $f(x) - q_n^*(x)$ احتمالاً صفرهای تقریبی یک چندجمله‌یی از درجه بالاتر معمولاً $T_{n+2}(x)$ چبیشف را خواهد داشت. اگر $f(x)$ در $[-1, 1]$ یا فرد باشد یا زوج این حالت ممکن است اتفاق بیفتد (مسئله ۲۹ را ببینید). تابع $f(x)$ را بر بازه $[-a, a]$ زوج گویند اگر به ازای تمام مقادیر x در $[-a, a]$ ، $f(-x) = f(x)$ ، فرد گویند اگر $f(-x) = -f(x)$. مثال قبلی با $f(x) = \tan^{-1}(x)$ یک نمایش خوبی از موضوع است. انگیزه دیگر در نظر گرفتن $I_n(x)$ به عنوان یک تقریب نزدیک مینیمکس مبتنی بر قضیه مهم زیر درباره چندجمله‌یهای چبیشف است.

قضیه ۱۲.۴ مسئله مینیم‌سازی

$$\tau_n = \inf_{\deg(Q) \leq n-1} \left[\text{Max}_{-1 \leq x \leq 1} |x^n + Q(x)| \right] \quad (۱۶.۷.۴)$$

را که در آن $Q(x)$ یک چندجمله‌یی است برای یک عدد صحیح ثابت $n > 0$ در نظر می‌گیریم. اگر $Q(x)$ را به طور ضمنی با

$$x^n + Q(x) = \frac{1}{\sqrt{n-1}} T_n(x) \quad (۱۷.۷.۴)$$

تعریف کنیم مینیم τ_n به طور یکتا حاصل می‌شود. این مینیم برابر است با

$$\tau_n = \frac{1}{\sqrt{n-1}} \quad (۱۸.۷.۴)$$

برهان نخست بعضی واقعیتها در مورد چندجمله‌یهای چبیشف را در نظر می‌گیریم. بنابه تعریف،

$$T_0(x) \equiv 1 \text{ و } T_1(x) = x. \text{ رابطه بازگشتی سه‌گانه}$$

$$T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x)$$

پایه یک برهان استقرایی است برای

$$T_n(x) = 2^{n-1}x^n + \text{جملات از درجهٔ پایینتر} \quad n \geq 1 \quad (۱۹.۷.۴)$$

بنابراین

$$\frac{1}{2^{n-1}}T_n(x) = x^n + \text{جملات از درجات پایینتر} \quad n \geq 1 \quad (۲۰.۷.۴)$$

چون $T_n(x) = \cos n\theta$, $x = \cos \theta$, $0 \leq \theta \leq \pi$ ، چندجمله‌یی $T_n(x)$ در $n+1$ نقطه در $[-1, 1]$ ماکسیمم و مینیمم می‌شود:

$$x_j = \cos\left(\frac{j\pi}{n}\right) \quad j = 0, 1, \dots, n \quad (۲۱.۷.۴)$$

چون تابع کسینوس دوره‌یی است، برای مقادیر دیگر j مقادیر جدیدی برای x_j به دست نمی‌آید. برای این نقاط

$$T_n(x_j) = (-1)^j \quad j = 0, 1, \dots, n \quad (۲۲.۷.۴)$$

و

$$-1 = x_n < x_{n-1} < \dots < x_1 < x_0 = 1$$

ضریب جملهٔ بزرگترین درجهٔ چندجمله‌یی $T_n(x)/2^{n-1}$ برابر یک است و

$$\text{Max}_{-1 \leq x \leq 1} \left| \frac{1}{2^{n-1}}T_n(x) \right| = \frac{1}{2^{n-1}} \quad (۲۳.۷.۴)$$

بنابراین $\tau_n < 1/2^{n-1}$. فرض می‌کنیم که

$$\tau_n < \frac{1}{2^{n-1}} \quad (۲۴.۷.۴)$$

نشان می‌دهیم که این فرض به تناقض می‌انجامد. فرض (۲۴.۷.۴) و تعریف (۱۶.۷.۴) دلالت بر وجود یک چندجمله‌یی به صورت زیر دارند

$$M(x) = x^n + Q(x) \quad \text{درجهٔ } Q(x) \text{ نابزرگتر از } n-1$$

با

$$\tau_n \leq \text{Max}_{-1 \leq x \leq 1} |M(x)| < \frac{1}{2^{n-1}} \quad (۲۵.۷.۴)$$

تعریف می‌کنیم

$$R(x) = \frac{1}{2^{n-1}}T_n(x) - M(x)$$

که دارای درجه نایزگتر از $n - 1$ است. علامت $R(x_j)$ را در نقاط (۲۱.۷.۴) امتحان می‌کنیم. با استفاده از (۲۲.۷.۴) و (۲۴.۷.۴)

$$R(x_0) = R(1) = \frac{1}{2^{n-1}} - M(1) > 0$$

$$R(x_1) = -\frac{1}{2^{n-1}} - M(x_1) = -\left[\frac{1}{2^{n-1}} + M(x_1)\right] < 0$$

و علامت $R(x_j)$ برابر $(-1)^j$ است. چون R دارای $n + 1$ تغییر علامت است، باید حداقل n صفر داشته باشد. ولی در این صورت درجه R که کوچکتر از n است، ایجاب می‌کند که $R \equiv 0$ ؛ پس $M \equiv (1/2^{n-1})T_n$.

برای اثبات اینکه هیچ چندجمله‌یی دیگری غیر از $1/2^{n-1}T_n(x)$ ، (۱۶.۷.۴) را مینیمم نمی‌کند، یک شکل دیگری از برهان قبل به‌کار برده می‌شود. ما از ذکر این برهان خودداری می‌کنیم. ■

اکنون مسأله تعیین $n + 1$ نقطه گرهی x_j در $[-1, 1]$ را در نظر می‌گیریم، که باید در ساختن یک چندجمله‌یی درونیاب $p_n(x)$ که تابع مفروض $f(x)$ را در $[-1, 1]$ تقریب می‌زند، مورد استفاده قرار گیرد. خطا در $p_n(x)$ برابر است با

$$f(x) - p_n(x) = \frac{(x - x_0) \dots (x - x_n)}{(n + 1)!} f^{(n+1)}(\xi_x) \quad (26.7.4)$$

مقدار $f^{(n+1)}(\xi_x)$ به $\{x_j\}$ بستگی دارد، ولی این بستگی چنان نیست که بتوان به صراحت به آن پرداخت. پس برای کوچک ساختن $\|f - p_n\|_\infty$ در حد ممکن، فقط کمیت

$$\text{Max}_{-1 \leq x \leq 1} |(x - x_0) \dots (x - x_n)| \quad (27.7.4)$$

را در نظر می‌گیریم. $\{x_j\}$ را طوری انتخاب می‌کنیم که این کمیت را مینیمم سازد. چندجمله‌یی (۲۷.۷.۴) از درجه $n + 1$ و ضریب بزرگترین درجه x در آن برابر یک است. طبق قضیه قبل، (۲۷.۷.۴) وقتی مینیمم می‌شود که آن را چندجمله‌یی $T_{n+1}(x)/2^n$ انتخاب کنیم و مقدار مینیمم (۲۷.۷.۴) برابر $1/2^n$ است. گره‌های $\{x_j\}$ صفرهای $T_{n+1}(x)$ اند و این نقاط در (۱۴.۷.۴) داده شده‌اند. با این انتخاب گره‌ها، $p_n = I_n$ و

$$\|f - I_n\|_\infty \leq \frac{1}{(n + 1)! 2^n} \|f^{(n+1)}\|_\infty \quad (28.7.4)$$

مثال گیریم $f(x) = e^x$ و $x = 3$. صورت تفاضلات منقسم چندجمله‌یی درونیاب نیوتن را به‌کار می‌بریم. گره‌ها، مقادیر تابع، و تفاضلات منقسم لازم در جدول ۷.۴ داده شده‌اند. با محاسبه

جدول ۷.۴ داده‌های درونیابی برای $f(x) = e^x$

i	x_i	$f(x_i)$	$f[x_0, \dots, x_i]$
۰	۰٫۹۲۳۸۸۰	۲٫۵۱۹۰۴۴۲	۲٫۵۱۹۰۴۴۲
۱	۰٫۳۸۲۶۸۳	۱٫۴۶۶۲۱۳۸	۱٫۹۴۵۳۷۶۹
۲	-۰٫۳۸۲۶۸۳	۰٫۶۸۲۰۲۸۸	۰٫۷۰۴۷۴۲۰
۳	۰٫۹۲۳۸۸۰	۰٫۳۹۶۹۷۶۰	۰٫۱۷۵۱۷۵۷

مستقیم داریم

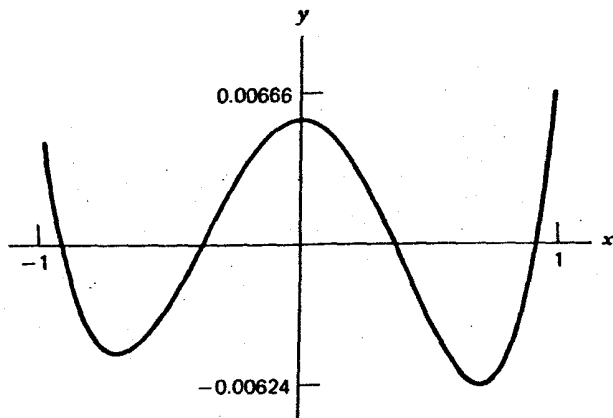
$$\text{Max}_{-1 \leq x \leq 1} |e^x - I_3(x)| \doteq 0.00666 \quad (29.7.4)$$

و حال آنکه کران در (۲۸.۷.۴) برابر 0.014 است. نمودار $e^x - I_3(x)$ در شکل ۷.۴ آمده است.

خطای $\|f - I_n\|_\infty$ معمولاً خیلی بدتر از $\rho_n(f)$ نیست. یک نتیجه دقیق چنین است

$$\|f - I_n\|_\infty \leq \left[\frac{2}{\pi} \log(n+1) + 2 \right] \rho_n(f) \quad n \geq 0 \quad (30.7.4)$$

برهان را در ریبولین (۱۹۷۰، ص ۱۳) ببینید. همان‌گونه که (۲۹.۷.۴) در مثالهای قبل $f(x) = e^x$ نشان می‌دهد، نتایج واقعی عددی معمولاً بهتر از آن است که با کران (۲۹.۷.۴) پیش‌بینی می‌شود.



شکل ۷.۴ $e^x - I_3(x)$

نوسان اجباری خطا گیریم $f(x) \in C[-1, 1]$ و تعریف می‌کنیم

$$F_n(x) = \sum_{k=0}^n c_{n,k} T_k(x) \quad -1 \leq x \leq 1 \quad (31.7.4)$$

گیریم

$$-1 \leq x_{n+1} < x_n < \dots < x_1 < x_0 \leq 1$$

گره‌هایی باشند که در مورد انتخاب آنها ذیلاً توضیح داده می‌شود. ضرایب $c_{n,k}$ را چنان اختیار می‌کنیم که خطای $f(x) - F_n(x)$ به طریقی نوسان کند که در قضیه ۱۰.۴ لازم فرض شده است:

$$f(x_i) - F_n(x_i) = (-1)^i E_n \quad i = 0, 1, \dots, n+1 \quad (32.7.4)$$

مجهول دیگر E_n را وارد کرده‌ایم که امیدواریم صفر نباشد. $n+2$ مجهول داریم: $c_{n,1}, c_{n,0}, \dots, c_{n,n}$ و E_n ، و $n+2$ معادله در (۳۲.۷.۴) موجودند. پس برای اینکه یک جواب داشته باشیم، شانس زیادی وجود دارد. اگر یک جواب موجود باشد، طبق قضیه ۹.۴،

$$|E_n| \leq \rho_n(f) \leq \|f - F_n\|_\infty \quad (33.7.4)$$

برای انتخاب نقاط گرهی، توجه کنید اگر $c_{n+1} T_{n+1}(x)$ در (۵.۷.۴) تقریباً خطای مینیماکس باشد، آنگاه مینیمم‌ها و ماکسیم‌های نسبی در خطای مینیماکس $f(x) - q_n^*(x)$ در مینیمم‌ها و ماکسیم‌های نسبی $T_{n+1}(x)$ واقع می‌شوند. این ماکسیم‌ها و مینیمم‌ها وقتی اتفاق می‌افتند که $T_{n+1}(x) = \pm 1$ و با رابطه زیر داده می‌شوند

$$x_i = \cos\left(\frac{i\pi}{n+1}\right) \quad i = 0, 1, \dots, n+1 \quad (34.7.4)$$

به نظر می‌آید که اگر $F_n(x)$ تقریب مینیماکس $q_n^*(x)$ اطلاق شود، نقاط فوق یک انتخاب عالی برای استفاده در (۳۲.۷.۴) باشند.

دستگاه (۳۲.۷.۴) چنین می‌شود

$$\sum_{k=0}^n c_{n,k} T_k(x_i) + (-1)^i E_n = f(x_i) \quad i = 0, 1, \dots, n+1 \quad (35.7.4)$$

توجه داشته باشید که

$$T_k(x_i) = \cos(k \cdot \cos^{-1}(x_i)) = \cos\left(\frac{ki\pi}{n+1}\right)$$

E_n را برابر $c_{n,n+1}/2$ اختیار می‌کنیم. دستگاه (۳۵.۷.۴) چنین می‌شود

$$\sum_{k=0}^{n+1} c_{n,k} \cos\left(\frac{ki\pi}{n+1}\right) = f(x_i) \quad i = 0, 1, \dots, n+1 \quad (36.7.4)$$

زیرا

$$\cos\left(\frac{ki\pi}{n+1}\right) = (-1)^i \quad \text{برای } k = n+1$$

نماد \sum بدین معناست که اولین و آخرین جمله، قبل از شروع به مجموعیابی باید نصف شوند. برای حل (۳۶.۷.۴) به روابط زیر نیاز داریم:

$$\sum_{i=0}^{n+1} \cos\left(\frac{ij\pi}{n+1}\right) \cos\left(\frac{ik\pi}{n+1}\right) = \begin{cases} n+1 & j=k=0 \text{ یا } n+1 \\ \frac{(n+1)}{2} & 0 < j=k < n+1 \\ 0 & j \neq k \quad 0 \leq j, k \leq n+1 \end{cases} \quad (37.7.4)$$

اثبات این روابط به اثبات روابط در مسأله ۴۲ فصل ۳، بستگی نزدیکی دارد.

معادله نام در (۳۶.۷.۴) را در $\cos\left[\frac{(ij\pi)}{(n+1)}\right]$ به ازای یک $0 \leq j \leq n+1$ ضرب و سپس روی i مجموعیابی می‌کنیم، با توجه به نصف کردن جملات اول و آخر، این نتیجه به دست می‌آید:

$$\sum_{k=0}^{n+1} c_{n,k} \sum_{i=0}^{n+1} \cos\left(\frac{ij\pi}{n+1}\right) \cos\left(\frac{ik\pi}{n+1}\right) = \sum_{i=0}^{n+1} f(x_i) \cos\left(\frac{ij\pi}{n+1}\right)$$

با استفاده از رابطه (۳۷.۷.۴)، تمام جملات غیر از یکی در مجموعیابی روی k صفر می‌شوند. با امتحان کردن دو حالت $0, n+1, j = 0, n+1, j < n+1$ همان فرمول به دست می‌آید:

$$c_{n,j} = \frac{2}{n+1} \sum_{i=0}^{n+1} f(x_i) \cos\left(\frac{ij\pi}{n+1}\right) \quad 0 \leq j \leq n+1 \quad (38.7.4)$$

فرمول برای E_n چنین است

$$E_n = \frac{1}{n+1} \sum_{i=0}^{n+1} (-1)^i f(x_i) \quad (39.7.4)$$

بین این تقریب $E_n(x)$ و بسط $C_n(x)$ چپیشف چند رابطه وجود دارد. مهمترین آنها، ضرایب $c_{n,j}$ تقریبه‌های c_j در $C_n(x)$ هستند. فرمول (۱۳.۷.۴) را برای c_j با استفاده از قاعده

دوزنقه‌یی (۵.۱.۵) با $n + 1$ زیر بازه محاسبه می‌کنیم:

$$c_j = \frac{2}{\pi} \int_0^\pi f(\cos \theta) \cos(j\theta) d\theta$$

$$\doteq \frac{2}{\pi} \sum_{i=0}^{n+1} f\left(\cos\left(\frac{i\pi}{n+1}\right)\right) \cdot \cos\left(\frac{ji\pi}{n+1}\right) \cdot \frac{\pi}{n+1} = c_{n,j} \quad (40.7.4)$$

به خوبی می‌دانیم که وقتی انتگرالده دوره‌یی باشد، قاعده انتگرالگیری عددی دوزنقه‌یی خیلی دقیق است (قضیه ۵.۵ بخش ۴.۵ را ببینید). به علاوه می‌توان نشان داد که

$$c_{nj} = c_j + c_{2(n+1)-j} + c_{2(n+1)+j} + c_{2(n+1)-j} + \dots \quad (41.7.4)$$

اگر ضرایب چبیشف در (۲.۷.۴) سریعاً کاهش یابند، تقریب $F_n(x)$ تقریباً برابر $C_n(x)$ می‌شود و محاسبه ساده‌تر از محاسبه $C_n(x)$ انجام می‌گیرد.

مثال مانند قبل، $f(x) = e^x$ ، $-1 \leq x \leq 1$ ، را به کار می‌بریم. به ازای $n = 1$ گرهما چنین‌اند

$$\{x_i\} = \{-1, 0, 1\}$$

$$E_1 = 0.272 \text{ و}$$

$$F_1(x) = 1.2715 + 1.1752x \quad (42.7.4)$$

برای خطا، ماکسیم‌های نسبی $e^x - F_1(x)$ در جدول ۸.۴ داده شده‌اند. برای $n = 3$ و $E_3 = 0.00547$

$$F_3(x) = 0.994526 + 0.995682x + 0.543081x^2 + 0.179519x^3 \quad (43.7.4)$$

نقاط خطای ماکسیم در جدول ۹.۴ داده شده‌اند. از قضیه ۹.۴ چنین به دست می‌آوریم:

$$0.00547 \leq \rho_3(f) \leq 0.00558$$

و این رابطه می‌گوید که $F_3(x)$ یک تقریب عالی برای $q_3^*(x)$ است.

برای آنکه ببینیم $F_3(x)$ یک تقریب $C_3(x)$ است، ضرایب $c_{3,j}$ را با ضرایب c_j در جدول ۳.۴ در مثال (۲۲.۵.۴) در پایان بخش ۵.۴ داده شده، مقایسه می‌کنیم. نتایج در جدول ۱۰.۴ داده شده‌اند.

جدول ۸.۴ ماکسیمهای نسبی $|e^x - F_1(x)|$

x	$e^x - F_1(x)$
-1.0	0.2772
0.1614	-0.286
1.0	0.2772

جدول ۹.۴ ماکسیمهای نسبی $|e^x - F_7(x)|$

x	$e^x - F_7(x)$
-1.0	0.00547
-0.6832	-0.00552
0.493	0.00558
0.7324	-0.00554
1.0	0.00547

مانند حالت چندجمله‌یی تقریب درونیاب $I_n(x)$ ، اگر $f(x)$ در $[-1, 1]$ فرد یا زوج باشد باید مواظب بود. در چنین حالتی، n را به شکل زیر انتخاب می‌کنیم.

$$(44.7.4) \quad \text{اگر } f \begin{cases} \text{زوج} \\ \text{فرد} \end{cases} \text{ باشد، آنگاه } n \text{ را } \begin{cases} \text{فرد} \\ \text{زوج} \end{cases} \text{ می‌گیریم.}$$

این انتخاب تضمین می‌کند که در (۵.۷.۴)، $c_{n+1} \neq 0$ ، و گره‌های انتخابی درست خواهند بود. یک تحلیل همگرایی $F_n(x)$ به $f(x)$ در شمپاین^۱ (۱۹۷۰) داده شده است، که کرانی مانند

جدول ۱۰.۴ ضرایب بسط $c_j(x)$ و $F_7(x)$ برای e^x

j	c_j	$c_{n,j}$
۰	۲,۵۳۲۱۳۱۷۶	۲,۵۳۲۱۳۲۱۵
۱	۱,۱۳۰۳۱۸۲۱	۱,۱۳۰۳۲۱۴۲
۲	۰,۲۷۱۴۹۵۳۴	۰,۲۷۱۵۴۰۳۲
۳	۰,۰۴۴۳۳۶۸۵	۰,۰۴۴۸۷۹۷۸
۴	۰,۰۰۵۴۷۴۲۴	$E_7 = 0.00547424$

کران (۳۰.۷.۴) برای $I_n(x)$ به دست می‌دهد:

$$\|f - F_n\|_\infty \leq w(n)\rho_n(f) \quad (۴۵.۷.۴)$$

که $w(n)$ از لحاظ تجربی تقریباً مساوی ضریب کرانی در (۳۰.۷.۴) است. $I_n(x)$ و $F_n(x)$ هر دو، تقریبهای عملی نزدیک مینیماکس هستند.

اکنون یک الگوریتم برای محاسبه $F_n(x)$ می‌دهیم که می‌توان آن را با استفاده از الگوریتم چپوال بخش ۵.۴ محاسبه کرد.

الگوریتم Approx (c, E, f, n)

۱. تبصره: با این الگوریتم ضرایب c_j در

$$F_n(x) = \sum_{j=0}^n c_j T_j(x) \quad -1 \leq x \leq 1$$

برطبق فرمول (۳۸.۷.۴) محاسبه می‌شود، و E از (۳۹.۷.۴) به دست می‌آید. جمله c قبل از استفاده از الگوریتم چپوال باید نصف شود.

۲. $x_i := \cos(i\pi/(n+1))$ را تولید کنید.

$$f_i = f(x_i) \quad i := 0, 1, \dots, n+1$$

۳. برای $i = 0, 1, \dots, n+1$ تا مرحله ۸ عمل را ادامه بدهید

$$\text{sum} := [f_0 + (-1)^i f_{n+1}] / 2 \quad ۴$$

۵. برای $i = 1, \dots, n$ تا مرحله ۶ عمل را انجام دهید

$$\text{sum} := \text{sum} + f_i \cos(ij\pi/(n+1)) \quad ۶$$

۷. حلقه را در i ختم کنید

$$c_j \doteq 2\text{sum}/(n+1) \quad ۸$$

۹. حلقه را در j ختم کنید.

$$E := c_{n+1}/2 \quad ۱۰$$

کسینوسهای موجود در مرحله ۶ را می‌توان با استفاده از فرمولهای جمع مثلثاتی برای توابع سینوس و کسینوس، به نحو کاراتری محاسبه کرد، ولی ما آنها را به علت سادگی آموزشی مطلب، انتخاب کرده‌ایم، زیرا زمان اجرای رایانه‌یی برای الگوریتم ما خیلی کم است. به همین دلیل، ما تکنیک FFTی بخش ۸.۳ را به کار نبرده‌ایم.

بحث در آثار خواندنی

نظریه تقریب از جنبه کلاسیک یک زمینه مهم ریاضیات است، و نیز ابزاری با اهمیت فزاینده در مطالعه مسائل جدید بسیار متنوع در ریاضیات کاربردی، مثلاً، در فیزیک ریاضی و ریاضیات ترکیبیاتی است. تنوع مسائل و روشهای نظریه تقریب را در کتابهای، اچیزر^۱ (۱۹۵۶)، اسکس^۲ (۱۹۷۵b)، دیویس (۱۹۶۳)، لورتس (۱۹۶۶)، میناردوس (۱۹۶۷)، پاول (۱۹۸۱)، رایس (۱۹۶۴)، (۱۹۶۸) می‌توان دید. کار کلاسیک در زمینه چندجمله‌بیهای متعامد در سگو (۱۹۶۷)، و یک بازنگری از کارهای تازه‌تر در اسکس (۱۹۷۵a) داده شده است. برای چندجمله‌بیهای چبیشف و کاربردهای زیاد آنها در همه ریاضیات کاربردی و آنالیز عددی، فاکس و پارکر (۱۹۶۸)، ریولین (۱۹۷۴) را ببینید. در این کتاب به موضوع مربوط به سریهای فوریه و تقریب با چندجمله‌بیهای مثلثاتی فقط اشاره شده، ولی این مبحث در عده زیادی از کاربردها از اهمیت کلیدی برخوردار است. یک مرجع کلاسیک زیگموند (۱۹۵۹) است. زمینه‌های زیادی در نظریه تقریب وجود دارند که ما حتی تعریف هم نکرده‌ایم. برای یک بازنگری عالی در این زمینه گاوچی^۳ (۱۹۷۵) را که شامل فهرست مراجع خوبی است ببینید. یک زمینه اصلی که در این کتاب حذف شده است، تقریب به وسیله توابع گویاست. برای زمینه کلی، میناردوس (۱۹۶۷)، فصل نهم) و رایس (۱۹۶۸)، فصل نهم) را ببینید. تعمیم چندجمله‌بیهای تیلر به توابع گویا را تقریب پاده^۴ می‌نامند؛ برای مقدمات بیکر (۱۹۷۵) و برزینسکی^۵ (۱۹۸۰) را ببینید. در زمینه مربوط به بسطهای وابسته بسط کسرهای مسلسل توابع وال (۱۹۴۸) را ببینید. بسیاری از توابع که از نظر کاربردی جالب هستند نمونه‌هایی هستند که توابع خاص در فیزیک ریاضی نامیده می‌شوند. این نمونه‌ها شامل توابع متعالی اساسی (سینوس، لگاریتم، نمایی، ریشه دوم) و علاوه بر آن، چندجمله‌بیهای متعامد، توابع بسل، تابع گاما و تابع ابرهندسی هستند. نوشته‌های مفصلی درباره توابع خاص وجود دارند و تقریبهای ویژه‌ای برای اغلب آنها ساخته شده است. مهمترین مراجع برای توابع خاص عبارت‌اند از: آبراموویتس و استگون (۱۹۶۴)، که یک کتاب راهنماست که با حمایت اداره استانداردهای ملی آمریکا تهیه شده و اردلی^۶ و همکاران (۱۹۵۳) و یک مجموعه سه جلدی که غالباً به آن با عنوان «پروژه بیتمن^۷» اشاره می‌شود. برای یک نگرش جامع و بازنگری روشهای تقریب توابع خاص گاوچی (۱۹۷۵) را ببینید. چکیده جامعی از قضایای نظری برای توابع خاص و روشهایی برای محاسبه عددی آنها در لوک^۸ (۱۹۶۹)، (۱۹۷۵)، (۱۹۷۷) آمده است. برای نمونه‌های نسبتاً متداولتر روندهای مطالعه توابع خاص گزارشهای سمپوزیوم اسکس (۱۹۷۵b) را ببینید.

1. Achieser

2. Askey

3. Gautschi, W

4. Padé

5. Beezinski

6. Erdélyi, S. W

7. Batemenn

8. Luke

از زمان استفاده وسیع از رایانه در دهه ۱۹۵۰، نیاز به تقریبهای با کیفیت بالای چندجمله‌یی یا تابع گویا برای توابع ریاضی اساسی و سایر توابع خاص به وجود آمده است. همان‌گونه که قبلاً اشاره شد تقریب این توابع نیاز به شناخت ویژگیهای آنها دارد. ولی در عین حال اطلاع دقیق از حساب رایانه‌های رقمی که در فصل اول بیان گردید لازم است. یک مطالعه کلی از روشهای عددی برای تولید تقریبهای چندجمله‌یی در فریزر (۱۹۶۵) داده شده است، که سازمان‌دهی این فصل را تحت تأثیر قرار داده است. برای یک بحث کامل از تقریب توابع مقدماتی همراه با الگوریتم‌های مشروح به کودی و ویت (۱۹۸۰) مراجعه کنید؛ بحثی از طرح برنامه‌ریزی مربوطه در کودی (۱۹۸۴) آمده است. برای بحث مشابهی از تقریبها، ولی بحثی که شامل بعضی از توابع خاص معمولی نیز باشد، هارت (۱۹۶۸) را ببینید. برای یک مجموعه وسیع از تقریبهای توابع خاص لوک (۱۹۷۵) و (۱۹۷۷) را ببینید. برای توابع کلی، یک برنامه موفق که برای تولید تقریبهای مینیمکس بسیار مورد استفاده قرار گرفته است، در کودی و همکاران (۱۹۶۸) داده شده است. برنامه‌های کلی برای محاسبه تقریبهای مینیمکس در کتابخانه‌های IMSL و NAG موجودند.

مراجع

- Abramowitz, M., and I. Stegun (eds.) (1964). *Handbook of Mathematical Functions*. National Bureau of Standards, U.S. Government Printing Office, Washington, D.C. (It is now published by Dover, New York.)
- Achieser, N. (1956). *Theory of Approximation* (transl. C. Hyman). Ungar, New York.
- Askey, R. (1975a). *Orthogonal Polynomials and Special Functions*. Society for Industrial and Applied Mathematics, Philadelphia.
- Askey, R. (ed.) (1975b). *Theory and Application of Special Functions*. Academic Press, New York.
- Baker, G., Jr. (1975). *Essentials of Padé Approximants*. Academic Press, New York.
- Brezinski, C. (1980). *Padé-Type Approximation and General Orthogonal Polynomials*. Birkhäuser, Basel.
- Cody, W. (1984). FUNPACK—A package of special function routines. In *Sources and Development of Mathematical Software*, W. Cowell (ed.), pp. 49–67. Prentice-Hall, Englewood Cliffs, N.J.
- Cody, W. and W. Waite (1980). *Software Manual for the Elementary Functions*. Prentice-Hall, Englewood Cliffs, N.J.
- Cody, W., W. Fraser, and J. Hart (1968). Rational Chebyshev approximation using linear equations. *Numer. Math.*, **12**, 242–251.
- Davis, P. (1963). *Interpolation and Approximation*. Ginn (Blaisdell), Boston.

- Erdélyi, A., W. Magnus, F. Oberhettinger, and F. Tricomi (1953). *Higher Transcendental Functions*, Vols. I, II, and III. McGraw-Hill, New York.
- Fox, L., and I. Parker (1968). *Chebyshev Polynomials in Numerical Analysis*. Oxford Univ. Press, Oxford, England.
- Fraser, W. (1965). A survey of methods of computing minimax and near-minimax polynomial approximations for functions of a single independent variable. *J. ACM*, 12, 295–314.
- Gautschi, W. (1975). Computational methods in special functions—A survey. In R. Askey (ed.), *Theory and Application of Special Functions*, pp. 1–98. Academic Press, New York.
- Hart, J., E. Cheney, C. Lawson, H. Maëhly, C. Mesztenyi, J. Rice, H. Thacher, and C. Witzgall (1968). *Computer Approximations*. Wiley, New York. (Reprinted in 1978, with corrections, by Krieger, Huntington, N.Y.)
- Isaacson, E., and H. Keller (1966). *Analysis of Numerical Methods*. Wiley, New York.
- Lorentz, G. (1966). *Approximation of Functions*. Holt, Rinehart & Winston, New York.
- Luke, Y. (1969). *The Special Functions and Their Applications*, Vols. I and II. Academic Press, New York.
- Luke, Y. (1975). *Mathematical Functions and Their Approximations*. Academic Press, New York.
- Luke, Y. (1977). *Algorithms for the Computation of Mathematical Functions*. Academic Press, New York.
- Meinardus, G. (1967). *Approximation of Functions: Theory and Numerical Methods* (transl. L. Schumaker). Springer-Verlag, New York.
- Powell, M. (1981). *Approximation Theory and Methods*. Cambridge Univ. Press, Cambridge, England.
- Rice, J. (1964). *The Approximation of Functions: Linear Theory*. Addison-Wesley, Reading, Mass.
- Rice, J. (1968). *The Approximation of Functions: Advanced Topics*. Addison-Wesley, Reading, Mass.
- Rivlin, T. (1974). *The Chebyshev Polynomials*. Wiley, New York.
- Shampine, L. (1970). Efficiency of a procedure for near-minimax approximation. *J. ACM*, 17, 655–660.
- Szego, G. (1967). *Orthogonal Polynomials*, 3rd ed. Amer. Math. Soc., Providence, R.I.
- Wall, H. (1948). *Analytic Theory of Continued Fractions*. Van Nostrand, New York.
- Zygmund, A. (1959). *Trigonometric Series*, Vols. I and II. Cambridge Univ. Press, Cambridge, England.

مسائل

۱. برای اینکه نشان دهید چندجمله‌بیهای $p_n(x)$ برنشتاین در قضیه ۱.۴، تقریبهای ضعیفی هستند، تقریب چندجمله‌یی درجه چهارم $p_n(x)$ را برای $f(x) = \sin(\pi x)$ ، $0 \leq x \leq 1$ ، حساب کنید. آن را با تقریب چندجمله‌یی درجه چهارم تیلر که حول $x = 1/2$ بسط داده شده مقایسه کنید.

۲. گیریم $S = \sum_{j=1}^{\infty} (-1)^j a_j$ یک سری همگرا باشد و فرض کنید که همه a_j ها ناکوچکتر از صفر باشند و

$$a_1 \geq a_2 \geq \dots \geq a_n \geq \dots$$

ثابت کنید که

$$\left| S - \sum_{j=1}^n (-1)^j a_j \right| \leq a_{n+1}$$

۳. با استفاده از مسئله ۲، همگرایی سریهای زیر را بررسی کنید. کران خطا را وقتی سریها پس از n جمله بریده شوند به دست آورید و توجه کنید که به x بستگی دارند. مقدار n را به گونه‌ای بیابید که خطا از 10^{-5} کمتر باشد. این مسئله روش معمول دیگری را برای کراندارکردن خطا در سری تیلر نشان می‌دهد،

$$J_1(x) = \sum_{j=0}^{\infty} \frac{(-1)^j \left(\frac{1}{4}x^2\right)^j}{(j!)^2} \quad (\text{الف})$$

$$\sum_{j=0}^{\infty} \frac{(-1)^j x^{2j}}{j^2} \quad (\text{ب})$$

۴. نمودار خطاهای تقریبهای $p_n(x)$ سری تیلر را برای $f(x) = \sin[(\pi/2)x]$ در $1 \leq x \leq -1$ ، به ازای $n = 1, 3, 5$ رسم کنید. به رفتار خطا، هم در نزدیکی مرکز و هم در نزدیکی نقاط انتهایی توجه نمایید.

۵. گیریم $f(x)$ بر بازه $[-\alpha, \alpha]$ به ازای مقداری از $\alpha > 0$ سه بار پیوسته مشتقپذیر باشد، و تقریب آن را توسط تابع گویای

$$R(x) = \frac{a + bx}{1 + cx}$$

در نظر می‌گیریم. برای تعمیم مفهوم سری تیلر، ثابتهای a و b و c را به گونه‌ای انتخاب کنید که

$$R^{(j)}(0) = f^{(j)}(0) \quad j = 0, 1, 2$$

آیا همیشه ممکن است چنین تابع تقریب $R(x)$ ی پیدا کرد؟ تابع $R(x)$ یک مثال از تقریب پاده است. بیکر (۱۹۸۰) و برزینسکی (۱۹۸۰) را ببینید.

۶. مسأله ۵ را برای حالت $f(x) = e^x$ به کار برید و $R(x)$ را پیدا کنید. خطای آن را بر $[-1, 1]$ تحلیل و آن را با خطای چندجمله‌یی درجه دوم تیلر مقایسه کنید.

۷. به وسیله اتحادهای گوناگون، اغلب می‌توان یک بازه‌ای را که لازم است تابعی در آن تقریب شود به یک بازه کوچکتر بدل کرد. نشان دهید که چگونه می‌توان هر یک از توابع زیر را از $-\infty < x < \infty$ به بازه داده شده تقلیل داد. معمولاً چند عمل اضافی، ولی ساده، لازم است.

$$\text{(الف) } e^x \quad 0 \leq x \leq 1$$

$$\text{(ب) } \cos(x) \quad 0 \leq x \leq \pi/4$$

$$\text{(ج) } \tan^{-1}(x) \quad 0 \leq x \leq 1$$

$$\text{(د) } \ln(x) \quad 1 \leq x \leq 2 \quad 0 < x < \infty \text{ را تقلیل دهید.}$$

۸. (الف) گیریم $f(x)$ بر $[a, b]$ پیوسته مشتقپذیر باشد. فرض کنید $p(x)$ یک چندجمله‌یی باشد که برای آن $\|f' - p\|_\infty \leq \varepsilon$ ، و تعریف می‌کنیم

$$q(x) = f(a) + \int_a^x p(t) dt \quad a \leq x \leq b$$

نشان دهید که $q(x)$ یک چندجمله‌یی است و در رابطه زیر صدق می‌کند

$$\|f - q\|_\infty \leq \varepsilon(b - a)$$

(ب) قسمت (الف) را به حالتی تعمیم دهید که $f(x)$ ، N بار پیوسته مشتقپذیر بر $[a, b]$ باشد، $N \geq 2$ و $p(x)$ یک چندجمله‌یی باشد که در رابطه زیر صدق کند

$$\|f^{(N)} - p\|_\infty \leq \varepsilon$$

فرمولی برای یک چندجمله‌یی $q(x)$ پیدا کنید که $f(x)$ را تقریب بزند و شامل تنها یک انتگرال باشد. (ج) فرض کنید که $f(x)$ بی‌نهایت بار مشتقپذیر بر $[a, b]$ باشد، یعنی، $f^{(j)}(x)$ در $[a, b]$ موجود و برای همه مقادیر $0 \leq j$ پیوسته باشد. [این فرض ایجاب نمی‌کند که $f(x)$ دارای بسط تیلر همگرا بر $[a, b]$ باشد.] ثابت کنید یک دنباله از چندجمله‌ییهای $\{p_n(x) \mid n \geq 1\}$ موجود است که برای آنها به ازای جميع مقادیر $0 \leq j$ ،

$$\lim_{n \rightarrow \infty} \|f^{(j)} - p_n^{(j)}\|_\infty = 0$$

راهنمایی: قضیه وایرستراس و قسمت (ب) را به کار برید.

۹. قضیه زیر را ثابت کنید: گیریم $f \in C^1[a, b]$ با $f''(x) > 0$ در $a \leq x \leq b$. اگر $q_1^*(x) = a_0 + a_1x$ تقریب خطی مینیمکس برای $f(x)$ بر $[a, b]$ باشد، آنگاه

$$a_1 = \frac{f(b) - f(a)}{b - a} \quad a_0 = \frac{f(a) + f(c)}{2} - \left(\frac{a+c}{2}\right) \left[\frac{f(b) - f(a)}{b - a}\right]$$

که c جواب یکتای

$$f'(c) = \frac{f(b) - f(a)}{b - a}$$

است. ρ چیست؟

۱۰. (الف) چندجمله‌بیهای خطی تیلر را برای $f(x) = \ln(x)$ در $1 \leq x \leq 2$ ، در حول $x_0 = \frac{3}{2}$ پیدا کنید. نمودار خطا را رسم کنید.

(ب) تقریب خطی مینیمکس برای $f(x) = \ln(x)$ را در $[1, 2]$ پیدا کنید. نمودار خطا را رسم و آن را با تقریب تیلر مقایسه کنید.

۱۱. (الف) نشان دهید که تقریب خطی مینیمکس $\sqrt{1+x^2}$ در $[0, 1]$ چنین است

$$q_1^*(x) = 0.955 + 0.414x$$

(ب) با استفاده از (الف)، تقریب

$$\sqrt{y^2 + z^2} \approx 0.955z + 0.414y \quad 0 \leq y \leq z$$

را پیدا و خطا را تعیین کنید.

۱۲. تقریب کمترین مربعات خطی برای $f(x) = \ln(x)$ را در $[1, 2]$ به دست آورید. خطا را با نتایج مسأله ۱۰ مقایسه کنید.

۱۳. مقداری از α را پیدا کنید که رابطه زیر را مینیم کند

$$\int_0^1 |e^x - \alpha| dx$$

مقدار مینیم چیست؟ این یک مثال ساده باز راه دیگری برای اندازه‌گیری خطا در یک تقریب و بهترین تقریب آن است.

۱۴. مسائل مینیم‌سازی زیر را حل و معلوم کنید که آیا مقداری یکتا برای α که این مینیم را به دست دهد وجود دارد؟ در هر یک از حالات، α می‌تواند در حوزه همه اعداد حقیقی تغییر کند. تابع $f(x) = x$ را با چندجمله‌بیهایی به شکل αx^2 تقریب می‌زنیم.

$$\text{Min}_{\alpha} \int_{-1}^1 [x - \alpha x^2]^2 dx \text{ (الف)}$$

$$\text{Min}_{\alpha} \int_{-1}^1 |x - \alpha x^2| dx \text{ (ب)}$$

$$\text{Min}_{\alpha} \text{Max}_{-1 \leq x \leq 1} |x - \alpha x^2| dx \text{ (ج)}$$

۱۵. با استفاده از (۱۰.۴.۴) نشان دهید که $\{p_n(x)\}$ یک خانواده متعامد است و

$$n \geq 0, \|p_n\|_2 = \sqrt{2/(2n+1)}$$

۱۶. تحقیق کنید که تابعهای

$$\varphi_n(x) = \frac{(-1)^n}{n!} e^x \frac{d^n}{dx^n} (x^n e^{-x})$$

به ازای $n \geq 0$ در بازه $[0, \infty)$ نسبت به تابع وزن $w(x) = e^{-x}$ متعامدند. (توجه کنید که

$$\int_0^{\infty} e^{-x} x^m dx = m! \text{ (برای } m = 0, 1, 2, \dots \text{)})$$

۱۷. (الف) ماکسیمم و مینیمم نسبی $T_n(x)$ را در $[-1, 1]$ پیدا کنید و (۲۱.۷.۴) را به دست

آورید.

(ب) صفرهای $T_n(x)$ را پیدا کنید و (۱۴.۷.۴) را به دست آورید.

۱۸. فرمولهایی را که برای b_n و c_n در رابطه بازگشتی سهگانه (۲۱.۴.۴) از قضیه ۵.۴ داده

شده‌اند پیدا کنید.

۱۹. روند قضیه ۲.۴-گرام-اشمیت را برای اجتناب از مرحله نرمال‌سازی $\|\psi_n\|_2$ $\varphi_n = \psi_n /$

تغییر دهید:

$$\psi_n(x) = x^n + b_{n,n-1} \psi_{n-1}(x) + \dots + b_{n,0} \psi_0(x)$$

و ضرایب $b_{n,j}$ $0 \leq j \leq n-1$ را بیابید.

۲۰. با استفاده از مسأله ۱۹، ψ_0, ψ_1, ψ_2 را برای توابع وزن $w(x)$ در بازه‌های $[a, b]$ که در زیر

تعیین شده‌اند، پیدا کنید.

$$0 \leq x \leq 1, w(x) = \ln(x) \text{ (الف)}$$

$$n \geq 0, \int_0^1 x^n \ln(x) dx = [-1/(n+1)]^2 \text{ راهنمایی:}$$

$$0 \leq x \leq 1, w(x) = x \text{ (ب)}$$

$$-1 \leq x \leq 1, w(x) = \sqrt{1-x^2} \text{ (ج)}$$

۲۱. گیریم $\{\varphi_n(x) \mid n \geq 1\}$ ها در (a, b) با تابع وزن $w(x) \geq 0$ متعامد باشند. صفرهای

$\varphi_n(x)$ را با

$$a < z_{n,n} < z_{n-1,n} < \dots < z_{1,n} < b$$

نشان دهید. ثابت کنید که صفرهای $\varphi_n(x)$ صفرهای $\varphi_{n+1}(x)$ را از هم جدا می‌کنند، یعنی

$$a < z_{n+1,n+1} < z_{n,n} < z_{n,n+1} < \dots < z_{2,n+1} < z_{1,n} < z_{n,n+1} < b$$

راهنمایی: از استقرا بر درجه n استفاده کنید. بنویسید $A_n > 0$, $\varphi_n(x) = A_n x^n + \dots$ و از رابطه بازگشتی سه‌گانه (۲۱.۴.۴) برای محاسبه چندجمله‌یها در صفرهای $\varphi_n(x)$ استفاده کنید. مشاهده می‌کنید که علامت برای $\varphi_{n+1}(x)$ و $\varphi_{n-1}(x)$ تغییر می‌کند.

۲۲. اتحاد (۲۷.۴.۴) کریستوفل - داربو را برای حالت $x = y$ بسط دهید و فرمول زیر را پیدا کنید.

$$\sum_{k=0}^n \frac{[\varphi_k(x)]^2}{\gamma_k}$$

راهنمایی: حد (۲۷.۴.۴) را وقتی $y \rightarrow x$ در نظر بگیرید.

۲۳. گیریم $f(x) = \cos^{-1}(x)$ برای $-1 \leq x \leq 1$ - (شاخه اصلی $0 \leq f \leq \pi$) باشد. چندجمله‌ی از درجه دو زیر را پیدا کنید

$$p(x) = a_0 + a_1 x + a_2 x^2$$

که عبارت زیر را مینیمم سازد.

$$\int_{-1}^1 \frac{[f(x) - p(x)]^2}{\sqrt{1-x^2}} dx$$

۲۴. $S_n(x)$ را با $\frac{1}{n+1} T'_{n+1}(x)$ ، که در آن $T_{n+1}(x)$ چندجمله‌ی چیشف از درجه $n+1$ است تعریف کنید. چندجمله‌ی $S_n(x)$ را چندجمله‌ی نوع دوم چیشف گویند. (الف) نشان دهید که $\{S_n(x) \mid n \geq 0\}$ یک خانواده متعامد در $[-1, 1]$ نسبت به تابع وزن $w(x) = \sqrt{1-x^2}$ است.

(ب) نشان دهید که خانواده $\{S_n(x)\}$ مانند $T_n(x)$ در رابطه بازگشتی سه‌گانه (۱۳.۴.۴) صدق می‌کند.

(ج) برای $f \in C[-1, 1]$ مفروض، مسأله زیر را حل کنید

$$\text{Min} \int_{-1}^1 \sqrt{1-x^2} [f(x) - p_n(x)]^2 dx$$

که در آن $p_n(x)$ می‌تواند هر چندجمله‌ی از درجه نابزرگتر از n باشد.

۲۵. تعداد عملیات محاسباتی الگوریتم چپوال بخش ۵.۴ را پیدا کنید. تعداد جمعها، تعداد ضربها را به دست آورید. سپس با الگوریتم ضرب معمولی تودرتو مقایسه کنید.

۲۶. نشان دهید که چارچوب بخشهای ۴.۴ و ۵.۴ برای چندجمله‌یهای مثلثاتی از درجه نایزگتر از n نیز به کار می‌رود. نشان دهید که خانواده $\{1, \sin(x), \cos(x), \dots, \sin(nx), \cos(nx)\}$ در $[0, 2\pi]$ متعامدند. تقریب کمترین مربعات برای $f(x)$ را در این بازه با استفاده از این گونه چندجمله‌یها به دست آورید. [اگر $n \rightarrow \infty$ ، سری معروف فوریه به دست می‌آید (زیگموند (۱۹۵۹) را ببینید].

۲۷. گیریم $f(x)$ تابعی پیوسته و زوج (فرد) بر $[-a, a]$ باشد. نشان دهید که تقریب مینیمکس $q_n^*(x)$ برای $f(x)$ یک تابع زوج (فرد) بر $[-a, a]$ خواهد بود خواه n زوج باشد خواه فرد.

راهنمایی: قضیه ۱۰.۴ و یکتایی در قضیه را به کار برید.

۲۸. با استفاده از (۱۰.۶.۴) کران $\rho_n(f)$ برای توابع $f(x)$ در زیر را در بازه‌های داده شده پیدا کنید $n = 1, 2, \dots, 10$.

$$\text{(الف) } \sin(x) \quad 0 \leq x \leq \pi/2$$

$$\text{(ب) } \ln(x) \quad 1 \leq x \leq e$$

$$\text{(ج) } \tan^{-1}(x) \quad 0 \leq x \leq \pi/4$$

$$\text{(د) } e^x \quad 0 \leq x \leq 1$$

۲۹. برای بسط چیشف (۲.۷.۴)، نشان دهید که اگر $f(x)$ در $[-1, 1]$ زوج (فرد) باشد، آنگاه $c_j = 0$ اگر j فرد (زوج) باشد.

۳۰. برای تابع $f(x) = \sin[(\pi/2)x]$ ، $-1 \leq x \leq 1$ ، تقریبات کمترین مربعات درجه ۳ی لژاندر و چیشف را برای $f(x)$ پیدا کنید. خطای هریک از تقریبات را معلوم و نمودار آنها را رسم کنید. قضیه ۹.۴ را برای کراندارکردن خطای مینیمکس $\rho_3(f)$ به کار برید.

راهنمایی: انتگرالگیری عددی را برای ساختن ضرایب کمترین مربعات به کار برید. به توضیحاتی که به دنبال (۱۳.۷.۴) برای تقریب کمترین مربعات چیشف آمده است توجه کنید.

۳۱. تقریب نزدیک مینیمکس درونیاب $I_n(x)$ را برای توابع اساسی ریاضی f که در زیر داده شده‌اند در بازه‌های داده شده به‌ازای $n = 1, 2, \dots, 8$ به دست آورید. با استفاده از روندهای استانده رایانه خود، خطا را محاسبه کنید. نمودار خطا را رسم کنید، و با استفاده از قضیه ۹.۴، کرانهای بالا و پایین $\rho_n(f)$ را به دست آورید.

$$\text{(الف) } e^x \quad 0 \leq x \leq 1$$

(ب) $0 \leq x \leq \pi/2$ $\sin(x)$

(ج) $0 \leq x \leq 1$ $\tan^{-1}(x)$

(د) $1 \leq x \leq 2$ $\ln(x)$

۳۲. مسأله ۳۱ را برای تقریب نزدیک مینیمکس $F_n(x)$ تکرار کنید.

۳۳. مسأله ۳۱ و ۳۲ را برای

$$f(x) = \frac{1}{x} \int_0^x \frac{\sin(t)}{t} dt \quad 0 \leq x \leq \frac{\pi}{2}$$

به کار برید.

راهنمایی: یک تقریب درجه بالای تیلر برای محاسبه $f(x)$ پیدا کنید، سپس تبدیل (۱۲.۵.۴) را به کار برید تا یک مسأله تقریب در $[-1, 1]$ به دست آید.

۳۴. برای $f(x) = e^x$ در $[-1, 1]$ ، ساختن $I_n(x)$ را در نظر بگیرید. رابطه زیر را برای خطای $I_n(x)$ با مقادیر مناسب α_n و β_n به دست آورید.

$$\alpha_n |T_{n+1}(x)| \leq |f(x) - I_n(x)| \leq \beta_n |T_{n+1}(x)| \quad -1 \leq x \leq 1$$

کرانه‌های ناصفر بالا و پایین $\rho_n(f)$ را پیدا کنید.

۳۵. (الف) تابع $\sin(x)$ در $x = 0$ صفر می‌شود. برای آنکه یک تقریب بهتری برای آن از لحاظ خطای نسبی پیدا کنید، تابع $f(x) = \sin(x)/x$ را در نظر بگیرید. تقریب نزدیک مینیمکس $I_n(x)$ را روی $0 \leq x \leq \pi/2$ برای $n = 1, 2, \dots, 7$ محاسبه و سپس $\sin(x) - xI_n(x)$ را با نتایج مسأله ۳۱ (ب) مقایسه کنید.

(ب) قسمت (الف) را برای تابع $f(x) = \tan^{-1}(x)$ ، $0 \leq x \leq 1$ ، تکرار کنید.

۳۶. گیریم $f(x) = a_n x^n + \dots + a_1 x + a_0$ و $a_n \neq 0$. تقریب مینیمکس $f(x)$ در $[-1, 1]$ را با یک چندجمله‌ای از درجه نایبشتر از $n - 1$ پیدا کرده و همچنین $\rho_{n-1}(f)$ را به دست آورید.

۳۷. گیریم $\alpha = \text{Min} \left[\text{Max}_{|x| \leq 1} |x^6 - x^2 - p_5(x)| \right]$ ، که مینیمم روی تمام چندجمله‌یهای از درجه نایبشتر از ۵ گرفته شده است.

(الف) α را پیدا کنید.

(ب) چندجمله‌ی $p_5(x)$ را که به ازای آن مینیمم حاصل می‌شود به دست آورید.

۳۸. قضیه (۱۰.۶.۴) را ثابت کنید. راهنمایی: تقریب نزدیک مینیمکس $I_n(x)$ را در نظر بگیرید.

۳۹. (الف) برای $f(x) = e^x$ در $[-1, 1]$ چندجمله‌یی $p_4(x)$ از درجهٔ چهار تیلر را که حول $x = 0$ بسط داده شده است پیدا کنید.

(ب) با استفاده از مسألهٔ ۳۶، چندجمله‌یی مینیماکس $m_{4,3}(x)$ از درجهٔ ۳ را که $p_4(x)$ را تقریب می‌زند پیدا کنید. نمودار خطای $e^x - m_{4,3}(x)$ را رسم و آن را با خطای $e^x - p_4(x)$ که در شکل ۱.۴ نشان داده شده مقایسه کنید. فرایند تبدیل چندجمله‌یی تیلر به یک چندجمله‌یی تیلر از درجهٔ پایینتر بدین طریق را، صرفه‌جویی یا تلسکوپی کردن گویند. این فرایند معمولاً چندین بار پی‌درپی به‌کار می‌رود تا یک چندجمله‌یی درجهٔ بالای تیلر را به یک چندجمله‌یی تقریب از درجهٔ به مراتب پایینتر تبدیل کند.

۴۰. با استفاده از برنامهٔ استاندارد در مرکز رایانه‌یی خود برای محاسبهٔ تقریبهای مینیماکس، تقریب مینیماکس $q_n^*(x)$ را برای توابع داده‌شدهٔ زیر در بازه‌های داده‌شده به‌دست آورید. این عمل را برای $n = 1, 2, 3, \dots, 8$ انجام دهید. نتایج را با نتایج مسألهٔ ۳۱ مقایسه کنید.

$$\text{(الف) } e^x \quad 0 \leq x \leq 1$$

$$\text{(ب) } \sin(x) \quad 0 \leq x \leq \pi/2$$

$$\text{(ج) } \tan^{-1}(x) \quad 0 \leq x \leq 1$$

$$\text{(د) } \ln(x) \quad 1 \leq x \leq 2$$

۴۱. تقریبهای مینیماکس $q_n^*(x)$ ، $n = 1, 3, 5, 7, 9$ را برای

$$f(x) = \frac{1}{x} \int_0^x e^{-t} dt \quad -1 \leq x \leq 1$$

به‌دست آورید.

راهنمایی: ابتدا یک تقریب تیلر با دقت زیاد تهیه کنید و سپس آن را با برنامهٔ مسألهٔ ۴۰ به‌کار

برید.

۴۲. مسألهٔ ۴۱ را برای

$$f(x) = \frac{1}{x} \int_0^x \frac{\sin(t)}{t} dt \quad |x| \leq \pi$$

تکرار کنید.

انتگرالگیری عددی

در این فصل روشهای عددی محاسبه انتگرالهای معین را پیدا و تحلیل می‌کنیم. این انتگرالها معمولاً به شکل زیرند

$$I(f) = \int_a^b f(x) dx \quad (۱.۰.۵)$$

که در آنها بازه $[a, b]$ متناهی است. چنین انتگرالهایی را اغلب نمی‌توان به طور صریح محاسبه نمود، و در بسیاری از آنها به جای محاسبه دقیق و استفاده از تابع اولیه پیچیده $f(x)$ ، سریعتر آن است که از راه انتگرالگیری عددی محاسبه شوند. به تقریب $I(f)$ معمولاً انتگرالگیری عددی یا تریس گویند. روشهای عددی بسیاری برای محاسبه (۱.۰.۵) وجود دارند، ولی بیشتر آنها را می‌توان به گونه‌ای ساخت که در چارچوب ساده زیر بگنجد. برای انتگرالده $f(x)$ ، یک خانواده تقریب‌زن $\{f_n(x), n \geq 1\}$ پیدا و تعریف می‌کنیم

$$I_n(f) = \int_a^b f_n(x) dx \doteq I(f_n) \quad (۲.۰.۵)$$

معمولاً می‌خواهیم تقریب $f_n(x)$ در رابطه زیر صدق کند

$$\|f - f_n\|_{\infty} \rightarrow 0 \quad \text{هرگاه} \quad n \rightarrow \infty \quad (۳.۰.۵)$$

و شکل هر $f_n(x)$ باید به گونه‌ای انتخاب شود که $I_n(f)$ به سادگی محاسبه شود. برای خطا داریم

$$E_n(f) = I(f) - I_n(f) = \int_a^b [f(x) - f_n(x)] dx$$

$$|E_n(f)| \leq \int_a^b |f(x) - f_n(x)| dx \leq (b-a) \|f - f_n\|_\infty \quad (۴.۰.۵)$$

بیشتر روشهای انتگرالگیری عددی را می‌توان در این چارچوب نگریست اگرچه بعضی از آنها از دیدگاههای دیگر بهتر مطالعه شده‌اند. یک دسته از روشهایی که در این چارچوب نمی‌گنجد آنهایی هستند که بر پایهٔ برونمایی با استفاده از برآوردهای مجانبی خطا، انجام می‌شوند. این روشها در بخش ۴.۵ مطالعه خواهند شد.

بیشتر انتگرالهای عددی $I_n(f)$ در محاسبه دارای شکل زیرند

$$I_n(f) = \sum_{j=1}^n w_{j,n} f(x_{j,n}) \quad n \geq 1 \quad (۵.۰.۵)$$

ضرایب $w_{j,n}$ ، وزنهای انتگرالگیری یا وزنهای تربیع خوانده می‌شوند؛ و نقاط $x_{j,n}$ را گره‌های انتگرالگیری می‌نامند که معمولاً در $[a, b]$ انتخاب می‌شوند. اندیس وابستگی وزنها و گره‌ها به n معمولاً با نوشتن به صورت w_j و x_j ، حذف می‌شود، زیرا به‌طور ضمنی وابستگی به n فهمیده می‌شود. روشهای استاندارد یا دارای وزنها و گره‌هایی هستند که فرمولهای ساده‌ای دارند، یا در غیر اینصورت در جداولی داده شده‌اند که در دسترس هستند. بنابراین نیازی به ساختن صریح توابع $f_n(x)$ (۲.۰.۵) نخواهد بود، اگرچه به‌یادداشتن نقش آنها در تعریف $I_n(f)$ ممکن است مفید باشد. مثال زیر یک مثال ساده از (۲.۰.۵) - (۴.۰.۵) است، ولی به شکل (۵.۰.۵) نیست.

مثال مطلوب است محاسبهٔ

$$I = \int_0^1 \frac{e^x - 1}{x} dx \quad (۶.۰.۵)$$

این انتگرالده یک نقطهٔ تکین برداشتنی در مبدأ دارد. برای تعریف $f_n(x)$ سری تیلر را برای e^x به‌کار می‌بریم (۴.۱.۱) [فصل ۱ را ببینید]. سپس تعریف می‌کنیم

$$\begin{aligned} I_n &= \int_0^1 \sum_{j=1}^n \frac{x^{j-1}}{j!} dx \\ &= \sum_{j=1}^n \frac{1}{(j!)(j)} \end{aligned} \quad (۷.۰.۵)$$

برای خطا در I_n ، فرمول (۴.۱.۱) تیلر را به کار می‌بریم تا به ازای مقداری چون $0 \leq \xi_x \leq x$ به دست آوریم

$$f(x) - f_n(x) = \frac{x^n}{(n+1)!} e^{\xi_x}$$

در این صورت

$$I - I_n = \int_0^1 \frac{x^n}{(n+1)!} e^{\xi_x} dx$$

$$\frac{1}{(n+1)!(n+1)} \leq I - I_n \leq \frac{e}{(n+1)!(n+1)} \quad (۸.۰.۵)$$

دنباله (۷.۰.۵) سریعاً همگراست و (۸.۰.۵) به ما اجازه می‌دهد که خطا را با دقت محاسبه کنیم. مثلاً، برای $n = 6$

$$I_6 = ۱,۳۱۷۸۷۰۳۷$$

و از (۸.۰.۵)

$$۲,۸۳ \times ۱۰^{-5} \leq I - I_6 \leq ۷,۷۰ \times ۱۰^{-5}$$

خطای واقعی برابر است با $۳,۱۸ \times ۱۰^{-5}$.

برای انتگرالهایی که انتگرالده در آنها یک نوع رفتار بدی داشته باشد، مثلاً، در نقطه‌ای بینهایت باشد، اغلب انتگرالده را به شکل زیر در نظر می‌گیریم

$$I(f) = \int_a^b w(x)f(x)dx \quad (۹.۰.۵)$$

فرض شده است که رفتار بد در $w(x)$ است، که تابع وزن خوانده می‌شود و تابع $f(x)$ خوش رفتار فرض شده است. برای مثال، محاسبه

$$\int_0^1 (\ln x)f(x)dx$$

را برای توابع پیوسته دلخواه $f(x)$ در نظر می‌گیریم. چارچوب (۲.۰.۵) - (۴.۰.۵) به سادگی برای پرداختن به (۹.۰.۵) تعمیم داده می‌شود. روشهایی برای این انتگرالها در بخش ۳.۵ و ۶.۵ در نظر گرفته شده‌اند.

بیشتر فرمولهای انتگرالگیری عددی برای تعریف $f_n(x)$ در (۲.۰.۵) با استفاده از درونیایی چند جمله‌یی یا درونیایی چند جمله‌یی تکه‌یی، قرار دارند. فرمولهایی را که در آنها از چنین درونیایی با نقاط

گرهی متساوی‌الفاصله، استفاده می‌شود، در بخشهای ۱.۵ و ۲.۵ به‌دست آورده و مورد بحث قرار داده‌ایم. فرمولهای انتگرالگیری گاوسی که از بعضی نظرها بهینه هستند و همگرایی خیلی سریع دارند، در بخش ۳.۵ داده شده‌اند. این فرمولها بر اساس تعریف $f_n(x)$ با استفاده از چند جمله‌یی درونیاب در نقاط گرهی که با دقت خاصی انتخاب شده‌اند و لازم نیست متساوی‌الفاصله باشند، استخراج شده‌اند. فرمولهای مجانبی خطا برای روشهای بخشهای ۱.۵ و ۲.۵ در بخش ۴.۵ داده شده و مورد بحث قرار گرفته‌اند و بعضی فرمولهای جدید، بر پایهٔ برونیابی با این فرمول‌های خطا، به‌دست آمده‌اند. روشهایی برای کنترل خطا به‌طور خودکار، با حفظ کارایی، در بخش ۵.۵ داده شده‌اند. در بخش ۶.۵ یک بررسی کلی از روشهایی برای انتگرالهایی که تکین یا به تعبیری بد رفتار هستند انجام شده است و در بخش ۷.۵ کار دشوار مشتقگیری عددی مورد بحث قرار می‌گیرد.

۱.۵ قاعدهٔ ذوزنقه‌یی و قاعدهٔ سیمپسون

مطالعهٔ انتگرالگیری عددی را با ارائهٔ دو روش عددی کاملاً معروف برای محاسبهٔ

$$I(f) = \int_a^b f(x) dx \quad (۱.۱.۵)$$

آغاز می‌کنیم. این روشها را به‌طور کامل شرح داده و تحلیل می‌نماییم و از آنها به‌عنوان مقدمات مطالب بخشهای بعد استفاده خواهیم کرد. بازهٔ $[a, b]$ در این بخش همیشه منتهای است.

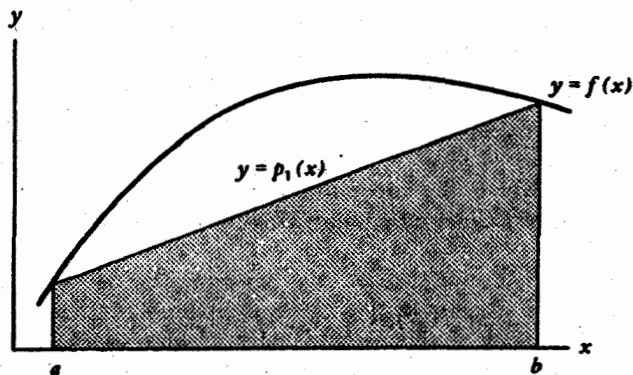
قاعدهٔ ذوزنقه‌یی قاعدهٔ ذوزنقه‌یی ساده بر پایهٔ تقریب‌زدن $f(x)$ با یک خط مستقیم، خط واصل بین دو نقطهٔ $(a, f(a))$ و $(b, f(b))$ قرار دارد. با انتگرالگیری این خط مستقیم به‌جای $f(x)$ ، تقریب زیر را به‌دست می‌آوریم

$$I_1(f) = \left(\frac{b-a}{2} \right) [f(a) + f(b)] \quad (۲.۱.۵)$$

این مقدار البته مساحت ذوزنقه‌ای است که در شکل ۱.۵ نمایش داده شده است. برای پیدا کردن فرمول خطا، فرمول خطای درونیابی (۱۰.۱.۳) را به‌کار می‌بریم

$$f(x) - \frac{(b-x)f(a) + (x-a)f(b)}{b-a} = (x-a)(x-b)f[a, b, x]$$

همچنین برای کار با خطا در قاعدهٔ ذوزنقه‌یی در این بخش فرض می‌کنیم که $f(x)$ در $[a, b]$



شکل ۱.۵ نمایش قاعده دوزنقه‌یی

دوبار پیوسته مشتقپذیر است. در این صورت

$$E_1(f) = \int_a^b f(x)dx - \frac{(b-a)}{2}[f(a) + f(b)]$$

$$= \int_a^b (x-a)(x-b)f[a, b, x]dx \quad (3.1.5)$$

با استفاده از قضیه مقدار میانگین برای انتگرالها (قضیه ۳.۱ فصل یک) و فرمول (۱۲.۲.۳) داریم

$$E_1(f) = f[a, b, \xi] \int_a^b (x-a)(x-b)dx \quad a \leq \xi \leq b$$

$$= \left[\frac{1}{2} f''(\eta) \right] \left[-\frac{1}{6}(b-a)^2 \right] \quad \eta \in [a, b]$$

بنابراین

$$E_1(f) = -\frac{(b-a)^2}{12} f''(\eta) \quad \eta \in [a, b] \quad (4.1.5)$$

اگر $b-a$ به اندازه کافی کوچک نباشد، قاعده دوزنقه‌یی (۲.۱.۵) خیلی مفید نیست. در چنین انتگرالی، آن را به مجموع انتگرالها روی زیربازه‌های کوچک تبدیل می‌کنیم، و بعد (۲.۱.۵) را برای هر یک از این انتگرالهای کوچکتر به کار می‌بریم. گیریم $n \geq 1$ و $h = (b-a)/n$ و $x_j = a + jh$ برای $j = 0, 1, \dots, n$ در این صورت

$$I(f) = \int_a^b f(x)dx = \sum_{j=1}^n \int_{x_{j-1}}^{x_j} f(x)dx$$

$$= \sum_{j=1}^n \left\{ \left(\frac{h}{2} \right) [f(x_{j-1}) + f(x_j)] - \frac{h^3}{12} f''(\eta_j) \right\}$$

با $x_{j-1} \leq \eta_j \leq x_j$. هیچ لزومی ندارد که زیربازه‌های $[x_{j-1}, x_j]$ همگی طولهای مساوی داشته باشند، ولی معمول است که ابتدا اصول کلی روش را بدین طریق معرفی نمایند. اگرچه در عمل هم این طریق متداول، ولی مواردی هست که در آن فاصلهٔ متغیر نقاط شبکه بیشتر مطلوب است. عبارات در مجموع را می‌توان ترکیب نموده قاعدهٔ ذوزنقه‌یی مرکب را به دست آورد،

$$I_n(f) = h \left[\frac{1}{2} f_0 + f_1 + f_2 + \dots + f_{n-1} + \frac{1}{2} f_n \right] \quad n \geq 1 \quad (5.1.5)$$

که در آن $f(x_j) \equiv f_j$ برای خطا در $I_n(f)$.

$$\begin{aligned} E_n(f) &= I(f) - I_n(f) = \sum_{j=1}^n -\frac{h^2}{12} f''(\eta_j) \\ &= -\frac{h^2 n}{12} \left[\frac{1}{n} \sum_{j=1}^n f''(\eta_j) \right] \end{aligned} \quad (6.1.5)$$

برای عبارت داخل کرشه

$$\text{Min}_{a \leq x \leq b} f''(x) \leq M \equiv \frac{1}{n} \sum_{j=1}^n f''(\eta_j) \leq \text{Max}_{a \leq x \leq b} f''(x)$$

چون $f''(x)$ در $a \leq x \leq b$ پیوسته است، و M بین ماکسیمم و مینیمم $f''(x)$ در $[a, b]$ است، در نقطه‌ای در $[a, b]$ مانند η ، $f''(\eta) = M$ پس می‌توانیم بنویسیم

$$E_n(f) = -\frac{(b-a)h^2}{12} f''(\eta) \quad \eta \in [a, b] \quad (7.1.5)$$

با استفاده از این تحلیل می‌توان برآورد دیگری برای خطا به دست آورد. از (۶.۱.۵) داریم

$$\begin{aligned} \text{Lim}_{n \rightarrow \infty} \frac{E_n(f)}{h^2} &= \text{Lim}_{n \rightarrow \infty} \left[-\frac{h}{12} \sum_{j=1}^n f''(\eta_j) \right] \\ &= -\frac{1}{12} \text{Lim}_{n \rightarrow \infty} \sum_{j=1}^n f''(\eta_j) h \end{aligned}$$

چون $x_{j-1} \leq \eta_j \leq x_j$ ، $j = 1, \dots, n$ ، آخرین مجموع، یک مجموع ریمان است؛ بنابراین

$$\text{Lim}_{n \rightarrow \infty} \frac{E_n(f)}{h^2} = -\frac{1}{12} \int_a^b f''(x) dx = -\frac{1}{12} [f'(b) - f'(a)] \quad (8.1.5)$$

$$E_n(f) \doteq -\frac{h^2}{12} [f'(b) - f'(a)] \equiv \tilde{E}_n(f) \quad (9.1.5)$$

عبارت $\tilde{E}_n(f)$ یک برآورد خطای مجانبی برای $E_n(f)$ نامیده می‌شود، و به تعبیر (۸.۱.۵) موّجه است.

تعریف گیریم $E_n(f)$ یک فرمول دقیق خطا و $\tilde{E}_n(f)$ برآوردی از آن باشد. گوییم که $\tilde{E}_n(f)$ برآورد خطای مجانبی برای $E_n(f)$ است اگر

$$\lim_{n \rightarrow \infty} \frac{\tilde{E}_n(f)}{E_n(f)} = 1 \quad (10.1.5)$$

یا هم‌ارز با آن

$$\lim_{n \rightarrow \infty} \frac{E_n(f) - \tilde{E}_n(f)}{E_n(f)} = 0.$$

برآورد (۹.۱.۵)، بر اساس (۸.۱.۵)، در این تعریف صدق می‌کند.

قاعدهٔ دوزنقه‌یی مرکب (۵.۱.۵) را می‌توانستیم از گذاردن تابع درونیاب خطی تکّه‌یی $f_n(x)$ به جای $f(x)$ در گره‌های x_0, x_1, \dots, x_n نیز به‌دست آوریم. از این پس قاعدهٔ دوزنقه‌یی مرکب را فقط قاعدهٔ دوزنقه‌یی می‌نامیم.

مثال قاعدهٔ دوزنقه‌یی (۹.۱.۵) را برای محاسبهٔ

$$I = \int_0^\pi e^x \cos(x) dx \quad (11.1.5)$$

به‌کار می‌بریم. مقدار صحیح آن برابر است با $12.0703463164 - (e^\pi + 1)/2 \approx 0$. مقادیر I_n در جدول ۱.۵ همراه با خطاهای درست E_n و برآوردهای مجانبی \tilde{E}_n که از (۹.۱.۵) به‌دست آمده، داده شده است. توجه نمایید که وقتی n دو برابر (و بنابراین h نصف) می‌شود خطا با مضرب ۴ کاهش می‌یابد. از عامل ضرب h^2 که در (۷.۱.۵) و (۹.۱.۵) وجود دارد این نتیجه قابل پیش‌بینی بود: وقتی h نصف می‌شود h^2 با مضرب ۴ کاهش می‌یابد. این مثال همچنین نشان می‌دهد که قاعدهٔ دوزنقه‌یی در مقایسه با سایر روشهایی که در این فصل بحث خواهد شد، چندان کارایی ندارد. با استفاده از برآورد خطای $\tilde{E}_n(f)$ ، می‌توانیم قاعدهٔ انتگرالگیری عددی بهتری را تعریف کنیم:

$$\begin{aligned} CT_n(f) &\equiv I_n(f) + \tilde{E}_n(f) \\ &= h \left[\frac{1}{3} f_0 + f_1 + \dots + f_{n-1} + \frac{1}{3} f_n \right] - \frac{h^2}{12} [f'(b) - f'(a)] \quad (12.1.5) \end{aligned}$$

این را قاعدهٔ دوزنقه‌یی تصحیح شده می‌نامند. دقت $\tilde{E}_n(f)$ موجب می‌شود که $CT_n(f)$ به مراتب دقیقتر از قاعدهٔ دوزنقه‌یی باشد. شکل دیگری برای به‌دست آوردن (۱۲.۱.۵) در مسألهٔ ۴

جدول ۱.۵ قاعده دوزنقه‌یی برای محاسبه (۱۱.۱.۵)

n	I_n	E_n	Ratio	\bar{E}
۲	-۱۷,۳۸۹,۲۵۹	۵,۳۲		۴,۹۶
۴	-۱۳,۳۳۶,۰۲۳	۱,۲۷	۴,۲۰	۱,۲۴
۸	-۱۲,۳۸۲,۱۶۲	۳,۱۲E-۱	۴,۰۶	۳,۱۰E-۱
۱۶	-۱۲,۱۴۸,۰۰۴	۷,۷۷E-۲	۴,۰۲	۷,۷۶E-۲
۳۲	-۱۲,۰۸۹۷۴۲	۱,۹۴E-۲	۴,۰۰	۱,۹۴E-۲
۶۴	۱۲,۰۷۵۱۹۴	۴,۸۵E-۳	۴,۰۰	۴,۵۸E-۳
۱۲۸	-۱۲,۰۷۱۵۵۸	۱,۲۱E-۳	۴,۰۰	۱,۲۱E-۳
۲۵۶	-۱۲,۰۷۰۶۴۹	۳,۰۳E-۴	۴,۰۰	۳,۰۳E-۴
۵۱۲	-۱۲,۰۷۰۴۲۲	۷,۵۷E-۵	۴,۰۰	۷,۵۷E-۵

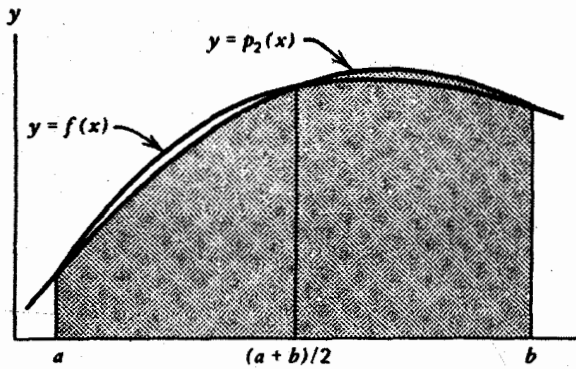
جدول ۲.۵ قاعده دوزنقه‌یی اصلاح شده برای (۱۱.۱.۵)

n	$CT_n(f)$	خطا	نسبت	خطای دوزنقه‌یی
۲	-۱۲,۴۲۵۵۲۸۳۶۷	۳,۵۵E-۱		۵,۳۲
۴	-۱۲,۰۹۵۰۹۰۱۰۶	۲,۴۷E-۲	۱۴,۴	۱,۲۷
۸	-۱۲,۰۷۱۹۲۹۲۴۵	۱,۵۸E-۳	۱۵,۶	۳,۱۲E-۱
۱۶	-۱۲,۰۷۰۴۴۵۸۰۴	۹,۹۵E-۵	۱۵,۹	۷,۷۷E-۲
۳۲	-۱۲,۰۷۰۳۵۲۵۴۳	۶,۲۳E-۶	۱۶,۰	۱,۹۴E-۲
۶۴	-۱۲,۰۷۰۳۴۶۷۰۶	۳,۸۹E-۷	۱۶,۰	۴,۸۵E-۳
۱۲۸	-۱۲,۰۷۰۳۴۶۳۴۱	۲,۴۳E-۸	۱۶,۰	۱,۲۱E-۳

توصیه شده است، که نشان می‌دهد (۱۲.۱.۵) در چارچوب تقریب نظری (۲.۰.۵) - (۴.۰.۵) درمی‌آید. مشکل اساسی در استفاده از $CT_n(f)$ نیاز به $f'(a)$ و $f'(b)$ است.

مثال $CT_n(f)$ را برای مثال قبلی (۱۱.۱.۵) به‌کار می‌بریم. نتایج همراه با خطاهای قاعده دوزنقه‌یی برای مقایسه در جدول ۲.۵ داده شده‌اند. به تجربه معلوم شده است که خطا در $CT_n(f)$ متناسب با h^4 است، در حالی که در قاعده دوزنقه‌یی متناسب با h^2 بود. یک اثبات این امر در مسأله ۴ داده شده است.

قاعده سیمپسون برای بهبود قاعده ساده دوزنقه‌یی (۲.۱.۵)، یک چندجمله‌یی درونیاب درجه دوم $p_2(x)$ را برای تقریب $f(x)$ روی $[a, b]$ به‌کار می‌بریم. گیریم $c = (a+b)/2$ ، و تعریف می‌کنیم



شکل ۲.۵ نمایش قاعده سیمپسون

$$I_r(f) = \int_a^b \left[\frac{(x-c)(x-b)}{(a-c)(a-b)} f(a) + \frac{(x-a)(x-b)}{(c-a)(c-b)} f(c) + \frac{(x-a)(x-c)}{(b-a)(b-c)} f(b) \right] dx$$

با انتگرالگیری، به دست می‌آوریم

$$I_r(f) = \frac{h}{3} \left[f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right] \quad h = \frac{b-a}{2} \quad (13.1.5)$$

این دستور قاعده سیمپسون خوانده می‌شود. نمایش آن در شکل ۲.۵ داده شده است، که قسمت سایه‌دار مساحت زیر نمودار $y = p_2(x)$ را نشان می‌دهد.

برای خطا با فرمول درونیابی خطای (۱۱.۲.۳) شروع می‌کنیم تا به دست آوریم

$$E_r(f) = I(f) - I_r(f) = \int_a^b (x-a)(x-c)(x-b)f[a, b, c, x]dx \quad (14.1.5)$$

ما نمی‌توانیم قضیه مقدار میانگین را به‌کار ببریم زیرا چندجمله‌یی انتگرالده در $x = c = (a+b)/2$ تغییر علامت می‌دهد. برای کار با قاعده سیمپسون در این فصل، فرض می‌کنیم که $f(x)$ در $[a, b]$ چهار بار پیوسته مشتقپذیر باشد. تعریف می‌کنیم

$$w(x) = \int_a^x (t-a)(t-c)(t-b)dt$$

به آسانی می‌توان نشان داد که برای $a < x < b$

$$w(x) > 0 \quad w(a) = w(b) = 0$$

با انتگرالگیری جزء به جزء داریم:

$$\begin{aligned} E_T(f) &= \int_a^b w'(x)f[a, b, c, x]dx \\ &= [w(x)f[a, b, c, x]]_{x=a}^{x=b} - \int_a^b w(x)\frac{d}{dx}f[a, b, c, x]dx \\ &= - \int_a^b w(x)f[a, b, c, x]dx \end{aligned}$$

در آخرین تساوی از رابطهٔ (۱۷.۲.۳) استفاده شده است. با استفاده از قضیهٔ مقدار میانگین انتگرالها و (۱۲.۲.۳)،

$$\begin{aligned} E_T(f) &= -f[a, b, c, \xi, \xi] \int_a^b w(x)dx & a \leq \xi \leq b \\ &= -\frac{f^{(\nu)}(\eta)}{24} \left[\frac{4}{15} h^5 \right] & h = \frac{b-a}{2} \quad \eta \in [a, b] \end{aligned}$$

بنابراین

$$E_T(f) = -\frac{h^5}{90} f^{(\nu)}(\eta) \quad \eta \in [a, b] \quad (۱۵.۱.۵)$$

از این رابطه می‌بینیم که اگر $f(x)$ یک چندجمله‌ی از درجهٔ نابزرگتر از ۳ باشد $E_T(f) = 0$ به رغم اینکه درونیابی درجهٔ دو وقتی دقیق است که $f(x)$ یک چندجمله‌ی حداکثر از درجهٔ ۲ باشد. حذف اتفاقی بیشتر خطاها در شکل ۲.۵ نموده شده است. این حذف موجب شده که قاعدهٔ سیمپسون بسیار دقیقتر از قاعدهٔ دوزنقه‌ی باشد.

بازهم یک قاعدهٔ ترکیبی می‌سازیم. برای n زوج و نا کوچکتر از ۲ تعریف می‌کنیم $x_j = a + jh, h = (b-a)/n$ و $j = 0, 1, \dots, n$. در این صورت

$$\begin{aligned} I(f) &= \int_a^b f(x)dx = \sum_{j=0}^{n/2} \int_{x_{2j-2}}^{x_{2j}} f(x)dx \\ &= \sum_{j=1}^{n/2} \left\{ \frac{h}{3} [f_{2j-2} + 4f_{2j-1} + f_{2j}] - \frac{h^5}{90} f^{(\nu)}(\eta_j) \right\} \end{aligned}$$

با $x_{2j-2} \leq \eta_j \leq x_{2j}$ جملات در مجموع را ساده می‌کنیم، قاعدهٔ ترکیبی سیمپسون به دست می‌آید:

$$\begin{aligned} I_n(f) &= \frac{h}{3} [f_0 + 4f_1 + 2f_2 + 4f_3 + 2f_4 \\ &\quad + \dots + 2f_{n-2} + 4f_{n-1} + f_n] \quad (۱۶.۱.۵) \end{aligned}$$

مانند قبل، این را فقط قاعدهٔ سیمپسون می‌نامیم. شاید بتوان گفت که این قاعدهٔ انتگرالگیری عددی است که بیشتر مورد استفاده قرار می‌گیرد. این روش ساده، به‌کار بردن آن آسان و برای طیف وسیعی از انتگرالها نسبتاً دقیق است.

برای خطا، مانند قاعدهٔ دوزنقه‌یی،

$$E_n(f) = I(f) - I_n(f) = -\frac{h^5(n/2)}{90} \cdot \frac{2}{n} \sum_{j=1}^{n/2} f^{(4)}(\eta_j)$$

$$E_n(f) = -\frac{h^5(b-a)}{180} f^{(4)}(\eta) \quad \eta \in [a, b] \quad (17.1.5)$$

می‌توانیم فرمول خطای مجانبی را نیز به‌دست آوریم

$$E_n(f) = -\frac{h^5}{180} [f^{(4)}(b) - f^{(4)}(a)] \equiv \tilde{E}_n(f) \quad (18.1.5)$$

اثبات اساساً مانند همان اثباتی است که برای پیدا کردن (۹.۱.۵) ارائه کردیم.

مثال قاعدهٔ سیمپسون (۱۶.۱.۵) را برای محاسبهٔ انتگرال (۱۱.۱.۵)، یعنی

$$I = \int_0^{\pi} e^x \cos(x) dx$$

را که قبلاً به‌عنوان یک مثال برای قاعدهٔ دوزنقه‌یی حساب کرده بودیم، به‌کار می‌بریم. نتایج عددی در جدول ۳.۵ داده شده‌اند. باز هم نرخ کاهش خطا، نتایج داده شده به وسیلهٔ (۱۷.۱.۵) و (۱۸.۱.۵) را تأیید می‌کند. در مقایسه با نتایج پیشین جدول ۱.۵ برای قاعدهٔ دوزنقه‌یی، روشن

جدول ۳.۵ قاعدهٔ سیمپسون برای محاسبهٔ (۱۱.۱.۵)

n	I_n	E_n	Ratio	\tilde{E}_n
۲	-۱۱,۵۹۲۸۳۹۵۵۳۴	-۴,۷۸E-۱		-۱,۶۳
۴	-۱۱,۹۸۴۹۴۴۰۱۹۸	-۸,۵۴E-۲	۵,۵۹	-۱,۰۲E-۱
۸	-۱۲,۰۶۴۲۰۸۹۵۷۲	-۶,۱۴E-۳	۱۴,۹	-۶,۳۸E-۳
۱۶	-۱۲,۰۶۹۹۵۱۳۲۳۳	-۳,۹۵E-۴	۱۵,۵	-۳,۹۹E-۴
۳۲	-۱۲,۰۷۰۳۲۱۴۵۶۱	-۲,۴۹E-۵	۱۵,۹	-۲,۴۹E-۵
۶۴	-۱۲,۰۷۰۳۴۴۷۵۹۹	-۱,۵۶E-۶	۱۶,۰	-۱,۵۶E-۶
۱۲۸	-۱۲,۰۷۰۳۴۶۲۱۹۱	-۹,۷۳E-۸	۱۶,۰	-۹,۷۳E-۸

است که قاعدهٔ سیمپسون بهتر است. در مقایسه با جدول ۲.۵، قاعدهٔ سیمپسون کمی بدتر است، ولی سرعت همگرایی یکی است. قاعدهٔ سیمپسون این مزیت را دارد که نیاز به مقادیر مشتق ندارد.

فرمولهای خطای هسته‌ی پتانو روش دیگری برای به دست آوردن فرمولهای خطا وجود دارد که از محاسبهٔ مشتق در یک نقطهٔ نامعلوم η نتیجه نمی‌شود. ابتدا قاعدهٔ دوزنقه‌ی را در نظر می‌گیریم. فرض کنید که $f' \in C[a, b]$ و $f''(x)$ در $[a, b]$ انتگرال‌پذیر باشند، در این صورت با استفاده از قضیهٔ تیلر [قضیهٔ ۴.۱ در فصل ۱]،

$$f(x) = p_1(x) + R_2(x) \quad p_1(x) = f(a) + (x - a)f'(a)$$

$$R_2(x) = \int_a^x (x - t)f''(t)dt$$

با توجه به (۳.۱.۵)، به ازای هر دو تابع $F, G \in C[a, b]$ داریم

$$E_1(F + G) = E_1(F) + E_1(G) \quad (۱۹.۱.۵)$$

بنابراین

$$E_1(f) = E_1(p_1) + E_1(R_2) = E_1(R_2)$$

زیرا با توجه به (۴.۱.۵) داریم $E_1(p_1) = 0$. اگر به جای $E_1(f)$ مقدار آن را بگذاریم خواهیم داشت

$$\begin{aligned} E_1(R_2) &= \int_a^b R_2(x)dx - \left(\frac{b-a}{2}\right)[R_2(a) + R_2(b)] \\ &= \int_a^b \int_a^x (x-t)f''(t)dt - \left(\frac{b-a}{2}\right) \int_a^b (b-t)f''(t)dt \end{aligned}$$

به طور کلی برای هر تابع انتگرال‌پذیر $G(x, t)$ ،

$$\int_a^b \int_a^x G(x, t)dt dx = \int_a^b \int_t^b G(x, t)dx dt \quad (۲۰.۱.۵)$$

بنابراین

$$E_1(R_2) = \int_a^b f''(t) \int_t^b (x-t)dx dt - \left(\frac{b-a}{2}\right) \int_a^b (b-t)f''(t)dt$$

اگر انتگرالها را ترکیب و نتایج را ساده کنیم، داریم

$$E_1(f) = \frac{1}{12} \int_a^b f''(t)(t-a)(t-b)dt \quad (۲۱.۱.۵)$$

برای قاعدهٔ دوزنقه‌ی مرکب (۵.۱.۵)

$$E_n(f) = \int_a^b K(t) f''(t) dt \quad (22.1.5)$$

$$K(t) = \frac{1}{4} (t - t_{j-1})(t - t_j) \quad t_{j-1} \leq t \leq t_j \quad j = 1, 2, \dots, n \quad (23.1.5)$$

فرمولهای (۲۱.۱.۵) و (۲۲.۱.۵) را عبارت خطای هسته‌ی پتانو و $K(t)$ را هستهٔ پتانو می‌نامند. برای ارائهٔ کلیتر موضوع، دیویس (۱۹۶۳، فصل ۳) را ببینید.

به‌عنوان یک مثال ساده از کاربرد آن، کرانه‌های (۲۲.۱.۵) را پیدا کرده به‌دست می‌آوریم

$$|E_n(f)| \leq \|K\|_\infty \int_a^b |f''(t)| dt = \frac{h^2}{8} \int_a^b |f''(t)| dt \quad (24.1.5)$$

اگر $f''(t)$ خیلی قلّه‌ی باشد، رابطهٔ اخیر ممکن است کرانی بهتر از (۷.۱.۵) برای خطا باشد، زیرا در (۷.۱.۵) معمولاً به‌جای $|f''(\eta)|$ باید $\|f''\|_\infty$ را بگذاریم. برای قاعدهٔ سیمپسون، قضیهٔ تیلر را به‌کار برده می‌نویسیم

$$f(x) = p_r(x) + R_r(x)$$

$$R_r(x) = \frac{1}{r!} \int_a^x (x-t)^r f^{(r)}(t) dt$$

مانند قبل

$$E_r(f) = E_r(p_r) + E_r(R_r) = E_r(R_r)$$

سپس $E_r(R_r)$ را با جایگذاری مستقیم و ساده‌کردن به‌دست می‌آوریم:

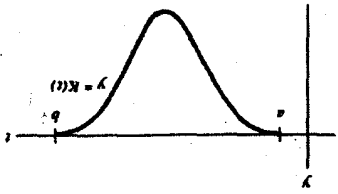
$$E_r(f) = \int_a^b R_r(x) dx - \frac{h}{3} \left[R_r(a) + 4R_r\left(\frac{a+b}{2}\right) + R_r(b) \right]$$

از اینجا نتیجه می‌شود

$$E_r(f) = \int_a^b K(t) f^{(r)}(t) dt \quad (25.1.5)$$

$$K(t) = \begin{cases} \frac{1}{42} (t-a)^2 (3t-a-2b) & a \leq t \leq \frac{a+b}{2} \\ \frac{1}{42} (b-t)^2 (b+2a-3t) & \frac{a+b}{2} \leq t \leq b \end{cases} \quad (26.1.5)$$

یک نمودار $K(t)$ در شکل ۳.۵ داده شده است. با محاسبهٔ مستقیم داریم:



شکل ۳.۵ هستهٔ پتانو برای قاعدهٔ سیمپسون

جدول ۴.۵ انتگرالگیری دوزنقه‌یی و سیمپسونی: مثال (۱)

n	قاعدهٔ دوزنقه‌یی		قاعدهٔ سیمپسون	
	خطا	نسبت	خطا	نسبت
۲	-۷,۱۹۷ - ۲		-۳,۳۷۰ - ۳	
۴	-۱,۸۱۷ - ۲	۳,۹۶	-۲,۳۱۵ - ۴	۱۴,۶
۸	-۴,۵۵۳ - ۳	۳,۹۹	-۱,۵۴۳ - ۵	۱۵,۰
۱۶	-۱,۱۳۹ - ۳	۴,۰۰	-۱,۰۰۸ - ۶	۱۵,۳
۳۲	-۲,۸۴۸ - ۴	۴,۰۰	-۶,۴۸۹ - ۸	۱۵,۵
۶۴	-۷,۱۲۱ - ۵	۴,۰۰	-۴,۱۴۱ - ۹	۱۵,۷
۱۲۸	-۱,۷۸۰ - ۵	۴,۰۰	-۲,۶۲۶ - ۱۰	۱۵,۸

$$\|K\|_{\infty} = \frac{h^2}{\sqrt{2}} \quad \int_a^b K(t) dt = -\frac{h^5}{90} \quad h = \frac{b-a}{2}$$

مانند قاعدهٔ دوزنقه‌یی مرکب، این نتایج به سادگی برای قاعدهٔ مرکب سیمپسون تعمیم داده می‌شوند. مثالهای زیر را برای توضیح کاملتر رفتار قاعده‌های سیمپسون و دوزنقه‌یی آورده‌ایم.

مثال ۱. $f(x) = x^2 \sqrt{x}, [a, b] = [0, 1], I = \frac{2}{9}$

جدول ۴.۵ خطا را برای مقادیر صعودی n به دست می‌دهد. مشتق $f^{(2)}(x)$ در $x = 0$ تکین است، بنابراین فرمول (۱۷.۱.۵) را نمی‌توان در این حالت به کار برد. به عنوان راه دیگر تعمیم (۲۵.۱.۵) را برای قاعدهٔ سیمپسون به کار می‌بریم، تا به دست آوریم

$$|E_n(f)| \leq \|K\|_{\infty} \int_a^b |f^{(2)}(t)| dt = \left(\frac{h^2}{\sqrt{2}}\right) \left(\frac{105}{8}\right) = \frac{35}{192} h^2$$

بنابراین وقتی که h نصف می‌شود (یعنی n دوبرابر می‌شود)، خطا باید ۱۶ بار کوچک شود. این روش نیز یک کران تا حدی واقع بینانه برای خطا به دست می‌دهد. توجه می‌کنید که چگونه

مقادیر تجربی در ستونهای نسبت جدول ۴.۵ و مقادیر پیش‌بینی شده نظری ۴ و ۱۶ به ترتیب، تطابق نزدیک دارند.

۲.

$$f(x) = \frac{1}{1 + (x - \pi)^2} \quad [a, b] = [0, 5] \quad I \doteq 2,33976628367$$

بر طبق نظریه، مشتق‌پذیری نامتناهی $f(x)$ ایجاب می‌کند که نسبتها در قاعده‌های دوزنقه‌یی و سیمپسون به ترتیب 4° و 16° باشد. ولی برای چند مقدار اولیه I_n ، همانگونه که جدول ۵.۵ نشان می‌دهد، لازم نیست که چنین باشد. انتگرالده نسبتاً قله‌یی است، بویژه مشتقات مراتب بالای آن، و این امر بر سرعت همگرایی اثر می‌گذارد.

۳.

$$f(x) = \sqrt{x} \quad [a, b] = [0, 1] \quad I = \frac{2}{3}$$

چون مقدار $f'(x)$ در $x = 0$ بینهایت می‌شود، هیچ یک از نتایج نظری را که قبلاً داده شده در این حالت نمی‌توان به‌کاربرد. نتایج عددی در جدول ۶.۵ داده شده‌اند. توجه می‌کنید که هنوز هم خطای یک رفتار منظمی دارد. در واقع، خطاهای دو روش، وقتی n دو برابر می‌شود به یک نسبت تنزل می‌کنند. این نسبت $2.83 \doteq 2^{1.5}$ در بخش ۴.۵، فرمول (۲۴.۴.۵)، توضیح داده شده است.

۴.

$$f(x) = e^{\cos(x)} \quad [a, b] = [0, 2\pi] \quad I \doteq 7,95492652101284$$

جدول ۵.۵ انتگرالگیری دوزنقه‌یی و سیمپسون: مثال ۲

n	قاعدهٔ دوزنقه‌یی		قاعدهٔ سیمپسون	
	خطا	نسبت	خطا	نسبت
۲	۱ - ۱,۷۳۱		۱ - ۲,۸۵۳	
۴	۲ - ۷,۱۱۰	۲,۴۳	۲ - ۳,۷۰۹	-۷,۶۹
۸	۳ - ۷,۴۹۶	۹,۴۸	۲ - ۱,۳۷۱	-۲,۷۱
۱۶	۳ - ۱,۹۵۳	۳,۸۴	۴ - ۱,۰۵۹	-۱۳۰
۳۲	۴ - ۲,۸۹۲	۳,۹۹	۶ - ۱,۰۸۰	۹,۸۱
۶۴	۴ - ۱,۲۲۳	۴,۰۰	۸ - ۶,۷۴۳	۱۶,۰
۱۲۸	۵ - ۳,۰۵۹	۴,۰۰	۹ - ۲,۲۱۷	۱۶,۰

جدول ۶.۵ انتگرالگیری دوزنقه‌یی و سیمپسون: مثال (۳)

n	قاعده دوزنقه		قاعده سیمپسون	
	خطا	نسبت	خطا	نسبت
۲	۲ - ۶,۳۱۱		۲ - ۲,۸۶۰	
۴	۲ - ۲,۳۳۸	۲,۷۰	۲ - ۱,۰۱۲	۲,۸۲
۸	۳ - ۸,۵۳۶	۲,۷۴	۳ - ۳,۵۸۷	۲,۸۳
۱۶	۳ - ۳,۰۸۵	۲,۷۷	۳ - ۱,۲۶۸	۲,۸۳
۳۲	۳ - ۱,۱۰۸	۲,۷۸	۴ - ۴,۴۸۵	۲,۸۳
۶۴	۴ - ۳,۹۵۹	۲,۸۰	۴ - ۱,۵۸۶	۲,۸۳
۱۲۸	۴ - ۱,۴۱۰	۲,۸۱	۵ - ۵,۶۰۶	۲,۸۳

جدول ۷.۵ انتگرالگیری دوزنقه‌یی و سیمپسون: مثال (۴)

n	قاعده دوزنقه‌یی		قاعده سیمپسون	
	خطا	نسبت	خطا	نسبت
۲	-۱,۷۴		۷,۲۱E - ۱	
۴	-۳,۴۴E - ۲	۵,۰۶E + ۱	۵,۳۴E - ۱	۱,۳۵
۸	-۱,۲۵E - ۶	۲,۷۵E + ۴	۱,۱۵E - ۲	۴,۶۴E + ۱
۱۶	< ۱,۰۰E - ۱۴	> ۱,۲۵E + ۸	۴,۱۷E - ۷	۲,۷۶E + ۴
۳۲			< ۱,۰۰E - ۱۴	> ۴,۱۷E + ۷

نتایج در جدول ۷.۵ نشان داده شده و بسیار خوب‌اند. هر دو روش همگرایی خیلی سریع هستند، بارچنان قاعده دوزنقه‌یی بر قاعده سیمپسون. این مثال همگرایی عالی قاعده دوزنقه‌یی را برای انتگرالده‌های دوره‌یی نشان می‌دهد، این موضوع در بخش ۴.۵ تحلیل شده است. نشانه‌ای از این رفتار را می‌توان در عبارات خطای مجانبی (۹.۱.۵) و (۱۸.۱.۵) ملاحظه کرد، زیرا هر دو برآورد $f(x)$ در چنین حالتی برابر صفرند.

۲.۵ فرمولهای انتگرالگیری نیوتن - کوتس

قاعده دوزنقه‌یی ساده (۲.۱.۵) و قاعده سیمپسون (۱۳.۱.۵) دو مثال اول فرمول انتگرالگیری نیوتن-کوتس^۱ هستند. برای $n \geq 1$ ، گیریم $n \geq 1$ ، $x_j = a + jh$ ، $h = (b-a)/n$ ، $j = 0, 1, \dots, n$.

برای تعریف $I_n(f)$ ، چندجمله‌یی درونیاب $f(x)$ در نقاط x_0, \dots, x_n را به جای $f(x)$ می‌گذاریم:

$$I(f) = \int_a^b f(x) dx \doteq I_n(f) = \int_a^b p_n(x) dx \quad (۱.۲.۵)$$

با استفاده از فرمول (۶.۱.۳) لاگرانژ برای $p_n(x)$

$$I_n(f) = \int_a^b \sum_{j=0}^n l_{j,n}(x) f(x_j) dx = \sum_{j=0}^n w_{j,n} f(x_j) \quad (۲.۲.۵)$$

با

$$w_{j,n} = \int_a^b l_{j,n}(x) dx \quad j = 0, 1, \dots, n \quad (۳.۲.۵)$$

معمولاً زیرنمایه n را حذف کرده فقط w_j می‌نویسیم. ماتاکون حالت‌های $n=1$ و $n=2$ را حساب کرده‌ایم. برای نشان دادن محاسبهٔ وزن‌ها، حالت w را برای $n=3$ مطرح می‌کنیم.

$$w_0 = \int_a^b l_0(x) dx = \int_{x_0}^{x_3} \frac{(x-x_1)(x-x_2)(x-x_3) dx}{(x_0-x_1)(x_0-x_2)(x_0-x_3)}$$

یک تبدیل متغیر محاسبات را آسان می‌کند. بگیریم $x = x_0 + \mu h$ ، $0 \leq \mu \leq 3$. در این صورت

$$\begin{aligned} w_0 &= -\frac{1}{6h^3} \int_{x_0}^{x_3} (x-x_1)(x-x_2)(x-x_3) dx \\ &= -\frac{1}{6h^3} \int_0^3 (\mu-1)h \cdot (\mu-2)h \cdot (\mu-3)h \cdot h d\mu \\ &= -\frac{h}{6} \int_0^3 (\mu-1)(\mu-2)(\mu-3) d\mu \\ w_0 &= \frac{3h}{8} \end{aligned}$$

فرمول کامل برای $n=3$ عبارت است از

$$I_3(f) = \frac{3h}{8} [f(x_0) + 3f(x_1) + 3f(x_2) + f(x_3)] \quad (۴.۲.۵)$$

که قاعدهٔ سه هاشتم نامیده می‌شود.

برای خطا، قضیه زیر را ارائه می‌دهیم.

قضیه ۱.۵ (الف) برای n زوج، فرض می‌کنیم $f(x)$ در $[a, b]$ ، $n + 2$ بار پیوسته مشتقپذیر باشد. در این صورت

$$I(f) - I_n(f) = C_n h^{n+2} f^{(n+2)}(\eta) \quad \eta \in [a, b] \quad (5.2.5)$$

با

$$C_n = \frac{1}{(n+2)!} \int_a^n \mu^2 (\mu - 1) \dots (\mu - n) d\mu \quad (6.2.5)$$

(ب) برای n فرد، فرض می‌کنیم $f(x)$ در $[a, b]$ ، $n + 1$ بار پیوسته - مشتقپذیر باشد، در این صورت

$$I(f) - I_n(f) = C_n h^{n+1} f^{(n+1)}(\eta) \quad \eta \in [a, b] \quad (7.2.5)$$

که در آن

$$C_n = \frac{1}{(n+1)!} \int_a^n \mu (\mu - 1) \dots (\mu - n) d\mu$$

برهان برهان را برای حالت (الف) که مهمترین حالت است به اختصار شرح می‌دهیم. برای اثبات کامل هر دو حالت (آیزکسون و کلر (۱۹۶۶)، صص ۳۰۸ - ۳۱۴) را ببینید. از (۱۱.۲.۳)،

$$E_n(f) = I(f) - I_n(f)$$

$$= \int_a^b (x - x_0)(x - x_1) \dots (x - x_n) f[x_0, x_1, \dots, x_n, x] dx$$

تعریف می‌کنیم

$$w(x) = \int_a^x (t - x_0) \dots (t - x_n) dt$$

پس

$$w(a) = w(b) = 0 \quad w(x) > 0 \quad \text{برای } a < x < b$$

اثبات $w(x) > 0$ را می‌توان در آیزکسون و کلر (۱۹۶۶)، صص ۳۰۹ دید. به سادگی می‌توان نشان داد که $w(b) = 0$ ، زیرا انتگرالده $(t - x_0) \dots (t - x_n)$ نسبت به گره میانی $(a+b)/2 = x_{n/2}$ تابع فرد است.

با استفاده از انتگرالگیری جزء به جزء و (۳-۲-۱۷) داریم،

$$\begin{aligned} E_n(f) &= \int_a^b w'(x) f[x_0, \dots, x_n, x] dx \\ &= [w(x) f[x_0, \dots, x_n, x]]_a^b - \int_a^b w(x) \frac{d}{dx} f[x_0, \dots, x_n, x] dx \\ E_n(f) &= - \int_a^b w(x) f[x_0, \dots, x_n, x, x] dx \end{aligned}$$

با استفاده از قضیه مقدار میانگین انتگرال و (۳-۲-۱۲)

$$\begin{aligned} E_n(f) &= -f[x_0, \dots, x_n, \xi, \xi] \int_a^b w(x) dx \\ &= -\frac{f^{(n+2)}(\eta)}{(n+2)!} \int_a^b \int_a^x (t-x_0) \dots (t-x_n) dt dx \quad (۸-۲-۵) \end{aligned}$$

ترتیب انتگرالگیری را عوض کرده و تبدیل متغیر $t = x_0 + \mu h$ ، $0 \leq \mu \leq n$ ، را به کار می‌بریم:

$$\begin{aligned} \int_a^b w(x) dx &= \int_a^b \int_t^b (t-x_0) \dots (t-x_n) dx dt \\ &= \int_{x_0}^{x_n} (t-x_0) \dots (t-x_{n-1})(t-x_n)(x_n-t) dt \\ &= -h^{n+2} \int_0^n \mu(\mu-1) \dots (\mu-n+1)(\mu-n)^2 d\mu \end{aligned}$$

تبدیل متغیر $\nu = n - \mu$ را به کار می‌بریم تا نتیجه زیر به دست آید

$$\int_a^b w(x) dx = -h^{n+2} \int_0^n (n-\nu) \dots (1-\nu) \nu^2 d\nu$$

با استفاده از این واقعیت که n زوج است و از ترکیب نتیجه آخر با (۵-۲-۵) نتایج (۵-۲-۵)

و (۶-۲-۵) به دست می‌آید. ■

به عنوان یک مرجع آسان، متداولترین فرمولهای نیوتن-کوتس در جدول ۸.۵ داده شده‌اند. برای $n = 4$ ، $I_4(f)$ را اغلب قاعده بول نامند. مانند گذشته، در جدول می‌گیریم $h = (b-a)/n$.

تعریف گویند یک فرمول انتگرالگیری $\tilde{I}(f)$ که $I(f)$ را تقریب می‌زند، دارای درجه دقت m است اگر

۱. برای همه چندجمله‌بیهای $f(x)$ از درجه نایزگتر از m ، $\tilde{I}(f) = I(f)$

۲. برای یک چندجمله‌بی f از درجه $m+1$ ، $\tilde{I}(f) \neq I(f)$

جدول ۸.۵ متداولترین فرمولهای نیوتن-کوتس

$$n = 1 \int_a^b f(x)dx = \frac{h}{2} [f(a) + f(b)] - \frac{h^3}{12} f''(\xi) \quad \text{قاعده دوزنقیه}$$

$$n = 2 \int_a^b f(x)dx = \frac{h}{3} \left[f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right] - \frac{h^5}{90} f^{(4)}(\xi) \quad \text{قاعده سیمپسون}$$

$$n = 3 \int_a^b f(x)dx = \frac{2h}{8} \left[f(a) + 3f(a+h) + 3f(b-h) + f(b) \right] - \frac{2h^7}{80} f^{(6)}(\xi)$$

$$n = 4 \int_a^b f(x)dx = \frac{2h}{45} \left[7f(a) + 32f(a+h) + 12f\left(\frac{a+b}{2}\right) + 32f(b-h) + 7f(b) \right] - \frac{8h^9}{945} f^{(8)}(\xi)$$

مثال $n = 1, 3$ در جدول ۸.۵، درجه‌های دقت نیز به ترتیب برابر است با $m = n = 1, 3$ ولی برای $n = 2, 4$ ، درجه‌های دقت به ترتیب برابر است با $m = n + 1 = 3, 5$. این روابط این نتیجه کلی را نشان می‌دهند که فرمولهای نیوتن-کوتس برای اندیس زوج n ، در مقایسه با فرمولهای با اندیس فرد n ، از درجه دقت بالاتری برخوردارند. [فرمولهای (۵.۲.۵) و (۷.۲.۵) را ملاحظه کنید]. هر یک از فرمولهای نیوتن-کوتس را می‌توان برای ساختن یک قاعده مرکب به‌کار برد. شاید سودمندترین آنها قاعده‌ای است که بر مبنای قاعده بول بنا شده است. (مسئله ۷ را ببینید). ما از ذکر جزئیات بیشتر صرف‌نظر می‌کنیم.

بحث همگرایی سؤال مورد توجه دیگر این است که وقتی $n \rightarrow \infty$ ، آیا $I_n(f)$ بسمت $I(f)$ همگرا می‌شود یا نه. با توجه به عدم همگرایی چند جمله‌بیهای درونیابی با نقاط متساوی‌الفاصله برای بعضی از انتخابهای $f(x)$ ۱۰.۵.۳ [را ببینید]، بایستی ما انتظار مشکلاتی را داشته باشیم. جدول ۹.۵ نتایج را برای یک مثال معروف نشان می‌دهد،

$$I = \int_{-4}^4 \frac{dx}{1+x^2} = 2 \cdot \tan^{-1}(4) \doteq 2,49516 \quad (9.2.5)$$

جدول ۹.۵ مثال (۹.۲.۵) نیوتن-کوتس

n	I_n
۲	۵,۴۹۰۲
۴	۲,۲۷۷۶
۶	۳,۳۲۸۸
۸	۱,۹۴۱۱
۱۰	۳,۵۹۵۶

این انتگرالهای عددی نیوتن-کوتس واگرا هستند؛ و این بیانگر آن است که فرمولهای انتگرالگیری نیوتن-کوتس $I_n(f)$ در (۲.۲.۵) لازم نیست که به $I(f)$ همگرا باشند.

برای فهمیدن دلایل عدم همگرایی فرمولهای نیوتن-کوتس برای (۹.۲.۵)، ما ابتدا یک بحث کلی همگرایی روشهای انتگرالگیری را عرضه می‌کنیم.

تعریف گیریم \mathcal{F} یک خانواده از توابع پیوسته در بازه داده شده $[a, b]$ باشد. \mathcal{F} را در $[a, b]$ چگال گوئیم اگر برای هر $f \in C[a, b]$ و هر $\epsilon > 0$ یک تابع f_ϵ در \mathcal{F} وجود داشته باشد که برای آن

$$\max_{a \leq x \leq b} |f(x) - f_\epsilon(x)| \leq \epsilon \quad (10.2.5)$$

مثال ۱. با توجه به قضیه وایرستراس [قضیه ۱.۴ را ببینید]، مجموعه همه چند جمله‌یها در $C[a, b]$ چگال است.

۲. گیریم $n \geq 1$ ، $h = (b-a)/n$ ، $x_j = a + jh$ ، $0 \leq j \leq n$. گیریم $f(x)$ بر هر زیر بازه $[x_{j-1}, x_j]$ خطی باشد. \mathcal{F} را مجموعه همه این گونه توابع خطی تکه‌یی $f(x)$ برای همه n ها تعریف می‌کنیم. اثبات چگال بودن \mathcal{F} در $C[a, b]$ را به مسأله ۱۱ واگذار می‌کنیم.

قضیه ۲.۵ گیریم

$$I_n(f) = \sum_{j=0}^n w_{j,n} f(x_{j,n}) \quad n \geq 1$$

یک دنباله از فرمولهای انتگرالگیری عددی باشد که

$$I(f) = \int_a^b f(x) dx$$

را تقریب می‌زند. گیریم \mathcal{F} یک خانواده چگال در $C[a, b]$ باشد. در این صورت

$$I_n(f) \rightarrow I(f) \quad f \in C[a, b] \quad \text{برای هر} \quad (11.2.5)$$

اگر و تنها اگر

۱.

$$I_n(f) \rightarrow I(f) \quad f \in \mathcal{F} \quad \text{برای همه} \quad (12.2.5)$$

$$B \equiv \sup_{n \geq 1} \sum_{j=0}^n |w_{j,n}| < \infty \quad (۱۳.۲.۵)$$

برهان (الف) نمایان است که، (۱۲.۲.۵) از (۱۱.۲.۵) نتیجه می‌شود. ولی اثبات نتیجه شدن (۱۳.۲.۵) از (۱۱.۲.۵) بسیار مشکل است. این یک مثال از اصل کراننداری یکنواخت است و تقریباً در همه کتابهای درسی آنالیز تابعی یافت می‌شود؛ برای مثال کرلیر (۱۹۸۲، ص ۱۲۱) را ببینید. (ب) اکنون ثابت می‌کنیم که از (۱۲.۲.۵) و (۱۳.۲.۵)، (۱۱.۲.۵) نتیجه می‌شود. گیریم $f \in C[a, b]$ داده شده باشد و $\epsilon > 0$ دلخواه باشد. با استفاده از این فرض که \mathcal{F} در $C[a, b]$ چگال است، $f_\epsilon \in \mathcal{F}$ را به گونه‌ای می‌گیریم که

$$\max_{a \leq x \leq b} |f(x) - f_\epsilon(x)| \leq \frac{\epsilon}{[2(b-a+B)]} \quad (۱۴.۲.۵)$$

سپس می‌نویسیم

$$\begin{aligned} I(f) - I_n(f) &= [I(f) - I(f_\epsilon)] + [I(f_\epsilon) - I_n(f_\epsilon)] \\ &\quad + [I_n(f_\epsilon) - I_n(f)] \end{aligned}$$

با استفاده از (۱۳.۲.۵) و (۱۴.۲.۵) به سادگی می‌توان به دست آورد

$$\begin{aligned} |I(f) - I_n(f)| &\leq |I(f) - I(f_\epsilon)| + |I(f_\epsilon) - I_n(f_\epsilon)| \\ &\quad + |I_n(f_\epsilon) - I_n(f)| \\ &\leq \frac{\epsilon}{2} + |I(f_\epsilon) - I_n(f_\epsilon)| \end{aligned}$$

با استفاده از (۱۲.۲.۵) وقتی $n \rightarrow \infty$ ، $I_n(f_\epsilon) \rightarrow I(f_\epsilon)$ ، بنابراین برای همه مقادیر n به اندازه کافی بزرگ مثلاً $n \geq n_\epsilon$

$$|I(f) - I_n(f)| \leq \epsilon$$

چون ϵ دلخواه بود، این نشان می‌دهد که وقتی $n \rightarrow \infty$ ، $I_n(f) \rightarrow I(f)$

چون انتگرالهای عددی $I_n(f)$ نیوتن-کوتس برای $f(x) = 1/(1+x^2)$ بر بازه $[-4, 4]$ به $I(f)$ همگرا نمی‌شود، باید نتیجه شود که یا شرط (۱۲.۲.۵) یا (۱۳.۲.۵) نقض شده است. اگر \mathcal{F} را مجموعه چندجمله‌یها انتخاب کنیم آنگاه (۱۲.۲.۵) برقرار است، زیرا برای هر چندجمله‌یی

p از درجه نایزگتر از n ، $I_n(p) = I(p)$. بنابراین (۱۳.۲.۵) باید نادرست باشد. برای فرمولهای (۲.۲.۵) نیوتن-کوتس

$$\sup_n \sum_{j=0}^n |w_{j,n}| = \infty \quad (15.2.5)$$

چون برای حالت خاص $I(f) = I_n(f)$ ، $f(x) \equiv 1$ برای هر n داریم

$$\sum_{j=0}^n w_{j,n} = b - a \quad n \leq 1 \quad (16.2.5)$$

از ترکیب این نتایج، علامت وزنه‌های $w_{j,n}$ وقتی n به اندازه کافی بزرگ می‌شود باید تغییر یابد. برای مثال، برای $n = 8$

$$\int_{x_0}^{x_8} f(x) dx \doteq I_8(f) = \frac{4h}{14175} [989(f_0 + f_8) + 5888(f_1 + f_7) - 928(f_2 + f_6) + 10496(f_3 + f_5) - 4540f_4]$$

چنین فرمولهایی ممکن است موجب خطاهای کاهش ارقام با معنی شوند، اگرچه بعید است تا وقتی n بزرگ نیست این یک مسئله جدی باشد. ولی به علت همین مسئله، افراد معمولاً از به‌کار بردن فرمولهای نیوتن-کوتس برای $n \geq 8$ احتراز می‌کنند، حتی از تشکیل فرمولهای مرکب. مهمترین مسئله (۲.۲.۵) در روش نیوتن-کوتس، آن است که برای انتگرالده‌های کاملاً خوشرفتار مانند (۹.۲.۵)، ممکن است انتگرال همگرا نباشد.

قاعده میانگاهی فرمول‌های نیوتن-کوتس دیگری وجود دارند که در آنها یک یا هر دو حد انتگرالگیری از نقاط شبکه درونیابی (و انتگرالگیری) حذف می‌شوند. بهترین آنها، که ساده‌ترین آنها نیز هست، قاعده میانگاهی است. این قاعده بر پایه درونیابی انتگرالده $f(x)$ با مقدار ثابت $f((a+b)/2)$ قرار دارد؛ و فرمول انتگرالگیری حاصل چنین است

$$\int_a^b f(x) dx = (b-a)f\left(\frac{a+b}{2}\right) + \frac{(b-a)^3}{24} f''(\eta) \quad \eta \in [a, b] \quad (17.2.5)$$

برای شکل مرکب آن، تعریف می‌کنیم

$$x_j = a + \left(j - \frac{1}{2}\right)h, \quad j = 1, 2, \dots, n$$

که میانگانه بازه‌های $[a + (j - 1)h, a + jh]$ است. در این صورت

$$\int_a^b f(x)dx = I_n(f) + E_n(f)$$

$$I_n(f) = h[f_1 + f_2 + \dots + f_n] \quad (18.2.5)$$

$$E_n(f) = \frac{h^2(b-a)}{24} f''(\eta) \quad \eta \in [a, b] \quad (19.2.5)$$

اثبات این نتایج به صورت مسأله ۱۰ واگذار شده است.

این فرمولهای انتگرالگیری که در آنها یک یا هر دو حد انتگرالگیری حذف شده‌اند، فرمولهای باز نیوتن-کوتس خوانده می‌شوند، و فرمولهای قبلی فرمولهای بسته نامیده می‌شوند. فرمولهای باز از مرتبه بالاتر برای استخراج فرمولهای عددی در حل معادلات دیفرانسیل معمولی به‌کار برده شده‌اند.

۳.۵ انتگرالگیری گاوسی

قاعده‌های دوزنقه‌ی مرکب و سیمپسون بر پایه استفاده از تقریب چندجمله‌ی از درجه پایین برای انتگرالده $f(x)$ در زیر بازه‌هایی با اندازه نزولی، قرار دارند. در این بخش رده‌ای از روش‌ها را بررسی می‌کنیم که در آنها از تقریبهای چندجمله‌ی $f(x)$ از درجه صعودی استفاده می‌شود. فرمولهای انتگرالگیری به دست آمده، در بیشتر حالات بسیار دقیق‌اند، و اگر فردی با محاسبه انتگرالهای زیادی روبرو باشد باید به این فرمولها جداً توجه کند.

برای عمومیت بیشتر، فرمولهای زیر را مورد توجه قرار می‌دهیم

$$I_n(f) = \sum_{j=1}^n w_{j,n} f(x_{j,n}) \doteq \int_a^b w(x) f(x) dx = I(f) \quad (1.3.5)$$

تابع وزن $w(x)$ بر $[a, b]$ غیرمنفی و انتگرالپذیر فرض می‌شود و همچنین باید در مفروضات (۸.۳.۴) و (۹.۳.۴) از بخش ۳.۴ نیز صدق کند. نقاط گرهی $\{x_{j,n}\}$ و وزنه‌های $\{w_{j,n}\}$ باید طوری انتخاب شوند که $I_n(f)$ ، برای چندجمله‌ی $f(x)$ از درجه در حد امکان بالا، برابر $I(f)$ شود. امید می‌رود که این شرایط به فرمول $I_n(f)$ ای بینجامد، که برای انتگرالده‌های $f(x)$ که توسط چندجمله‌یها به خوبی تقریب‌پذیرند، تقریباً دقیق باشد. در بخش ۲.۵ فرمولهای نیوتن-کوتس، وقتی n بزرگ می‌شود، درجه دقت افزایش می‌یابد، ولی با وجود این برای بسیاری از انتگرالده‌های خوشرفتار همگرا نیستند. مشکل در فرمولهای نیوتن-کوتس این است که نقاط شبکه $\{x_{j,n}\}$

باید متساوی الفاصله باشند. با حذف این محدودیت ما قادریم فرمولهای جدید $I_n(f)$ ای به دست آوریم که برای هر تابع $f \in C[a, b]$ همگرا باشند.

برای آنکه درک مستقیمی از پیدا کردن $I_n(f)$ به دست آوریم، حالت خاص زیر را در نظر می‌گیریم

$$\int_{-1}^1 f(x) dx \doteq \sum_{j=1}^n w_j f(x_j) \quad (۲.۳.۵)$$

که در آن $w(x) \equiv 1$ و وابستگی صریح $\{w_j\}$ و $\{x_j\}$ به n حذف شده است. وزنها $\{w_j\}$ و نقاط شبکه‌یی $\{x_j\}$ باید تعیین شوند، که خطای

$$E_n(f) = \int_{-1}^1 f(x) dx - \sum_{j=1}^n w_j f(x_j) \quad (۳.۳.۵)$$

برای چندجمله‌یی $f(x)$ از درجهٔ هرچه ممکن بالا، برابر صفر شود. برای پیدا کردن معادلات گرهها و وزنها، ابتدا توجه می‌کنیم که

$$E_n(a_0 + a_1 x + \dots + a_m x^m) = a_0 E_n(1) + a_1 E_n(x) + \dots + a_m E_n(x^m) \quad (۴.۳.۵)$$

بنابراین $E_n(f)$ برای چندجمله‌یی از درجهٔ نابزرگتر از m برابر صفر است اگر و تنها اگر

$$E_n(x^i) = 0 \quad i = 0, 1, \dots, m \quad (۵.۳.۵)$$

حالت ۱. $n = 1$. چون دو پارامتر w_1 و x_1 وجود دارند، نیاز داریم

$$E_n(1) = 0 \quad E_n(x) = 0$$

که نتیجه می‌دهد

$$\int_{-1}^1 1 dx - w_1 = 0 \quad \int_{-1}^1 x dx - w_1 x_1 = 0$$

از اینجا لازم می‌آید $w_1 = 2$ و $x_1 = 0$. پس فرمول (۲.۳.۵) چنین می‌شود

$$\int_{-1}^1 f(x) dx \doteq 2f(0)$$

که قاعدهٔ میانگاهی است.

حالت ۲. $n = 2$. چهار پارامتر موجودند، w_1, x_1, w_2, x_2 و بنابراین چهار شرط برای این پارامترها می‌گذاریم

$$E_n(x^i) = \int_{-1}^1 x^i dx - [w_1 x_1^i + w_2 x_2^i] = 0 \quad i = 0, 1, 2, 3$$

یا

$$w_1 + w_2 = 2$$

$$w_1 x_1 + w_2 x_2 = 0$$

$$w_1 x_1^2 + w_2 x_2^2 = \frac{2}{3}$$

$$w_1 x_1^3 + w_2 x_2^3 = 0$$

این معادلات به یک فرمول یکتای زیر می‌انجامند

$$\int_{-1}^1 f(x) dx \doteq f\left(-\frac{\sqrt{3}}{3}\right) + f\left(\frac{\sqrt{3}}{3}\right) \quad (6.3.5)$$

که درجه دقت آن ۳ است. این را با قاعده سیمپسون (۱۳.۱.۵) مقایسه کنید که برای رسیدن به این دقت از سه نقطه گرهی استفاده می‌کند.

حالت ۳. برای حالت کلی $2n$ پارامتر آزاد $\{x_i\}$ و $\{w_i\}$ موجودند و باید حدس بزنیم که یک فرمول (۲.۳.۵) وجود دارد که با استفاده از n گره درجه دقتی برابر $2n - 1$ خواهد داشت. معادلاتی که باید حل شوند عبارت‌اند از

$$E_n(x^i) = 0 \quad i = 0, 1, \dots, 2n - 1$$

یا

$$\sum_{j=1}^n w_j x_j^i = \begin{cases} 0 & i = 1, 3, \dots, 2n - 1 \\ \frac{2}{i+1} & i = 0, 2, \dots, 2n - 2 \end{cases} \quad (7.3.5)$$

این معادلات، غیرخطی‌اند و حلپذیری آنها به هیچ وجه ساده نیست. به علت دشواری استفاده از این دستگاه غیرخطی، از روش دیگری برای این نظریه برای (۲.۳.۵) استفاده می‌کنیم، که تا اندازه‌ای پیچ در پیچ است.

گیریم $\{\varphi_n(x) \mid n \geq 0\}$ نسبت به تابع وزن $w(x) \geq 0$ بر (a, b) چند جمله‌بیهای متعام باشند. صفرهای $\varphi_n(x)$ را با

$$a < x_1 < \dots < x_n < b \quad (8.3.5)$$

نشان می‌دهیم. همچنین از نمادگذاری $(۴-۴-۱۸)-(۴-۴-۲۰)$ داریم

$$\varphi_n(x) = A_n x^n + \dots \quad a_n = \frac{A_{n+1}}{A_n}$$

$$\gamma_n = \int_a^b w(x) [\varphi_n(x)]^2 dx \quad (۹.۳.۵)$$

قضیه ۳.۵ برای هر $n \geq 1$ ، یک فرمول انتگرالگیری عددی یکتای $(۱.۳.۵)$ از درجه دقت $2n - 1$ وجود دارد. با فرض اینکه $f(x)$ در بازه $[a, b]$ ، $2n$ بار پیوسته مشتقپذیر باشد، فرمول $I_n(f)$ و خطای آن با عبارت زیر داده می‌شود

$$\int_a^b w(x) f(x) dx = \sum_{j=1}^n w_j f(x_j) + \frac{\gamma_n}{A_n^2 (2n)!} f^{(2n)}(\eta) \quad (۱۰.۳.۵)$$

بمازای یک $a < \eta < b$. نقاط گرهی $\{x_j\}$ صفرهای $\varphi_n(x)$ هستند و وزنه‌های $\{w_j\}$ چنین‌اند؛

$$w_j = \frac{-a_n \gamma_n}{\varphi'_n(x_j) \varphi_{n+1}(x_j)} \quad j = 1, \dots, n \quad (۱۱.۳.۵)$$

برهان برهان به سه قسمت تقسیم شده است. ابتدا با استفاده از نقاط گرهی $(۸.۳.۵)$ یک فرمول با درجه دقت $2n - 1$ پیدا می‌کنیم. سپس نشان می‌دهیم که این فرمول یکتاست. سرانجام پیدا کردن فرمول خطا و وزنها را مطرح می‌کنیم.

(الف) ساختن فرمول. درونیایی ارمیت به‌عنوان ابزار برای ساختن به‌کار برده شده است (بخش ۶.۳) را برای یادآوری نمادها و نتایج ببینید. برای نقاط گرهی $(۸.۳.۵)$ ، چندجمله‌ی درونیاب ارمیت $f(x)$ و $f'(x)$ عبارت است از

$$H_n(x) = \sum_{j=1}^n f(x_j) h_j(x) + \sum_{j=1}^n f'(x_j) \tilde{h}_j(x) \quad (۱۲.۳.۵)$$

با $h_j(x)$ و $\tilde{h}_j(x)$ که در $(۲.۶.۳)$ بخش ۶.۳ تعریف شده‌اند. خطای درونیایی چنین است

$$\begin{aligned} \mathcal{E}_n(x) &= f(x) - H_n(x) = [\psi_n(x)]^2 f[x_1, x_1, \dots, x_n, x_n, x] \\ &= \frac{[\psi_n(x)]^2}{(2n)!} f^{(2n)}(\xi) \quad \xi \in [a, b] \end{aligned} \quad (۱۳.۳.۵)$$

$$\psi_n(x) = (x - x_1) \dots (x - x_n)$$

توجه کنید که

$$\psi_n(x) = \frac{\varphi_n(x)}{A_n} \quad (۱۴.۳.۵)$$

زیرا $\psi_n(x)$ و $\varphi_n(x)$ از درجه n و صفرهای آنها یکی هستند. با استفاده از (۱۲.۳.۵)، اگر $f(x)$ پیوسته مشتقپذیر باشد، آنگاه

$$\begin{aligned} \int_a^b w(x)f(x)dx &= \int_a^b w(x)H_n(x)dx + \int_a^b w(x)\mathcal{E}_n(x)dx \\ &\equiv I_n(f) + E_n(f) \end{aligned} \quad (۱۵.۳.۵)$$

درجه دقت حداقل برابر $2n - 1$ است، زیرا از (۱۳.۳.۵) نتیجه می‌شود که اگر $f(x)$ چندجمله‌یی از درجه کوچکتر از $2n$ باشد، $\mathcal{E}_n(x) = 0$. همچنین از (۱۳.۳.۵)،

$$E_n(x^{2n}) = \int_a^b w(x)\mathcal{E}_n(x)dx = \int_a^b w(x)[\psi_n(x)]^2 dx > 0 \quad (۱۶.۳.۵)$$

بنابراین درجه دقت $I_n(f)$ دقیقاً برابر $2n - 1$ است.

برای پیدا کردن یک فرمول ساده‌تر برای $I_n(f)$

$$I_n(f) = \sum_{j=1}^n f(x_j) \int_a^b w(x)h_j(x)dx + \sum_{j=1}^n f'(x_j) \int_a^b w(x)\tilde{h}_j(x)dx \quad (۱۷.۳.۵)$$

نشان می‌دهیم که تمام انتگرالها در دومین مجموع صفرند. یادآوری می‌کنیم که از (۲.۶.۳)،

$$\begin{aligned} \tilde{h}_j(x) &= (x - x_j)[l_j(x)]^2 \\ l_j(x) &= \frac{\psi_n(x)}{(x - x_j)\psi'_n(x_j)} = \frac{\varphi_n(x)}{(x - x_j)\varphi'_n(x_j)} \end{aligned}$$

در آخرین مرحله از (۱۴.۳.۵) استفاده می‌کنیم. بنابراین

$$\tilde{h}_j(x) = (x - x_j)l_j(x)l_j(x) = \frac{\varphi_n(x)l_j(x)}{\varphi'_n(x_j)} \quad (۱۸.۳.۵)$$

چون درجه (l_j) برابر $n - 1$ است و چون $\varphi_n(x)$ بر همه چندجمله‌یهای با درجه کوچکتر از n متعامد است، داریم

$$\int_a^b w(x)\tilde{h}_j(x)dx = \frac{1}{\varphi'_n(x_j)} \int_a^b w(x)\varphi_n(x)l_j(x)dx = 0 \quad j = 1, \dots, n \quad (۱۹.۳.۵)$$

فرمول انتگرالگیری (۱۵.۳.۵) به صورت زیر در می آید

$$\int_a^b w(x)f(x)dx = \sum_{j=1}^n w_j f(x_j) + E_n(f)$$

$$w_j = \int_a^b w(x)h_j(x)dx \quad j = 1, \dots, n \quad (20.3.5)$$

(ب) یکتایی فرمول (۱۹.۳.۵). فرض می کنیم که یک فرمول انتگرالگیری عددی

$$\int_a^b w(x)f(x)dx \doteq \sum_{j=1}^n v_j f(z_j) \quad (21.3.5)$$

داریم با درجه دقت نا کوچکتر از $2n - 1$. فرمول درونیایی ارمیت را برای $f(x)$ در گره های z_1, \dots, z_n می سازیم. در این صورت برای هر چند جمله ای $f(x)$ از درجه نابزرگتر از $2n - 1$,

$$f(x) = \sum_{j=1}^n f(z_j)h_j(x) + \sum_{j=1}^n f'(z_j)\tilde{h}_j(x) \quad \deg(f) \leq 2n - 1 \quad (22.3.5)$$

که $h_j(x)$ و $\tilde{h}_j(x)$ با استفاده از $\{z_j\}$ تعریف شده اند. رابطه (۲۲.۳.۵) را در $w(x)$ ضرب کرده و از فرض درباره درجه دقت (۲۱.۳.۵) استفاده می کنیم و انتگرال می گیریم تا به دست آوریم

$$\sum_{j=1}^n v_j f(z_j) = \sum_{j=1}^n f(z_j) \int_a^b w(x)h_j(x)dx + \sum_{j=1}^n f'(z_j) \int_a^b w(x)\tilde{h}_j(x)dx \quad (23.3.5)$$

برای هر چند جمله ای $f(x)$ از درجه نابزرگتر از $2n - 1$.

گیریم $f(x) = \tilde{h}_i(x)$. ویژگی (۳.۶.۳) از $\tilde{h}_i(x)$ را به کار می بریم تا از (۲۳.۳.۵) پیدا کنیم که

$$0 = \int_a^b w(x)\tilde{h}_i(x)dx \quad i = 1, \dots, n \quad (24.3.5)$$

مانند (۱۸.۳.۵)، می توانیم بنویسیم

$$\tilde{h}_i(x) = (x - z_i)[l_i(x)]^2 = \frac{l_i(x)\omega_n(x)}{\omega'_n(z_i)}$$

$$\omega(x) = (x - z_1) \dots (x - z_n)$$

در این صورت (۲۴.۳.۵) چنین می شود

$$\int_a^b w(x)\omega_n(x)l_i(x)dx = 0 \quad i = 1, 2, \dots, n$$

چون همه چندجمله‌بیهای از درجه نایزگتر از $n-1$ را می‌توان به صورت ترکیبی از $l_1(x), \dots, l_n(x)$ نوشت، $w_n(x)$ بر همه چندجمله‌بیهای از درجه نایزگتر از $n-1$ متعامد است. با استفاده از یکتایی چندجمله‌بیهای متعامد [از قضیه ۲.۴]، $w_n(x)$ مضرب ثابتی از $\varphi_n(x)$ است. بنابراین باید صفرهای آنها یکی باشند، و

$$z_i = x_i \quad i = 1, \dots, n$$

برای کامل نمودن اثبات یکتایی، باید نشان دهیم که $w_i = v_i$ که v_i وزن در (۲۱.۳.۵) و w_i وزن در (۱۰.۳.۵) است. از (۲۳.۳.۵) با (۲۴.۳.۵) و $f(x) = h_i(x)$ استفاده می‌کنیم. نتیجه بلافاصله به دست می‌آید، زیرا اکنون $h_i(x)$ با استفاده از $\{x_i\}$ ساخته شده است.

(ج) فرمول خطا. با استنتاج بعضی از ویژگیهای سودمند دیگر وزنه‌های $\{w_i\}$ در (۱۰.۳.۵) آغاز می‌کنیم. با استفاده از تعریف (۲.۶.۳) ی $h_i(x)$

$$\begin{aligned} w_i &= \int_a^b w(x)h_i(x)dx = \int_a^b w(x)[1 - 2l'_i(x_i)(x - x_i)][l_i(x)]^2 dx \\ &= \int_a^b w(x)[l_i(x)]^2 dx - 2l'_i(x_i) \int_a^b w(x)(x - x_i)[l_i(x)]^2 dx \end{aligned}$$

با توجه به (۱۹.۳.۵)، آخرین انتگرال صفر است، زیرا $\tilde{h}_i(x) = (x - x_i)[l_i(x)]^2$. بنابراین

$$w_i = \int_a^b w(x)[l_i(x)]^2 dx > 0 \quad i = 1, 2, \dots, n \quad (25.3.5)$$

و همه وزنها برای همه مقادیر n مثبت‌اند.

برای ساختن w_i ؛ با گذاشتن $f(x) = l_i(x)$ در (۲۰.۳.۵)، کار را شروع و توجه می‌کنیم که $E_n(f) = 0$ ، زیرا درجه (l_i) برابر $n-1$ است، با استفاده از $l_i(x_j) = \delta_{ij}$ خواهیم داشت

$$w_i = \int_a^b w(x)l_i(x)dx \quad i = 1, \dots, n \quad (26.3.5)$$

برای ساده‌کردن بیشتر این فرمول، اتحاد کریستوفل - داربو (قضیه ۶.۴) را می‌توان به‌کار برد، و با یک رشته عملیات فرمول (۱۱.۳.۵) را به‌دست آورد. برای شرح جزئیات آیرکسون و کلر (۱۹۶۶ صص ۳۳۳ - ۳۳۴) را ببینید.

برای خطای انتگرالگیری اگر $f(x)$, νn بار در $[a, b]$ پیوسته مشتقپذیر باشد، آنگاه

$$\begin{aligned} E_n(f) &= \int_a^b w(x) \mathcal{E}_n(x) dx \\ &= \int_a^b w(x) [\psi_n(x)]^2 f[x_1, x_1, \dots, x_n, x_n, x] dx \\ &= f[x_1, x_1, \dots, x_n, x_n, \xi] \int_a^b w(x) [\psi_n(x)]^2 dx \quad \xi \in [a, b] \end{aligned}$$

که در آخرین مرحله از قضیه مقدار میانگین انتگرال استفاده شده است. با استفاده از (۱۴.۳.۵) در آخرین انتگرال و گذاشتن مشتق به جای تفاضلات منقسم خواهیم داشت

$$E_n(f) = \frac{f^{(\nu n)}(\eta)}{(\nu n)!} \int_a^b w(x) \frac{[\varphi_n(x)]^2}{A_n^2} dx$$

که فرمول خطای (۱۰.۳.۵) را به دست می دهد.

انتگرالگیری گاوس-لژاندر برای $w(x) \equiv 1$ فرمول گاوس روی $[-1, 1]$ چنین است

$$\int_{-1}^1 f(x) dx \doteq \sum_{j=1}^n w_j f(x_j) \quad (27.3.5)$$

که گره‌ها صفرهای $P_n(x)$ ، چند جمله‌یی درجه n لژاندر، در بازه $[-1, 1]$ هستند. وزن‌ها عبارت‌اند از

$$w_i = \frac{-2}{(n+1)P'_n(x_i)P_{n+1}(x_i)} \quad i = 1, 2, \dots, n \quad (28.3.5)$$

و

$$E_n(f) = \frac{2^{2n+1}(n!)^2}{(2n+1)[(2n)!]^2} \cdot \frac{f^{(\nu n)}(\eta)}{(2n)!} = e_n \frac{f^{(\nu n)}(\eta)}{(2n)!} \quad (29.3.5)$$

برای $-1 < \eta < 1$.

برای انتگرالها در بازه‌های متناهی دیگر و تابع وزن $w(x) \equiv 1$ ، تبدیل متغیرهای خطی زیر را به کار می بریم:

$$\int_a^b f(t) dt = \left(\frac{b-a}{2} \right) \int_{-1}^1 f \left(\frac{a+b+x(b-a)}{2} \right) dx \quad (30.3.5)$$

که انتگرال را به بازه استاندارد $[-1, 1]$ بدل می کند.

برای راحتی، جدول ۱۰.۵ را که گره‌ها و وزن‌ها را برای فرمول (۲۷.۳.۵) برای مقادیر کوچک

جدول ۱۰.۵ گرههای لژاندار-گاوس و وزنها

n	x_i	ω_i
۲	± 0.5773502692	1_0
۳	± 0.7745966692 0_0	0.5555555556 0.8888888889
۴	± 0.8611363116 ± 0.3399810436	0.3478546451 0.6521451549
۵	± 0.9061798459 ± 0.5384693101 0_0	0.2369268851 0.4786286705 0.5688888889
۶	± 0.9324695142 ± 0.6612093865 ± 0.2386191861	0.1713244924 0.3607615730 0.4679139346
۷	± 0.9491079123 ± 0.7415311856 ± 0.4058451514 0_0	0.1294849662 0.2797053915 0.3818300505 0.4179591837
۸	± 0.9602898565 ± 0.7966664774 ± 0.5255324099 ± 0.1834346425	0.1012285363 0.2223810345 0.3137066459 0.3626837834

n می‌دهد، داده‌ایم. برای مقادیر بزرگ n جدول بسیار کامل در استرود و سکرست^۱ (۱۹۶۶) را ببینید، که تا $n = 512$ پیش رفته است.

مثال انتگرال (۱۱.۱.۵) را محاسبه کنید،

$$I = \int_0^{\pi} e^x \cos(x) dx = -12.0703463164$$

که قبلاً در بخش ۱.۵ به عنوان یک مثال برای قاعدهٔ دوزنقه‌یی (جدول ۱.۵) و قاعدهٔ سیمپسون (جدول ۳.۵) به کار برده شده‌است. نتایجی که در جدول ۱۱.۵ داده شده‌اند مزیت محسوس انتگرالگیری گاوسی را نشان می‌دهند.

جدول ۱۱.۵ انتگرالگیری گاوسی برای (۱۱.۱.۵)

n	I_n	$I - I_n$
۲	-۱۲٫۳۳۶۲۱۰۴۶۵۷۰	۲٫۶۶E - ۱
۳	-۱۲٫۱۲۷۴۲۰۴۵۰۱۷	۵٫۷۱E - ۲
۴	-۱۲٫۰۷۰۱۸۹۴۹۰۲۹	-۱٫۵۷E - ۴
۵	-۱۲٫۰۷۰۳۲۸۵۳۵۸۹	-۱٫۷۸E - ۵
۶	-۱۲٫۰۷۰۳۴۶۳۳۱۱۰	۱٫۴۷E - ۸
۷	-۱۲٫۰۷۰۳۴۶۳۱۷۵۳	۱٫۱۴E - ۹
۸	-۱۲٫۰۷۰۳۴۶۳۱۶۳۹	-۴٫۲۵E - ۱۳

یک قضیه کلی خطا. در تلاش برای توضیح همگرایی بسیار خوب انتگرالگیری گاوسی، قضیه مفیدی را می‌آوریم. در زیر بخش بعد، خطا در انتگرالگیری گاوس-لژاندر را با تفصیل بیشتر بررسی خواهیم کرد.

قضیه ۴.۵ فرض کنید $[a, b]$ متناهی باشد. در این صورت خطا در انتگرالگیری گاوسی،

$$E_n(f) = \int_a^b w(x)f(x)dx - \sum_{j=1}^n w_j f(x_j)$$

در رابطه زیر صدق می‌کند

$$|E_n(f)| \leq 2 \left[\int_a^b w(x)dx \right] \rho_{2n-1}(f) \quad n \geq 1 \quad (31.3.5)$$

که $\rho_{2n-1}(f)$ خطای مینیمکس (۱.۲.۴) است.

برهان برای هر چندجمله‌یی $p(x)$ از درجه نایزگتر از $2n-1$ ، $E_n(p) = 0$. همچنین تابع خطا به ازای هر $F, G \in C[a, b]$ در رابطه زیر صدق می‌کند

$$E_n(F+G) = E_n(F) + E_n(G)$$

گیریم $p(x) = q_{2n-1}^*(x)$ تقریب مینیمکس از درجه نایزگتر از $2n-1$ برای $f(x)$ بر $[a, b]$

باشد. در این صورت

$$E_n(f) = E_n(f) - E_n(q_{r_{n-1}}^*) = E_n(f - q_{r_{n-1}}^*) \\ = \int_a^b w(x)[f(x) - q_{r_{n-1}}^*(x)]dx - \sum_{j=1}^n w_j[f(x_j) - q_{r_{n-1}}^*(x_j)]$$

$$|E_n(f)| \leq \|f - q_{r_{n-1}}^*\|_\infty \left[\int_a^b w(x)dx + \sum_{j=1}^n |w_j| \right]$$

با توجه به (۲۵.۳.۵)، همه w_j ها بزرگتر از صفرند. همچنین چون $1 \equiv p(x)$ از درجه صفر است، پس

$$\sum_{j=1}^n w_j = \int_a^b w(x)dx$$

و برهان (۳۱.۳.۵) کامل می شود. ■

از نتایج بخش ۶.۴ و ۷.۴، سرعت همگرایی به صفر $\rho_m(f)$ با همواری انتگرالده، افزایش می یابد. از (۲۱.۳.۵)، همین نتیجه برای انتگرالگیری گاوسی درست است. در مقابل، قاعده دوزنقه‌ی مرکب، بدون توجه به همواری $f(x)$ ، معمولاً سریعتر از مرتبه h^2 همگرا نخواهد بود [به ویژه، اگر $f'(b) - f'(a) \neq 0$]. انتگرالگیری گاوسی، برخلاف بسیاری از قواعد انتگرالگیری مرکب از همواری بیشتر انتگرالده بهره می گیرد.

مثال قاعده انتگرالگیری گاوس - لژاندر را برای انتگرالگیری زیر در نظر می گیریم

$$I = \int_0^1 e^{-x^2} dx \approx 0.74681241322812 \quad (32.3.5)$$

جدول ۱۲.۵ شامل کرانهای خطا بر اساس فرمول (۳۱.۳.۵)،

$$|E_n(f)| \leq 2\rho_{r_{n-1}}(f) \quad (33.3.5)$$

همراه با خطای واقعی است. کران خطا تقریباً به اندازه خطای واقعی است.

بحث انتگرالگیری گاوس - لژاندر نخست سعی می کنیم عبارت خطای (۲۹.۳.۵) را بیشتر قابل فهم کنیم. ابتدا تعریف می کنیم

$$M_m = \max_{-1 \leq x \leq 1} \frac{|f^{(m)}(x)|}{m!} \quad m \geq 0 \quad (34.3.5)$$

جدول ۱۲.۵ انتگرالگیری گاوسی (۳۲.۳.۵)

n	$E_n(f)$	(۳۲.۳.۵)
۱	$-۳,۲۰E - ۲$	$۱,۰۶E - ۱$
۲	$۲,۲۹E - ۴$	$۱,۳۳E - ۳$
۳	$۹,۵۵E - ۶$	$۳,۲۴E - ۵$
۴	$-۳,۳۵E - ۷$	$۹,۲۴E - ۷$
۵	$۶,۰۵E - ۹$	$۱,۶۱E - ۸$

برای ردهٔ بزرگی از توابع نامتناهی مشتقپذیر f بر $[-۱, ۱]$ ، داریم $\text{Sup}_{m \geq 0} M_m < \infty$.
برای مثال، اگر ناحیهٔ R از صفحهٔ مختلط با

$$R = \{z : |z - x| \leq ۱, -۱ \leq x \leq ۱\}$$

تعریف شده باشد، و $f(z)$ در R تحلیلی باشد، این امر درست است. در بسیاری از توابع، هرگاه $m \rightarrow \infty$ آنگاه $M_m \rightarrow 0$ ، برای مثال $\cos x$ و e^x ، $f(x) = e^x$ از ترکیب (۲۹.۳.۵) و (۳۴.۳.۵) به دست می‌آوریم

$$|E_n(f)| \leq e_n M_{2n} \quad n \geq ۱ \quad (۳۵.۳.۵)$$

و اندازهٔ e_n در بررسی سرعت همگرایی اساسی است.

برای آشنایی بیشتر با عامل e_n ، می‌توان آن را با استفاده از فرمول استرلینگ

$$n! \doteq e^{-n} n^n \sqrt{2\pi n}$$

که به مفهوم خطای نسبی وقتی $n \rightarrow \infty$ ، درست است برآورد نمود. در این صورت چنین به دست می‌آوریم

$$e_n \doteq \frac{\pi}{4n} \quad \text{وقتی} \quad n \rightarrow \infty \quad (۳۶.۳.۵)$$

این یک برآورد کاملاً خوبی است. برای مثال، $e_5 = 0.00293$ و $e_{307} = 0.000000307$ برآورد (۳۶.۳.۵) را می‌دهد. چنانچه این برآورد با (۳۵.۳.۵) ترکیب شود، نتیجه می‌دهد:

$$|E_n(f)| \leq \frac{\pi}{4n} \cdot M_{2n} \quad (۳۷.۳.۵)$$

که وقتی $n \rightarrow \infty$ ، به مفهوم مجانبی یک کران صحیح است. این رابطه نشان می‌دهد که $E_n(f) \rightarrow 0$ یعنی بایک نرخ نمایی به صورت تابعی از n کاهش می‌یابد. این نرخ را با نرخ چندجمله‌بیهای $1/(n^2)$ و $1/(n^4)$ به ترتیب برای قاعده‌های دوزنقه‌یی و سیمپسون مقایسه نمایید. برای مطالعه انتگرالده‌هایی که بینهایت مشتق‌پذیر نیستند، می‌توانیم شکل خطای هسته‌یی پتانو را درست مانند بخش ۱.۵ برای قاعده‌های سیمپسون و دوزنقه‌یی به کار ببریم. اگر $f(x)$ در $[-1, 1]$ ، r بار مشتق‌پذیر و $f^{(r)}(x)$ در $[-1, 1]$ انتگرال‌پذیر باشد، آنگاه برای $K_{n,r}(t)$ ، یک هسته مناسب پتانو، داریم

$$E_n(f) = \int_{-1}^1 K_{n,r}(t) f^{(r)}(t) dt \quad n > \frac{r}{2} \quad (38.3.5)$$

شیوه ساختن $K_{n,r}(t)$ درست مانند ساختن هسته‌های پتانو (۲۱.۱.۵) و (۲۵.۱.۵) در بخش (۱.۵) است. با توجه به (۳۸.۳.۵) داریم

$$|E_n(f)| \leq e_{n,r} M_r$$

$$e_{n,r} = r! \int_{-1}^1 |K_{n,r}(t)| dt \quad (39.3.5)$$

مقادیر $e_{n,r}$ که در جدول ۱۳.۵ داده شده‌اند از استرود-سکرست (۱۹۶۶، صص ۱۵۲-۱۵۳) گرفته شده‌است. جدول ۱۳.۵ نشان می‌دهد که اگر f دوبار پیوسته مشتق‌پذیر باشد، انتگرالگیری گاوسی حداقل به همان سرعت قاعده دوزنقه‌یی (۵.۱.۵) همگراست. با استفاده از (۳۹.۳.۵) و جدول، می‌توانیم یک کران تجربی بسازیم

$$|E_n(f)| \leq \frac{r!}{n^2} \left[\text{Max}_{-1 \leq x \leq 1} |f''(x)| \right] \quad (40.3.5)$$

جدول ۱۳.۵ مقادیر ثابت خطای $e_{n,r}$ برای (۳۹.۳.۵)

n	$e_{n,2}$	نسبت	$e_{n,4}$	نسبت
۲	0.162		0.178	
۴	$0.437E-1$	۳٫۷	$0.647E-2$	۲۷٫۵
۸	$0.118E-1$	۳٫۷	$0.417E-3$	۱۵٫۵
۱۶	$0.311E-2$	۳٫۸	$0.279E-4$	۱۴٫۹
۳۲	$0.800E-3$	۳٫۹	$0.183E-5$	۱۵٫۳
۶۴	$0.203E-3$	۳٫۹		
۱۲۸	$0.511E-4$	۴٫۰		

از فرمول متناظر با (۷.۱.۵) برای قاعدهٔ دوزنقه‌یی در $[-1, 1]$ چنین نتیجه می‌شود

$$| \text{خطای دوزنقه‌یی} | \leq \frac{r^6}{n^2} \left[\text{Max}_{-1 \leq x \leq 1} | f'''(x) | \right]$$

که اندکی بزرگتر از (۴۰.۳.۵) است. در محاسبات واقعی، انتگرالگیری گاوسی ظاهراً همیشه بر قاعدهٔ دوزنقه‌یی برتری دارد، بجز مواردی که انتگرالده‌ها دوره‌یی باشند و بازهٔ انتگرالگیری مضرب صحیحی از دورهٔ انتگرالده باشد، مانند جدول ۷.۵. بحث مشابهی را با استفاده از جدول ۱۳.۵ با e_n برای انتگرالده‌های $f(x)$ که چهار بار مشتق‌پذیر باشند می‌توان انجام داد. (مسئلهٔ ۲۰ را ببینید).

مثال سه مثال دیگر می‌آوریم که به خوشرفتاری مثالهای جدولهای ۱۱.۵ و ۱۲.۵ نیستند. مثالهای زیر را در نظر می‌گیریم

$$I^{(1)} = \int_0^1 \sqrt{x} dx \frac{2}{3} \quad I^{(2)} = \int_0^{\pi} \frac{dx}{1 + (x - \pi)^2} = 2.339766628367$$

$$I^{(3)} = \int_0^{2\pi} e^{-x} \sin(5x) dx = 0.199954669278 \quad (41.3.5)$$

مقادیر جدول ۱۴.۵ نشان می‌دهند که انتگرالگیری گاوسی، علی‌رغم رفتار بد انتگرالده، هنوز خیلی کاراست.

نتایج $I^{(1)}$ را با نتایج جدول ۶.۵ برای قاعده‌های دوزنقه‌یی و سیمپسون مقایسه کنید. انتگرالگیری گاوسی با خطایی متناسب با $1/n^2$ همگراست و حال آنکه در جدول ۶.۵، خطا با نرخ متناسب با $1/n^{1.5}$ همگراست. انتگرال زیر را که در آن $\alpha > -1$ و غیر صحیح، $f(x)$ هموار و $f(0) \neq 0$

جدول ۱۴.۵ مثالهای انتگرالگیری گاوسی (۴۱.۳.۵)

n	$I^{(1)} - I_n^{(1)}$	نسبت	$I^{(2)} - I_n^{(2)}$	$I^{(3)} - I_n^{(3)}$
۲	$-7.22E-3$		$3.50E-1$	$3.48E-1$
۴	$-1.16E-3$	۶٫۲	$-9.19E-2$	$-1.04E-1$
۸	$-1.69E-4$	۶٫۹	$-4.03E-3$	$-1.80E-2$
۱۶	$-2.30E-5$	۷٫۴	$-6.24E-7$	$-3.34E-1$
۳۲	$-3.00E-6$	۷٫۶	$-2.98E-11$	$1.16E-1$
۶۴	$-3.84E-7$	۷٫۸	-	$1.53E-1$
۱۲۸	$-4.84E-8$	۷٫۹	-	$6.69E-15$

در نظر می‌گیریم

$$I(f) = \int_{-1}^1 x^\alpha f(x) dx \quad (42.3.5)$$

دونالدسن^۱ و الیوت^۲ (۱۹۷۲) نشان داده‌اند که خطا در انتگرالگیری گاوس - لژاندر در (۴۲.۳.۵) دارای برآورد مجانبی

$$E_n(f) \doteq \frac{c(f, \alpha)}{n^{\nu(1+\alpha)}} \quad (43.3.5)$$

است. این نتیجه با $I_n^{(1)}$ در جدول ۱۴.۵ با $\alpha = 1/2$ تطابق دارد. نتایج مهم دیگری از انتگرالگیری گاوس - لژاندر، در مقاله دونالدسن و الیوت داده شده است.

همگرایی اولیه $I_n^{(2)}$ به $I^{(2)}$ خیلی آهسته است ولی وقتی n افزایش می‌یابد، سرعت همگرایی به طور چشمگیری افزایش می‌یابد. برای $n \geq 64$ ، در محدوده حساب ماشین، $I_n^{(2)} = I^2$. این نتایج را همچنین با نتایج جدول ۵.۵ برای قاعده‌های ذوزنقه‌یی و سیمپسون مقایسه کنید.

تقریبه‌های جدول ۱۴.۵ برای $I^{(2)}$ بسیار کم است، چون نوسان انتگرالده زیاد است. انتگرالده ۱۰۱ صفر در بازه انتگرالگیری دارد. برای پیدا کردن یک مقدار دقیق $I_n^{(2)}$ ، درجه چندجمله‌یی تقریب زن در انتگرالگیری زیربنائی گاوسی باید بسیار بزرگ باشد. $I_n^{(2)}$ ، با $n = 128$ ، یک تقریب بسیار دقیق برای $I^{(2)}$ است.

توضیحات کلی انتگرالگیریهای گاوسی تعدادی نقاط قوت و ضعف دارند.

۱. به علت شکل گره‌ها و وزن‌ها و در نتیجه نیاز به استفاده از جدول، بسیاری از افراد یک فرمول ساده‌تری مانند قاعده سیمپسون، را ترجیح می‌دهند. وقتی که برای انتگرالگیری از کامپیوتر استفاده می‌شود این امر مسأله‌ای نیست. برنامه‌هایی باید نوشته شوند که شامل این وزن‌ها و نقاط گرهی برای مقادیر استاندارد n باشند، مثلاً $n = 2, 4, 6, 8, \dots, 512$ [که از استرود و سکرست (۱۹۶۶) گرفته شده‌اند]. به علاوه تعداد زیادی برنامه‌های سریع برای محاسبه نقاط گره‌ها و وزن‌ها، برای بعضی از توابع وزن متداول، وجود دارند. در بین الگوریتم‌های شناخته شده، بهتر از همه الگوریتمی است که در گلوب^۳ و ولش^۴ (۱۹۶۹) آمده است.

۲. برآورد خطا مشکل است، بنابراین معمولاً

$$I - I_n \doteq I_m - I_n \quad (44.3.5)$$

را برای یک $m > n$ انتخاب می‌کنند. مثلاً $m = n + 2$ را برای انتگرالده‌های خوشرفتار، و در غیر این صورت $m = 2n$. این نتیجه‌ای بهتر از آنچه لازم است به دست می‌دهد، حتی با افزایش تعداد ارزشیابیهای تابعی، انتگرالگیری گاوسی باز از خیلی از روشهای دیگر سریعتر است.

۳. گره‌ها برای هر فرمول I_n از گرههای فرمولهای قبلی I_m متمایزند و این امر موجب یک عدم کارایی می‌شود. اگر I_n براساس برآورد خطایی مانند (۴۴.۳.۵)، دقیق نباشد باید مقدار جدیدی برای I_n حساب کنیم. این باعث اتلاف انرژی خواهد بود، چون از هیچ یک از مقادیر قبلی انتگرالده نمی‌توان در محاسبه مجدد I_n ، استفاده کرد. این موضوع خیلی کاملتر در آخرین قسمت این بخش توضیح داده شده است و به روشهای جدیدی منجر می‌شود که این نقطه ضعف را ندارند. مع‌هذا، در بسیاری از حالات، عدم کارایی انتگرالگیری گاوسی، به علت نرخ سریع همگرایی معمولاً مهم نخواهد بود.

۴. اگر قرار باشد که رده وسیعی از انتگرالهای با ماهیت مشابه محاسبه شوند، مانند زیر عمل می‌کنیم. تعدادی انتگرال نماینده شامل یک انتگرال که احتمالاً انتگرالده آن بدترین رفتار را دارد انتخاب می‌کنیم. n را مقداری تعیین می‌کنیم که $I_n(f)$ در بین مجموعه نماینده‌ها دارای دقت کافی باشد. سپس این مقدار n را ثابت نگه می‌داریم و $I_n(f)$ را به‌عنوان انتگرال عددی تمام اعضای رده اصلی انتگرالها به‌کار می‌بریم.

۵. انتگرالگیری گاوسی برای خیلی از انتگرالده‌های نزدیک به تکین، به‌طور خیلی مؤثر عمل می‌کند، همان‌گونه که در (۴۳.۳.۵) برای (۴۲.۳.۵) نشان داده شده است. ولی تمام نقاط با رفتار تکین باید در نقاط انتهایی بازه انتگرالگیری قرار گیرند. انتگرالگیری گاوسی برای انتگرالی چون

$$\int_0^1 \sqrt{|x - 0.7|} dx$$

که دارای یک نقطه تکین در بازه انتگرالگیری است، بسیار ضعیف است. (بسیاری از روشهای انتگرالگیری عددی دیگر نیز برای این انتگرال ضعیف عمل می‌کنند). این انتگرال باید به شکل زیر شکسته و محاسبه شود

$$\int_0^{0.7} \sqrt{0.7 - x} dx + \int_{0.7}^1 \sqrt{x - 0.7} dx$$

تعمیمهایی که از نقاط گرهی مجدداً استفاده می‌کنند فرض کنید که فرمول انتگرالگیری زیر را داریم:

$$I_n(f) = \sum_{k=1}^n w_k f(x_k) \doteq \int_a^b w(x) f(x) dx \quad (45.3.5)$$

می‌خواهیم یک فرمول انتگرالگیری جدیدی بسازیم که در آن از n نقطه x_1, \dots, x_n و m نقطه جدید x_{n+1}, \dots, x_{n+m} استفاده شود:

$$I_{n+m}(f) = \sum_{k=1}^{n+m} v_k f(x_k) \doteq \int_a^b w(x) f(x) dx \quad (۴۶.۳.۵)$$

این $n + 2m$ پارامتر نامعین، یعنی گره‌های x_{n+1}, \dots, x_{n+m} و وزنهای u_1, \dots, u_{n+m} باید طوری انتخاب شوند که درجه دقت (۴۵.۳.۵) در حد امکان زیاد باشد. ما دنبال یک فرمول با درجه دقت $n + 2m - 1$ هستیم. اینکه آیا چنین فرمولی را می‌توان با گره‌های جدید x_{n+1}, \dots, x_{n+m} که در $[a, b]$ قرار دارند، تعیین کرد معلوم نیست.

برای جالنتی که (۴۵.۳.۵) یک فرمول گاوسی است، کرونرود تعمیمهای (۴۶.۳.۵) را برای $m = n + 1$ مطالعه نموده است. این نوع فرمولهای زوج، راه کم هزینه‌تری (در مقایسه با استفاده از یک قاعده گاوسی با $n + 1$ نقطه گرهی) برای تولید برآورد خطا برای قاعده گاوسی به دست می‌دهد و درجه دقت آنقدر بالاست که نوع دقت متناظر با قاعده‌های گاوسی را ایجاد می‌کند.

شکل دیگری از موضوع اخیر در پاترسن^۱ (۱۹۶۸) ارائه شده است. برای $w(x) \equiv 1$ ، او با یک قاعده گاوس-لژاندر $I_n(f)$ شروع می‌نماید. سپس با تکرار فرمولهای (۴۶.۳.۵) دنباله‌ای از فرمولها را از عناصر قبلی دنباله با $m = n + 1$ به وجود می‌آورد. یکی از مقاله‌های پاترسن (۱۹۷۳) شامل الگوریتمی است بر پایه دنباله قواعد $I_2, I_4, I_6, I_8, I_{10}, I_{12}, I_{14}, I_{16}, I_{18}, I_{20}, I_{22}, I_{24}, I_{26}, I_{28}, I_{30}$ ؛ فرمول I_2 قاعده گاوس سه نقطه‌ای است. دنباله دیگری از این نوع $\{I_{10}, I_{21}, I_{43}, I_{87}\}$ توسط پیسنز^۲ و همکاران (۱۹۸۳، صص ۱۹، ۲۶، ۲۷) داده شده است، که I_{10} قاعده گاوس ده نقطه‌ای است. تمام فرمولهای این چینی پاترسن تاکنون، همه دارای نقاط گرهی داخل بازه انتگرالگیری با وزنهای مثبت بوده‌اند.

درجه دقت قاعده‌های پاترسن با تعداد نقاط گرهی افزایش می‌یابد. برای دنباله $I_2, I_4, I_6, \dots, I_{20}, I_{22}, \dots$ که قبلاً اشاره شد، درجه دقت به ترتیب چنین‌اند ۳، ۵، ۷، ۹، ۱۱، ۱۳، ۱۵، ۱۷، ۱۹، ۲۱، ۲۳، ۲۵، ۲۷، ۲۹، ۳۱، ۳۳، ۳۵، ۳۷، ۳۹، ۴۱، ۴۳، ۴۵، ۴۷، ۴۹، ۵۱، ۵۳، ۵۵، ۵۷، ۵۹، ۶۱، ۶۳، ۶۵، ۶۷، ۶۹، ۷۱، ۷۳، ۷۵، ۷۷، ۷۹، ۸۱، ۸۳، ۸۵، ۸۷، ۸۹، ۹۱، ۹۳، ۹۵، ۹۷، ۹۹، ۱۰۱، ۱۰۳، ۱۰۵، ۱۰۷، ۱۰۹، ۱۱۱، ۱۱۳، ۱۱۵، ۱۱۷، ۱۱۹، ۱۲۱، ۱۲۳، ۱۲۵، ۱۲۷، ۱۲۹، ۱۳۱، ۱۳۳، ۱۳۵، ۱۳۷، ۱۳۹، ۱۴۱، ۱۴۳، ۱۴۵، ۱۴۷، ۱۴۹، ۱۵۱، ۱۵۳، ۱۵۵، ۱۵۷، ۱۵۹، ۱۶۱، ۱۶۳، ۱۶۵، ۱۶۷، ۱۶۹، ۱۷۱، ۱۷۳، ۱۷۵، ۱۷۷، ۱۷۹، ۱۸۱، ۱۸۳، ۱۸۵، ۱۸۷، ۱۸۹، ۱۹۱، ۱۹۳، ۱۹۵، ۱۹۷، ۱۹۹، ۲۰۱، ۲۰۳، ۲۰۵، ۲۰۷، ۲۰۹، ۲۱۱، ۲۱۳، ۲۱۵، ۲۱۷، ۲۱۹، ۲۲۱، ۲۲۳، ۲۲۵، ۲۲۷، ۲۲۹، ۲۳۱، ۲۳۳، ۲۳۵، ۲۳۷، ۲۳۹، ۲۴۱، ۲۴۳، ۲۴۵، ۲۴۷، ۲۴۹، ۲۵۱، ۲۵۳، ۲۵۵، ۲۵۷، ۲۵۹، ۲۶۱، ۲۶۳، ۲۶۵، ۲۶۷، ۲۶۹، ۲۷۱، ۲۷۳، ۲۷۵، ۲۷۷، ۲۷۹، ۲۸۱، ۲۸۳، ۲۸۵، ۲۸۷، ۲۸۹، ۲۹۱، ۲۹۳، ۲۹۵، ۲۹۷، ۲۹۹، ۳۰۱، ۳۰۳، ۳۰۵، ۳۰۷، ۳۰۹، ۳۱۱، ۳۱۳، ۳۱۵، ۳۱۷، ۳۱۹، ۳۲۱، ۳۲۳، ۳۲۵، ۳۲۷، ۳۲۹، ۳۳۱، ۳۳۳، ۳۳۵، ۳۳۷، ۳۳۹، ۳۴۱، ۳۴۳، ۳۴۵، ۳۴۷، ۳۴۹، ۳۵۱، ۳۵۳، ۳۵۵، ۳۵۷، ۳۵۹، ۳۶۱، ۳۶۳، ۳۶۵، ۳۶۷، ۳۶۹، ۳۷۱، ۳۷۳، ۳۷۵، ۳۷۷، ۳۷۹، ۳۸۱، ۳۸۳، ۳۸۵، ۳۸۷، ۳۸۹، ۳۹۱، ۳۹۳، ۳۹۵، ۳۹۷، ۳۹۹، ۴۰۱، ۴۰۳، ۴۰۵، ۴۰۷، ۴۰۹، ۴۱۱، ۴۱۳، ۴۱۵، ۴۱۷، ۴۱۹، ۴۲۱، ۴۲۳، ۴۲۵، ۴۲۷، ۴۲۹، ۴۳۱، ۴۳۳، ۴۳۵، ۴۳۷، ۴۳۹، ۴۴۱، ۴۴۳، ۴۴۵، ۴۴۷، ۴۴۹، ۴۵۱، ۴۵۳، ۴۵۵، ۴۵۷، ۴۵۹، ۴۶۱، ۴۶۳، ۴۶۵، ۴۶۷، ۴۶۹، ۴۷۱، ۴۷۳، ۴۷۵، ۴۷۷، ۴۷۹، ۴۸۱، ۴۸۳، ۴۸۵، ۴۸۷، ۴۸۹، ۴۹۱، ۴۹۳، ۴۹۵، ۴۹۷، ۴۹۹، ۵۰۱، ۵۰۳، ۵۰۵، ۵۰۷، ۵۰۹، ۵۱۱، ۵۱۳، ۵۱۵، ۵۱۷، ۵۱۹، ۵۲۱، ۵۲۳، ۵۲۵، ۵۲۷، ۵۲۹، ۵۳۱، ۵۳۳، ۵۳۵، ۵۳۷، ۵۳۹، ۵۴۱، ۵۴۳، ۵۴۵، ۵۴۷، ۵۴۹، ۵۵۱، ۵۵۳، ۵۵۵، ۵۵۷، ۵۵۹، ۵۶۱، ۵۶۳، ۵۶۵، ۵۶۷، ۵۶۹، ۵۷۱، ۵۷۳، ۵۷۵، ۵۷۷، ۵۷۹، ۵۸۱، ۵۸۳، ۵۸۵، ۵۸۷، ۵۸۹، ۵۹۱، ۵۹۳، ۵۹۵، ۵۹۷، ۵۹۹، ۶۰۱، ۶۰۳، ۶۰۵، ۶۰۷، ۶۰۹، ۶۱۱، ۶۱۳، ۶۱۵، ۶۱۷، ۶۱۹، ۶۲۱، ۶۲۳، ۶۲۵، ۶۲۷، ۶۲۹، ۶۳۱، ۶۳۳، ۶۳۵، ۶۳۷، ۶۳۹، ۶۴۱، ۶۴۳، ۶۴۵، ۶۴۷، ۶۴۹، ۶۵۱، ۶۵۳، ۶۵۵، ۶۵۷، ۶۵۹، ۶۶۱، ۶۶۳، ۶۶۵، ۶۶۷، ۶۶۹، ۶۷۱، ۶۷۳، ۶۷۵، ۶۷۷، ۶۷۹، ۶۸۱، ۶۸۳، ۶۸۵، ۶۸۷، ۶۸۹، ۶۹۱، ۶۹۳، ۶۹۵، ۶۹۷، ۶۹۹، ۷۰۱، ۷۰۳، ۷۰۵، ۷۰۷، ۷۰۹، ۷۱۱، ۷۱۳، ۷۱۵، ۷۱۷، ۷۱۹، ۷۲۱، ۷۲۳، ۷۲۵، ۷۲۷، ۷۲۹، ۷۳۱، ۷۳۳، ۷۳۵، ۷۳۷، ۷۳۹، ۷۴۱، ۷۴۳، ۷۴۵، ۷۴۷، ۷۴۹، ۷۵۱، ۷۵۳، ۷۵۵، ۷۵۷، ۷۵۹، ۷۶۱، ۷۶۳، ۷۶۵، ۷۶۷، ۷۶۹، ۷۷۱، ۷۷۳، ۷۷۵، ۷۷۷، ۷۷۹، ۷۸۱، ۷۸۳، ۷۸۵، ۷۸۷، ۷۸۹، ۷۹۱، ۷۹۳، ۷۹۵، ۷۹۷، ۷۹۹، ۸۰۱، ۸۰۳، ۸۰۵، ۸۰۷، ۸۰۹، ۸۱۱، ۸۱۳، ۸۱۵، ۸۱۷، ۸۱۹، ۸۲۱، ۸۲۳، ۸۲۵، ۸۲۷، ۸۲۹، ۸۳۱، ۸۳۳، ۸۳۵، ۸۳۷، ۸۳۹، ۸۴۱، ۸۴۳، ۸۴۵، ۸۴۷، ۸۴۹، ۸۵۱، ۸۵۳، ۸۵۵، ۸۵۷، ۸۵۹، ۸۶۱، ۸۶۳، ۸۶۵، ۸۶۷، ۸۶۹، ۸۷۱، ۸۷۳، ۸۷۵، ۸۷۷، ۸۷۹، ۸۸۱، ۸۸۳، ۸۸۵، ۸۸۷، ۸۸۹، ۸۹۱، ۸۹۳، ۸۹۵، ۸۹۷، ۸۹۹، ۹۰۱، ۹۰۳، ۹۰۵، ۹۰۷، ۹۰۹، ۹۱۱، ۹۱۳، ۹۱۵، ۹۱۷، ۹۱۹، ۹۲۱، ۹۲۳، ۹۲۵، ۹۲۷، ۹۲۹، ۹۳۱، ۹۳۳، ۹۳۵، ۹۳۷، ۹۳۹، ۹۴۱، ۹۴۳، ۹۴۵، ۹۴۷، ۹۴۹، ۹۵۱، ۹۵۳، ۹۵۵، ۹۵۷، ۹۵۹، ۹۶۱، ۹۶۳، ۹۶۵، ۹۶۷، ۹۶۹، ۹۷۱، ۹۷۳، ۹۷۵، ۹۷۷، ۹۷۹، ۹۸۱، ۹۸۳، ۹۸۵، ۹۸۷، ۹۸۹، ۹۹۱، ۹۹۳، ۹۹۵، ۹۹۷، ۹۹۹، ۱۰۰۱، ۱۰۰۳، ۱۰۰۵، ۱۰۰۷، ۱۰۰۹، ۱۰۱۱، ۱۰۱۳، ۱۰۱۵، ۱۰۱۷، ۱۰۱۹، ۱۰۲۱، ۱۰۲۳، ۱۰۲۵، ۱۰۲۷، ۱۰۲۹، ۱۰۳۱، ۱۰۳۳، ۱۰۳۵، ۱۰۳۷، ۱۰۳۹، ۱۰۴۱، ۱۰۴۳، ۱۰۴۵، ۱۰۴۷، ۱۰۴۹، ۱۰۵۱، ۱۰۵۳، ۱۰۵۵، ۱۰۵۷، ۱۰۵۹، ۱۰۶۱، ۱۰۶۳، ۱۰۶۵، ۱۰۶۷، ۱۰۶۹، ۱۰۷۱، ۱۰۷۳، ۱۰۷۵، ۱۰۷۷، ۱۰۷۹، ۱۰۸۱، ۱۰۸۳، ۱۰۸۵، ۱۰۸۷، ۱۰۸۹، ۱۰۹۱، ۱۰۹۳، ۱۰۹۵، ۱۰۹۷، ۱۰۹۹، ۱۱۰۱، ۱۱۰۳، ۱۱۰۵، ۱۱۰۷، ۱۱۰۹، ۱۱۱۱، ۱۱۱۳، ۱۱۱۵، ۱۱۱۷، ۱۱۱۹، ۱۱۲۱، ۱۱۲۳، ۱۱۲۵، ۱۱۲۷، ۱۱۲۹، ۱۱۳۱، ۱۱۳۳، ۱۱۳۵، ۱۱۳۷، ۱۱۳۹، ۱۱۴۱، ۱۱۴۳، ۱۱۴۵، ۱۱۴۷، ۱۱۴۹، ۱۱۵۱، ۱۱۵۳، ۱۱۵۵، ۱۱۵۷، ۱۱۵۹، ۱۱۶۱، ۱۱۶۳، ۱۱۶۵، ۱۱۶۷، ۱۱۶۹، ۱۱۷۱، ۱۱۷۳، ۱۱۷۵، ۱۱۷۷، ۱۱۷۹، ۱۱۸۱، ۱۱۸۳، ۱۱۸۵، ۱۱۸۷، ۱۱۸۹، ۱۱۹۱، ۱۱۹۳، ۱۱۹۵، ۱۱۹۷، ۱۱۹۹، ۱۲۰۱، ۱۲۰۳، ۱۲۰۵، ۱۲۰۷، ۱۲۰۹، ۱۲۱۱، ۱۲۱۳، ۱۲۱۵، ۱۲۱۷، ۱۲۱۹، ۱۲۲۱، ۱۲۲۳، ۱۲۲۵، ۱۲۲۷، ۱۲۲۹، ۱۲۳۱، ۱۲۳۳، ۱۲۳۵، ۱۲۳۷، ۱۲۳۹، ۱۲۴۱، ۱۲۴۳، ۱۲۴۵، ۱۲۴۷، ۱۲۴۹، ۱۲۵۱، ۱۲۵۳، ۱۲۵۵، ۱۲۵۷، ۱۲۵۹، ۱۲۶۱، ۱۲۶۳، ۱۲۶۵، ۱۲۶۷، ۱۲۶۹، ۱۲۷۱، ۱۲۷۳، ۱۲۷۵، ۱۲۷۷، ۱۲۷۹، ۱۲۸۱، ۱۲۸۳، ۱۲۸۵، ۱۲۸۷، ۱۲۸۹، ۱۲۹۱، ۱۲۹۳، ۱۲۹۵، ۱۲۹۷، ۱۲۹۹، ۱۳۰۱، ۱۳۰۳، ۱۳۰۵، ۱۳۰۷، ۱۳۰۹، ۱۳۱۱، ۱۳۱۳، ۱۳۱۵، ۱۳۱۷، ۱۳۱۹، ۱۳۲۱، ۱۳۲۳، ۱۳۲۵، ۱۳۲۷، ۱۳۲۹، ۱۳۳۱، ۱۳۳۳، ۱۳۳۵، ۱۳۳۷، ۱۳۳۹، ۱۳۴۱، ۱۳۴۳، ۱۳۴۵، ۱۳۴۷، ۱۳۴۹، ۱۳۵۱، ۱۳۵۳، ۱۳۵۵، ۱۳۵۷، ۱۳۵۹، ۱۳۶۱، ۱۳۶۳، ۱۳۶۵، ۱۳۶۷، ۱۳۶۹، ۱۳۷۱، ۱۳۷۳، ۱۳۷۵، ۱۳۷۷، ۱۳۷۹، ۱۳۸۱، ۱۳۸۳، ۱۳۸۵، ۱۳۸۷، ۱۳۸۹، ۱۳۹۱، ۱۳۹۳، ۱۳۹۵، ۱۳۹۷، ۱۳۹۹، ۱۴۰۱، ۱۴۰۳، ۱۴۰۵، ۱۴۰۷، ۱۴۰۹، ۱۴۱۱، ۱۴۱۳، ۱۴۱۵، ۱۴۱۷، ۱۴۱۹، ۱۴۲۱، ۱۴۲۳، ۱۴۲۵، ۱۴۲۷، ۱۴۲۹، ۱۴۳۱، ۱۴۳۳، ۱۴۳۵، ۱۴۳۷، ۱۴۳۹، ۱۴۴۱، ۱۴۴۳، ۱۴۴۵، ۱۴۴۷، ۱۴۴۹، ۱۴۵۱، ۱۴۵۳، ۱۴۵۵، ۱۴۵۷، ۱۴۵۹، ۱۴۶۱، ۱۴۶۳، ۱۴۶۵، ۱۴۶۷، ۱۴۶۹، ۱۴۷۱، ۱۴۷۳، ۱۴۷۵، ۱۴۷۷، ۱۴۷۹، ۱۴۸۱، ۱۴۸۳، ۱۴۸۵، ۱۴۸۷، ۱۴۸۹، ۱۴۹۱، ۱۴۹۳، ۱۴۹۵، ۱۴۹۷، ۱۴۹۹، ۱۵۰۱، ۱۵۰۳، ۱۵۰۵، ۱۵۰۷، ۱۵۰۹، ۱۵۱۱، ۱۵۱۳، ۱۵۱۵، ۱۵۱۷، ۱۵۱۹، ۱۵۲۱، ۱۵۲۳، ۱۵۲۵، ۱۵۲۷، ۱۵۲۹، ۱۵۳۱، ۱۵۳۳، ۱۵۳۵، ۱۵۳۷، ۱۵۳۹، ۱۵۴۱، ۱۵۴۳، ۱۵۴۵، ۱۵۴۷، ۱۵۴۹، ۱۵۵۱، ۱۵۵۳، ۱۵۵۵، ۱۵۵۷، ۱۵۵۹، ۱۵۶۱، ۱۵۶۳، ۱۵۶۵، ۱۵۶۷، ۱۵۶۹، ۱۵۷۱، ۱۵۷۳، ۱۵۷۵، ۱۵۷۷، ۱۵۷۹، ۱۵۸۱، ۱۵۸۳، ۱۵۸۵، ۱۵۸۷، ۱۵۸۹، ۱۵۹۱، ۱۵۹۳، ۱۵۹۵، ۱۵۹۷، ۱۵۹۹، ۱۶۰۱، ۱۶۰۳، ۱۶۰۵، ۱۶۰۷، ۱۶۰۹، ۱۶۱۱، ۱۶۱۳، ۱۶۱۵، ۱۶۱۷، ۱۶۱۹، ۱۶۲۱، ۱۶۲۳، ۱۶۲۵، ۱۶۲۷، ۱۶۲۹، ۱۶۳۱، ۱۶۳۳، ۱۶۳۵، ۱۶۳۷، ۱۶۳۹، ۱۶۴۱، ۱۶۴۳، ۱۶۴۵، ۱۶۴۷، ۱۶۴۹، ۱۶۵۱، ۱۶۵۳، ۱۶۵۵، ۱۶۵۷، ۱۶۵۹، ۱۶۶۱، ۱۶۶۳، ۱۶۶۵، ۱۶۶۷، ۱۶۶۹، ۱۶۷۱، ۱۶۷۳، ۱۶۷۵، ۱۶۷۷، ۱۶۷۹، ۱۶۸۱، ۱۶۸۳، ۱۶۸۵، ۱۶۸۷، ۱۶۸۹، ۱۶۹۱، ۱۶۹۳، ۱۶۹۵، ۱۶۹۷، ۱۶۹۹، ۱۷۰۱، ۱۷۰۳، ۱۷۰۵، ۱۷۰۷، ۱۷۰۹، ۱۷۱۱، ۱۷۱۳، ۱۷۱۵، ۱۷۱۷، ۱۷۱۹، ۱۷۲۱، ۱۷۲۳، ۱۷۲۵، ۱۷۲۷، ۱۷۲۹، ۱۷۳۱، ۱۷۳۳، ۱۷۳۵، ۱۷۳۷، ۱۷۳۹، ۱۷۴۱، ۱۷۴۳، ۱۷۴۵، ۱۷۴۷، ۱۷۴۹، ۱۷۵۱، ۱۷۵۳، ۱۷۵۵، ۱۷۵۷، ۱۷۵۹، ۱۷۶۱، ۱۷۶۳، ۱۷۶۵، ۱۷۶۷، ۱۷۶۹، ۱۷۷۱، ۱۷۷۳، ۱۷۷۵، ۱۷۷۷، ۱۷۷۹، ۱۷۸۱، ۱۷۸۳، ۱۷۸۵، ۱۷۸۷، ۱۷۸۹، ۱۷۹۱، ۱۷۹۳، ۱۷۹۵، ۱۷۹۷، ۱۷۹۹، ۱۸۰۱، ۱۸۰۳، ۱۸۰۵، ۱۸۰۷، ۱۸۰۹، ۱۸۱۱، ۱۸۱۳، ۱۸۱۵، ۱۸۱۷، ۱۸۱۹، ۱۸۲۱، ۱۸۲۳، ۱۸۲۵، ۱۸۲۷، ۱۸۲۹، ۱۸۳۱، ۱۸۳۳، ۱۸۳۵، ۱۸۳۷، ۱۸۳۹، ۱۸۴۱، ۱۸۴۳، ۱۸۴۵، ۱۸۴۷، ۱۸۴۹، ۱۸۵۱، ۱۸۵۳، ۱۸۵۵، ۱۸۵۷، ۱۸۵۹، ۱۸۶۱، ۱۸۶۳، ۱۸۶۵، ۱۸۶۷، ۱۸۶۹، ۱۸۷۱، ۱۸۷۳، ۱۸۷۵، ۱۸۷۷، ۱۸۷۹، ۱۸۸۱، ۱۸۸۳، ۱۸۸۵، ۱۸۸۷، ۱۸۸۹، ۱۸۹۱، ۱۸۹۳، ۱۸۹۵، ۱۸۹۷، ۱۸۹۹، ۱۹۰۱، ۱۹۰۳، ۱۹۰۵، ۱۹۰۷، ۱۹۰۹، ۱۹۱۱، ۱۹۱۳، ۱۹۱۵، ۱۹۱۷، ۱۹۱۹، ۱۹۲۱، ۱۹۲۳، ۱۹۲۵، ۱۹۲۷، ۱۹۲۹، ۱۹۳۱، ۱۹۳۳، ۱۹۳۵، ۱۹۳۷، ۱۹۳۹، ۱۹۴۱، ۱۹۴۳، ۱۹۴۵، ۱۹۴۷، ۱۹۴۹، ۱۹۵۱، ۱۹۵۳، ۱۹۵۵، ۱۹۵۷، ۱۹۵۹، ۱۹۶۱، ۱۹۶۳، ۱۹۶۵، ۱۹۶۷، ۱۹۶۹، ۱۹۷۱، ۱۹۷۳، ۱۹۷۵، ۱۹۷۷، ۱۹۷۹، ۱۹۸۱، ۱۹۸۳، ۱۹۸۵، ۱۹۸۷، ۱۹۸۹، ۱۹۹۱، ۱۹۹۳، ۱۹۹۵، ۱۹۹۷، ۱۹۹۹، ۲۰۰۱، ۲۰۰۳، ۲۰۰۵، ۲۰۰۷، ۲۰۰۹، ۲۰۱۱، ۲۰۱۳، ۲۰۱۵، ۲۰۱۷، ۲۰۱۹، ۲۰۲۱، ۲۰۲۳، ۲۰۲۵، ۲۰۲۷، ۲۰۲۹، ۲۰۳۱، ۲۰۳۳، ۲۰۳۵، ۲۰۳۷، ۲۰۳۹، ۲۰۴۱، ۲۰۴۳، ۲۰۴۵، ۲۰۴۷، ۲۰۴۹، ۲۰۵۱، ۲۰۵۳، ۲۰۵۵، ۲۰۵۷، ۲۰۵۹، ۲۰۶۱، ۲۰۶۳، ۲۰۶۵، ۲۰۶۷، ۲۰۶۹، ۲۰۷۱، ۲۰۷۳، ۲۰۷۵، ۲۰۷۷، ۲۰۷۹، ۲۰۸۱، ۲۰۸۳، ۲۰۸۵، ۲۰۸۷، ۲۰۸۹، ۲۰۹۱، ۲۰۹۳، ۲۰۹۵، ۲۰۹۷، ۲۰۹۹، ۲۱۰۱، ۲۱۰۳، ۲۱۰۵، ۲۱۰۷، ۲۱۰۹، ۲۱۱۱، ۲۱۱۳، ۲۱۱۵، ۲۱۱۷، ۲۱۱۹، ۲۱۲۱، ۲۱۲۳، ۲۱۲۵، ۲۱۲۷، ۲۱۲۹، ۲۱۳۱، ۲۱۳۳، ۲۱۳۵، ۲۱۳۷، ۲۱۳۹، ۲۱۴۱، ۲۱۴۳، ۲۱۴۵، ۲۱۴۷، ۲۱۴۹، ۲۱۵۱، ۲۱۵۳، ۲۱۵۵، ۲۱۵۷، ۲۱۵۹، ۲۱۶۱، ۲۱۶۳، ۲۱۶۵، ۲۱۶۷، ۲۱۶۹، ۲۱۷۱، ۲۱۷۳، ۲۱۷۵، ۲۱۷۷، ۲۱۷۹، ۲۱۸۱، ۲۱۸۳، ۲۱۸۵، ۲۱۸۷، ۲۱۸۹، ۲۱۹۱، ۲۱۹۳، ۲۱۹۵، ۲۱۹۷، ۲۱۹۹، ۲۲۰۱، ۲۲۰۳، ۲۲۰۵، ۲۲۰۷، ۲۲۰۹، ۲۲۱۱، ۲۲۱۳، ۲۲۱۵، ۲۲۱۷، ۲۲۱۹، ۲۲۲۱، ۲۲۲۳، ۲۲۲۵، ۲۲۲۷، ۲۲۲۹، ۲۲۳۱، ۲۲۳۳، ۲۲۳۵، ۲۲۳۷، ۲۲۳۹، ۲۲۴۱، ۲۲۴۳، ۲۲۴۵، ۲۲۴۷، ۲۲۴۹، ۲۲۵۱، ۲۲۵۳، ۲۲۵۵، ۲۲۵۷، ۲۲۵۹، ۲۲۶۱، ۲۲۶۳، ۲۲۶۵، ۲۲۶۷، ۲۲۶۹، ۲۲۷۱، ۲۲۷۳، ۲۲۷۵، ۲۲۷۷، ۲۲۷۹، ۲۲۸۱، ۲۲۸۳، ۲۲۸۵، ۲۲۸۷، ۲۲۸۹، ۲۲۹۱، ۲۲۹۳، ۲۲۹۵، ۲۲۹۷، ۲۲۹۹، ۲۳۰۱، ۲۳۰۳، ۲۳۰۵، ۲۳۰۷، ۲۳۰۹، ۲۳۱۱، ۲۳۱۳، ۲۳۱۵، ۲۳۱۷، ۲۳۱۹، ۲۳۲۱، ۲۳۲۳، ۲۳۲۵، ۲۳۲۷، ۲۳۲۹، ۲۳۳۱، ۲۳۳۳، ۲۳۳۵، ۲۳۳۷، ۲۳۳۹، ۲۳۴۱، ۲۳۴۳، ۲۳۴۵، ۲۳۴۷، ۲۳۴۹، ۲۳۵۱، ۲۳۵۳، ۲۳۵۵، ۲۳۵۷، ۲۳۵۹، ۲۳۶۱، ۲۳۶۳، ۲۳۶۵، ۲۳۶۷، ۲۳۶۹، ۲۳۷۱، ۲۳۷۳، ۲۳۷۵، ۲۳۷۷، ۲۳۷۹، ۲۳۸۱، ۲۳۸۳، ۲۳۸۵، ۲۳۸۷، ۲۳۸۹، ۲۳۹۱، ۲۳۹۳، ۲۳۹۵، ۲۳۹۷، ۲۳۹۹، ۲۴۰۱، ۲۴۰۳، ۲۴۰۵، ۲۴۰۷، ۲۴۰۹، ۲۴۱۱، ۲۴۱۳، ۲۴۱۵، ۲۴۱۷، ۲۴۱۹، ۲۴۲۱، ۲۴۲۳، ۲۴۲۵، ۲۴۲۷، ۲۴۲۹، ۲۴۳۱، ۲۴۳۳، ۲۴۳۵، ۲۴۳۷، ۲۴۳۹، ۲۴۴۱، ۲۴۴۳، ۲۴۴۵، ۲۴۴۷، ۲۴۴۹، ۲۴۵۱، ۲۴۵۳، ۲۴۵۵، ۲۴۵۷، ۲۴۵۹، ۲۴۶۱، ۲۴

مثال گیریم (۴۵.۳.۵) قاعده سه نقطه‌یی گاوسی بر $[-1, 1]$ باشد:

$$I_3(f) = \frac{1}{9} f(0) + \frac{5}{9} [f(-\sqrt{0.6}) + f(\sqrt{0.6})] \quad (47.3.5)$$

قاعده کرونرود برای آن چنین است

$$I_4(f) = \alpha_0 f(0) + \alpha_1 [f(-\sqrt{0.6}) + f(\sqrt{0.6})] \\ + \alpha_2 [f(-\beta_1) + f(\beta_1)] + \alpha_3 [f(-\beta_2) + f(\beta_2)] \quad (48.3.5)$$

که β_1^2 و β_2^2 به ترتیب ریشه‌های کوچکتر و بزرگتر معادله زیر هستند.

$$x^2 - \frac{10}{9}x + \frac{155}{191} = 0$$

وزنه‌های $\alpha_0, \alpha_1, \alpha_2, \alpha_3$ از انتگرالگیری چندجمله‌یی $p_4(x)$ لاگرانژ بر $[-1, 1]$ که $f(x)$ را در نقاط $\{0, \pm\sqrt{0.6}, \pm\beta_1, \pm\beta_2\}$ درونیابی می‌نماید، به دست آمده‌اند. مقادیر تقریبی عبارتند از

$$\alpha_0 = 0.450916538658 \quad \alpha_1 = 0.268488089868$$

$$\alpha_2 = 0.401397414776 \quad \alpha_3 = 0.104656226026$$

۴.۵ فرمولهای خطای مجانبی و کاربردهای آنها

تعریف فرمول خطای مجانبی (۱۰.۱.۵) را برای یک فرمول انتگرالگیری عددی به یاد می‌آوریم:

$\tilde{E}_n(f)$ یک فرمول خطای مجانبی برای $I_n(f)$ است $E_n(f) = I(f) - I_n(f)$

$$\lim_{n \rightarrow \infty} \frac{\tilde{E}_n(f)}{E_n(f)} = 1 \quad (1.4.5)$$

یا هم‌ارز با آن

$$\lim_{n \rightarrow \infty} \frac{E_n(f) - \tilde{E}_n(f)}{E_n(f)} = 0$$

(۹.۱.۵) و (۱۸.۱.۵) از بخش ۱.۵ مثالهایی در این مورد هستند.

با به دست آوردن فرمول خطای مجانبی، شکل یا ساختار خطا را به دست می‌آوریم. با این اطلاعات، می‌توانیم یا خطا در $I_n(f)$ را مثل در جدولهای ۱.۵ و ۳.۵ برآورد نماییم، یا یک فرمول جدید و دقیقتری مانند قاعده دوزنقه‌یی اصلاح شده در (۱۲.۱.۵) پیدا کنیم. در این بخش، این هر دو حالت بعداً نشان داده شده و به روش انتگرالگیری همگرای سریع رامبرگ ختم شده‌است. مطلب را با بسط بیشتر فرمولهای خطای مجانبی آغاز می‌کنیم.

چندجمله‌بیهای برنولی برای استفاده در قضیه بعد، چندجمله‌بیهای برنولی $B_n(x)$ را $n \geq 0$ معرفی می‌کنیم. این چندجمله‌بیها به‌طور ضمنی با تابع مولد زیر تعریف شده‌اند

$$\frac{t(e^{xt} - 1)}{e^t - 1} = \sum_{j=1}^{\infty} B_j(x) \frac{t^j}{j!} \quad (۲.۴.۵)$$

چند تایی اول این چند جمله‌بیها به قرار زیرند:

$$\begin{aligned} B_0(x) &= 1, & B_1(x) &= x, & B_2(x) &= x^2 - x, \\ B_3(x) &= x^3 - \frac{3x^2}{2} + \frac{x}{2}, & B_4(x) &= x^4(1-x)^2 \end{aligned} \quad (۳.۴.۵)$$

در این چندجمله‌بیها

$$B_k(0) = 0 \quad k \geq 1 \quad (۴.۴.۵)$$

برای محاسبه این چندجمله‌بیها روابط بازگشتی که به سادگی قابل محاسبه‌اند وجود دارند (مسئله ۲۳ را ببینید).

اعداد برنولی نیز، که به‌طور ضمنی با رابطه

$$\frac{t}{e^t - 1} = \sum_{j=0}^{\infty} B_j \cdot \frac{t^j}{j!} \quad (۵.۴.۵)$$

تعریف می‌شوند، مورد توجه ما هستند. چند تایی اول این اعداد چنین‌اند:

$$\begin{aligned} B_0 &= 1 & B_1 &= -\frac{1}{2} & B_2 &= \frac{1}{6} & B_3 &= -\frac{1}{30} \\ B_4 &= \frac{1}{42} & B_5 &= -\frac{1}{42} \end{aligned} \quad (۶.۴.۵)$$

و برای تمام اعداد صحیح فرد $j \geq 3$ ، $B_j = 0$. برای پیدا کردن یک رابطه بین چندجمله‌بیهای $B_j(x)$ برنولی، از (۲.۴.۵) نسبت به x در $[0, 1]$ انتگرال می‌گیریم. در این صورت

$$1 - \frac{t}{e^t - 1} = \sum_{j=1}^{\infty} \frac{t^j}{j!} \int_0^1 B_j(x) dx$$

و بنابراین

$$B_j = - \int_0^1 B_j(x) dx \quad j \geq 1 \quad (۷.۴.۵)$$

ما همچنین به یک تعریف توسیع دوره‌یی از $B_j(x)$ نیاز داریم

$$\bar{B}_j(x) = \begin{cases} B_j(x) & 0 \leq x \leq 1 \\ \bar{B}_j(x-1) & x \geq 1 \end{cases} \quad (۸.۴.۵)$$

فرمول اوایلر-مک لورن قضیه زیر یک فرمول خیلی مفصلی برای مقدار مجانبی خطا در قاعده دوزنقه‌یی به دست می‌دهد. این قضیه اساس بسیاری از تحلیلهای مجانبی خطا، در این بخش، است. ارتباط آن با بعضی از فرمولهای دیگر انتگرالگیری، بعداً در این بخش ظاهر خواهد شد.

قضیه ۵.۵ (فرمول اوایلر-مک لورن) گیریم $n \geq 1, m \geq 0$ ، و تعریف می‌کنیم $h = (b-a)/n$. همچنین فرض می‌کنیم $f(x)$ $2m+2$ بار پیوسته مشتقپذیر بر $[a, b]$ برای مقداری از $m \geq 0$ باشد. در این صورت برای خطا در قاعده دوزنقه‌یی،

$$\begin{aligned} E_m(f) &= \int_a^b f(x) dx - h \sum_{j=0}^n f(x_j) \\ &= - \sum_{i=1}^m \frac{B_{2i}}{(2i)!} h^{2i} [f^{(2i-1)}(b) - f^{(2i-1)}(a)] \\ &\quad + \frac{h^{2m+2}}{(2m+2)!} \int_a^b \bar{B}_{2m+2}\left(\frac{x-a}{h}\right) f^{(2m+2)}(x) dx \quad (۹.۴.۵) \end{aligned}$$

توجه: پریم دوگانه روی نماد مجموعیابی بدین معناست که اولین و آخرین جمله قبل از جمع شدن باید نصف شوند.

برهان یک برهان کامل در رلستن^۱ (۱۹۶۵، صص ۱۳۱-۱۳۳) داده شده و یک مطالعه کلتر در (بخش ۲) لاینس^۲ و پیوری^۳ (۱۹۷۳) آمده است. برهان رلستن کوتاه و دقیق است و از ویژگیهای خاص چندجمله‌بیهای برنولی کمال استفاده را می‌نماید. ما برهان ساده‌تر ولی با کلیت کمتری را با نشان دادن اینکه اثبات بر پایه انتگرالگیری جزء به جزء و کمی عملیات جبری استادانه مبتنی است ارائه می‌کنیم.

برای اثبات (۹.۴.۵) در حالت کلی $n > 1$ ، ابتدا اثبات را برای $n = 1$ انجام می‌دهیم. پس

توجه خود را بر رابطه زیر متمرکز می‌نماییم

$$\begin{aligned} E_{\setminus}(f) &= \int_0^h f(x) dx - \frac{h}{2} [f(0) + f(h)] \\ &= \frac{1}{2} \int_0^h f''(x) x(x-h) dx \end{aligned} \quad (10.4.5)$$

فرمول اخیر از (۲۱.۱.۵) به دست می‌آید. چون فرمول مجانبی زیر را می‌دانیم

$$E_{\setminus}(f) \doteq -\frac{h^2}{12} [f'(h) - f'(0)]$$

کوشش می‌کنیم با عملیاتی روی (۱۰.۴.۵) این نتیجه را به دست آوریم. می‌نویسیم

$$E_{\setminus}(f) = \int_0^h f''(x) \left[-\frac{h^2}{12} \right] dx + \int_0^h f''(x) \left[\frac{x(x-h)}{2} + \frac{h^2}{12} \right] dx$$

پس

$$E_{\setminus}(f) = -\frac{h^2}{12} [f'(h) - f'(0)] + \int_0^h f''(x) \left[\frac{x^2}{2} - \frac{xh}{2} + \frac{h^2}{12} \right] dx$$

با استفاده از انتگرالگیری جزء به جزء

$$\begin{aligned} E_{\setminus}(f) &= -\frac{h^2}{12} [f'(h) - f'(0)] + \left[f''(x) \left(\frac{x^3}{6} - \frac{x^2 h}{4} + \frac{h^2 x}{12} \right) \right]_0^h \\ &\quad - \int_0^h f^{(3)}(x) \left[\frac{x^3}{6} - \frac{x^2 h}{4} + \frac{h^2 x}{12} \right] dx \end{aligned}$$

محاسبه کمیت داخل کروشه‌ها، در $x = h$ و $x = 0$ ، صفر به ما می‌دهد. با انتگرالگیری جزء به جزء مجدد، قسمت خارج انتگرال باز هم صفر می‌شود. در نتیجه

$$E_{\setminus}(f) = -\frac{h^2}{12} [f'(h) - f'(0)] + \frac{1}{24} \int_0^h f^{(3)}(x) x^2 (x-h)^2 dx \quad (11.4.5)$$

که همان (۹.۴.۵) با $m = 1$ است. برای به دست آوردن حالت $m = 2$ ، ابتدا توجه می‌کنیم که

$$\frac{1}{24} \int_0^h x^2 (x-h)^2 dx = \frac{h^5}{720}$$

سپس، مانند قبل می‌نویسیم

$$\begin{aligned} \frac{1}{\Gamma(\nu)} \int_0^h f^{(\nu)}(x) x^{\nu-1} (x-h)^{\nu-1} dx &= \int_0^h f^{(\nu)}(x) \left[\frac{h^{\nu}}{\Gamma(\nu)} \right] dx \\ &+ \int_0^h f^{(\nu)}(x) \left[\frac{x^{\nu} (x-h)^{\nu-1}}{\Gamma(\nu)} - \frac{h^{\nu}}{\Gamma(\nu)} \right] dx \\ &= \frac{h^{\nu}}{\Gamma(\nu)} [f^{(\nu)}(h) - f^{(\nu)}(0)] \\ &+ \int_0^h f^{(\nu)}(x) \left[\frac{x^{\nu} - \nu x^{\nu-1} h + x^{\nu} h^{\nu-1}}{\Gamma(\nu)} - \frac{h^{\nu}}{\Gamma(\nu)} \right] dx \end{aligned}$$

اگر دوبار انتگرالگیری جزء به جزء را انجام دهیم، فرمول (۹.۴.۵) را برای حالت $m = 2$ به دست می‌آوریم. این عمل را می‌توان به‌طور نامتناهی ادامه داد. در برهانی که رلستن ارائه داده از انتگرالگیری جزء به جزء و بهره‌گیری از روابط خاص چندجمله‌بیهای برنولی استفاده شده است (مسئله ۲۳ را ببینید).

برای اثبات حالت کلی $m > 1$ ، می‌نویسیم

$$E_n(f) = \sum_{j=1}^n \left\{ \int_{x_{j-1}}^{x_j} f(x) dx - \frac{h}{\nu} [f(x_{j-1}) + f(x_j)] \right\}$$

برای حالت $m = 1$ ، با به‌کار بردن (۱۱.۴.۵)،

$$\begin{aligned} E_n(f) &= \sum_{j=1}^n \left\{ -\frac{h^{\nu}}{\Gamma(\nu)} [f'(x_j) - f'(x_{j-1})] \right\} \\ &+ \sum_{j=1}^n \frac{1}{\Gamma(\nu)} \int_{x_{j-1}}^{x_j} f^{(\nu)}(x) (x-x_{j-1})^{\nu-1} (x-x_j)^{\nu-1} dx \\ &= -\frac{h^{\nu}}{\Gamma(\nu)} [f'(b) - f'(a)] + \frac{h^{\nu}}{\Gamma(\nu)} \int_a^b f^{(\nu)}(x) \overline{B}_{\nu} \left(\frac{x-a}{h} \right) dx \quad (12.4.5) \end{aligned}$$

اثبات برای $m > 1$ ، اساساً مشابه همین است. ■

جمله خطا در (۹.۴.۵) با استفاده از قضیه مقدار میانگین انتگرال ساده می‌شود. می‌توان نشان داد که

$$B_{\nu_j}(x) > 0 \quad 0 < x < 1 \quad (13.4.5)$$

و در نتیجه جمله خطا در رابطه زیر صدق می‌کند

$$\begin{aligned} \int_a^b \overline{B}_{r_{m+r}}\left(\frac{x-a}{h}\right) f^{(r_{m+r})}(x) dx &= f^{(r_{m+r})}(\xi) \int_a^b \overline{B}_{r_{m+r}}\left(\frac{b-a}{h}\right) dx \\ &= n f^{(r_{m+r})}(\xi) \int_a^{a+h} B_{r_{m+r}}\left(\frac{x-a}{n}\right) dx \\ &= n h f^{(r_{m+r})}(\xi) \int_0^1 B_{r_{m+r}}(u) du \\ &= -(b-a) B_{r_{m+r}} f^{(r_{m+r})}(\xi) \end{aligned}$$

برای مقداری چون $a \leq \xi \leq b$

بنابراین (۹.۴.۵) چنین می‌شود

$$\begin{aligned} E_n(f) &= - \sum_{i=1}^m \frac{B_{r_i}}{(r_i)!} h^{r_i} [f^{(r_{i-1})}(b) - f^{(r_{i-1})}(a)] \\ &\quad - \frac{h^{r_{m+2}}(b-a) B_{r_{m+2}}}{(r_{m+2})!} f^{(r_{m+2})}(\xi) \end{aligned} \quad (۱۴.۴.۵)$$

برای مقداری چون $a \leq \xi \leq b$

به‌عنوان یک فرع مهم (۹.۴.۵)، می‌توانیم نشان دهیم که قاعده دوزنقه‌ی وقتی برای توابع دوره‌ی به‌کار برده شود، بسیار خوب عمل می‌کند.

فرع ۱ گیریم $f(x)$ برای $a \leq x \leq b$ به‌طور نامتناهی مشتق‌پذیر باشد و فرض می‌کنیم که تمام مشتقات مرتبه فرد آن دوره‌ی و $a - b$ مضرب صحیحی از این دوره باشد. آنگاه مرتبه همگرایی قاعده دوزنقه‌ی $I_n(f)$ برای

$$I(f) = \int_a^b f(x) dx$$

به‌کار گرفته شده است از هر توانی از h بزرگتر است.

برهان مستقیماً از فرضها درباره $f(x)$

$$f^{(r_{j-1})}(b) = f^{(r_{j-1})}(a) \quad j \geq 1 \quad (۱۵.۴.۵)$$

نتیجتاً برای هر $m \geq 0$ با $h = (b-a)/n$ ، (۱۴.۴.۵) نتیجه می‌دهد که

$$I(f) - I_n(f) = - \frac{h^{r_{m+2}}}{(r_{m+2})!} (b-a) B_{r_{m+2}} f^{(r_{m+2})}(\xi), \quad a \leq \xi \leq b \quad (۱۶.۴.۵)$$

بنابراین وقتی $n \rightarrow \infty$ (و $h \rightarrow 0$) نرخ همگرایی متناسب می‌شود با h^{2m+2} . ولی چون m دلخواه بود، پس نتیجه مطلوب به دست آمده است.

این نتیجه برای $f(x) = \exp(\cos(x))$ در جدول ۷.۵ نشان داده شده است. قاعدهٔ دوزنقه‌یی اغلب بهترین قاعدهٔ انتگرالگیری برای انتگرالدهای هموار دوره‌یی از نوعی است که در فرع قبلی مشخص شدند. برای مقایسهٔ فرمول گاوس-لژاندر و قاعدهٔ دوزنقه‌یی برای خانوادهٔ یک پارامتری از توابع دوره‌یی، با رفتارهای گوناگون، دونالدسن و الیوت (۱۹۷۲، ص ۵۹۲) را ببینید. اینان توصیه می‌کنند که حتی برای انتگرالدهای خیلی قله‌یی، قاعدهٔ دوزنقه‌یی برتر است. در تحلیل قبلی به نظر می‌آمد که انتگرالگیری گاوسی برای انتگرالدهای قله‌یی بهتر است. این نتیجه‌گیری تحلیل قبلی را اصلاح می‌کند.

فرمول مجموعیابی اوپلر-مک لورن یکی از کاربردهای مهم (۹.۴.۵) یا (۱۴.۴.۵)، اگر چه انتگرالگیری عددی را شامل نمی‌شود، در مجموعیابی سری‌هاست.

فرع ۲ (فرمول مجموعیابی اوپلر-مک لورن) فرض کنید $f(x)$ ، در $0 \leq x < \infty$ برای مقداری از $0 \leq m \leq 2m+2$ بار پیوسته مشتقپذیر باشد. در این صورت برای تمام مقادیر $n \geq 1$

$$\begin{aligned} \sum_{j=0}^n f(j) &= \int_0^n f(x) dx + \frac{1}{2} [f(0) + f(n)] \\ &+ \sum_{i=1}^m \frac{B_{2i}}{(2i)!} [f^{(2i-1)}(n) - f^{(2i-1)}(0)] \\ &- \frac{1}{(2m+2)!} \int_0^n \bar{B}_{2m+2}(x) f^{(2m+2)}(x) dx \end{aligned} \quad (17.4.5)$$

برهان فقط $a = 0$ و $b = n$ را در (۹.۴.۵) بگذارید و جملات را به شکل مناسبی مرتب کنید.

مثال در یک مثال بعدی به مجموع

$$S = \sum_{n=1}^{\infty} \frac{1}{n^{3/2}} \quad (18.4.5)$$

نیاز داریم. اگر انتخاب واضح $f(x) = (x+1)^{-3/2}$ را در (۱۷.۴.۵) به کار ببریم، نتایج مایوس

کننده‌اند. ولی اگر $n \rightarrow \infty$ ، به دست می‌آوریم

$$S = \int_0^{\infty} \frac{dx}{(x+1)^{r/2}} + \frac{1}{2} - \sum_1^m \frac{B_{2i}}{(2i)!} f^{(2i-1)}(0) - \frac{1}{(2m+2)!} \int_0^{\infty} \bar{B}_{2m+2}(x) f^{(2m+2)}(x) dx$$

و جمله خطا برای هیچ انتخابی از m ، کوچک نمی‌شود. ولی اگر سری را به دو قسمت تقسیم کنیم، می‌توانیم با آن به گونه‌ای خیلی دقیق برخورد کنیم. گیریم $f(x) = (x+10)^{-3/2}$ ، پس با $m=1$

$$\sum_1^{\infty} \frac{1}{n^{3/2}} = \sum_1^{\infty} f(j) = \int_0^{\infty} \frac{dx}{(x+10)^{3/2}} + \frac{1}{2(10)^{3/2}} - \frac{(1/6) \cdot (-3/2)}{2} \cdot \frac{(10)^{-5/2}}{(10)^{5/2}} + E$$

$$E = -\frac{1}{24} \int_0^{\infty} \bar{B}_2(x) f^{(2)}(x) dx$$

چون $\bar{B}_2(x) \geq 0$ ، $f^{(2)}(x) > 0$ داریم $E < 0$. همچنین

$$0 < -E < \frac{1}{24} \int_0^{\infty} \left(\frac{1}{16}\right) f^{(2)}(x) dx = \frac{35}{(10 \cdot 24)(10)^{1/2}} \doteq 1.08 \times 10^{-6}$$

بنابراین

$$\sum_1^{\infty} \frac{1}{n^{3/2}} = 0.64493406685 + E$$

از مجموعیابی به طور مستقیم $\sum_1^{\infty} (1/n^{3/2}) \doteq 1.963713717$ به دست می‌آوریم

$$\sum_1^{\infty} \frac{1}{n^{3/2}} = 2.6123759 + E \quad 0 < -E < 1.08 \times 10^{-6} \quad (19.4.5)$$

برای اطلاعات بیشتر در مورد تکنیکهای مجموعیابی، به رلستن (۱۹۶۵، صص ۱۳۴-۱۳۸) مراجعه کنید. برای اینکه به اهمیت روش مجموعیابی بالایی ببرید، برای به دست آوردن دقت مشابه برای S لازم است تعداد 3.43×10^{12} جمله به S علاوه کنید.

تعمیم فرمول اولر-مک لورن برای انتگرالهایی که در آنها انتگرالده در بعضی نقاط مشتقپذیر نباشد، باز هم می‌توان یک بسط خطای مجانبی پیدا کرد. برای قاعده دوزنقه‌یی و قاعده‌های دیگر انتگرالگیری عددی که برای انتگرالده‌های با نقاط تکین جبری و [یا] لگاریتمی به کار رفته‌اند، مقاله

لاینس و نینهام^۱ (۱۹۷۶) را ببینید. در پاراگرافهای بعد نتایج آنها را برای انتگرال

$$I = \int_0^1 x^\alpha f(x) dx \quad \alpha > 0 \quad (20.4.5)$$

با $f \in C^{m+1}[0, 1]$ و با استفاده از قاعده انتگرالگیری عددی ذوزنقه‌یی مشخص می‌کنیم. فرض کنید α عدد صحیح نباشد،

$$E_n(f) = \sum_{j=0}^{m-1} \frac{c_j}{n^{\alpha+j+1}} + \sum_{j=1}^{m-1} \frac{d_j}{n^{j+1}} + O\left(\frac{1}{n^{m+1}}\right) \quad (21.4.5)$$

جمله $O(1/n^{m+1})$ کمیته است با مقداری متناسب با $1/(n^{m+1})$ یا شاید کوچکتر از آن. ثابتها در زیر داده شده‌اند،

$$c_j = \frac{2\Gamma(\alpha + j + 1) \sin[(\pi/2)(\alpha + j)] \zeta(\alpha + j + 1) f^{(j)}(0)}{(2\pi)^{\alpha+j+1} j!}$$

$$d_j = 0 \quad j \text{ زوج}$$

$$d_j = (-1)^{(j-1)/2} \frac{2\zeta(j+1)}{(2\pi)^{j+1}} g^{(j)}(1) \quad j \text{ فرد}$$

با $g(x) = x^\alpha f(x)$ ، $\Gamma(x)$ تابع گاما و $\zeta(p)$ تابع زتا،

$$\zeta(p) = \sum_{j=1}^{\infty} \frac{1}{j^p} \quad p > 1 \quad (22.4.5)$$

برای $0 < \alpha < 1$ با $m = 1$ ، برآورد خطای مجانبی را به دست می‌آوریم

$$E_n(f) = \frac{2\Gamma(\alpha + 1) \sin[(\pi/2)\alpha] \zeta(\alpha + 1) f(0)}{(2\pi)^{\alpha+1} n^{\alpha+1}} + O\left(\frac{1}{n^2}\right) \quad f \in C^1[0, 1] \quad (23.4.5)$$

برای مثال، با $I = \int_0^1 \sqrt{x} f(x) dx$ و استفاده از (۱۹.۴.۵) در محاسبه $\zeta(3/2)$ ،

$$E_n(f) = \frac{c}{n\sqrt{n}} f(0) + O\left(\frac{1}{n^2}\right) \quad c = \frac{\zeta(3/2)}{2\pi} \doteq 0.208 \quad (24.4.5)$$

این مطلب به طور عددی در مثال جدول ۶.۵، در بخش ۱.۵، تأیید شده است.

برای تکینگی لگاریتمی در نقاط انتهایی، از کارهای لاینس و نینهام (۱۹۷۶) یک فرمول خطای مجانبی به شکل زیر نتیجه می‌شود

$$E_n(f) \doteq \frac{c \cdot \log n}{n^p} \quad (25.4.5)$$

برای یک مقدار $p > 0$ و یک مقدار ثابت c ، برای منظورهای عددی، این مقدار اصولاً $O(1/n^p)$ است. برای توجیه آن، حد زیر را با استفاده از قاعده هوییتال برای $p > q$ محاسبه می‌کنیم:

$$\lim_{n \rightarrow \infty} \frac{\log(n)/n^p}{1/n^q} = \lim_{n \rightarrow \infty} \frac{\log(n)}{n^{p-q}} = 0$$

این بدان معناست که $\log(n)/n^p$ برای هر $q < p$ سریعتر از $1/n^q$ کاهش می‌یابد. و روشن است که با سرعتی کمی کمتر از $1/n^p$ کاهش می‌یابد.

برای محاسبات عملی، (۲۵.۴.۵) اصولاً $O(1/n^p)$ است. برای مثال، حد خطاهای پی‌درپی زیر را محاسبه می‌کنیم.

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{I - I_n}{I - I_{2n}} &= \lim_{n \rightarrow \infty} \frac{c \cdot \log(n)/n^p}{c \cdot \log(2n)/2^p n^p} = \lim_{n \rightarrow \infty} 2^p \cdot \frac{\log(n)}{\log(2n)} \\ &= \lim_{n \rightarrow \infty} 2^p \cdot \frac{1}{1 + (\log 2 / \log n)} = 2^p \end{aligned}$$

این همان نسبت حدی است که پیدا می‌شود وقتی که خطا درست $O(1/n^p)$ است.

برونیاپی ایتکین^۱ با انگیزه‌ای از مطالب قبل، فرض می‌کنیم که فرمول انتگرالگیری دارای یک فرمول خطای مجانبی

$$I - I_n \doteq \frac{c}{n^p} \quad p > 0 \quad (26.4.5)$$

باشد. این رابطه همیشه برقرار نیست. برای مثال، انتگرالگیری گاوسی معمولاً در (۲۶.۴.۵) صدق نمی‌کند و قاعده دوزنقه‌یی وقتی برای انتگرالده‌های دوره‌یی به‌کار رفته باشد در آن صدق نخواهد کرد. معذالک، بسیاری از قاعده‌های انتگرالگیری عددی، برای تعداد وسیعی از انتگرالدها در آن صدق می‌کنند. با استفاده از این شکل مفروض برای خطا، سعی می‌کنیم که خطا را برآورد نماییم. مشابه این کار، کار برونیاپی ایتکین در بخش ۶.۲ است.

ابتدا p را برآورد می‌نماییم. با استفاده از (۲۶.۴.۵)،

$$\frac{I_{2n} - I_n}{I_{4n} - I_{2n}} = \frac{(I - I_n) - (I - I_{2n})}{(I - I_{2n}) - (I - I_{4n})} = \frac{(c/n^p) - (c/2^p n^p)}{(c/2^p n^p) - (c/4^p n^p)} = 2^p$$

$$R_{2n} \equiv \frac{I_{2n} - I_n}{I_{4n} - I_{2n}} \doteq 2^p \quad (27.4.5)$$

این روش یک راه ساده برای محاسبه p به دست می‌دهد.

مثال استفاده از قاعده سیمپسون را برای $\int_0^1 x\sqrt{x}dx = 0.4$ در نظر می‌گیریم. در جدول ۱۵.۵، ستون R_n باید به $2^{2.5} \doteq 5.66$ میل کند، که یک نتیجه نظری از لاینس و نینهام (۱۹۶۷) برای مرتبه همگرایی است. روشن است که نتایج عددی نظریه را تأیید می‌کند.

برای برآورد انتگرال I با دقت اضافی، I_n ، I_{2n} و I_{4n} محاسبه شده‌اند. با استفاده از (۲۶.۴.۵)

$$\frac{I - I_n}{I - I_{2n}} \doteq 2^p \doteq \frac{I - I_{2n}}{I - I_{4n}}$$

و بنابراین

$$(I - I_n)(I - I_{4n}) \doteq (I - I_{2n})^2$$

جدول ۱۵.۵ خطاهای انتگرالگیری سیمپسون برای $\int_0^1 x\sqrt{x}dx$

n	I_n	$I - I_n$	$I_n - I_{n/2}$	R_n
۲	۰.۴۰۲۳۶۸۹۲۷۰۶۲	-۲.۳۶۹ - ۳		
۴	۰.۴۰۰۴۳۱۹۱۶۰۴۵	-۴.۳۱۹ - ۴	-۱.۹۳۷ - ۳	
۸	۰.۴۰۰۰۷۷۲۴۹۴۴۷	-۷.۷۲۵ - ۵	-۳.۵۴۷ - ۴	۵.۴۶
۱۶	۰.۴۰۰۰۱۳۷۱۳۴۶۹	-۱.۳۷۱ - ۵	-۶.۳۵۴ - ۵	۵.۵۸
۳۲	۰.۴۰۰۰۰۲۴۲۷۸۴۶	-۲.۴۲۸ - ۶	-۱.۱۲۹ - ۵	۵.۶۳
۶۴	۰.۴۰۰۰۰۰۴۲۹۴۱۳	-۴.۲۹۴ - ۷	-۱.۹۹۸ - ۶	۵.۶۵
۱۲۸	۰.۴۰۰۰۰۰۰۷۵۹۲۴	-۷.۵۹۲ - ۸	-۳.۵۳۵ - ۷	۵.۶۵
۲۵۶	۰.۴۰۰۰۰۰۰۰۱۳۴۲۳	-۱.۳۴۲ - ۸	-۶.۲۵۰ - ۸	۵.۶۶
۵۱۲	۰.۴۰۰۰۰۰۰۰۰۲۳۷۳	-۲.۳۷۳ - ۹	-۱.۱۰۵ - ۸	۵.۶۶

جدول ۱۶.۵ جدول تفاضلی برای انتگرالگیری سیمپسون

m	I_m	$\Delta I_m = I_{2m} - I_m$	$\Delta^2 I_m$
۱۶	۰.۴۰۰۰۱۳۷۱۳۴۶۹	-۱,۱۲۸۵۶۲۳E - ۵	
۳۲	۰.۴۰۰۰۰۰۲۴۲۷۸۴۶	-۱,۹۹۸۴۳۳E - ۶	۹,۲۸۷۱۹E - ۶
۶۴	۰.۴۰۰۰۰۰۰۴۲۹۴۱۳		

حال آن را نسبت به I حل می‌کنیم و عملیات جبری انجام می‌دهیم تا شکل مطلوبی به دست آید

$$I \doteq \tilde{I}_{2n} \equiv I_{2n} - \frac{(I_{2n} - I_{2n})^2}{(I_{2n} - I_{2n}) - (I_{2n} - I_n)} \quad (28.4.5)$$

مثال با استفاده از مسأله قبل برای $f(x) = x\sqrt{x}$ و جدول ۱۵.۵، جدول تفاضلی را در جدول ۱۶.۵ به دست می‌آوریم. در این صورت

$$I \doteq \tilde{I}_{64} = 0.3999999999387$$

$$I - \tilde{I}_{64} = 6.13 \times 10^{-10} \quad I - I_{64} = -4.29 \times 10^{-7}$$

بنابراین \tilde{I}_{64} یک بهبود قابل توجه در I_{64} است. ضمناً توجه کنید که $\tilde{I}_{64} - I_{64}$ یک تقریب عالی برای $I - I_{64}$ است.

جمع‌بندی کنیم؛ یک قاعده انتگرالگیری عددی که در (۲۶.۴.۵) صدق می‌کند و سه مقدار I_n ، I_{2n} و I_{4n} داده شده‌اند، \tilde{I}_{4n} ، برونمایی ایتکین (۲۸.۴.۵) را به دست می‌آوریم. این معمولاً یک بهبود کلی در I_{4n} به عنوان تقریب I خواهد بود؛ و بر این اساس

$$I - I_{4n} \doteq \tilde{I}_{4n} - I_{4n} \quad (29.4.5)$$

با قاعده سیمپسون، یا هر فرمول بسته مرکب دیگر نیوتن-کوتس، هزینه محاسبه I_n ، I_{2n} و I_{4n} بیشتر از محاسبه I_{4n} به تنهایی، یعنی $4n + 1$ محاسبه تابع، نخواهد بود. و هرگاه الگوریتم به صورت صحیحی طراحی شده باشد، نیازی به ذخیره موقت برای تعداد زیادی مقدار تابع $f(x_j)$ نیست. به این دلیل هرگز نباید قاعده سیمپسون را فقط برای یک مقدار اندیس n به کار برد. بدون هیچگونه وقت اضافی، و بایک الگوریتم کمی پیچیده‌تر، می‌توان یک برونمایی ایتکین و یک برآورد خطا پیدا کرد.

برونیابی ریچاردسن^۱ اگر فرض کنیم انتگرالده $f(x)$ در انتگرال $I(f)$ به قدر کافی هموار است، آنگاه می‌توانیم عبارت خطای (۹.۴.۵) را به شکل زیر بنویسیم

$$I - I_n = \frac{d_r^{(0)}}{n^2} + \frac{d_r^{(0)}}{n^2} + \dots + \frac{d_{2m}^{(0)}}{n^{2m}} + F_{n,m} \quad (30.4.5)$$

که I_n معرف قاعدهٔ ذوزنقه‌بی است و

$$F_{n,m} = \frac{(b-a)^{2m+2}}{(2m+2)!n^{2m+2}} \int_a^b \bar{B}_{2m+2} \left(\frac{x-a}{h} \right) f^{(2m+2)}(x) dx$$

$$d_{2j}^{(0)} = -\frac{B_{2j}}{(2j)!} (b-a)^{2j} [f^{(2j-1)}(b) - f^{(2j-1)}(a)] \quad (31.4.5)$$

اگر چه سربهایی که با آنها کار می‌کنیم همیشه متناهی‌اند و یک جملهٔ خطا دارند، معمولاً به جملهٔ خطا توجه نمی‌کنیم
وقتی که n زوج است

$$I - I_{n/2} = \frac{4d_r^{(0)}}{n^2} + \frac{16d_r^{(0)}}{n^4} + \frac{64d_r^{(0)}}{n^6} + \dots \quad (32.4.5)$$

(۳۰.۴.۵) را در ۴ ضرب و (۳۲.۴.۵) را از آن کم می‌کنیم:

$$4(I - I_n) - (I - I_{n/2}) = \frac{-12d_r^{(0)}}{n^2} - \frac{60d_r^{(0)}}{n^4} - \dots$$

$$I = \frac{4I_n - I_{n/2}}{3} - \frac{4d_r^{(0)}}{n^2} - \frac{20d_r^{(0)}}{n^4} - \dots$$

تعریف می‌کنیم

$$I_n^{(1)} = \frac{1}{3} [4I_n^{(0)} - I_{n/2}^{(0)}] \quad \text{زوج } n \geq 2 \quad (33.4.5)$$

و $I_m^{(0)} \equiv I_m$. ما $\{I_n^{(1)}\}$ را برونیاب ریچاردسن برای $\{I_n^{(0)}\}$ می‌نامیم.

دنبالهٔ

$$I_2^{(1)}, I_4^{(1)}, I_6^{(1)}, \dots$$

یک قاعدهٔ انتگرالگیری عددی جدید است. برای خطا.

$$I - I_n^{(1)} = \frac{d_r^{(1)}}{n^2} + \frac{d_r^{(1)}}{n^4} + \dots \quad (34.4.5)$$

$$d_r^{(1)} = -4d_r^{(0)}, d_r^{(1)} = -20d_r^{(0)}, \dots \quad (35.4.5)$$

برای آنکه فرمول صریحی برای $I_n^{(1)}$ پیدا کنیم، می‌گیریم $h = (b - a)/n$ و $x_j = a + jh$ و سپس با استفاده از (۳۳.۴.۵) و تعریف قاعدهٔ دوزنقه‌بندی

$$I_n^{(1)} = \frac{4h}{3} \left[\frac{1}{4}f_0 + f_1 + \dots + f_{n-1} + \frac{1}{4}f_n \right] - \frac{2h}{3} \left[\frac{1}{4}f_0 + f_1 + f_2 + \dots + f_{n-2} + \frac{1}{4}f_n \right]$$

$$I_n^{(1)} = \frac{h}{3} [f_0 + 4f_1 + 2f_2 + 4f_3 + \dots + 2f_{n-2} + 4f_{n-1} + f_n] \quad (۳۶.۴.۵)$$

که قاعدهٔ سیمپسون با n تقسیم جزئی است. برای خطا، با استفاده از (۳۵.۴.۵) و (۳۱.۴.۵)،

$$I - I_n^{(1)} = -\frac{h^4}{180} [f^{(4)}(b) - f^{(4)}(a)] + \frac{h^6}{1512} [f^{(6)}(b) - f^{(6)}(a)] + \dots \quad (۳۷.۴.۵)$$

معنی این عبارت این است که کار با فرمول اویلر-مک لورن با عملیات جبری ساده به قاعدهٔ سیمپسون تبدیل می‌شود. ما مثال عددی ذکر نمی‌کنیم، زیرا طبق (۳۶.۴.۵)، $I_n^{(1)}$ درست، همان قاعدهٔ سادهٔ سیمپسون است.

استدلال فوق را که به $I_n^{(1)}$ منتهی شد، می‌توان ادامه داد و فرمولهای جدیدی به دست آورد. مانند گذشته، اگر n مضربی از ۴ باشد، آنگاه

$$I - I_{n/2}^{(1)} = \frac{16d_4^{(1)}}{n^4} + \frac{64d_6^{(1)}}{n^6} + \dots$$

$$16(I - I_n^{(1)}) - (I - I_{n/2}^{(1)}) = \frac{-48d_6^{(1)}}{n^6} + \dots$$

$$I = \frac{16I_n^{(1)} - I_{n/2}^{(1)}}{15} - \frac{48d_6^{(1)}}{15n^6} + \dots \quad (۳۸.۴.۵)$$

در این صورت

$$I - I_n^{(2)} = \frac{d_6^{(2)}}{n^6} + \frac{d_8^{(2)}}{n^8} + \dots \quad (۳۹.۴.۵)$$

با

$$I_n^{(2)} = \frac{16I_n^{(1)} - I_{n/2}^{(1)}}{15} \quad n \geq 4 \quad (۴۰.۴.۵)$$

و n بر ۴ بخشپذیر است. $\{I_n^{(2)}\}$ را برونیاب ریچاردسن $\{I_n^{(1)}\}$ می‌نامیم. اگر وزنه‌های واقعی انتگرالگیری $I_n^{(2)}$ را پیدا کنیم، مشابه با (۳۶.۴.۵)، خواهیم دید که $I_n^{(2)}$ همان قاعدهٔ مرکب بول^۱ است.

با به‌کار بردن فرمولهای قبلی، می‌توانیم، با استفاده از (۳۹.۴.۵)، برآوردهای سودمندی برای خطا به دست آوریم.

$$\begin{aligned} I - I_n^{(1)} &= \frac{16I_n^{(1)} - I_{n/2}^{(1)}}{15} - I_n^{(1)} + \frac{d_\varphi^{(2)}}{n^6} + \dots \\ &= \frac{I_n^{(1)} - I_{n/2}^{(1)}}{15} + \frac{d_\varphi^{(2)}}{n^6} + \dots \end{aligned}$$

با استفاده از $h = (b - a)/n$

$$I - I_n^{(1)} = \frac{1}{15} [I_n^{(1)} - I_{n/2}^{(1)}] + O(h^6) \quad (41.4.5)$$

و بنابراین

$$I - I_n^{(1)} \doteq \frac{1}{15} [I_n^{(1)} - I_{n/2}^{(1)}] \quad (42.4.5)$$

چون هر دو عبارت $O(h^6)$ هستند و جمله باقیمانده $O(h^6)$ است. این را برآورد خطای ریچاردسن برای قاعده سیمپسون می‌نامند.

این روند برونابی را می‌توان به طریق استقرایی ادامه داد. تعریف می‌کنیم

$$I_n^{(k)} = \frac{4^k I_n^{(k-1)} - I_{n/2}^{(k-1)}}{4^k - 1} \quad n \geq 2^k \quad (43.4.5)$$

با n مضربی از 2^k ، $k \geq 1$ می‌توان نشان داد که خطا شکل زیر را دارد

$$\begin{aligned} I - I_n^{(k)} &= \frac{d_{\varphi_{k+2}}^{(k)}}{n^{\varphi_{k+2}}} + \dots \\ &= A_k (b - a) h^{\varphi_{k+2}} f^{(\varphi_{k+2})}(\zeta_n) \quad a < \zeta_n < b \quad (44.4.5) \end{aligned}$$

که A_k ثابتی مستقل از f و h است و

$$d_{\varphi_{k+2}}^{(k)} = A_k (b - a)^{\varphi_{k+2}} [f^{(\varphi_{k+2})}(b) - f^{(\varphi_{k+2})}(a)]$$

بالاخره، می‌توان نشان داد که برای هر $f \in C[a, b]$

$$\lim_{n \rightarrow \infty} I_n^{(k)}(f) = I(f) \quad (45.4.5)$$

قاعده‌های $I_n^{(k)}(f)$ برای $k \geq 2$ ، هیچ رابطه مستقیمی با قاعده‌های مرکب نیوتن-کوتس ندارند. برای جزئیات کامل باؤتر^۱ و دیگران (۱۹۶۳) را ببینید.

$$\begin{array}{cccccc}
 I_1^{(0)} & & & & & \\
 I_2^{(0)} & I_2^{(1)} & & & & \\
 I_4^{(0)} & I_4^{(1)} & I_4^{(2)} & & & \\
 I_8^{(0)} & I_8^{(1)} & I_8^{(2)} & I_8^{(3)} & & \\
 I_{16}^{(0)} & I_{16}^{(1)} & I_{16}^{(2)} & I_{16}^{(3)} & I_{16}^{(4)} & \\
 \vdots & \vdots & \vdots & \vdots & \vdots &
 \end{array}$$

شکل ۴.۵ جدول انتگرالگیری رامبرگ

انتگرالگیری رامبرگ^۱ تعریف می‌کنیم:

$$J_k(f) = I_{\nu_k}^{(k)} \quad k = 0, 1, 2, \dots \quad (۴۶.۴.۵)$$

این قاعده، قاعده انتگرالگیری رامبرگ است. نمودار شکل ۴.۵ برای برونیا‌بهای ریچاردسن در مورد قاعده دوزنقه‌یی را در نظر می‌گیریم، که شماره تقسیمهای جزئی توانی از ۲ است. ستون اول قاعده دوزنقه‌یی را نشان می‌دهد، ستون دوم قاعده سیمپسون را، و غیره. به موجب (۴۵.۴.۵)، هر ستون به $I(f)$ همگراست. قاعده رامبرگ استفاده از قطر است. چون هر ستون از ستون قبلی سریعتر همگرا می‌شود، با فرض اینکه $f(x)$ به‌طور نامتناهی مشتق‌پذیر است، می‌توان انتظار داشت که $I_k(f)$ برای هر k سریعتر از $\{I_n^{(k)}\}$ همگرا شود. معمولاً چنین است و نتیجتاً از اواخر دهه ۱۹۵۰، این روش بسیار عمومیت یافته است. در مقایسه با روشهای انتگرالگیری گاوسی، انتگرالگیری رامبرگ این مزیت را دارد که از طولهایی به فاصله‌های مساوی استفاده می‌کند. برای تحلیل کاملتر انتگرالگیری رامبرگ باؤتر و دیگران (۱۹۶۳) را ببینید.

مثال با استفاده از انتگرالگیری رامبرگ، انتگرال زیر را محاسبه کنید

$$I = \int_0^{\pi} e^x \cos(x) dx = -\frac{1}{4}(e^{\pi} + 1)$$

از این انتگرال به‌عنوان یک مثال قبلاً در جدولهای ۱.۵، ۳.۵ و ۱۱.۵ به ترتیب برای قاعده‌های دوزنقه‌یی، سیمپسون، و گاوس-لژاندر استفاده شد. نتایج رامبرگ در جدول ۱۷.۵ داده شده‌اند. این اعداد نشان می‌دهند که انتگرالگیری رامبرگ نسبت به قاعده سیمپسون برتری دارد، ولی انتگرالگیری گاوسی باز همگرایی سریعتری دارد.

جدول ۱۷.۵ مثال انتگرالگیری رامبرگ

k	گره‌ها	$J_k(f)$	خطا
۰	۲	-۳۴٫۷۷۸۵۱۸۶۶۰۲۶	$۲٫۲۷E + ۱$
۱	۳	-۱۱٫۵۹۲۸۳۹۵۵۳۴۲	$-۴٫۷۸E - ۱$
۲	۵	-۱۲٫۰۱۱۰۸۴۳۱۷۵۴	$-۵٫۹۳E - ۲$
۳	۹	-۱۲٫۰۷۰۴۲۰۴۱۲۸۷	$۷٫۴۱E - ۵$
۴	۱۷	-۱۲٫۰۷۰۳۴۷۲۰۸۷۳	$۸٫۹۲E - ۷$
۵	۳۳	-۱۲٫۰۷۰۳۴۶۳۱۶۳۲	$-۶٫۸۲E - ۱۱$
۶	۶۵	-۱۲٫۰۷۰۳۴۶۳۱۶۳۹	$< ۵٫۰۰E - ۱۲$

برای محاسبه $J_k(f)$ برای یک مقدار خاص k ، وقتی $J_1(f), \dots, J_{k-1}(f)$ قبلاً محاسبه شده باشند، سطر

$$I_n^{(0)}, I_n^{(1)}, \dots, I_n^{(k-1)} \quad n = 2^{k-1} \quad (۴۷.۴.۵)$$

بایستی در یک حافظه موقتی نگهداری شود. سپس $I_{2^{k-1}}^{(0)}$ از $I_{2^{k-1}}^{(0)}$ و مقدار جدید تابع، محاسبه خواهد شد. با استفاده از (۳۳.۴.۵)، (۴۰.۴.۵) و (۴۳.۴.۵) سطر بعدی که در جدول جمله $J_k(f)$ است، محاسبه خواهد شد. $J_k(f)$ با $J_{k-1}(f)$ باید مقایسه شود تا ملاحظه شود آیا دقت کافی وجود دارد تا $J_k(f)$ به عنوان یک تقریب دقیق برای $I(f)$ پذیرفته شود، یا نه.

ما این شیوه را به یک طریق صوری در الگوریتم زیر می آوریم. این الگوریتم به دلایل آموزشی گنجانده شده است و نیابستی به عنوان یک برنامه جدی تلقی شود، مگر آنکه بعضی اصلاحات در آن صورت گیرد. برای مثال، آزمون خطای خیلی مقدماتی و قدیمی است و یک بررسی مطمئنی برای انتگرالهای عددی متناظر با مقادیر کوچک k ، وقتی مقادیر تابع به قدر کافی نمونه‌گیری نشده است، بایستی ضمیمه شود،

الگوریتم Romberg (f,a,b,ε,int)

۱. تبصره: انتگرالگیری رامبرگ را برای محاسبه int به عنوان برابری از انتگرال

$$I = \int_a^b f(x) dx$$

به کار برید. عمل را متوقف کنید وقتی که $|I - \text{int}| \leq \epsilon$.

۲. مقادیر اولیه:

$$k := 0, n := 1,$$

$$T_0 := R_0 := \alpha_0 := (b-a)[f(a) + f(b)]/2$$

۳. حلقه اصلی را شروع کنید: $n := 2n, k := k + 1, h = (b-a)/n$

$$\text{sum} := \sum_{j=1}^{n/2} f(a + (2j-1)h) \quad ۴$$

$$T_k := h \cdot \text{sum} + \frac{1}{2} T_{k-1} \quad ۵$$

$$\beta_j := \alpha_j \quad j = 0, 1, \dots, k-1 \quad ۶$$

$$\alpha_0 := T_k \quad m := 1 \quad ۷$$

۸. تا مرحله ۱۰ برای $j = 1, 2, \dots, k$ عمل را انجام دهید.

$$m := 2m \quad ۹$$

$$\alpha_j := \frac{m \cdot \alpha_{j-1} - \beta_{j-1}}{m-1} \quad ۱۰$$

$$R_k := \alpha_k \quad ۱۱$$

۱۲. اگر $|R_k - R_{k-1}| > \epsilon$ به مرحله ۳ بروید.

۱۳. چون $\epsilon \leq |R_k - R_{k-1}|$ ، $\text{int} = R_{k-1}$ را بپذیرید و الگوریتم را خاتمه دهید.

شکلهای متفاوتی از انتگرالگیری رامبرگ موجود است. برای مثال، راههای دیگر افزایش تعداد گرهها مطالعه شده است. برای یک بررسی خیلی کامل از نوشته‌هایی درباره انتگرالگیری رامبرگ، دیویس^۱ و رابینوویتس^۲ (۱۹۸۴، صص ۴۳۴-۴۴۶) را ببینید. یک برنامه فورترن برای یک انتگرالگیری رامبرگ نیز در آنجا داده شده است.

۵.۵ انتگرالگیری عددی خودکار

در برنامه انتگرالگیری عددی خودکار، مقدار تقریبی انتگرال با دقتی که کاربر مشخص می‌کند محاسبه می‌شود. کاربر نیازی به تعیین روش یا تعداد نقاط گرهی ندارد. برنامه‌های انتگرالگیری خودکار بسیار خوبی وجود دارند که بسیاری افراد از آنها استفاده می‌نمایند. با چنین برنامه‌ای، وقت کاربر برای برنامه‌نویسی صرف نخواهد شد، و نیازی به آشنایی با نظریه انتگرالگیری عددی نخواهد داشت. معذالک، تقریباً همیشه ممکن است برنامه خودکار را بهتر نمود، گرچه این امر به

شناخت کافی از انتگرالگیری عددی که برای مسأله به کار رفته، نیاز دارد. وقتی فقط به دفعات کمی انتگرالگیری عددی انجام می‌شود، انتگرالگیری خودکار اغلب راه خوبی برای صرفه‌جویی در وقت است. ولی برای مسائلی که مستلزم انتگرالگیری‌های زیاد است، شاید بهتر باشد که با صرف وقت روند انتگرالگیری عددی کم‌خرجتری پیدا کنیم.

یک برنامه انتگرالگیری عددی خودکار، مانند یک «جعبه سیاه» عمل می‌کند، بدون اینکه کاربر قادر باشد مراحل میانی محاسبات را ببیند. به همین علت، مهمترین ویژگی لزوم چنین برنامه‌ای این است که قابل اعتماد باشد: مقدار تقریبی انتگرال که توسط برنامه به دست می‌آید باید دقت ادعا شده را داشته باشد. از جنبه نظری، همان‌گونه که در پاراگراف بعد توضیح داده می‌شود، چنین الگوریتمی وجود ندارد. ولی برای آن نوع انتگرالده‌هایی که معمولاً در عمل در نظر گرفته می‌شوند، برنامه‌هایی هستند که تا درجه بالایی قابل اعتمادند. این قابلیت اعتماد می‌تواند بهبود یابد اگر کاربر شرح برنامه را بخواند و به محدودیتها و فرضهای برنامه توجه کند.

برای پی بردن به عدم امکان برنامه انتگرالگیری خودکار کاملاً اعتمادپذیر از لحاظ نظری، باید توجه داشت که این برنامه مقادیر انتگرالده $f(x)$ را در تعداد متناهی نقطه چون، x_n, \dots, x_1 محاسبه می‌کند. در این صورت تعداد نامتناهی تابع پیوسته $\hat{f}(x)$ موجودند که برای آنها

$$\hat{f}(x_i) = f(x_i) \quad i = 1, \dots, n$$

و

$$\int_a^b \hat{f}(x) dx \neq \int_a^b f(x) dx$$

در واقع تعداد نامتناهی از این گونه توابع $\hat{f}(x)$ وجود دارند که به طور نامتناهی مشتق‌پذیرند. در مسائل عملی، نامحتمل می‌نماید که یک برنامه انتگرالگیری عددی خوش‌ساخت، اعتمادپذیر نباشد، ولی این شدنی است. اعتمادپذیری برنامه انتگرالگیری خودکار را می‌توان با سخت‌تر کردن آزمونهای خطا بالا برد، ولی این نیز باعث از دست رفتن کارایی برنامه می‌شود. معمولاً بین اعتمادپذیری و کارایی یک ارتباطی وجود دارد. برای بحث بیشتر درباره مسائل اعتمادپذیری و کارایی برنامه‌های انتگرالگیری خودکار لاینس^۱ و کاکانوف^۲ (۱۹۷۶) را ببینید.

انتگرالگیری تطبیقی برنامه‌های خودکار را می‌توان به دو دسته تقسیم کرد. (۱) آنهایی که در آنها از یک قاعده فراگیر مانند قاعده گاوس یا قاعده دوزنقه‌یی با فاصله‌گذاری مساوی استفاده می‌کنند، و (۲) آنهایی که در آنها یک خط مشی تطبیقی به کار می‌برند، که در آن، قاعده انتگرالگیری جای

نقاط گرهی و حتی تعریف آن را تغییر می‌دهد تا تغییر رفتار موضعی انتگرالده را منعکس نماید. در خط مشی فراگیر همان نوع برآورد خطا که در بخش قبل بحث شد به کار برده می‌شود. اکنون مفهوم و طرز کار خط مشی تطبیقی را مورد بحث قرار می‌دهیم.

همواری و مشتقپذیری بسیاری از انتگرالده‌ها در نقاط مختلف بازه انتگرالگیری $[a, b]$ تغییر می‌کنند. برای مثال، در

$$I = \int_0^1 \sqrt{x} dx$$

انتگرالده در $x = 0$ داری شیب بی‌نهایت است ولی در نقاط نزدیک به ۱ خوش‌رفتار است. در بیشتر روشهای عددی یک شبکه یکنواخت از نقاط گرهی را به کار می‌برند، یعنی، چگالی نقاط گرهی را در سرتاسر بازه انتگرالگیری تقریباً مساوی می‌گیرند. این روش شامل فرمولهای مرکب نیوتن-کوتس، انتگرالگیری گاوسی و انتگرالگیری رامبرگ می‌شود. هرگاه انتگرالده در نقطه‌ای چون α از بازه $[a, b]$ بدرفتار باشد، نقاط گرهی بسیاری باید نزدیک α جا داده شوند تا این رفتار را تعدیل کنند. ولی این امر موجب می‌شود که نقاط گرهی بیشتر از حد نیاز، در سایر قسمتهای $[a, b]$ به کار روند. در انتگرالگیری تطبیقی سعی می‌شود تا نقاط گرهی را مطابق با رفتار انتگرالده جای دهند، به گونه‌ای که چگالی نقاط گرهی در نزدیکی نقاط بدرفتار انتگرالده بیشتر باشد.

اکنون مفهوم اساسی انتگرالگیری تطبیقی را با استفاده از شکل ساده شده قاعده تطبیقی سیمپسون شرح می‌دهیم. برای آنکه ببینیم چرا فاصله‌گذاری متغیر لازم است، قاعده سیمپسون را با چنین فاصله‌گذاری نقاط گرهی در نظر می‌گیریم:

$$I(f) = \sum_{j=1}^{n/2} \int_{x_{2j-2}}^{x_{2j}} f(x) dx \doteq I_n(f)$$

$$= \sum_{j=1}^{n/2} \left(\frac{x_{2j} - x_{2j-2}}{6} \right) (f_{2j-2} + 4f_{2j-1} + f_{2j}) \quad (1.5.5)$$

با $x_{2j-1} = (x_{2j-2} + x_{2j})/2$. با استفاده از (۱.۵.۵)

$$I(f) - I_n(f) = -\frac{1}{2880} \sum_{j=1}^{n/2} (x_{2j} - x_{2j-2})^5 f^{(4)}(\xi_j) \quad (2.5.5)$$

با $x_{2j-2} \leq \xi_j \leq x_{2j}$. واضح است که می‌خواهید $x_{2j} - x_{2j-2}$ را مطابق با اندازه $f^{(4)}(\xi)$ انتخاب کنید که معمولاً نامعلوم است. اگر قدر مطلق $f^{(4)}(x)$ بسیار متغیر باشد، نمی‌خواهید نقاط گرهی متساوی‌فاصله باشند.

نمادگذاری، زیر را وارد می‌کنیم

$$I_{\alpha,\beta} = \int_{\alpha}^{\beta} f(x) dx$$

$$I_{\alpha,\beta}^{(1)} = \frac{h}{3} \left[f(\alpha) + 4f\left(\frac{\alpha+\beta}{2}\right) + f(\beta) \right] \quad h = \frac{\beta-\alpha}{2} \quad (3.5.5)$$

$$I_{\alpha,\beta}^{(2)} = I_{\alpha,\gamma}^{(1)} + I_{\gamma,\beta}^{(1)} \quad \gamma = \frac{\beta+\alpha}{2}$$

برای توضیح محاسبه

$$I = \int_a^b f(x) dx$$

با الگوریتم تطبیقی، از یک تعریف بازگشتی استفاده می‌کنیم. فرض می‌کنیم $\varepsilon > 0$ داده شده است و می‌خواهیم یک انتگرال تقریبی \bar{I} پیدا کنیم که

$$|I - \bar{I}| \leq \varepsilon \quad (4.5.5)$$

با $\alpha = a$ و $\beta = b$ آغاز و $I_{\alpha,\beta}^{(1)}$ و $I_{\alpha,\beta}^{(2)}$ را حساب می‌کنیم. اگر

$$|I_{\alpha,\beta}^{(2)} - I_{\alpha,\beta}^{(1)}| \leq \varepsilon \quad (5.5.5)$$

آنگاه $I_{\alpha,\beta}^{(2)}$ را به عنوان تقریب تطبیقی انتگرال $I_{\alpha,\beta}$ می‌پذیریم. در غیر این صورت می‌نویسیم $\varepsilon := \varepsilon/2$ و انتگرال تطبیقی برای $I_{\alpha,\beta}$ را برابر مجموع انتگرالهای تطبیقی برای $I_{\alpha,\gamma}$ و $I_{\gamma,\beta}$ قرار می‌دهیم که $\gamma = \alpha + \beta/2$ و هر کدام از انتگرالها باید با یک تحمل خطای ε محاسبه شوند.

در یک کاربرد عملی مانند کاربرد در یک برنامه رایانه‌ای، محدودیتهای خیلی بیشتری به عنوان وسایل ایمنی، ضمیمه می‌شوند و برآورد خطا معمولاً بسیار پیچیده‌تر است. و در همه محاسبه‌های تابع، دقت می‌شود که اطمینان حاصل شود که انتگرالده در یک نقطه دوبار محاسبه نشود. این امر نیاز به یک شیوه انباشتن ماهرانه مقادیر $f(x)$ دارد که باید موقتاً نگهداری شوند چون برای محاسبات بعدی مورد احتیاج‌اند. تغییرات بسیار جزئی زیادی می‌توان به‌کار برد تا اجرای برنامه بهبود یابد، ولی معمولاً در مرحله اول مقدار زیادی تجربه و بررسی عملی لازم است. بدین دلیل و دلایل دیگر، توصیه شده است از شیوه‌های تطبیقی آزمایش شده استفاده شود. [مثلاً، دوبر (۱۹۷۱) پیسنز و همکاران (۱۹۸۳) را ببینید]. این موضوع در آخر بخش بیشتر توضیح داده شده است.

جدول ۱۸.۵ مثال تطبیقی (۶.۵.۵) سیمپسون

$[\alpha, \beta]$	$I^{(2)}$	$I - I^{(2)}$	$I - I^{(1)}$	$ I^{(2)} - I^{(1)} $	ε
$[0, 0.0625]$	0.10258	$1.6E-4$	$4.5E-4$	$2.9E-4$	0.0003125
$[0.0625, 0.125]$	0.19046	$1.2E-7$	$1.1E-6$	$1.0E-6$	0.0003125
$[0.125, 0.25]$	0.53871	$4.5E-7$	$3.6E-6$	$4.0E-6$	0.000625
$[0.25, 0.5]$	0.152368	$9.3E-7$	$1.1E-5$	$1.0E-5$	0.00125
$[0.5, 1.0]$	0.430962	$2.4E-6$	$3.0E-5$	$2.8E-5$	0.0025

مثال استفاده از شیوه تطبیقی ساده قبلی سیمپسون را برای محاسبه انتگرال

$$I = \int_0^1 \sqrt{x} dx \quad (6.5.5)$$

با $\varepsilon = 0.05$ بر بازه $[0, 1]$ در نظر می‌گیریم. بازه‌های نهایی $[\alpha, \beta]$ و انتگرالهای $I_{\alpha, \beta}^{(2)}$ در جدول ۱۸.۵ داده شده‌اند. ستونی که با ε مشخص شده، تحمل خطا را که در آزمون (۵.۵.۵) برای برآورد خطا در $I_{\alpha, \beta}^{(1)}$ به‌کار رفته نشان می‌دهد، خطایی که برای $I_{\alpha, \beta}^{(1)}$ برآورده شده بر بازه $[0, 0.0625]$ نادقیق ولی برای سایر زیربازه‌ها دقیق است. مقداری که برای برآورد $I_{\alpha, \beta}$ به‌کار برده شده عملاً برابر $I_{\alpha, \beta}^{(2)}$ و در همه زیربازه‌ها به قدر کافی دقیق است. انتگرال کل که از مجموع تمام $I_{\alpha, \beta}^{(1)}$ ها به‌دست می‌آید، عبارت است از

$$\tilde{I} = 0.666505 \quad I - \tilde{I} = 1.6E - 4$$

و کران محاسبه شده چنین است

$$|I - \tilde{I}| \leq 3.3E - 4$$

که از مجموع مقادیر ستونی که با $|I^{(2)} - I^{(1)}|$ نشان داده شده به‌دست آمده است. توجه دارید که خطا بیشتر در اولین زیرفاصله متمرکز است، همان‌گونه که از رفتار انتگرالده در نزدیکی $x = 0$ می‌توانستیم پیش‌بینی نماییم. برای یک مثال که آزمون (۵.۵.۵) برای آن مناسب نیست، مسئله ۳۲ را ببینید.

چند برنامه انتگرالگیری خودکار یکی از معروفترین برنامه‌های انتگرالگیری خودکار، برنامه تطبیقی CADRE است که در دوبر (۱۹۷۱) آمده و شامل وسیله شناسایی تکنیکهای جبری در نقاط انتهایی بازه انتگرالگیری است. فرمولهای خطای مجانبی لاینس و نیهام (۱۹۶۷) که برای یک حالت خاص در (۲۱.۴.۵) داده شده، به‌کار گرفته شده است تا روش انتگرالگیری همگرای

سرعتی را باز هم بر پایه برونمایی مکرر ریچاردسون، فراهم نماید. به تجربه دریافته‌اند که برنامه CADRE هم کاملاً اعتمادپذیر و هم کارآست.

یک بسته نرم‌افزار که اخیراً ساخته شده QUADPACK است، که بعضی از برنامه‌های آن برای منظورهای کلی و بعضی دیگر برای رده‌های خاصی از انتگرالها تدوین شده است. این بسته نرم‌افزار نتیجه یک همکاری دسته‌جمعی است و شرح کامل آن در پیسنز و همکاران (۱۹۸۳) داده شده است. این بسته به خوبی آزمایش شده و به نظر می‌آید که یک مجموعه عالی از برنامه‌ها باشد. بحث قبلی را با محاسبه تقریبهای عددی انتگرالهای زیر روشن می‌سازیم:

$$I_1 = \int_0^1 \frac{x dx}{1 + 256(x - 0.375)^2} = \frac{1}{4} [\tan^{-1}(10) + \tan^{-1}(6)]$$

$$I_2 = \int_0^1 x^2 \sqrt{x} dx = \frac{2}{7}$$

$$I_3 = \int_0^1 \sqrt{x} dx = \frac{2}{3}$$

$$I_4 = \int_0^1 \log(x) dx = -1$$

$$I_5 = \int_0^1 \log |x - 0.7| dx = 0.3 \log(0.3) + 0.7 \log(0.7) - 1$$

$$I_6 = \int_0^{2\pi} e^{-x} \sin(50x) dx = \frac{50}{2501} (1 - e^{-2\pi})$$

$$I_7 = \int_{-1}^{1} \frac{dx}{\sqrt{|x|}} = 2.06$$

از QUADPACK، برنامه DQAGP را انتخاب می‌کنیم. این برنامه نیز برای تشخیص تکنیکهای جبری در دو سر بازه، و ترمیم حضور آنها، راههایی دارد. برای اجرای بهتر، برنامه به‌کاربر اجازه می‌دهد که نقاط داخل بازه انتگرالگیری را که انتگرالده در آن نقاط تکین است، مشخص نماید.

ما هر دو برنامه CADRE و DQAGP را برای محاسبه انتگرالهای فوق، با تحمل خطای 10^{-2} ، 10^{-5} ، 10^{-8} به‌کار برده‌ایم. نتایج در جداول ۱۹.۵ و ۲۰.۵ داده شده‌اند. برای مقایسه منصفانه‌تر CADRE و DQAGP، برای دو انتگرال در هر یک از انتگرالهای I_7 و I_5 به‌کار برده‌ایم تا تکنیکها در نقاط انتهایی بازه‌ها قرار گیرند. برای مثال، CADRE را برای هر یک از انتگرالهای موجود در

$$I_7 = \int_{-1}^0 \frac{dx}{\sqrt{-x}} + \int_0^{1} \frac{dx}{\sqrt{x}} \quad (7.5.5)$$

جدول ۱۹.۵ مثالهای انتگرالگیری برای CADRE

خطای مطلوب							
		۱۰ ^{-۲}			۱۰ ^{-۵}		
انتگرال	خطا	N	خطا	N	خطا	N	
I_1	$A = ۲,۴۹E - ۶$ $P = ۵,۳۰E - ۴$	۷۶	$A = ۱,۴۰E - ۷$ $P = ۴,۴۵E - ۶$	۷۳	$A = ۴,۶۰E - ۱۱$ $P = ۲,۴۸E - ۹$	۲۲۵	
$I_۲$	$A = ۱,۱۸E - ۵$ $P = ۳,۲۷E - ۳$	۹	$A = ۳,۹۶E - ۷$ $P = ۳,۵۶E - ۶$	۱۷	$P = ۲,۷۳E - ۱۰$ $P = ۲,۸۱E - ۹$	۱۲۹	
$I_۳$	$A = ۱,۰۳E - ۴$ $P = ۲,۹۸E - ۳$	۱۷	$A = ۳,۲۳E - ۸$ $P = ۴,۴۳E - ۸$	۳۳	$A = ۱,۹۸E - ۹$ $P = ۲,۸۶E - ۹$	۶۵	
$I_۴$	$A = ۶,۵۷E - ۵$ $P = ۴,۹۸E - ۳$	۱۰۵	$A = ۶,۴۵E - ۸$ $P = ۹,۲۲E - ۶$	۲۰۹	$A = ۴,۸۰E - ۹$ $P = ۱,۵۵E - ۸$	۲۸۱	
$I_۵$	$A = ۲,۷۷E - ۵$ $P = ۳,۰۲E - ۳$	۲۲۶	$A = ۷,۴۱E - ۸$ $P = ۱,۰۰E - ۵$	۴۱۸	$A = ۵,۸۹E - ۹$ $P = ۱,۱۱E - ۸$	۵۶۲	
$I_۶$	$A = ۸,۴۹E - ۶$ $P = ۸,۴۸E - ۳$	۹۵۵	$A = ۲,۳۷E - ۸$ $P = ۱,۶۷E - ۵$	۱۱۷۱	$A = ۴,۳۰E - ۱۱$ $P = ۲,۰۷E - ۸$	۲۵۷۷	
$I_۷$	$A = ۴,۵۴E - ۴$ $P = ۱,۳۰E - ۳$	۹۸	$A = ۷,۷۲E - ۷$ $P = ۸,۰۲E - ۶$	۴۱۸	$A = **$ $P = **$	۱۵۰۶	

به کار برده ایم. در جدولها، P معرّف کران خطای پیش‌بینی شده در برنامه، و A خطای مطلق واقعی در جواب محاسبه شده است. ستون N تعداد محاسبات انتگرالده را می‌دهد. در تمام نقاطی که انتگرالده نامعین بوده، به طور دلخواه مقدار تابع صفر گرفته شده است. مثالها با درجه دقت مضاعف در یک Prime 850 با واحد گرد کردن $۱۰^{-۱۴} \times ۱,۴ \approx ۲^{-۴۶}$ محاسبه شده‌اند.

در جدول ۱۹.۵، CADRE برای $I_۷$ با تحمل خطای $۱۰^{-۸}$ نتیجه‌ای نداده است، با اینکه از (۷.۵.۵) استفاده شده است. در غیر از این حالت نتیجه کاملاً خوب بوده است. وقتی از تجزیه به دو انتگرال (۷.۵.۵) استفاده نشود و CADRE فقط یک بار برای بازه $[-۹, ۱۰۰۰۰]$ به کار رود، هر سه تحمل خطا مردود است.

در جدول ۲۰.۵ خطای پیش‌بینی شده، در بعضی حالات از خطای واقعی کوچکتر است. به نظر می‌آید این اشکال به علت کار در حد نهایی درجه دقت حساب ماشین باشد، و در تمام حالات، خطای نهایی در حدود خواسته شده خوب بوده است.

در مقایسه، برنامه‌های DQAGP و CADRE هر دو اعتمادپذیر و کارا هستند. همچنین، هر دو برنامه، به طور نسبتاً ضعیفی برای انتگرال خیلی نوسانی $I_۶$ عمل می‌کنند. این امر نشان

جدول ۲۰.۵ مثالهای انتگرالگیری برای DQAGP

خطای مطلوب						
		۱۰-۲	۱۰-۵		۱۰-۸	
انتگرال	خطا	<i>N</i>	خطا	<i>N</i>	خطا	<i>N</i>
I_1	$A = 2,88E - 9$ $P = 2,96E - 3$	۱۰۵	$A = 5,40E - 13$ $P = 5,21E - 10$	۱۴۷	$A = 5,40E - 13$ $P = 5,21E - 10$	۱۴۷
I_2	$A = 1,17E - 11$ $P = 7,46E - 9$	۲۱	$A = 1,17E - 11$ $P = 7,46E - 9$	۲۱	$A = 1,17E - 11$ $P = 7,46E - 9$	۲۱
I_3	$A = 4,79E - 6$ $P = 4,95E - 3$	۲۱	$A = 4,62E - 13$ $P = 4,77E - 14$	۱۸۹	$A = 4,62E - 13$ $P = 4,77E - 14$	۱۸۹
I_4	$A = 5,97E - 13$ $P = 7,15E - 14$	۲۳۱	$A = 5,97E - 13$ $P = 7,15E - 14$	۲۳۱	$A = 5,97E - 13$ $P = 7,15E - 14$	۲۳۱
I_5	$A = 8,67E - 13$ $P = 1,15E - 13$	۴۶۲	$A = 8,67E - 13$ $P = 1,15E - 13$	۴۶۲	$A = 8,67E - 13$ $P = 1,15E - 13$	۴۶۲
I_6	$A = 1,00E - 3$ $P = 4,36E - 3$	۵۲۵	$A = 6,33E - 14$ $P = 8,13E - 6$	۸۶۱	$A = 5,33E - 14$ $P = 7,12E - 9$	۱۲۳۹
I_7	$A = 1,67E - 10$ $P = 1,16E - 10$	۴۶۲	$A = 1,67E - 10$ $P = 1,16E - 10$	۴۶۲	$A = 1,67E - 10$ $P = 1,16E - 10$	۴۶۲

می‌دهد که I_6 بایستی با برنامه‌ای محاسبه شود که برای انتگرالهای نوسانی (مانند DQAWO در QUADPACK برای محاسبه ضرایب فوریه) نوشته شده‌اند. از جدولها چنین برمی‌آید که، DAQGP تا اندازه‌ای در برخورد با انتگرالهای مشکل‌توان‌تر است در حالی که برای سایر انتگرالها، در مقایسه با CADRE، کارایی یکسان دارد. مثالهای مفصل‌تر برای CADRE در رابینسون^۱ (۱۹۷۹) داده شده است.

در محاسبات پرحجم، خطر با استفاده کردن از برنامه‌های انتگرالگیری خود کار زیاد است، و ممکن است موجب نتایج اشتباه و عدم کارایی شود. برای توضیح در استفاده از این نوع برنامه‌ها در محاسبات وسیع و توصیه‌ها برای اینکه چه وقت آنها را انتخاب کنید، لاینس (۱۹۸۳) را ببیند. آنچه در زیر می‌آید از نکاتی است که او نتیجه گرفته است.

قاعده انتگرالگیری خودکار (AQR) یک بخش عملی و مؤثر در نرم‌افزار عددی است. مزیت اصلی آن برای کاربر راحتی آن است. کاربر می‌تواند آن را از قفسه کتابخانه برداشته، به جریان اندازه و به کار آن اطمینان داشته باشد. به دلیل این سهولت، معمولاً یک هزینه جزئی در زمان CPU وجود دارد، این

هزینه اضافی دارای ضربی در حدود ۳ است. روال قاعده محاسبه انتگرالگیری (REQR) [قاعده انتگرالگیری غیرخودکار] این هزینه اضافی را ندارد، ولی کدگذاری و کنترل REQR ممکن است چند ساعت از وقت کاربر را بگیرد. بنابراین اگر زمان پیش‌بینی شده CPU زیاد نباشد، بسیاری از کاربرها با طیب خاطر هزینه اضافی را می‌پردازند تا در وقت صرفه‌جویی کنند و گرفتار در دسر نشوند. مع‌ذالک، بعضی مسائل مشخص - معمولاً با مقیاس بزرگ - وجود دارند که AQR برای آنها طراحی نشده است و استفاده حساب نشده از آن به زمان CPU اضافی می‌انجامد که ضریب آن ۱۰۰ یا بیشتر خواهد بود. این مسائل در شرایطی پیش می‌آیند که در آنها تعداد زیادی انتگرالگیرهای مجزا وجود دارند و نتایج این انتگرالگیرها بعداً در بعضی از فرایندهای عددی دیگر، به عنوان ورودی مورد استفاده قرار می‌گیرند. برای شناختن چنین وضعیتی، لازم است که فرایند عددی بعدی امتحان شود تا ببینند که آیا به تابع ورودی هموار نیاز دارد ... برای بعضی از این مسائل یک REQR کاملاً مناسب است، و حال آنکه یک AQR ممکن است به یک ناکامی عددی بینجامد.

۶.۵ انتگرالهای تکین

در این بخش محاسبه تقریبی انتگرالهایی را مورد بحث قرار می‌دهیم که برای آنها روشهای بحث شده در بخشهای ۱.۵ تا ۴.۵، خوب عمل نمی‌کنند: این روشها شامل قواعد مرکب نیوتن - کوتس (مثلاً، قاعده ذوزنقه‌یی)، تربیع گاوس - لژاندر و انتگرالگیری رامبرگ هستند. انتگرالهایی که در اینجا بحث می‌شوند، اگر با قواعدی که اخیراً نامبرده شد حساب شوند، به دلایل مختلف، به همگرایی عددی ضعیف می‌انجامند. در اینجا در موارد زیر بحث خواهیم کرد انتگرالهایی که (۱) انتگرالده آنها در بازه انتگرالگیری (a, b) دارای تکینی باشند، (۲) انتگرالهایی که بازه انتگرالگیری نامتناهی دارند. روشهای انتگرالگیری تطبیقی را می‌توان برای این انتگرالها به‌کار گرفت، ولی اگر ماهیت رفتار تکین آنها به دقت بررسی و ترمیم شوند معمولاً ممکن است بتوان یک تقریب با همگرایی سریعتر پیدا کرد.

تبدیل متغیر انتگرالگیری اهمیت این روش را با چند مثال نشان خواهیم داد. برای

$$I = \int_a^b \frac{f(x)}{\sqrt{x}} dx \quad (1.6.5)$$

که $f(x)$ تابعی چندین بار پیوسته مشتقپذیر است، می‌گیریم $x = u^2$ و $0 \leq u \leq \sqrt{b}$ در این صورت

$$I = 2 \int_0^{\sqrt{b}} f(u^2) du$$

این انتگرال دارای انتگرالده هموار است و روشهای استاندارد را می‌توان برای آن به‌کار برد.

همچنین با تبدیل متغیر $u = \sqrt{1-x}$

$$\int_0^1 \sin(x) \sqrt{1-x^2} dx = 2 \int_0^1 u^2 \sqrt{2-u^2} \sin(1-u^2) du$$

انتگرالده سمت راست تعداد نامتناهی مشتق پیوسته در $[0, 1]$ دارد، حال آنکه مشتق انتگرالده نخست، در $x = 1$ تکین بود.

برای یک باره انتگرالگیری نامتناهی، روش تبدیل متغیر باز هم مفید خواهد بود. فرض کنید

$$I = \int_1^{\infty} \frac{f(x)}{x^p} dx \quad p > 1 \quad (2.6.5)$$

و $\lim_{x \rightarrow \infty} f(x)$ وجود دارد. همچنین فرض کنید $f(x)$ در بازه $[1, \infty)$ هموار است. با استفاده از تبدیل متغیر

$$x = \frac{1}{u^\alpha} \quad dx = \frac{-\alpha}{u^{1+\alpha}} \quad \alpha > 0 \quad \text{برای}$$

داریم

$$I = \alpha \int_0^1 u^{p\alpha} f\left(\frac{1}{u^\alpha}\right) \frac{du}{u^{1+\alpha}} = \alpha \int_0^1 u^{(p-1)\alpha-1} f\left(\frac{1}{u^\alpha}\right) du \quad (3.6.5)$$

با انتخاب α به گونه‌ای که نمای $(p-1)\alpha-1$ بزرگ شود، همواری انتگرالده جدید را در $u = 0$ ماکسیم می‌نماییم. برای مثال با

$$I = \int_1^{\infty} \frac{f(x)}{x\sqrt{x}} dx$$

تبدیل متغیر $x = 1/u^4$ به انتگرال زیر می‌انجامد

$$I = 4 \int_0^1 u f\left(\frac{1}{u^4}\right) du \quad (4.6.5)$$

اگر رفتار $f(x)$ را در $x = \infty$

$$f(x) = c_0 + \frac{c_1}{x} + \frac{c_2}{x^2} + \dots$$

فرض کنیم، آنگاه

$$u f\left(\frac{1}{u^4}\right) = c_0 u + c_1 u^5 + c_2 u^9 + \dots$$

و (۴.۶.۵) یک انتگرالده هموار در $u = 0$ خواهد داشت.
در برخورد با تکینی‌های انتگرالده در نقاط انتهایی انتگرال

$$I = \int_a^b f(x) dx \quad (5.6.5)$$

آیری^۱ و همکاران (۱۹۷۰) یک فکر جالب دارند: تعریف می‌کنیم

$$\psi(t) = \exp\left(\frac{-c}{1-t^2}\right) \quad (6.6.5)$$

$$\varphi(t) = a + \frac{b-a}{\gamma} \int_{-1}^t \psi(u) du \quad -1 \leq t \leq 1 \quad (7.6.5)$$

که در آن c یک ثابت مثبت است و

$$\gamma = \int_{-1}^1 \psi(u) du$$

وقتی t از -1 تا 1 تغییر می‌کند، $\varphi(t)$ از a تا b تغییر می‌کند. با استفاده از $x = \varphi(t)$ به‌عنوان تبدیل متغیر در (۵.۶.۵)، به‌دست می‌آوریم

$$I = \int_{-1}^1 f(\varphi(t)) \varphi'(t) dt \quad (8.6.5)$$

تابع $\varphi'(t) = ((b-a)/\gamma)\psi(t)$ در $[-1, 1]$ بینهایت بار مشتق‌پذیر است و تابع و تمام مشتقات آن در $t = \pm 1$ صفرند. در (۸.۶.۵) انتگرالده و تمام مشتقات آن در $t = \pm 1$ ، برای تقریباً تمام توابع قابل توجه صفر می‌شوند. با استفاده از فرمول خطای (۹.۴.۵) برای قاعدهٔ ذوزنقه‌یی در $[-1, 1]$ ، می‌توان دید که قاعدهٔ ذوزنقه‌یی وقتی برای (۸.۶.۵) به‌کار برده شود، با سرعت زیاد همگراست. ما این روش را روش IMT می‌نامیم.

این روش در یدونکر و پیسنز^۲ (۱۹۷۶) و در مقایسه‌های کلی رایبسون (۱۹۷۹)، عرضه شده است و به‌عنوان یک راه بسیار اعتمادپذیر و کارا در مواجهه با انتگرالهای (۴.۶.۵) که در نقاط انتهایی تکینی دارند، شناخته شده است. یدونکر و پیسنز (۱۹۷۶) همچنین انتگرالها را روی $[0, \infty)$ نیز بررسی کرده‌اند، بدین طریق که ابتدا تبدیل متغیر $x = (1+u)/(1-u)$ را به‌کار برده‌اند. $-1 \leq u < 1$ و سپس تبدیل متغیر $u = \varphi(t)$ را به‌کار برده‌اند.

مثال روش قبلی (۵.۶.۵) - (۸.۶.۵) را با قاعدهٔ ذوزنقه‌یی برای محاسبهٔ انتگرال زیر به‌کار برید

$$I = \int_0^1 \sqrt{-\ln x} dx = \frac{\sqrt{\pi}}{4} \doteq 0.8862269 \quad (9.6.5)$$

جدول ۲۱.۵ مثال برای روش IMT

گره‌ها	خطا
۲	$-۶.۵۴E-۲$
۴	$۵.۸۲E-۳$
۸	$-۱.۳۰E-۴$
۱۶	$۷.۴۲E-۶$
۳۲	$۱.۱۷E-۸$
۶۴	$۱.۱۸E-۱۲$

توجه داشته باشید که انتگرالده در هر دو نقطه انتهایی دارای رفتار تکین است، اگرچه در این دو حالت رفتار متفاوت است. ثابت c در (۶.۶.۵) برابر ۴ است و محاسبه (۷.۶.۵) از دِدوکر و رابینسون اخذ شده است. نتایج در جدول ۲۱.۵ نشان داده شده‌اند. ستونی که با گره‌ها مشخص شده تعداد نقاط گرهی را در بازه $[۰, ۱]$ مشخص می‌نماید.

انتگرالگیری گاوسی در بخش ۳.۵ یک نظریه کلی فرمولهای انتگرالگیری گاوسی

$$\int_a^b w(x)f(x)dx \doteq \sum_{j=1}^n w_{j,n}f(x_{j,n}) \quad n \geq 1$$

را که درجه دقت $2n - 1$ دارند، توضیح دادیم. پیدا کردن گره‌ها و وزنها و شکل خطا در قضیه ۳.۵ داده شده است. برای کار در این بخش ملاحظه می‌کنیم که: (۱) بازه (a, b) ممکن است نامتناهی باشد، و (۲) $w(x)$ ممکن است در (a, b) تکنیهایی داشته باشد، به شرطی که منفی نباشد و در فرضهای (۸.۳.۴) و (۹.۳.۴) از بخش ۳.۴ صدق کند. برای همگرایی سریع، همچنین پیش‌بینی می‌کنیم که لازم است $f(x)$ تابعی هموار باشد، همان گونه که در انتگرالگیری گاوس-لژاندر در بخش ۳.۵ نشان داده شد.

وزنها و گره‌ها برای تعداد زیادی از توابع وزن $w(x)$ و بازه‌های (a, b) معلوم‌اند. جداول استرود و سکرست (۱۹۶۶) شامل انتگرالهای

$$\int_0^{\infty} x^{\alpha} e^{-x} f(x) dx \quad \int_{-\infty}^{\infty} e^{-x^2} f(x) dx \quad \int_0^1 \ln\left(\frac{1}{x}\right) f(x) dx \quad (۱۰.۶.۵)$$

و بعضی انتگرالهای دیگر است. ثابت α بزرگتر از -1 است. کتابهای دیگری هستند که جداولی برای انتگرالهای غیر از انتگرالهای (۱۰.۶.۵) دارند. همچنین مقاله گلوب^۱ و وِلش^۲ (۱۹۶۹)

شیوه‌ای برای پیدا کردن گره‌ها و وزنهای (۱۰.۶.۵) بر مبنای حلّ مسألهٔ ویژه مقدار ماتریسی، به دست می‌دهد. یک برنامه داده شده است که بیشتر شامل انتگرالهای وزین متداولتری است که در انتگرالگیری گاوسی به کار برده می‌شوند. برای یک بحث اضافی از انتگرالگیری گاوس، با ارجاع به نوشتارها (از جمله جداول و برنامه‌ها)، دیویس و رابینویش (۱۹۸۴)، صص ۹۵-۱۳۲، (۲۲۲-۲۲۹) را ببینید.

مثال استفاده از انتگرالگیری گاوسی را برای محاسبهٔ انتگرال زیر نشان می‌دهیم:

$$I = \int_0^{\infty} g(x) dx$$

انتگرالگیری گاوس-لاگر را که در آن $w(x) = e^{-x}$ اختیار شده به کار می‌بریم. سپس I را به صورت زیر می‌نویسیم

$$I = \int_0^{\infty} e^{-x} [e^x g(x)] dx = \int_0^{\infty} e^{-x} f(x) dx \quad (۱۱.۶.۵)$$

نتایج را برای سه انتگرال زیر می‌دهیم

$$I^{(۱)} = \int_0^{\infty} \frac{x dx}{e^x - 1} = \frac{\pi^2}{6}$$

$$I^{(۲)} = \int_0^{\infty} \frac{x dx}{(1+x^2)^5} = \frac{1}{8}$$

$$I^{(۳)} = \int_0^{\infty} \frac{dx}{1+x^2} = \frac{\pi}{2}$$

انتگرالگیری گاوس-لاگر برای انتگرالدهایی که وقتی $x \rightarrow \infty$ ، به طور نمایی کاهش می‌یابند، از همه بهتر است. برای انتگرالدهایی که وقتی $x \rightarrow \infty$ ، $O(\frac{1}{x^p})$ ، $p > 1$ هستند، نرخ همگرایی وقتی $p \rightarrow 1$ بسیار ضعیف می‌شود. این توضیحات در جدول ۲۲.۵ نمایش داده شده است. برای یک بحث صوری در همگرایی انتگرالگیری گاوس-لاگر، دیویس و رابینویش (۱۹۸۴)، صص (۲۲۲) را ببینید.

یک مورد به ویژه سادهٔ انتگرالگیری گاوسی، انتگرال تکین زیر است

$$I(f) = \int_{-1}^1 \frac{f(x) dx}{\sqrt{1-x^2}} \quad (۱۲.۶.۵)$$

جدول ۲۲.۵ مثالهایی از انتگرالگیری گاوس - لاگر

گره‌ها	$I^{(1)} - I_n^{(1)}$	$I^{(2)} - I_n^{(2)}$	$I^{(3)} - I_n^{(3)}$
۲	$۱,۰۱E - ۴$	$-۸,۰۵E - ۲$	$۷,۷۵E - ۲$
۴	$۱,۲۸E - ۵$	$-۴,۲۰E - ۲$	$۶,۶۹E - ۲$
۸	$-۹,۴۸E - ۸$	$۱,۲۷E - ۲$	$۳,۷۰E - ۲$
۱۶	$۳,۱۶E - ۱۱$	$-۱,۳۹E - ۳$	$۱,۷۱E - ۲$
۳۲	$۷,۱۱E - ۱۴$	$۳,۰۵E - ۵$	$۸,۳۱E - ۳$
۶۴	-	$۱,۰۶E - ۷$	$۴,۰۷E - ۳$

با این تابع وزن، چند جمله‌بندی‌های متعامد چند جمله‌بندی‌های چیشف $\{T_n(x), n \geq 0\}$ هستند. بنابراین گره‌های انتگرالگیری در (۱۰.۶.۵) چنین داده می‌شوند

$$x_{j,n} = \cos\left(\frac{2j-1}{2n}\pi\right) \quad j = 1, \dots, n \quad (۱۳.۶.۵)$$

و با توجه به (۱۱.۳.۵)، وزنها عبارت‌اند از

$$w_{j,n} = \frac{\pi}{n} \quad j = 1, \dots, n$$

با استفاده از فرمول (۱۰.۳.۵) برای خطا، فرمول انتگرالگیری گاوسی برای (۱۲.۶.۵) با رابطه زیر داده می‌شود

$$\int_{-1}^1 \frac{f(x)dx}{\sqrt{1-x^2}} = \frac{\pi}{n} \sum_{j=1}^n f(x_{j,n}) + \frac{2\pi}{2^{2n}(2n)!} f^{(2n)}(\eta) \quad (۱۴.۶.۵)$$

برای $-1 < \eta < 1$

این فرمول بستگی به قاعده میانگاهی مرکب (۱۸.۲.۵) دارد. اگر تبدیل متغیر $x = \cos \theta$ در (۱۴.۶.۵) داده شود، چنین به دست خواهد آمد

$$\int_0^\pi f(\cos \theta) d\theta = \frac{\pi}{n} \sum_{j=1}^n f(\cos \theta_{j,n}) + E \quad (۱۵.۶.۵)$$

که $\theta_{j,n} = (2j-1)\pi/2n$. بنابراین انتگرالگیری گاوسی برای (۱۲.۶.۵) هم‌ارز است با قاعده میانگاهی مرکب که برای انتگرال سمت چپ (۱۵.۶.۵) به‌کار برده شده است. مانند قاعده ذوزنقه‌یی، قاعده میانگاهی دارای بسط خطایی است بسیار شبیه به بسط خطایی که در (۹.۴.۵)، با استفاده از فرمول اوپلر-مکلورن داده شده است. فرع ۱ مربوط به قضیه ۵.۵ نیز معتبر است، که نشان

می دهد قاعده میانگاهی مرکب برای توابع دوره‌یی بسیار دقیق است. این دقت زیاد در (۱۴.۶.۵) منعکس شده است. بنابراین انتگرالگیری گاوسی برای (۱۲.۶.۵) به فرمولی می انجامد که مناسب بودن آن از بسط خطای مجانبی برای قاعده میانگاهی مرکب وقتی که برای انتگرال سمت چپ (۱۵.۶.۵) به کار برده شود، به دست می آید.

برخورد تحلیلی با تکینی بازه انتگرالگیری را به دو بخش تقسیم می کنیم، یک بخش شامل نقطه تکینی است که بایستی به گونه ای تحلیلی به آن پرداخته شود. برای مثال، انتگرال زیر را در نظر می گیریم

$$I = \int_0^b f(x) \log(x) dx = \left[\int_0^\varepsilon + \int_\varepsilon^b \right] f(x) \log(x) dx \equiv I_1 + I_2 \quad (۱۶.۶.۵)$$

با فرض اینکه $f(x)$ در $[\varepsilon, b]$ هموار است، یک روش استاندارد برای محاسبه I_2 به کار می بریم. برای I_1 فرض می کنیم $f(x)$ حول صفر دارای یک سری همگرایی تیلر در $[\varepsilon, 0]$ باشد. در این صورت

$$\begin{aligned} I_1 &= \int_0^\varepsilon f(x) \log(x) dx = \int_0^\varepsilon \left(\sum_{j=0}^{\infty} a_j x^j \right) \log(x) dx \\ &= \sum_{j=0}^{\infty} a_j \frac{\varepsilon^{j+1}}{j+1} \left[\log(\varepsilon) - \frac{1}{j+1} \right] \end{aligned} \quad (۱۷.۶.۵)$$

برای مثال، در مورد

$$I = \int_0^{2\pi} \cos(x) \log(x) dx$$

تعریف می کنیم

$$\begin{aligned} I_1 &= \int_0^{0.1} \cos(x) \log(x) dx \quad \varepsilon = 0.1 \\ &= \varepsilon [\log(\varepsilon) - 1] - \frac{\varepsilon^3}{6} \left[\log(\varepsilon) - \frac{1}{3} \right] + \frac{\varepsilon^5}{600} \left[\log(\varepsilon) - \frac{1}{5} \right] - \dots \end{aligned}$$

که یک سری متناوب است و بنابراین روشن است که با سه جمله اول یک تقریب دقیق برای I_1 به دست خواهد آمد. یک روش استاندارد می توان برای I_2 در $[0.1, 2\pi]$ به کار برد. تکنیکهای مشابهی می توان برای بازه های نامتناهی انتگرالگیری $[a, \infty)$ به کار برد که انتگرال را در بازه $[b, \infty)$ برای مقدار بزرگی از b حذف می کند. این روش در اینجا بحث نشده است.

انتگرالگیری حاصلضرب بگیریم $I(f) = \int_a^b w(x) f(x) dx$ که در آن تابع وزن $w(x)$ نزدیک به تکین یا تکین، و $f(x)$ یک تابع هموار است. هدف اصلی این است که دنباله ای از توابع $f_n(x)$ تولید کنیم به طوری که

$$\|f - f_n\|_\infty = \max_{a \leq x \leq b} |f(x) - f_n(x)| \rightarrow 0 \quad \text{وقتی } n \rightarrow \infty \quad ۱$$

۲. انتگرالهای

$$I_n(f) \equiv \int_a^b w(x) f_n(x) dx \quad (۱۸.۶.۵)$$

به سادگی قابل محاسبه باشند.

این هدف تعمیم طرح (۲۰.۵) است که در مقدمه آمده است. برای خطا،

$$\begin{aligned} |I(f) - I_n(f)| &\leq \int_a^b |w(x)| |f(x) - f_n(x)| dx \\ &\leq \|f - f_n\|_\infty \int_a^b |w(x)| dx \end{aligned} \quad (۱۹.۶.۵)$$

بنابراین وقتی $I_n(f) \rightarrow I(f)$, $n \rightarrow \infty$ ، و نرخ همگرایی حداقل به همان نرخ همگرایی $f_n(x)$ به $f(x)$ در $[a, b]$ است.

در چارچوب فوق، برای تعریف روشهای انتگرالگیری حاصلضرب معمولاً $f_n(x)$ را از $f(x)$ با استفاده از درونیابی چندجمله‌یی تکه‌یی تعریف می‌کنند. برای نشان دادن طرح اصلی، با جبری ساده روش ذوزنقه‌یی حاصلضرب را برای محاسبه انتگرال زیر تعریف می‌کنیم

$$I(f) = \int_a^b f(x) \log(x) dx \quad (۲۰.۶.۵)$$

گیریم $n \geq 1$ ، $h = b/n$ و $x_j = jh$ برای $j = 0, 1, \dots, n$. $f_n(x)$ را تابع خطی تکه‌یی که $f(x)$ را در گره‌های x_0, x_1, \dots, x_n درونیابی می‌نماید چنین تعریف می‌کنیم.

برای $j = 1, 2, \dots, n$ ، $x_{j-1} \leq x \leq x_j$ تعریف می‌کنیم

$$f_n(x) = \frac{1}{h} [(x_j - x)f(x_{j-1}) + (x - x_{j-1})f(x_j)] \quad (۲۱.۶.۵)$$

با توجه به (۱۰.۱.۳) می‌توان به سادگی نشان داد

$$\|f - f_n\|_\infty \leq \frac{h^2}{8} \|f''\|_\infty \quad (۲۲.۶.۵)$$

مشروط بر آنکه $f(x)$ دوبار پیوسته مشتق‌پذیر در $a \leq x \leq b$ باشد. از (۱۹.۶.۵) کران خطا را به دست می‌آوریم

$$|I(f) - I_n(f)| \leq \frac{h^2}{8} \|f''\|_\infty \int_a^b |\log(x)| dx \quad (۲۳.۶.۵)$$

این روش تعریف $I_n(f)$ مشابه قاعدهٔ دوزنقه‌ی (۵.۱.۵) است. قاعدهٔ (۵.۱.۵) را می‌توانستیم از انتگرال‌گیری تابع $f_n(x)$ قبلی نیز به دست آوریم. در این صورت تابع وزن $w(x) \equiv 1$ می‌شد. به سادگی می‌توانیم مطالب فوق را با به کار بردن درونیایی چندجمله‌ی تکه‌ی از درجهٔ بالاتر تعمیم دهیم. استفاده از درونیایی مرتبهٔ دوم تکه‌ی برای تعریف $f_n(x)$ ، به یک فرمول $I_n(f)$ می‌انجامد که قاعدهٔ حاصلضرب سیمپسون خوانده می‌شود. و با استفاده از استدلالی که به (۲۳.۶.۵) رسیدیم، می‌توان نشان داد که

$$|I(f) - I_n(f)| \leq \frac{\sqrt{3}}{2\sqrt{7}} h^2 \|f^{(2)}\|_\infty \int_a^b |\log(x)| dx \quad (24.6.5)$$

فرمولهای مراتب بالاتر را می‌توان با استفاده از درونیایها درجهٔ بالاتر به دست آورد. برای محاسبهٔ $I_n(f)$ با استفاده از (۲۱.۶.۵) داریم

$$I_n(f) = \sum_{j=1}^n \int_{x_{j-1}}^{x_j} \log(x) \left[\frac{(x_j - x)f(x_{j-1}) + (x - x_{j-1})f(x_j)}{h} \right] dx$$

$$= \sum_{k=0}^n w_k f(x_k) \quad (25.6.5)$$

$$w_0 = \frac{1}{h} \int_{x_0}^{x_1} (x_1 - x) \log(x) dx \quad w_n = \frac{1}{h} \int_{x_{n-1}}^{x_n} (x - x_{n-1}) \log(x) dx,$$

$$w_j = \frac{1}{h} \int_{x_{j-1}}^{x_j} (x - x_{j-1}) \log(x) dx$$

$$+ \frac{1}{h} \int_{x_j}^{x_{j+1}} (x_{j+1} - x) \log(x) dx \quad j = 1, \dots, n-1 \quad (26.6.5)$$

محاسبهٔ این وزنها را می‌توان خیلی ساده کرد. با تبدیل متغیر $x - x_{j-1} = uh$ ، $0 \leq u \leq 1$ داریم

$$\frac{1}{h} \int_{x_{j-1}}^{x_j} (x - x_{j-1}) \log(x) dx = h \int_0^1 u \log[(j-1+u)h] du$$

$$= \frac{h}{\sqrt{e}} \log(h) + h \int_0^1 u \log(j-1+u) du$$

و

$$\frac{1}{h} \int_{x_{j-1}}^{x_j} (x_j - x) \log(x) dx = h \int_0^1 (1-u) \log[(j-1+u)h] du$$

$$= \frac{h}{\sqrt{e}} \log(h) + h \int_0^1 (1-u) \log(j-1+u) du$$

تعریف می‌کنیم

جدول ۲۳.۵ وزنها برای قاعده دوزنقه‌ی حاصلضرب

k	$\psi_1(k)$	$\psi_2(k)$
۰	-۰٫۲۵۰	-۰٫۷۵۰
۱	۰٫۲۵۰	۰٫۱۳۶۲۹۴۳۶۱۱
۲	۰٫۴۸۸۳۷۵۹۲۸۱	۰٫۴۲۱۱۶۶۵۷۶۸
۳	۰٫۶۴۸۵۷۷۸۵۴۵	۰٫۶۰۰۷۶۲۷۲۳۹
۴	۰٫۷۶۹۵۷۰۵۴۵۷	۰٫۷۳۲۴۴۱۵۷۲۰
۵	۰٫۸۶۶۸۶۰۲۷۴۷	۰٫۸۳۶۵۰۶۹۷۸۵
۶	۰٫۹۴۸۲۴۲۸۳۷۶	۰٫۹۲۲۵۷۱۳۹۰۴
۷	۱٫۰۱۸۲۰۱۶۵۲	۰٫۹۹۵۹۵۹۶۳۸۵

$$\psi_1(k) = \int_0^1 u \log(u+k) du$$

$$\psi_2(k) = \int_0^1 (1-u) \log(u+k) du \quad (27.6.5)$$

برای $k = 0, 1, 2, \dots$ در این صورت

$$w_0 = \frac{h}{\gamma} \log(h) + h\psi_2(0) \quad w_n = \frac{h}{\gamma} \log(h) + h\psi_1(n-1)$$

$$w_j = h \log(h) + h[\psi_1(j-1) + \psi_2(j)] \quad j = 1, 2, \dots, n-1 \quad (28.6.5)$$

توابع $\psi_1(k)$ و $\psi_2(k)$ به h, b یا n بستگی ندارند. آنها را می‌توان محاسبه و در یک جدول نگهداری کرد و برای مقادیر متفاوت b و n به‌کار برد. برای مثال، یک جدول $\psi_1(k)$ و $\psi_2(k)$ به‌ازای $k = 0, 1, \dots, 99$ را می‌توان برای هر مقدار $b > 0$ و هر $n \leq 100$ به‌کار برد. وقتی که جدول مقادیر $\psi_1(k)$ و $\psi_2(k)$ محاسبه شده باشد، هزینه استفاده از قاعده دوزنقه‌ی حاصلضرب بیش از هزینه هیچ قاعده انتگرالگیری دیگری نخواهد بود.

انتگرالهای $\psi_1(k)$ و $\psi_2(k)$ در (۲۷.۶.۵) را می‌توان به‌طور صریح محاسبه کرد؛ بعضی از مقادیر آنها در جدول ۲۳.۵ داده شده است.

مثال ۴۴۸۴۱۳۷ $I = \int_0^1 (1/(x+2)) \log(x) dx = -0.4484137$ را محاسبه کنید. مقادیر محاسبه شده در جدول ۲۴.۵ داده شده‌اند. آهنگ همگرایی محاسبه شده با مرتبه همگرایی (۲۳.۶.۵) مطابقت دارد. انواع بسیاری درونیابی را می‌توان برای تعریف $f_n(x)$ به‌کار برد، ولی تاکنون در بیشتر کاربردها، چندجمله‌بیهای درونیابی تکه‌یی با نقاط گرهی متساوی‌الفاصله به‌کار رفته‌اند. تابعهای وزن دیگری

جدول ۲۴.۵ مثال برای قاعدهٔ دوزنقیبی حاصلضرب

n	I_n	$I - I_n$	نسبت
۱	-۰٫۴۵۸۳۳۳۳	۰٫۰۰۰۹۹۲	
۲	-۰٫۴۵۱۶۰۹۶	۰٫۰۰۰۳۲۰	۳٫۱۰
۴	-۰٫۴۴۹۳۰۱۱	۰٫۰۰۰۰۸۸۷	۲٫۶۱
۸	-۰٫۴۴۸۶۴۶۰	۰٫۰۰۰۰۲۳۲	۳٫۸۲

را نیز می‌توان به‌کار برد. برای مثال

$$w(x) = x^\alpha \quad \alpha > -1 \quad x \geq 0 \quad (29.6.5)$$

و باز هم تابعهای وزن را می‌توان به فرمولهای نسبتاً ساده‌تری شبیه (۲۸.۶.۵) بدل کرد. برای یک مقدار گنگ α ، چون

$$w(x) = \frac{1}{x\sqrt{x-1}}$$

برای از بین بردن تکینی در انتگرال، دیگر تبدیل متغیر را نمی‌توان به‌کار برد. همچنین یکی دیگر از کاربردهای اصلی انتگرالگیری حاصلضرب در معادلات انتگرالی است که در آنها تابع هسته تکینی جبری و/یا لگاریتمی دارد. برای چنین معادلاتی حتی با تکینی ریشهٔ دوم دیگر تبدیل متغیر ممکن نیست. برای مثال، معادلهٔ

$$\lambda \varphi(x) - \int_a^b \frac{\varphi(y)dy}{|x-y|^{1/2}} = f(x) \quad a \leq x \leq b$$

با مقادیر داده شدهٔ λ ، a ، b ، f و تابع مجهول خواسته شدهٔ φ را در نظر می‌گیریم. انتگرالگیری حاصلضرب برای چنین معادلاتی به روندهای کارا می‌انجامد، به شرطی که $\varphi(y)$ یک تابع هموار باشد [ایتکین^۱ ۱۹۷۶، ص ۱۰۶ را ببینید].

برای توابع وزن پیچیده که در آنها وزنه‌های w_j را دیگر نمی‌توان محاسبه نمود، اغلب امکان دارد مسأله را به مسألهٔ دیگری برگرداند که انتگرالگیری عددی برای آن باز به‌سادگی قابل اجرا باشد. این مطلب را با یک مثال بررسی خواهیم کرد.

مثال انتگرال $I = \int_0^\pi f(x) \log(\sin x) dx$ را در نظر می‌گیریم. انتگرالده در هر دو نقطهٔ $x = 0$

و $x = \pi$ تکینی دارد. با استفاده از

$$\log(\sin x) = \log \left[\frac{\sin x}{x(\pi - x)} \right] + \log(x) + \log(\pi - x)$$

داریم

$$\begin{aligned} I &= \int_0^\pi f(x) \log \left[\frac{\sin x}{x(\pi - x)} \right] dx + \int_0^\pi f(x) \log(x) dx \\ &\quad + \int_0^\pi f(x) \log(\pi - x) dx \\ &\equiv I_1 + I_2 + I_3 \end{aligned} \quad (30.6.5)$$

انتگرال I_1 ، انتگرالده بینهایت مشتقپذیر دارد و هر روش عددی استاندارد برای آن به خوبی اجرا می‌شود. انتگرال I_2 قبلاً مورد بحث قرار گرفته است، با $w(x) = \log(x)$. برای I_3 ، یک تبدیل متغیر به‌کار می‌بریم و می‌نویسیم

$$I_3 = \int_0^\pi f(x) \log(\pi - x) dx = \int_0^\pi f(\pi - z) \log(z) dz$$

از ترکیب با I_2 داریم

$$I_2 + I_3 = \int_0^\pi \log(x) [f(x) + f(\pi - x)] dx$$

که برای آن همان عمل قبلی به‌کار می‌رود. با این دستکاریها حالت‌های $w(x) = \log(x)$ و $w(x) = x^\alpha$ بسیار بیشتر از آن که در ابتدا تصور می‌شد، قابل استفاده است.

برای تحلیل خطای مجانبی انتگرالگیری حاصلضرب، اثر دهوگ^۱ و وایس^۲ (۱۹۷۳) را ببینید، که در آن بعضی از تعمیمهای بسط اوپلر-مکلورن داده شده است. با استفاده از نتایجی که به‌دست آورده‌اند، می‌توان نشان داد که خطا در قاعده حاصلضرب سیمپسون برابر $O(h^r \log(h))$ است. بنابراین کران (۲۴.۶.۵) بر پایه خطای درونیابی $f(x) - f_n(x)$ ، نرخ همگرایی صحیحی را پیش‌بینی نمی‌کند. این نتیجه مشابه با نتیجه (۱۷.۱.۵) برای خطای قاعده سیمپسون است که در آنجا خطا کوچکتر از آن بود که با استفاده از درونیابی درجه دوم می‌توانستیم نتیجه بگیریم.

۷.۵ مشتقگیری عددی

تقریبهای عددی برای مشتقات عمدتاً به دو منظور مورد استفاده واقع می‌شوند. نخست، علاقه‌مندیم مشتقات داده‌هایی را حساب کنیم که اغلب به‌طور تجربی به‌دست آمده‌اند. دوم، فرمولهای مشتقگیری

عددی را برای پیدا کردن روشهای عددی در حل معادلات دیفرانسیل معمولی و معادلات دیفرانسیل جزئی به کار می‌بریم. ما این بخش را با به دست آوردن فرمولهایی که بیشترین کاربرد را برای مشتقگیری عددی دارند، آغاز می‌کنیم.

مسئله مشتقگیری عددی از بعضی جهات مشکلتر از انتگرالگیری عددی است. وقتی از مقادیر تابع که به صورت تجربی معین شده‌اند، استفاده می‌کنیم، خطای این مقادیر، معمولاً به ناپایداری در مشتقگیری عددی تابع می‌انجامد. برعکس، انتگرالگیری عددی وقتی با چنین خطاهایی مواجه می‌شود پایدار است (مسئله ۱۳ را ببینید).

فرمولهای کلاسیک یکی از راههای اساسی پیدا کردن تقریب عددی برای $f'(x)$ ، استفاده از مشتق چندجمله‌یی $p_n(x)$ است که $f(x)$ را در یک مجموعه نقاط گزینی معین درونیابی می‌کند. گیریم x_0, x_1, \dots, x_n داده شده باشد و $p_n(x)$ تابع $f(x)$ را در این گره‌ها درونیابی کند. معمولاً $\{x_i\}$ ها متساوی‌الفاصله‌اند. در این صورت از رابطه زیر استفاده می‌کنیم

$$f'(x) \doteq p'_n(x) \quad (۱.۷.۵)$$

از روابط (۶.۱.۳)، (۴.۲.۳) و (۱۱.۲.۳) داریم:

$$p_n(x) = \sum_{j=0}^n f(x_j) l_j(x)$$

$$l_j(x) = \frac{\Psi_n(x)}{(x - x_j) \Psi'_n(x_j)}$$

$$= \frac{(x - x_0) \dots (x - x_{j-1})(x - x_{j+1}) \dots (x - x_n)}{(x_j - x_0) \dots (x_j - x_{j-1})(x_j - x_{j+1}) \dots (x_j - x_n)}$$

$$\Psi_n(x) = (x - x_0) \dots (x - x_n)$$

$$f(x) - p_n(x) = \Psi_n(x) f[x_0, \dots, x_n, x] \quad (۲.۷.۵)$$

بنابراین

$$f'(x) \doteq p'_n(x) = \sum_{j=0}^n f(x_j) l'_j(x) \equiv D_n f(x) \quad (۳.۷.۵)$$

$$f'(x) - D_h f(x) = \Psi'_n(x) f[x_0, \dots, x_n, x] + \Psi_n(x) f[x_0, \dots, x_n, x, x] \quad (۴.۷.۵)$$

که در مرحله آخر از رابطه (۱۷.۲.۳) استفاده شده است. با استفاده از (۱۲.۲.۳)،

$$f'(x) - D_h f(x) = \Psi'_n(x) \frac{f^{(n+1)}(\xi_1)}{(n+1)!} + \Psi_n(x) \frac{f^{(n+2)}(\xi_2)}{(n+2)!} \quad (5.7.5)$$

که $\xi_1, \xi_2 \in \mathcal{H}\{x_0, \dots, x_n, x\}$. فرمولهای مشتقگیری مراتب بالاتر و خطای آنها را می‌توان از مشتقگیری بیشتر از (۳.۷.۵) و (۴.۷.۵) به دست آورد.

عادیترین کاربردهای مطالب فوق در نقاط گرهی متساوی الفاصله $\{x_i\}$ است. در این حالت گیریم $x_1 = x_0 + ih$ و $i \geq 0$ و $h > 0$. به سادگی می‌توان نشان داد که

$$\Psi_n(x) = O(h^{n+1}) \quad \Psi'_n(x) = O(h^n) \quad (6.7.5)$$

بنابراین

$$f'(x) - p'_n(x) = \begin{cases} O(h^n) & \Psi'_n(x) \neq 0 \\ O(h^{n+1}) & \Psi'_n(x) = 0 \end{cases} \quad (7.7.5)$$

اکنون مثالی برای هر حالت به دست می‌آوریم.

گیریم $n = 1$ ، پس $p_n(x)$ درست همان درونیایی خطی $(x_0, f(x_0))$ و $(x_1, f(x_1))$ خواهد بود. در این صورت رابطه (۳.۷.۵) چنین به دست می‌دهد

$$f'(x_0) \doteq D_h f(x_0) \equiv \frac{1}{h} [f(x_0 + h) - f(x_0)] \quad (8.7.5)$$

از (۵.۷.۵)،

$$f'(x_0) - D_h f(x_0) = \frac{h}{2} f''(\xi_1) \quad x_0 \leq \xi_1 \leq x_1 \quad (9.7.5)$$

زیرا $\Psi(x_0) = 0$.

برای بهبود این وضعیت با درونیایی خطی، انتخاب می‌کنیم $x = m \equiv (x_0 + x_1)/2$. در این صورت

$$f'(m) \doteq \frac{1}{h} [f(x_1) - f(x_0)]$$

این فرمول را با قراردادن $\delta = h/2$ بازنویسی می‌کنیم، تا به دست آوریم

$$f'(m) \doteq D_\delta f(m) = \frac{1}{2\delta} [f(m + \delta) - f(m - \delta)] \quad (10.7.5)$$

برای خطا، با استفاده از (۵.۷.۵) و $\Psi_{\nu}(m) = 0$

$$f'(m) - D_{\delta}f(m) = \frac{-\delta^2}{6} f^{(3)}(\xi_2) \quad m - \delta \leq \xi_2 \leq m + \delta \quad (۱۱.۷.۵)$$

به طور کلی، برای به دست آوردن مرتبه بالاتر حالت (۷.۷.۵)، می‌خواهیم نقاط گرهی $\{x_i\}$ را طوری انتخاب کنیم که $\Psi'_n(x) = 0$. این درست خواهد بود اگر n فرد باشد و گره‌ها مانند (۱۰.۷.۵)، نسبت به x قرینه باشند.

برای پیدا کردن فرمولهایی از مراتب بالاتر وقتی که نقاط گرهی تماماً در یک طرف x باشند مقادیر بزرگتری از n را در (۳.۷.۵) به کار می‌بریم. برای مثال با $x = x_0$ و $n = 2$ داریم،

$$f'(x_0) \doteq D_h f(x_0) = \frac{1}{\sqrt{h}} [-3f(x_0) + 4f(x_1) - f(x_2)] \quad (۱۲.۷.۵)$$

$$f'(x_0) - D_h f(x_0) = \frac{h^2}{3} f^{(3)}(\xi_1) \quad x_0 \leq \xi_1 \leq x_2 \quad (۱۳.۷.۵)$$

روش ضرایب نامعین یک روش دیگر برای پیدا کردن فرمولهایی برای انتگرالگیری عددی، مشتقگیری عددی و درونیابی، روش موسوم به ضرایب نامعین است. این روش اغلب هم‌ارز است با فرمولهای حاصل از فرمول درونیابی یک چندجمله‌یی، ولی گاهی به دست آوردن آن ساده‌تر است. ما این روش را با پیدا کردن فرمولی برای $f''(x)$ نشان خواهیم داد.

فرض می‌کنیم

$$f''(x) \doteq D_h^{(2)} f(x) = Af(x+h) + Bf(x) + Cf(x-h) \quad (۱۴.۷.۵)$$

که در آن A و B و C نامشخص‌اند. به جای $f(x+h)$ و $f(x-h)$ بسطهای تیلر آنها را می‌گذاریم

$$f(x \pm h) = f(x) \pm hf'(x) + \frac{h^2}{2} f''(x) \pm \frac{h^3}{6} f^{(3)}(x) + \frac{h^4}{24} f^{(4)}(\xi_{\pm})$$

که در آن $x-h \leq \xi_- \leq x \leq \xi_+ \leq x+h$. از قرار دادن این مقادیر در فرمول (۱۴.۷.۵) و مرتب کردن برحسب توان h :

$$Af(x+h) + Bf(x) + Cf(x-h)$$

$$= (A+B+C)f(x) + h(A-C)f'(x) + \frac{h^2}{2}(A+C)f''(x)$$

$$+ \frac{h^3}{6}(A-C)f^{(3)}(x) + \frac{h^4}{24}[Af^{(4)}(\xi_+) + Bf^{(4)}(\xi_-)] \quad (۱۵.۷.۵)$$

برای آنکه این عبارت برابر $f''(x)$ شود، قرار می‌دهیم

$$A + B + C = 0 \quad A - C = 0 \quad A + C = \frac{2}{h^2}$$

جواب این دستگاه چنین است

$$A = C = \frac{1}{h^2} \quad B = -\frac{2}{h^2} \quad (۱۶.۷.۵)$$

که فرمول زیر را به دست می‌دهد

$$D_h^{(2)} f(x) = \frac{f(x+h) - 2f(x) + f(x-h)}{h^2} \quad (۱۷.۷.۵)$$

برای خطا، (۱۶.۷.۵) را در (۱۵.۷.۵) قرار می‌دهیم و از (۱۷.۷.۵) استفاده می‌کنیم تا به دست آوریم

$$f''(x) - D_h^{(2)} f(x) = \frac{h^2}{24} [f^{(4)}(\xi_+) + f^{(4)}(\xi_-)]$$

با استفاده از مسأله ۱ فصل ۱ و با فرض اینکه $f(x)$ چهار بار پیوسته مشتقپذیر است

$$f''(x) - D_h^{(2)} f(x) = -\frac{h^2}{12} f^{(4)}(\xi) \quad (۱۸.۷.۵)$$

به‌ازای $x-h \leq \xi \leq x+h$. فرمولهای (۱۷.۷.۵) و (۱۸.۷.۵) را می‌توانستیم از محاسبه $p''(x)$ درونیاب چندجمله‌یی درجه دوم $f(x)$ در $x-h, x, x+h$ به دست آوریم ولی روش فوق احتمالاً ساده‌تر است.

نظر کلی در روش ضرایب نامعین آن است که ضرایب بسط تیلر را بر حسب h به‌گونه‌ای پیدا کنیم که مشتق (یا انتگرال) مورد نظر تا حد ممکن به مقدار حقیقی نزدیک باشد.

اثر خطا در مقادیر تابع فرمولهای قبلی برای پیدا کردن روشهای حل معادلات دیفرانسیل معمولی و معادلات دیفرانسیل جزئی مفیدند، ولی اگر برای مقادیر تابع که به‌صورت تجربی به دست آمده‌اند به‌کار روند ممکن است خطاهای جدی ایجاد نمایند. برای نشان دادن روشی که اثرات چنین خطاهایی را تحلیل نماید تقریب مشتق دوم (۱۷.۷.۵) را در نظر می‌گیریم.

مطلب را با بازنویسی (۱۷.۷.۵) به شکل زیر آغاز می‌کنیم

$$f''(x_1) \doteq D_h^{(2)} f(x_1) = \frac{f(x_2) - 2f(x_1) + f(x_0)}{h^2}$$

که $x_j = x_0 + jh$. x_j گیریم که مقادیر واقعی که به کار رفته‌اند \tilde{f}_i باشند با

$$f(x_i) = \tilde{f}_i + \varepsilon_i \quad i = 0, 1, 2 \quad (19.7.5)$$

مشتق عددی واقعی محاسبه شده چنین است

$$\tilde{D}_h^{(2)} f(x_1) = \frac{\tilde{f}_2 - 2\tilde{f}_1 + \tilde{f}_0}{h^2} \quad (20.7.5)$$

برای خطای آن (۱۹.۷.۵) را در (۲۰.۷.۵) می‌گذاریم و به دست می‌آوریم

$$\begin{aligned} f''(x_1) - \tilde{D}_h^{(2)} f(x_1) &= f''(x_1) - \frac{f(x_2) - 2f(x_1) + f(x_0)}{h^2} + \frac{\varepsilon_2 - 2\varepsilon_1 + \varepsilon_0}{h^2} \\ &= \frac{-h^2}{12} f^{(4)}(\xi) + \frac{\varepsilon_2 - 2\varepsilon_1 + \varepsilon_0}{h^2} \end{aligned} \quad (21.7.5)$$

برای عبارت شامل $\{\varepsilon_i\}$ فرض کنید که این خطاها در بازه $-E \leq \varepsilon \leq E$ ، تصادفی باشند.

در این صورت

$$|f''(x_1) - \tilde{D}_h^{(2)} f(x_1)| \leq \frac{h^2}{12} |f^{(4)}(\xi)| + \frac{4E}{h^2} \quad (22.7.5)$$

و این کران اخیر در بسیاری از مواد به دست می‌آید. یک مثال برای چنین خطاهایی، خطاهای گرد کردن است که E کران اندازه آنهاست.

کران خطا در (۲۲.۷.۵) در ابتدا با کاهش h کوچکتر می‌شود ولی برای h به اندازه کافی نزدیک به صفر، خطا دوباره شروع به افزایش می‌کند. یک مقدار بهینه برای h وجود دارد که آن را h^* می‌نامیم که سمت راست (۲۲.۷.۵) را مینیمم می‌کند، و احتمالاً مقدار مشابهی نیز برای خطای واقعی $f''(x_1) - \tilde{D}_h^{(2)} f(x_1)$ وجود خواهد داشت.

مثال گیریم $f(x) = -\cos(x)$ ، با استفاده از تقریب عددی (۱۷.۷.۵)، $f''(0)$ را محاسبه می‌کنیم. در جدول ۲۵.۵، دو ستون خطا داده‌ایم، خطاهای (۱) $D_h^{(2)} f(0)$ که دقیقاً محاسبه شده است؛ (۲) خطاهای $\tilde{D}_h^{(2)} f(0)$ که با حساب اعشاری هشت رقمی محاسبه و گرد شده است. در این حالت اخیر،

$$|f''(0) - \tilde{D}_h^{(2)} f(0)| \leq \frac{h^2}{12} + \frac{2 \times 10^{-8}}{h^2} \quad (23.7.5)$$

کران (۲۳.۷.۵) در $h^* = 0.022$ مینیمم می‌شود که با خطاهای $f''(0) - \tilde{D}_h^{(2)} f(0)$ که در جدول داده شده‌اند سازگار است. برای $D_h^{(2)} f(0)$ که دقیقاً محاسبه شده است، توجه کنید که وقتی h نصف می‌شود، خطاها با نسبت $1/4$ کوچک می‌شوند که با فرمول خطای (۱۸.۷.۵) سازگار است.

جدول ۲۵.۵ مثال $\tilde{D}_h^{(r)}$ و $D_h^{(r)} f(\circ)$

$f''(\circ) - \tilde{D}_h^{(r)} f(\circ)$	نسبت	$f''(\circ) - D_h^{(r)} f(\circ)$	h
$۲,۰۷E - ۲$		$۲,۰۷E - ۲$	$۰,۵$
$۵,۲۰E - ۳$	$۳,۹۸$	$۵,۲۰E - ۳$	$۰,۲۵$
$۱,۳۰E - ۳$	$۳,۹۹$	$۱,۳۰E - ۳$	$۰,۱۲۵$
$۳,۲۵E - ۴$	$۴,۰۰$	$۳,۲۵E - ۴$	$۰,۰۶۲۵$
$۸,۴۵E - ۵$	$۴,۰۰$	$۸,۱۴E - ۵$	$۰,۰۳۱۲۵$
$۲,۵۶E - ۶$	$۴,۰۰$	$۲,۰۳E - ۵$	$۰,۰۱۵۶۲۵$
$-۷,۹۴E - ۵$	$۴,۰۰$	$۵,۰۹E - ۶$	$۰,۰۰۷۸۱۲۵$
$-۷,۹۴E - ۵$	$۴,۰۰$	$۱,۲۷E - ۶$	$۰,۰۰۳۹۰۶۲۵$
$-۱,۳۹E - ۳$	$۴,۰۰$	$۳,۱۸E - ۷$	$۰,۰۰۱۹۵۳۱۲۵$

بحث در آثار خواندنی

اگر چه مبحث انتگرالگیری عددی یکی از قدیمیترین مباحث آنالیز عددی است و نوشته‌های زیادی در مورد آنها موجود است، مقالات جدیدی به مقدار نسبتاً زیاد مرتباً منتشر می‌شوند. بسیاری از این مقاله‌ها روشهایی را برای رده‌هایی خاص از مسائل، مثل انتگرالهای نوسانی، مطرح می‌کنند، و عده دیگر واکنشی به تغییرات در رایانه‌ها، مثل استفاده از معماریهای خط لوله برداری هستند. بهترین بررسی انتگرالگیری عددی کار وسیع و مشروح دیویس و رابینوویس (۱۹۸۴) است. این اثر شامل بررسی فراگیر از بیشتر روشهای انتگرالگیری، فهرست منابع جامع، مجموعه‌ای از برنامه‌های رایانه‌یی و فهرست برنامه‌های انتگرالگیری منتشر شده، است. همچنین شامل مقاله آبراموویس درباره «محاسبه عملی انتگرالها» است، که توصیه‌های بسیار عالی در روشهای تحلیلی انتگرالگیری به دست می‌دهد. کتابهای درسی مهم دیگر در مورد انتگرالگیری عددی عبارت‌اند از، انگلس (۱۹۸۰)، کریلوف (۱۹۶۲)، استرود (۱۹۷۱). برای تاریخچه روشهای انتگرالگیری عددی کلاسیک، گولدستاین (۱۹۷۷) را ببینید.

به علت کمی جا، مجبور بوده‌ایم بعضی از مفاهیم مهم را حذف کنیم. مهمترین آنها عبارت‌اند از (۱) انتگرالگیری کلنشار-کرتیس^۱ و (۲) انتگرالگیری چند متغیره. اولی مبتنی است بر انتگرالگیری بسط چیشیف انتگرالده: که به تجربه ثابت شده است که برای انواع بسیاری از انتگرالها روشی عالی است. این روش ابتکاری در کلنشار-کرتیس (۱۹۶۰) عرضه شده است. یک شرح معمولی از آن در (صفحات ۲۸ تا ۲۹) پیسنز و همکاران (۱۹۸۳) داده شده است. زمینه انتگرالگیری

چند متغیره یک زمینه بسیار فعال پژوهشی است، و کتابهای درسی انگلس (۱۹۸۰) و استرود (۱۹۷۱) بهترین مدخل بر این زمینه‌اند. به لحاظ استفاده گسترده از انتگرالگیری چند متغیره در روش عنصر متناهی برای حل معادلات دیفرانسیل جزئی، کتابهای روش عنصر متناهی اغلب شامل فرمولهای انتگرالگیری برای ناحیه‌های مثلثی و مستطیلی هستند.

انتگرالگیری عددی خودکار-وقتی احساس شد که اغلب انتگرالگیریهای عددی را می‌توان از این طریق محاسبه کرد- یک زمینه بسیار فعال پژوهشی در دهه‌های ۱۹۶۰ و ۱۹۷۰ بود. اخیراً یک بازگشت به سوی استفاده وسیع از انتگرالگیریهای غیر خودکار، به ویژه برای انتگرالهای دم‌دستی، صورت گرفته است. یک بحث عالی درباره مزایا و معایب نسبی انتگرالگیری خودکار در لاینس (۱۹۸۳) داده شده است. تواناترین و انعطاف‌پذیرترین برنامه‌های خودکار امروزی شاید آتھایی باشند که در QUADPACK داده شده‌اند، این برنامه‌ها را پیسنز و همکاران (۱۹۸۳) بحث و روشن کرده‌اند. نسخه‌هایی از QUADPACK در IMSL و NAG موجود است.

برای ریزرایانه‌ها و محاسبات دستی، قاعده سیمپسون، به دلیل سادگی، هنوز خیلی طرفدار دارد. مع‌هذا عنایت خاصی بایستی به روش گاوس به دلیل دقت بالاتر آن بشود. گره‌ها و وزنها در آبرموویتس و اشتگون (۱۹۶۴) و در استرود و سكرت (۱۹۶۶) اکنون در دسترس قرار داده شده‌اند و برنامه‌هایی برای محاسبه آنها مهیا شده‌اند.

مشتقگیری عددی به تعبیر بخش ۶.۱ یک مسأله بد حالت است. روندهای مشتقگیری عددی که در آنها این نکته در نظر گرفته شده است در ده پانزده سال گذشته توسعه یافته‌اند. به ویژه، اندرسون و بلومفیلد (۱۹۷۴a) (۱۹۷۴b)، کالوم^۱ (۱۹۷۲) و ابا^۲ (۱۹۸۰)، ولترینگ^۳ (۱۹۸۶) را ببینید.

مراجع

- Abramowitz, M., and I. Stegun, Eds. (1964). *Handbook of Mathematical Functions*. National Bureau of Standards, U.S. Government Printing Office, Washington, D.C.
- Anderssen, R., and P. Bloomfield (1974a). Numerical differentiation procedures for non-exact data, *Numer. Math.*, **22**, 157-182.
- Anderssen, R., and P. Bloomfield (1974b). A time series approach to numerical differentiation, *Technometrics* **16**, 69-75.
- Atkinson, K. (1976). *A Survey of Numerical Methods for the Solution of Fredholm Integral Equations of the Second Kind*. Society for Industrial and Applied Mathematics, Philadelphia.

- Atkinson, K. (1982). Numerical integration on the sphere, *J. Austr. Math. Soc. (Ser. B)* **23**, 332–347.
- Bauer, F., H. Rutishauser, and E. Stiefel (1963). New aspects in numerical quadrature. In *Experimental Arithmetic, High Speed Computing, and Mathematics*, pp. 199–218. Amer. Math. Soc., Providence, R.I.
- de Boor, C. (1971). CADRE: An algorithm for numerical quadrature. In *Mathematical Software*, pp. 201–209. Academic Press, New York.
- Clenshaw, C., and A. Curtis (1960). A method for numerical integration on an automatic computer, *Numer. Math.* **2**, 197–205.
- Cryer, C. (1982). *Numerical Functional Analysis*. Oxford Univ. Press (Clarendon), Oxford, England.
- Cullum, J. (1971). Numerical differentiation and regularization, *SIAM J. Numer. Anal.* **8**, 254–265.
- Davis, P. (1963). *Interpolation and Approximation*. Ginn (Blaisdell), Boston.
- Davis, P., and P. Rabinowitz (1984). *Methods of Numerical Integration*, 2nd ed. Academic Press, New York.
- Dixon, V. (1974). Numerical quadrature: A survey of the available algorithms. In *Software for Numerical Mathematics*, D. Evans, Ed., pp. 105–137. Academic Press, London.
- Donaldson, J., and D. Elliott (1972). A unified approach to quadrature rules with asymptotic estimates of their remainders, *SIAM J. Numer. Anal.* **9**, 573–602.
- de Doncker, E., and R. Piessens (1976). Algorithm 32: Automatic computation of integrals with singular integrand, over a finite or an infinite interval, *Computing* **17**, 265–279.
- Engels, H. (1980). *Numerical Quadrature and Cubature*. Academic Press, New York.
- Goldstine, H. (1977). *A History of Numerical Analysis*. Springer-Verlag, New York.
- Golub, G., and J. Welsch (1969). Calculation of Gauss quadrature rules, *Math. Comput.* **23**, 221–230.
- de Hoog, F., and R. Weiss (1973). Asymptotic expansions for product integration, *Math. Comput.* **27**, 295–306.
- Iri, M., S. Moriguti, and Y. Takasawa (1970). On a numerical integration formula (in Japanese), *J. Res. Inst. Math. Sci.*, **91**, 82. Kyoto Univ., Kyoto, Japan.
- Isaacson, E., and H. Keller (1966). *Analysis of Numerical Methods*. Wiley, New York.
- Kronrod, A. (1965). *Nodes and Weights of Quadrature Formulas*, Consultants Bureau, New York.

- Krylov, V. (1962). *Approximate Calculation of Integrals*. Macmillan, New York.
- Lyness, J. (1983). When not to use an automatic quadrature routine, *SIAM Rev.* **25**, 63-88.
- Lyness, J., and C. Moler (1967). Numerical differentiation of analytic functions, *SIAM J. Num. Anal.* **4**, 202-210.
- Lyness, J., and J. Kaganove (1976). Comments on the nature of automatic quadratic routines, *ACM Trans. Math. Softw.* **2**, 65-81.
- Lyness, J., and K. Puri (1973). The Euler-MacLaurin expansion for the simplex, *Math. Comput.* **27**, 273-293.
- Lyness, J., and B. Ninham (1967). Numerical quadrature and asymptotic expansions, *Math. Comput.* **21**, 162-178.
- Patterson, T. (1968). The optimum addition of points to quadrature formulae, *Math. Comput.* **22**, 847-856. (Includes a microfiche enclosure; a microfiche correction is in **23**, 1969.)
- Patterson, T. (1973) Algorithm 468: Algorithm for numerical integration over a finite interval. *Commun. ACM* **16**, 694-699.
- Piessens, R., E. deDoncker-Kapenga, C. Überhuber, and D. Kahaner (1983). *QUADPACK: A Subroutine Package for Automatic Integration*. Springer-Verlag, New York.
- Ralston, A. (1965). *A First Course in Numerical Analysis*. McGraw-Hill, New York.
- Robinson, I. (1979). A comparison of numerical integration programs, *J. Comput. Appl. Math.* **5**, 207-223.
- Robinson, I., and E. de Doncker (1981). Automatic computation of improper integrals over a bounded or unbounded planar region, *Computing* **27**, 253-284.
- Stenger, F. (1981). Numerical methods based on Whittaker cardinal or sinc functions, *SIAM Rev.* **23**, 165-224.
- Stroud, A. (1971). *Approximate Calculation of Multiple Integrals*. Prentice-Hall, Englewood Cliffs, N.J.
- Stroud, A., and D. Secrest (1966). *Gaussian Quadrature Formulas*. Prentice-Hall, Englewood Cliffs, N.J.
- Wahba, G. (1980). Ill-posed problems: Numerical and statistical methods for mildly, moderately, and severely ill-posed problems with noisy data, Tech. Rep. #595, Statistics Dept., Univ. of Wisconsin, Madison. (Prepared for the Proc. Int. Symp. on Ill-posed Problems, Newark, Del., 1979.)
- Woltring, H. (1986). A Fortran package for generalized, cross-validatorspline smoothing and differentiation, *Adv. Eng. Softw.* **8**, 104-113.

مسائل

۱. برنامه‌ای بنویسید که با استفاده از قاعدهٔ دوزنقه‌ی و n تقسیم جزئی، $I = \int_a^b f(x) dx$ را محاسبه نماید. نتیجه را I_n بنامید. این برنامه را برای محاسبهٔ انتگرالهای زیر با $n = 2, 4, 8, 16, \dots, 512$ به کار برید

$$\int_{-2}^2 \frac{dx}{1+x^2} \quad (\text{ج}) \quad \int_0^1 x^{5/2} dx \quad (\text{ب}) \quad \int_0^1 e^{-x^2} dx \quad (\text{الف})$$

$$\int_0^\pi e^x \cos(4x) dx \quad (\text{د}) \quad \int_0^{2\pi} \frac{dx}{2 + \cos(x)}$$

نرخ همگرایی I_n به I را به صورت تجربی با محاسبهٔ نسبت‌های (۲۷.۴.۵) تحلیل کنید:

$$R_n = \frac{I_{2n} - I_n}{I_{4n} - I_{2n}}$$

۲. مسألهٔ ۱ را با استفاده از قاعدهٔ سیمپسون تکرار کنید.

۳. قاعدهٔ دوزنقه‌ی تصحیح شده (۱۲.۱.۵) را برای انتگرالهای مسألهٔ ۱ به کار برید. نتایج را با نتایج مسألهٔ ۲ که با قاعدهٔ سیمپسون به دست آمده است مقایسه کنید.

۴. به عنوان روش دیگر برای قاعدهٔ دوزنقه‌ی تصحیح شده (۱۲.۱.۵)، چند جمله‌ی درجهٔ سوم درونیابی ارمیت را برای $f(x)$ به کار برید تا

$$\int_a^b f(x) dx \doteq \left(\frac{b-a}{2} \right) [f(a) + f(b)] - \frac{(b-a)^2}{12} [f'(b) - f'(a)]$$

را به دست آورید. فرمول خطا برای درونیابی ارمیت را به کار برید تا یک فرمول خطا برای

تقریب پیشین پیدا شود. این نتایج را برای (۱۲.۱.۵)، با تقسیم $[a, b]$ به n جزء، تعمیم دهید.

۵. (الف) فرض کنید که در بازهٔ $[0, b]$ ، $f(x)$ پیوسته و $f'(x)$ انتگرالپذیر باشد. نشان دهید که خطا در قاعدهٔ دوزنقه‌ی برای محاسبهٔ $\int_0^1 f(x) dx$ به شکل زیر است

$$E_n(f) = \int_0^1 K(t) f'(t) dt$$

$$K(t) = \frac{t_{j-1} + t_j}{2} - t \quad t_{j-1} \leq t \leq t_j \quad j = 1, \dots, n$$

این نتایج را با (۲۲.۱.۵) که در آن $f'(x)$ پیوسته و $f''(x)$ انتگرالپذیر است، مقایسه کنید.

(ب) نتیجه را برای $f(x) = x^\alpha \log(x)$ و $f(x) = x^\alpha$ ، $0 < \alpha < 1$ به کار برید. این یک مرتبهٔ همگرایی به دست می‌دهد، اگرچه کوچکتر از مرتبهٔ واقعی است [مسألهٔ ۶ و (۲۳.۴.۵) را ببینید].

۶. با استفاده از برنامهٔ مسألهٔ ۱ نرخ همگرایی قاعدهٔ دوزنقه‌ی را به طور تجربی برای $\int_0^1 x^\alpha \ln(x) dx$ در حالت $0 \leq \alpha \leq 1$ برای $0, 1, 25, 50, 75, 100$ به دست آورید.

۷. شکل مرکب قاعده بول را که در جدول ۸.۵ در درایه $n = 4$ داده شده است به دست آورید. فرمولهای خطایی مشابه فرمولهایی که در (۱۷.۱.۵) و (۱۸.۱.۵) برای قاعده سیمپسون داده شده اند پیدا کنید.

۸. مسأله ۱ را با استفاده از قاعده بول که در مسأله ۷ به دست آمده تکرار کنید.

۹. گیریم $p_2(x)$ چندجمله‌یی درجه دومی باشد که $f(x)$ را در $x = 0, h, 2h$ درونیابی می‌کند. با استفاده از این درونیاب یک فرمول انتگرالگیری عددی I_h برای $I = \int_0^{2h} f(x) dx$ پیدا کنید بسط یک سری تیلر $f(x)$ را به کار برید تا نشان دهید که

$$I - I_h = \frac{3}{8} h^4 f^{(3)}(0) + O(h^5)$$

۱۰. برای فرمول انتگرالگیری میانگامی (۱۷.۲.۵)، فرمول خطای هسته‌یی پتانوزیر را به دست آورید

$$\int_0^h f(x) dx - hf\left(\frac{h}{2}\right) = \int_0^h K(t) f''(t) dt$$

$$K(t) = \begin{cases} \frac{1}{2} t^2 & 0 \leq t \leq \frac{h}{2} \\ \frac{1}{2} (h-t)^2 & \frac{h}{2} \leq t \leq h \end{cases}$$

با استفاده از این فرمول جمله خطا در (۱۷.۲.۵) را پیدا کنید.

۱۱. توابعی را که به طریق زیر تعریف شده اند در نظر می‌گیریم. گیریم $n > 0$ و $h = (b-a)/n$ $t_j = a + jh$ برای $j = 0, 1, \dots, n$. در هر زیر بازه $[t_{j-1}, t_j]$ به ازای $j = 1, \dots, n$ خطی باشد. نشان دهید که مجموعه همه چنین f هایی، برای مقادیر $n \geq 1$ در $C[a, b]$ چگال است.

۱۲. گیریم $w(x)$ یک تابع انتگرالپذیر بر بازه $[a, b]$ باشد،

$$\int_a^b |w(x)| dx < \infty$$

و گیریم

$$\int_a^b w(x) f(x) dx \doteq \sum_{j=1}^n w_{j,n} f(x_{j,n})$$

دنباله‌ای از قاعده‌های انتگرالگیری عددی باشد. قضیه ۲.۵ را برای این حالت تعمیم دهید.

۱۳. فرض کنید قاعده انتگرالگیری عددی

$$\int_a^b f(x) dx \doteq \sum_{j=1}^n w_{j,n} f(x_{j,n})$$

برای کلیه توابع پیوسته همگرا باشد. اثر خطاها را در مقادیر تابع در نظر بگیرید. فرض کنید که از

$$\tilde{f}_i = f(x_i) \text{ با}$$

$$|f(x_i) - \tilde{f}_i| \leq \varepsilon \quad 1 \leq i \leq n$$

استفاده می‌کنیم. تأثیر این خطاها در مقادیر تابع، بر این انتگرالگیری عددی چه خواهد بود؟

۱۴. انتگرالگیری گاوس-لژاندر را برای انتگرالهای مسئله ۱ به‌کار برید. نتایج را با نتایج قاعده‌های ذوزنقی و سیمپسون مقایسه کنید.

۱۵. انتگرالگیری گاوس-لژاندر را برای محاسبه $\int_{-1}^1 dx/(1+x^2)$ با فرمول $n = 2, 4, 6, 8$ نقطه‌گرهی به‌کار برید. نتایج را با نتایج جدول ۹.۵ که از فرمول نیوتن-کوتس حاصل شده است مقایسه کنید.

۱۶. ثابت کنید که گره‌ها و وزنه‌های گاوس-لژاندر در $[-1, 1]$ به‌طور متقارن حول $x = 0$ توزیع شده‌اند.

۱۷. فرمول: انتگرالگیری گاوسی را برای

$$I(f) = \int_0^1 f(x) \log\left(\frac{1}{x}\right) dx$$

که در آن تابع وزن $w(x) = \log(1/x)$ است پیدا کنید. راهنمایی: مسئله ۲۰ (الف) از فصل ۴ را ببینید. همچنین از مشابه (۷.۳.۵) برای محاسبه وزنها استفاده کنید نه از طریق فرمول (۱۱.۳.۵).

۱۸. فرمولهای انتگرالگیری گاوسی یک نقطه‌یی و دو نقطه‌یی را برای

$$I = \int_0^1 xf(x)dx = \sum_{j=1}^n w_j f(x_j)$$

با تابع وزن $w(x) = x$ به دست آورید.

۱۹. برای انتگرال $I = \int_{-1}^1 \sqrt{1-x^2} f(x) dx$ با وزن $w(x) = \sqrt{1-x^2}$ فرمولهای صریح برای وزنها و گره‌های فرمول انتگرالگیری گاوسی به دست آورید. همچنین فرمول خطا را پیدا کنید. راهنمایی: مسئله ۲۴ فصل ۴ را ببینید.

۲۰. با استفاده از ستون $e_{n,4}$ جدول ۱۳.۵، فرمول خطای مرتبه ۴ مشابه فرمول خطای مرتبه ۲ (۴۰.۳.۵) را برای $e_{n,2}$ پیدا کنید. آن را با فرمول خطای مرتبه چهار (۱۷.۱.۵) سیمپسون مقایسه کنید.

۲۱. در فرمول (۴۸.۳.۵) کروند^۱ وزنها را از حل دستگاه چهار معادله خطی، می‌توان محاسبه کرد. این دستگاه را پیدا و حل کنید تا مقادیر داده شده بعد از (۴۸.۳.۵) تأیید شود. راهنمایی: طریقی را که به (۷.۳.۵) انجامید به‌کار برید.

۲۲. فرمول هفت نقطه‌ی گاوس - لژاندر را با فرمول هفت نقطه‌ی (۴۸.۳.۵) کرونرود مقایسه کنید. هر یک از آنها را برای انتگرالهای گوناگون به کار برید و خطاهای متناظر آنها را مقایسه کنید.
۲۳. (الف) رابطه (۷.۴.۵) را برای چند جمله‌یهای $B_j(x)$ برنولی و اعداد B_j برنولی پیدا کنید. نشان دهید که برای تمام اعداد صحیح فرد $z \geq 3$, $B_j = 0$.
- (ب) اتحادهای زیر را به دست آورید

$$B'_j(x) = jB_{j-1}(x) \quad j \geq 2 \text{ زوج}$$

$$B'_j(x) = j[B_{j-1}(x) + B_{j-1}] \quad j \geq 3 \text{ فرد}$$

- این اتحادها را می‌توان برای یک اثبات کلی فرمول اوایلر - مک لورن (۹.۴.۵) به کار گرفت.
۲۴. با استفاده از فرمول مجموعیابی اوایلر - مک لورن (۱۷.۴.۵)، برآوردی برای $\zeta(\frac{5}{4})$ با دقتی تا سه رقم اعشار پیدا کنید. $\zeta(p)$ تابع زتا در (۲۲.۴.۵) تعریف شده است.
۲۵. فرمول خطای مجانبی برای قاعده دوزنقه‌ی را که برای $\int_0^1 \sqrt{x} f(x) dx$ به کار رفته پیدا کنید. برآورد در مسأله ۲۴ را به کار برید.
۲۶. جدول انتگرالهای تقریبی I_n در زیر را که با استفاده از قاعده سیمپسون به دست آمده‌اند در نظر می‌گیریم. مرتبه همگرایی I_n به I را پیش‌بینی کنید.

n	I_n
۲	۰٫۲۸۴۵۱۷۷۹۶۸۶
۴	۰٫۲۸۵۵۹۲۵۴۵۷۶
۸	۰٫۲۸۵۷۰۲۴۸۷۴۸
۱۶	۰٫۲۸۵۷۱۳۱۷۷۳۱
۳۲	۰٫۲۸۵۷۱۴۱۸۳۶۳
۶۴	۰٫۲۸۵۷۱۴۲۷۶۴۳

- یعنی اگر $I - I_n = c/n^p$ چیست؟ آیا به نظر می‌آید که این شکل معتبری برای خطای این داده‌ها باشد؟ مقداری برای c و خطای I_{64} پیش‌بینی کنید. n به چه بزرگی باید انتخاب شود اگر بخواهیم I_n خطایی کوچکتر از 10^{-11} داشته باشد؟
۲۷. فرض کنید که خطا در یک فرمول انتگرالگیری بسط مجانبی زیر را داشته باشد

$$I - I_n = \frac{C_1}{n\sqrt{n}} + \frac{C_2}{n^2} + \frac{C_3}{n^2\sqrt{n}} + \frac{C_4}{n^3} + \dots$$

فرآیند برونمایی ریچاردسن از بخش ۵.۴ را برای پیدا کردن فرمول‌هایی برای C_1 و C_2 تعمیم دهید. فرض کنیم سه مقدار I_n , I_{2n} , I_{4n} محاسبه شده باشند. این مقادیر را برای محاسبه C_1 و C_2 و برآوردی از I با خطایی از مرتبه $1/n^2 \sqrt{n}$ به کار برید.

۲۸. برای قاعدهٔ دوزنقه‌یی (که با $I_n^{(T)}$ نمایش داده شده) در محاسبهٔ $I = \int_a^b f(x) dx$ فرمول خطای مجانبی را به شکل زیر داریم

$$I - I_n^{(T)} = -\frac{h^2}{12} [f'(b) - f'(a)] + O(h^4)$$

و برای فرمول میانگاهی $I_n^{(M)}$ داریم

$$I - I_n^{(M)} = \frac{h^2}{24} [f'(b) - f'(a)] + O(h^4)$$

به شرط آنکه f به اندازهٔ کافی در $[a, b]$ مشتقپذیر باشد. با استفاده از این نتایج، از ترکیب $I_n^{(M)}$ و $I_n^{(T)}$ یک فرمول انتگرالگیری عددی جدید \tilde{I}_n با درجهٔ همگرایی بالاتر به دست آورید. وزنهای فرمول جدید \tilde{I}_n را بنویسید.

۲۹. یک فرمول خطای مجانبی برای قاعدهٔ سیمپسون، قابل مقایسه با فرمول اوایلر-مک لورن (۹.۴.۵) برای قاعدهٔ دوزنقه‌یی، به دست آورید. از (۹.۴.۵)، (۳۳.۴.۵) و (۳۶.۴.۵) به همان‌گونه که در (۳۷.۴.۵) عمل کردیم، استفاده کنید.

۳۰. نشان دهید که فرمول (۴۰.۴.۵) برای $I_n^{(2)}$ قاعدهٔ مرکب بول^۱ است. مسألهٔ ۷ را ببینید.
۳۱. الگوریتم رامبرگ بخش ۴.۵ را اجرا کنید و سپس آن را برای انتگرالهای مسألهٔ ۱ به کار برید. نتایج را با نتایج قاعده‌های سیمپسون و دوزنقه‌یی مقایسه کنید.

۳۲. انتگرال $I = \int_0^1 x^\alpha dx$ را با استفاده از قاعدهٔ تطبیقی سیمپسون که به دنبال (۳.۵.۵) شرح داده شده بود محاسبه کنید. برای اینکه ببینید آزمون (۵.۵.۵) ممکن است به شکست انجامد حالت $0 < \alpha < 1$ را در نظر بگیرید که انتگرالده به طور دلخواه در $x = 0$ برابر صفر گرفته شده است. نشان دهید که برای مقادیر به اندازهٔ کافی کوچک ε ، آزمون (۵.۵.۵) در زیر بازهٔ $[0, h]$ هرگز صدق نمی‌کند. توجه داشته باشید که برای ε مشخص در $[0, 1]$ ، تحمل خطا برای $[0, h]$ برابر εh خواهد بود.

۳۳. قاعدهٔ سیمپسون با فاصله‌گذاری برابر را برای انتگرالگیری $\int_0^1 \log(x) dx$ به کار برید. به ازای $x = 0$ ، انتگرالده را برابر صفر بگیرید. نتایج را با آنچه در جداول (۱۹.۵) و (۲۰.۵) برای I_4 داده شده مقایسه نمایید.

۳۴. از یک برنامه انتگرالگیری تطبیقی [مثل DQAGP از پیسنز و همکاران (۱۹۸۳)] برای محاسبه انتگرالهای مسأله ۱ استفاده کنید. نتایج را با نتایج مسائل ۱، ۲، ۱۴، ۳۱ مقایسه نمایید.

۳۵. رفتار تکین انتگرالده در

$$I = \int_0^1 f(x) \log(x) dx$$

را با استفاده از تبدیل متغیر $x = t^r$ ، $r > 0$ ، کاهش دهید. همواری انتگرالده نتیجه را تحلیل کنید. همچنین رفتار تجربی قاعده‌های دوزنقه‌ی و سیمپسون را برای مقادیر گوناگون r بررسی نمایید.

۳۶. روش IMT را، که به دنبال (۵.۶.۵) توضیح داده شده‌است، برای محاسبه $I = \int_0^\infty f(x) dx$ ، با تبدیل متغیری از $[0, \infty)$ به یک بازه متناهی، به‌کار برید. از همین روش برای محاسبه انتگرالهایی که به دنبال (۱۱.۶.۵) آمده‌اند استفاده کنید.

۳۷. انتگرالگیری گاوس-لاگر با $n = 2, 4, 6, 8$ نقطه گرهی را برای محاسبه انتگرالهای زیر به‌کار برید. از فرمول (۱۱.۶.۵) استفاده کنید تا انتگرالها به شکل مناسب درآیند.

$$\int_0^\infty \frac{x dx}{(1+x^2)^2} = \frac{1}{4} \quad (\text{ب}) \quad \int_0^\infty e^{-x^2} dx = \frac{\sqrt{\pi}}{2} \quad (\text{الف})$$

$$\int_0^\infty \frac{\sin(x)}{x} dx = \frac{\pi}{2} \quad (\text{ج})$$

۳۸. انتگرالهای تکین زیر را در حدود تحمل خطای مشخص شده محاسبه کنید. از برنامه خودکار استفاده نکنید.

$$\int_0^{\epsilon} \pi \cos(u) \log(u) du \quad \epsilon = 10^{-8} \quad (\text{الف})$$

$$\int_0^{\epsilon} x^2 \sin\left(\frac{1}{x}\right) dx \quad \epsilon = 10^{-2} \quad (\text{ب})$$

$$\int_0^1 \frac{\cos(x)}{\sqrt{x^2}} dx \quad \epsilon = 10^{-8} \quad (\text{ج})$$

$$\int_0^1 \frac{\sqrt{1-x^2}}{x^\alpha} dx \quad \epsilon = 10^{-5}, \quad \alpha = 1 - \frac{1}{\pi} \quad (\text{د})$$

۳۹. برای محاسبه انتگرالهای مسأله ۳۸ از یک برنامه تطبیقی (مثلاً DQAGP) استفاده کنید.

۴۰. (الف) قاعده انتگرالگیری دوزنقه‌ی حاصلضرب را برای محاسبه

$$I(f) = \int_0^b \frac{f(x)}{x^\alpha} dx \quad 0 < \alpha \leq 1$$

بسط دهید. وزن‌ها را به شکل مناسبی، مشابه با (۲۸.۶.۵) برای قاعده دوزنقه‌ی حاصلضرب وقتی $w(x) = \log(x)$ ، انتخاب کنید. همچنین یک فرمول خطا مشابه با (۲۴.۶.۵) به‌دست آورید.

(ب) با استفاده از نتایج قسمت (الف) برای محاسبه انتگرالهای زیر یک برنامه ساده بنویسید

$$\int_0^1 \frac{dx}{\sin(x^2/\pi)} \quad (\text{ii}) \qquad \int_0^1 \frac{e^x}{x^{1/\pi}} dx \quad (\text{i})$$

راهنمایی: برای قسمت (ii) ابتدا $u = x^{2/\pi}$ بگیرد.

۴۱. برای نشان دادن بدحالتی مشتقگیری از یک تابع $y(t)$ در یک بازه $[0, 1]$ ، محاسبه

$$x(t) = y'(t) \quad \text{و} \quad x_n(t) = y'_n(t), \quad \text{با}$$

$$y_n(t) = y(t) + \frac{1}{n} t^n \qquad n \geq 1$$

در نظر می‌گیریم. تعریف بدحالت و خوش حالت را از بخش ۶.۱ به یاد آورید. با استفاده از ساخت بالا نشان دهید که محاسبه $x(t)$ نسبت به تغییرات y ناپایدار است. برای اندازه‌گیری تغییرات در x و y از نرم ماکسیمم (۸.۱.۱) و (۴.۱.۴) استفاده کنید.

۴۲. با استفاده از رایانه یا ماشین حساب خود، مثال مشتقگیری عددی (۲۳.۷.۵) و جدول ۲۵.۵ را تکرار کنید.

۴۳. نتایج خطا برای $D_\delta f(m)$ در (۱۰.۷.۵) را مربوط به اثرات خطاهای گرد کردن، شبیه به نتایج (۲۱.۷.۵) - (۲۲.۷.۵) برای $\tilde{D}_h^{(2)} f(x)$ به دست آورید. آن را برای $f(x) = e^x$ در $x = 0$ به کار برید و نتایج را با نتایج واقعی محاسبات دستی خود مقایسه نمایید.

روشهای عددی برای معادلات دیفرانسیل معمولی

معادلات دیفرانسیل یکی از مهمترین ابزارهای ریاضی هستند که در مسائل مدلسازی در علوم فیزیکی به کار برده می شوند. در این فصل ما روشهای عددی حل مسائل معادلات دیفرانسیل معمولی را به دست آورده و تحلیل می نماییم. شکل اصلی مسأله ای که ما مطالعه می کنیم، مسأله مقدار اولیه است:

$$y' = f(x, y) \quad y(x_0) = Y_0. \quad (1.0.6)$$

تابع $f(x, y)$ باید به ازای همه مقادیر (x, y) در حوزه ای چون D از صفحه xy ، پیوسته باشد و (x_0, Y_0) نقطه ای در D است. نتایجی که برای (1.0.6) به دست می آید، هم برای دستگاه معادلات دیفرانسیل و هم برای معادلات از مراتب بالاتر، به آسانی تعمیم داده می شوند، به شرطی که نمادهای برداری و ماتریسی مناسب به کار گرفته شوند. این تعمیمها در دو بخش آینده بحث و نشان داده شده اند.

گوییم $Y(x)$ یک جواب (1.0.6) در $[a, b]$ است اگر برای همه مقادیر $a \leq x \leq b$

$$(x, Y(x)) \in D. \quad 1.$$

$$Y(x_0) = Y_0. \quad 2.$$

$$Y'(x) = f(x, Y(x)) \text{ موجود و} \quad 3.$$

در سرتاسر این فصل، نماد $Y(x)$ معرف جواب درست هر مسأله معادله دیفرانسیلی است که در نظر گرفته می شود.

مثال ۱. معادله کلی دیفرانسیل مرتبه اول خطی به شکل زیر است

$$y' = a_0(x)y + g(x) \quad a \leq x \leq b$$

که در آن ضرایب $a_0(x)$ و $g(x)$ بنا به فرض در $[a, b]$ پیوسته اند. حوزه D برای این مسأله چنین است

$$D = \{(x, y) \mid a \leq x \leq b, \quad -\infty < y < \infty\}$$

جواب دقیق این معادله را در هر کتاب درسی مقدماتی معادلات دیفرانسیل [مثلاً بویس^۱ و دپیرما^۲ (۱۹۸۶)] می توان یافت. به عنوان یک حالت خاص، معادله زیر را در نظر می گیریم

$$y' = \lambda y + g(x) \quad 0 \leq x < \infty \quad (۲.۰.۶)$$

با $g(x)$ پیوسته در $[0, \infty)$. جوابی که در $Y(0) = Y_0$ صدق می کند با رابطه زیر داده می شود

$$Y(x) = Y_0 e^{\lambda x} + \int_0^x e^{\lambda(x-t)} g(t) dt \quad 0 \leq x < \infty \quad (۳.۰.۶)$$

این رابطه بعداً برای نشان دادن نتایج نظری گوناگون به کار خواهد رفت.

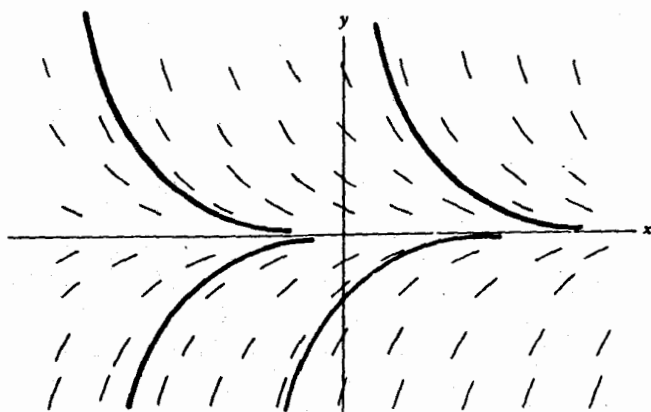
۲. معادله $y' = -y^2$ غیرخطی است. $Y(x) \equiv 0$ یکی از جوابهای آن است، بقیه جوابها به شکل

$$Y(x) = \frac{1}{x+c}$$

هستند، ثابت c اختیاری است. توجه کنید که $|Y(-c)| = \infty$. بنابراین همواری کلی $f(x, y) = -y^2$ همواری جوابها را تضمین نمی کند.

برای آنکه شناختی هندسی از جوابهای یک معادله دیفرانسیل مرتبه اول به دست آوریم، می توانیم به میدان سوهایی که این معادله در صفحه xy القا می کند نظر افکنیم. اگر $Y(x)$ جوابی باشد که از (x_0, Y_0) بگذرد، آنگاه شیب $Y(x)$ در (x_0, Y_0) برابر است با $Y'(x_0) = f(x_0, Y_0)$. در حوزه D $f(x, y)$ یک مجموعه نقاط (x, y) را گرفته یک پاره خط کوتاه با شیب $f(x, y)$ از هریک از نقاط (x, y) رسم می کنیم.

مثال معادله $y' = -y$ را در نظر بگیرید. میدان سوها در شکل ۱.۶ داده شده است، و چندین



شکل ۱.۶ میدان سوها برای $y' = -y$

جواب نمونه در آن رسم شده است. از نمودار روشن است که تمام جوابها در رابطه زیر صدق می‌کنند

$$\lim_{x \rightarrow \infty} Y(x) = 0$$

برای آنکه ترسیم میدان سوهای $y' = f(x, y)$ را ساده‌تر کنیم، دنبال آن خمهایی از صفحه xy می‌گردیم که در امتداد آنها $f(x, y)$ ثابت باشد. رابطه

$$f(x, y) = c$$

را به‌ازای مقادیر مختلف c حل می‌کنیم. هر جواب $Y(x)$ معادله $y' = f(x, y)$ که منحنی $f(x, y) = c$ را قطع کند، در نقطه برخورد در رابطه $Y'(x) = c$ صدق می‌نماید. خمهای $f(x, y) = c$ را خمهای تراز این معادله دیفرانسیل می‌نامند.

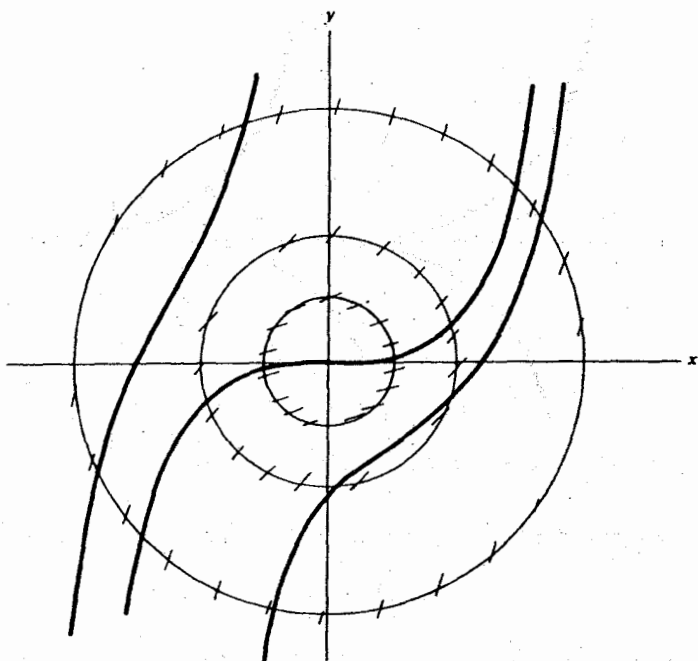
مثال نمایش رفتار کیفی جوابهای

$$y' = x^2 + y^2$$

را رسم کنید. خمهای تراز معادله با

$$x^2 + y^2 = c \quad c > 0$$

داده می‌شوند که دایره‌هایی به مرکز $(0, 0)$ و به شعاع \sqrt{c} هستند. میدان سوها با بعضی جوابهای نمونه در شکل ۲.۶ داده شده است. سه دایره با شعاعهای ۱ و $\frac{1}{2}$ و $c = \frac{1}{4}$ رسم شده است. این ترسیم یک اطلاع کیفی درباره جواب $Y(x)$ در اختیار ما می‌گذارد و گاهی اوقات می‌تواند مفید باشد.



شکل ۲.۶ میدان سوها برای $y' = x^2 + y^2$

۱.۶ وجود، یکتایی و نظریه پایداری

از مثالهای مقدمه باید مشهود باشد که در بیشتر حالات، مسأله مقدار اولیه (۱.۰.۶) یک جواب یکتا دارد. پیش از شروع به آنالیز عددی (۱.۰.۶) بعضی قضایای نظری آن را ارائه می‌دهیم. شرایطی داده شده است تا یکتایی جواب را تضمین نماید، و ما پایداری جواب را هنگامی که داده Y و مشتق $f(x, y)$ به اندازه کوچکی تغییر نمایند، بررسی خواهیم کرد. این قضایا برای بهتر فهمیدن روشهای عددی که بعداً ارائه خواهند شد لازم‌اند. از این پس فرض می‌شود که حوزه D در شرط فنی فرعی زیر صدق می‌نماید: اگر دو نقطه (x, y_1) و (x, y_2) هر دو به D تعلق داشته باشند، آنگاه پاره خط قائمی که این دو نقطه را به هم وصل می‌کند نیز به D تعلق دارد.

قضیه ۱.۶ گیریم $f(x, y)$ یک تابع پیوسته از x و y به ازای جمیع مقادیر (x, y) در D ، و (x_0, Y_0) یک نقطه داخلی در D باشد. فرض می‌کنیم $f(x, y)$ در شرط لیبشیتس زیر

$$|f(x, y_1) - f(x, y_2)| \leq K |y_1 - y_2| \quad (x, y_1), (x, y_2) \in D$$

(۱.۱.۶)

وجود، یکتایی و نظریهٔ پایداری ۳۷۷

برای مقداری از $K \geq 0$ ، صدق می‌کند. در این صورت در یک بازهٔ $I = [x_0 - \alpha, x_0 + \alpha]$ که مناسب انتخاب شده باشد، (۱.۰.۶) جواب یکتای $Y(x)$ دارد.

برهان برهان این قضیه را می‌توانید در اغلب کتابهای نظریهٔ معادلات دیفرانسیل معمولی [مثلاً بویس و دیپیرما (۱۹۸۶، ص ۹۵)] ببینید. به همین جهت ما آن را حذف می‌کنیم. ■

نشان می‌دهیم که شرط لیشیتس برقرار است اگر $\partial f(x, y)/\partial y$ در D موجود و کراندار باشد. می‌گیریم

$$K = \text{Max}_{(x,y) \in D} \left| \frac{\partial f(x, y)}{\partial y} \right| \quad (۲.۱.۶)$$

سپس با استفاده از قضیهٔ مقدار میانگین

$$f(x, y_1) - f(x, y_2) = \frac{\partial f(x, \xi)}{\partial y} (y_1 - y_2)$$

برای مقداری از ξ بین y_1 و y_2 . با استفاده از (۲.۱.۶) نتیجهٔ (۱.۱.۶) فوراً به‌دست می‌آید.

مثال دو مثال زیر این قضیه را روشن می‌سازند

۱. معادلهٔ $y' = 1 + \sin(xy)$ را با

$$D = \{(x, y) \mid 0 \leq x \leq 1, -\infty < y < \infty\}$$

در نظر می‌گیریم. برای محاسبهٔ K ، ثابت لیشیتس، (۲.۱.۶) را به‌کار می‌بریم. پس

$$\frac{\partial f(x, y)}{\partial y} = x \cdot \cos(xy) \quad K = 1$$

بنابراین برای هر (x_0, Y_0) ، با $0 < x_0 < 1$ ، یک جواب $Y(x)$ مربوط به مسألهٔ مقدار اولیه در بازه‌ای چون $[0, 1] \subset [x_0 - \alpha, x_0 + \alpha]$ وجود دارد.

۲. مسألهٔ

$$y' = \frac{2x}{a^2} y^2 \quad y(0) = 1$$

را برای هر ثابت $a > 0$ در نظر می‌گیریم. جواب چنین است

$$Y(x) = \frac{a^2}{a^2 - x^2} \quad -a < x < a$$

اگر a کوچک باشد، جواب فقط در بازه کوچکی وجود دارد. برای تعیین ثابت لپشیتس باید

$$\frac{\partial f(x, y)}{\partial y} = \frac{2xy}{a^2}$$

را حساب کنیم. برای آنکه یک ثابت متناهی لپشیتس روی D داشته باشیم، ناحیه D باید نسبت به x و y کراندار باشد، مثلاً $-c \leq x \leq c$ و $-b \leq y \leq b$. در این صورت، قضیه قبلی بیان می‌کند که یک جواب $Y(x)$ در بازه‌ای چون $-\alpha \leq x \leq \alpha$ با $\alpha \leq c$ وجود دارد.

مبادا چنین پنداشت که می‌توانیم مشتق جزئی در نقطه آغازی $(x_0, Y_0) = (0, 1)$ را حساب کرده اطلاعات کافی برای برآورد ثابت لپشیتس K به دست آوریم، توجه داشته باشید که برای هر مقدار $a > 0$

$$\frac{\partial f(0, 1)}{\partial y} = 0.$$

پایداری جواب پایداری جواب $Y(x)$ هنگامی بررسی می‌شود که مسأله مقدار اولیه به اندازه کوچکی تغییر یابد. این مطلب به بحث پایداری بخش ۶.۱ فصل ۱ مرتبط است. ما مسأله اختلال‌یافته زیر را، با همان مفروضات $f(x, y)$ در قضیه ۱.۶ در نظر می‌گیریم

$$\begin{aligned} y' &= f(x, y) + \delta(x) \\ y(x_0) &= Y_0 + \varepsilon \end{aligned} \quad (3.1.6)$$

وانگهی فرض می‌کنیم به ازای تمام مقادیر x هایی که به ازای یک y ، $(x, y) \in D$ باشد، $\delta(x)$ پیوسته است. در این صورت می‌توان نشان داد که مسأله (۳.۱.۶) یک جواب یکتا دارد که با $Y(x; \delta, \varepsilon)$ نشان داده می‌شود.

قضیه ۲.۶ همان فرضهای قضیه ۱.۶ را می‌پذیریم. در این صورت مسأله (۲.۱.۶) در بازه‌ای چون $[x_0 - \alpha, x_0 + \alpha]$ ، به ازای مقداری از $\alpha > 0$ ، یک جواب یکتای $Y(x; \delta, \varepsilon)$ یکنواخت دارد برای تمام اختلالهای ε و $\delta(x)$ که $\delta(x)$ برای مقادیر به اندازه کافی کوچک ε در

$$\|\varepsilon\| \leq \varepsilon, \quad \|\delta\|_{\infty} \leq \varepsilon.$$

صدق می‌کند. به علاوه، اگر $Y(x)$ جواب مسأله اختلال نیافته باشد، آنگاه

$$\text{Max}_{|x-x_0| \leq \alpha} |Y(x) - Y(x; \delta, \varepsilon)| \leq k[|\varepsilon| + \alpha \|\delta\|_\infty] \quad (۴.۱.۶)$$

که در آن، با استفاده از ثابت لیشیتس K از (۱.۱.۶)، $k = 1/(1 - \alpha K)$.

برهان به دست آوردن (۴.۱.۶) بسیار شبیه به اثبات قضیه ۱.۶ است و می‌توان آن را در بسیاری از کتابهای درسی کارشناسی ارشد در معادلات دیفرانسیل معمولی یافت. ■

با استفاده از این قضیه، می‌توانیم بگوییم که مسأله مقدار اولیه (۱.۰.۶)، به تعبیر بخش ۶.۱ فصل ۱، خوش حالت یا پایدار است. اگر تغییرات کوچکی در معادله دیفرانسیل یا در مقدار اولیه داده شود جواب نیز به اندازه کوچکی تغییر می‌کند. جواب Y به داده‌های مسأله، یعنی تابع f و مقدار اولیه Y_0 ، بستگی پیوسته دارد.

در بخش ۶.۱ اشاره شد که مسأله می‌تواند پایدار ولی نسبت به محاسبات عددی بد-وضع باشد. این وضعیت برای معادلات دیفرانسیل مصداق دارد، اگرچه در عمل زیاد پیش نمی‌آید. برای اینکه بهتر بفهمیم که چه موقع ممکن است چنین اتفاقی بیفتد، ما اختلال Y را که از اختلالهای مسأله نتیجه می‌شود برآورد می‌کنیم. برای ساده‌کردن بحث، فقط اختلال ε در مقدار اولیه Y_0 را در نظر می‌گیریم. اختلالات $\delta(x)$ در معادله تقریباً به همان‌گونه که در (۴.۱.۶) نشان داده شده است در جواب نهایی مسأله وارد می‌شوند.

مانند (۳.۱.۶)، اختلالی در مقدار آغازی Y_0 ایجاد می‌کنیم، گیریم $Y(x; \varepsilon)$ معرف جواب اختلال یافته باشد. پس

$$Y'(x; \varepsilon) = f(x, Y(x; \varepsilon)) \quad x_0 - \alpha \leq x \leq x_0 + \alpha \quad (۵.۱.۶)$$

$$Y(x_0; \varepsilon) = Y_0 + \varepsilon$$

معادلات (۱.۰.۶) برای $Y(x)$ را از معادلات فوق کم می‌کنیم و قرار می‌دهیم $Z(x) = Y(x; \varepsilon) - Y(x)$ در این صورت

$$\begin{aligned} Z'(x; \varepsilon) &= f(x, Y(x; \varepsilon)) - f(x, Y(x)) \\ &\doteq \frac{\partial f(x, Y(x))}{\partial y} Z(x; \varepsilon) \end{aligned} \quad (۶.۱.۶)$$

و $Z(x_0; \varepsilon) = \varepsilon$. تقریب (۶.۱.۶) وقتی معتبر است که $Y(x; \varepsilon)$ به اندازه کافی به $Y(x)$ نزدیک باشد، که برای مقادیر کوچک ε و بازه کوچک $[x_0 - \alpha, x_0 + \alpha]$ چنین است. به سادگی می‌توانیم معادله دیفرانسیل تقریبی (۶.۱.۶) را حل کنیم و به دست آوریم

$$Z(x; \varepsilon) \doteq \varepsilon \cdot \exp \left[\int_{x_0}^x \frac{\partial f(t, Y(t))}{\partial y} dt \right] \quad (7.1.6)$$

اگر مشتق جزئی در رابطه

$$\frac{\partial f(t, Y(t))}{\partial y} \leq \sigma \quad |x_0 - t| \leq \alpha \quad (8.1.6)$$

صدق کند آنگاه، هرگاه x افزایش یابد احتمال دارد $Z(x, \varepsilon)$ به وسیله ε کراندار بماند. در این حالت گوییم مسأله مقدار اولیه خوش-وضع است.

به عنوان یک مثال از رفتار معکوس، مسأله زیر را در نظر می‌گیریم

$$y' = \lambda y + g(x) \quad y(0) = Y. \quad (9.1.6)$$

با $\lambda > 0$. پس $\partial f / \partial y = \lambda$ و می‌توانیم دقیقاً حساب کنیم

$$Z(x; \varepsilon) = \varepsilon e^{\lambda x}$$

پس تغییر در $Y(x)$ با افزایش x بسیار بزرگ می‌گردد.

مثال معادله

$$y' = 100y - 101e^{-x} \quad y(0) = 1 \quad (10.1.6)$$

دارای جواب $Y(x) = e^{-x}$ است. مسأله اختلال یافته

$$y' = 100y - 101e^{-x} \quad y(0) \equiv 1 + \varepsilon$$

دارای جواب

$$Y(x; \varepsilon) = e^{-x} + \varepsilon e^{100x}$$

است که به سرعت از جواب واقعی دور می‌شود. گوییم (۱۰.۱.۶) یک مسأله بد-وضع است. برای اینکه یک مسأله خوش-وضع باشد، نیاز داریم که انتگرال

$$\int_{x_0}^x \frac{\partial f(t, Y(t))}{\partial t} dt$$

با افزایش x ، کران بالای صفر یا یک عدد مثبت کوچک داشته باشد. پس اختلال $Z(x; \varepsilon)$ با مضرب ثابتی از ε کراندار می‌شود، که مقدار ثابت خیلی بزرگ نخواهد بود.

در حالتی که (۸.۱.۶) برقرار ولی قدرمطلق مشتق نسبی بزرگ باشد، با افزایش x ، اختلال $Z(x; \varepsilon)$ به سرعت به صفر می‌گراید. چنین معادلاتی خوش‌وضع به حساب می‌آیند، ولی آنها نیز می‌توانند برای بیشتر روشهای عددی این فصل مشکل‌آفرین باشند. این معادلات را معادلات دیفرانسیل سرسخت (stiff) خوانند و ما در بخش ۹.۶ به آنها باز می‌گردیم.

دستگاه معادلات دیفرانسیل مطالب این فصل به دستگاه m معادله مرتبه اول زیر تعمیم می‌یابد.

$$\begin{aligned} y_1' &= f_1(x, y_1, \dots, y_m) & y_1(x_0) &= Y_{1,0} \\ &\vdots & & \end{aligned} \quad (11.1.6)$$

$$y_m' = f_m(x, y_1, \dots, y_m) \quad y_m(x_0) = Y_{m,0}$$

این دستگاه اغلب با استفاده از نمادهای برداری به شکل یک معادله مرتبه اول نوشته می‌شود. قرار می‌دهیم

$$\mathbf{y}(x) = \begin{bmatrix} y_1(x) \\ \vdots \\ y_m(x) \end{bmatrix} \quad \mathbf{f}(x, \mathbf{y}) = \begin{bmatrix} f_1(x, \mathbf{y}) \\ \vdots \\ f_m(x, \mathbf{y}) \end{bmatrix} \quad \mathbf{Y}_0 = \begin{bmatrix} Y_{1,0} \\ \vdots \\ Y_{m,0} \end{bmatrix} \quad (12.1.6)$$

در این صورت دستگاه (۱۱.۱.۶) می‌تواند به شکل زیر نوشته شود

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}) \quad \mathbf{y}(x_0) = \mathbf{Y}_0. \quad (13.1.6)$$

با استفاده از نرمهای برداری به جای قدرمطلق، در واقع تمام نتایج قبلی برای این مسأله مقدار اولیه برداری تعمیم می‌یابد.

برای معادلات از مراتب بالاتر، مسأله مقدار اولیه

$$\begin{aligned} y^{(m)} &= f(x, y, y', \dots, y^{(m-1)}) \\ y(x_0) &= Y_0, \dots, y^{(m-1)}(x_0) = Y_0^{(m-1)} \end{aligned} \quad (14.1.6)$$

را با معرفی توابع مجهول جدید زیر، می‌توان به یک دستگاه معادله مرتبه اول برگرداند

$$y_1 = y, y_2 = y', \dots, y_m = y^{(m-1)}$$

این توابع در دستگاه زیر صدق می‌کنند

$$\begin{aligned}
 y_1' &= y_2 & y_1(x_0) &= Y_0 \\
 y_2' &= y_3 & y_2(x_0) &= Y_0' \\
 &\vdots & & \vdots \\
 y_{m-1}' &= y_m & & \\
 y_m' &= f(x, y_1, \dots, y_m) & y_m(x_0) &= Y_0^{(m-1)}
 \end{aligned}
 \tag{۱۵.۱.۶}$$

مثال معادله خطی مرتبه دوم زیر

$$y'' = a_1(x)y' + a_0(x)y + g(x) \quad y(x_0) = \alpha \quad y'(x_0) = \beta$$

به دستگاه زیر تبدیل می شود

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}' = \begin{bmatrix} 0 & 1 \\ a_0(x) & a_1(x) \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} + \begin{bmatrix} 0 \\ g(x) \end{bmatrix} \quad \begin{bmatrix} y_1(x_0) \\ y_2(x_0) \end{bmatrix} = \begin{bmatrix} \alpha \\ \beta \end{bmatrix}$$

در شکل برداری که در آن $A(x)$ ماتریس ضرایب است

$$y' = A(x)y + G(x) \quad y(x_0) = Y_0 = \begin{bmatrix} \alpha \\ \beta \end{bmatrix}$$

که یک دستگاه معادلات دیفرانسیل خطی مرتبه اول است.

روشهای ویژه‌ای برای معادلات مرتبه m ام وجود دارند، ولی این روشها فقط برای $m = 2$ که در کاربردهای مکانیک نیوتونی از قانون دوم مکانیک نیوتن ظاهر می شوند، به اندازه قابل توجهی بسط داده شده اند. بیشتر معادلات مراتب بالا پس از تبدیل به یک دستگاه معادلات مرتبه اول هم ارز، همان گونه که در بالا توضیح داده شد، حل می شوند.

۲.۶ روش اویلر

متداولترین روشهای عددی برای حل (۱.۰.۶)، روشهای تفاضلات متناهی نامیده می شوند. مقادیر تقریبی برای جواب در یک مجموعه از نقاط گرهی

$$x_0 < x_1 < x_2 < \dots < x_n < \dots \tag{۱.۲.۶}$$

به دست آمده اند و مقدار تقریبی در هر x_n با استفاده از مقادیری که در مراحل قبلی محاسبه شده، به دست آمده است. ما بحث را با یک روش ساده ولی از جنبه محاسبه غیرکارا که به لئونهارت اویلر

نسبت داده می‌شود آغاز می‌کنیم. تحلیل آن بسیاری از جنبه‌های آنالیز روشهای تفاضلات متناهی کاراتر را در بر دارد، بی‌آنکه پیچیدگی اضافی آنها را داشته باشد. ابتدا چندین شکل پیدا کردن روش اویلر را بیان و با تحلیل کامل همگرایی و پایداری موضوع را دنبال می‌کنیم. یک فرمول مجانبی خطا به دست آورده این بخش را با تعمیم نتایج قبلی برای یک دستگاه معادلات به پایان می‌رسانیم. مانند قبل، $Y(x)$ نمایانگر جواب درست (۱.۰.۶) است:

$$Y'(x) = f(x, Y(x)) \quad Y(x_0) = Y_0 \quad (2.2.6)$$

جواب تقریبی را با $y(x)$ و مقادیر $y(x_0), y(x_1), y(x_2), \dots, y(x_n)$ را اغلب با $y_0, y_1, \dots, y_n, \dots$ نشان می‌دهیم. اندازه گامهای مساوی $h > 0$ را برای تعریف نقاط گرهی به‌کار می‌بریم.

$$x_j = x_0 + jh \quad j = 0, 1, \dots$$

وقتی جوابهای عددی را به‌ازای مقادیر مختلف h مقایسه می‌کنیم، نماد $y_h(x)$ را نیز برای اشاره به $y(x)$ با طول گام h به‌کار می‌بریم. مسأله (۱.۰.۶) را در یک بازه متناهی ثابت که همیشه با نماد $[x_0, b]$ نشان داده می‌شود، حل خواهیم کرد. نماد $N(h)$ نشان دهنده بزرگترین اندیس N است که برای آن

$$x_N \leq b \quad x_{N+1} > b$$

در بخشهای بعدی، تغییر اندازه گام را در هر نقطه x_n مورد بحث قرار می‌دهیم تا خطا را کنترل کنیم.

پیدا کردن روش اویلر روش اویلر چنین تعریف شده است

$$y_{n+1} = y_n + hf(x_n, y_n) \quad n = 0, 1, 2, \dots \quad (3.2.6)$$

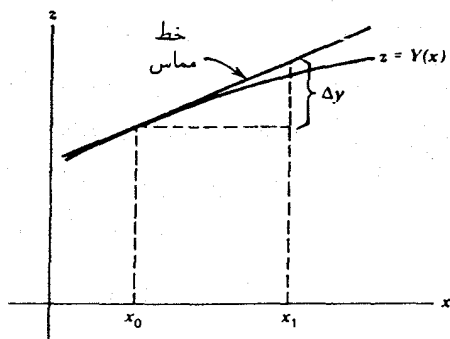
با $y_0 \doteq Y_0$. این روش از چهار دیدگاه نشان داده شده است.

۱. از دیدگاه هندسی. نمودار جواب $Y(x)$ در شکل ۳.۶ را در نظر می‌گیریم. خط مماس بر نمودار $Y(x)$ در x_0 را رسم می‌کنیم و این خط را به‌عنوان یک تقریب خم برای $x_0 \leq x \leq x_1$ به‌کار می‌بریم. پس

$$\frac{\Delta y}{h} = Y'(x_0) = f(x_0, Y_0)$$

$$Y(x_1) - Y(x_0) \doteq \Delta y = hY'(x_0)$$

$$Y(x_1) \doteq Y(x_0) + hf(x_0, Y(x_0))$$



شکل ۳.۶ تغییر هندسی روش اویلر

با تکرار این استدلال در $[x_1, x_2]$, $[x_2, x_3]$, ... فرمول کلی (۳.۲.۶) را به دست می آوریم.

۲. سری تیلر. $Y(x_{n+1})$ را حول x_n بسط می دهیم.

$$Y(x_{n+1}) = Y(x_n) + hY'(x_n) + \frac{h^2}{2}Y''(\xi_n) \quad x_n \leq \xi_n \leq x_{n+1} \quad (4.2.6)$$

با حذف عبارت خطا، روش (۳.۲.۶) ی اویلر را به دست می آوریم. جمله

$$T_n = \frac{h^2}{2}Y''(\xi_n) \quad (5.2.6)$$

خطای برشی یا خطای گسسته سازی در x_{n+1} خوانده می شود. ما خطای برشی را در این کتاب به کار می بریم.

۳. مشتقگیری عددی. از تعریف مشتق

$$\frac{Y(x_{n+1}) - Y(x_n)}{h} \doteq Y'(x_n) = f(x_n, Y(x_n))$$

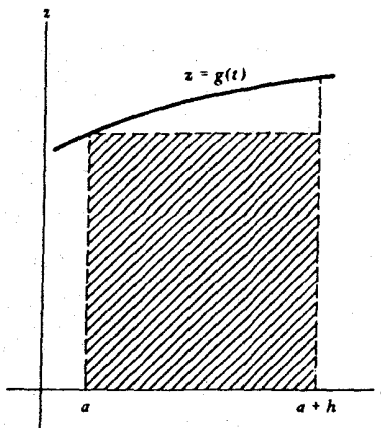
$$Y(x_{n+1}) \doteq Y(x_n) + hf(x_n, Y(x_n))$$

۴. انتگرالگیری عددی. از $Y'(t) = f(t, Y(t))$ در $[x_n, x_{n+1}]$ انتگرال می گیریم:

$$Y(x_{n+1}) = Y(x_n) + \int_{x_n}^{x_{n+1}} f(t, Y(t)) dt \quad (6.2.6)$$

روش ساده انتگرالگیری عددی را در نظر می گیریم

$$\int_a^{a+h} g(t) dt \doteq hg(a) \quad (7.2.6)$$



شکل ۴.۶ نمایش (۷.۲.۶)

که قاعده مستطیلی دست چپ خوانده می‌شود. قسمت هاشورزده شکل ۴.۶ انتگرال عددی را نشان می‌دهد. با استفاده از این انتگرال در (۶.۲.۶)، مانند گذشته خواهیم داشت:

$$Y(x_{n+1}) \doteq Y(x_n) + hf(x_n, Y(x_n))$$

از سه روش تحلیلی (۲) تا (۴)، هم (۲) و هم (۴) ساده‌ترین حالات یک مجموعه روشهای دقیق افزایشی هستند. روش (۲) به روشهای تک گامی، به ویژه فرمولهای رونگه - کوتا می‌انجامد. روش (۴) ما را به روشهای چند گامی، به ویژه روشهای مصحح - پیشگو می‌رساند. شاید شگفت‌آور باشد که روش (۳) اغلب به روشهای موفق دیگری نمی‌رسد. اولین مثال برای (۳)، روش میانگاهی در بخش ۴.۶ است که به مسائل ناپایداری عددی می‌انجامد. در مقابل، در بخش ۹.۶، مشتقگیری عددی به‌کار رفته است که یک رده از روشهایی را برای حل معادلات دیفرانسیل سرسخت به‌دست می‌دهد. قسمت عمده این فصل به روشهای چند مرحله‌یی اختصاص داده شده است، تا اندازه‌ای بدین علت که آنها معمولاً کاراترین رده از روشها هستند، و تا اندازه‌ای به این دلیل که تحلیل آنها پیچیده‌تر از روشهای رونگه - کوتاست. این روش اخیر را در بخش ۱۰.۶ ادامه خواهیم داد.

قبل از تحلیل روش اویلر، چند مثال عددی می‌آوریم. این مثالها همچنین برای توضیح قسمتی از نظریه‌ای که ارائه شده به‌کار گرفته می‌شوند.

مثال ۱. معادله $y' = y$ ، $y(0) = 1$ را در نظر می‌گیریم. جواب درست آن $Y(x) = e^x$ است. نتایج عددی در جدول ۱.۶ برای چند مقدار h داده شده‌اند. جوابهای $y_h(x_n)$ به جای آنکه در

جدول ۱.۶ روش اویلر برای مثال (۱)

	x	$y_h(x)$	$Y(x)$	$Y(x) - y_h(x)$
$h = 0.2$	0.40	۱,۴۴۰۰۰	۱,۴۹۱۸۲	0.05182
	0.80	۲,۰۷۳۶۰	۲,۲۲۵۵۴	0.15194
	۱.۲۰	۲,۹۸۵۹۸	۳,۳۲۰۱۲	0.۳۳۴۱۳
	۱.۶۰	۴,۲۹۹۸۲	۴,۹۵۳۰۳	0.۶۵۳۲۱
	۲.۰۰	۶,۱۹۱۷۴	۷,۳۸۹۰۶	۱,۱۹۷۳۲
$h = 0.1$	0.40	۱,۴۶۴۱۰	۱,۴۹۱۸۲	0.02772
	0.80	۲,۱۴۳۵۹	۲,۲۲۵۵۴	0.08195
	۱.۲۰	۳,۱۳۸۴۳	۳,۳۲۰۱۲	0.۱۸۱۶۹
	۱.۶۰	۴,۵۹۴۹۷	۴,۹۵۳۰۳	0.۳۵۸۰۶
	۲.۰۰	۶,۷۲۷۵۰	۷,۳۸۹۰۶	0.۶۶۱۵۶
$h = 0.05$	0.40	۱,۴۷۷۴۶	۱,۴۹۱۸۲	0.01437
	0.80	۲,۱۸۲۸۷	۲,۲۲۵۵۴	0.04267
	۱.۲۰	۳,۲۲۵۱۰	۳,۳۲۰۱۲	0.09502
	۱.۶۰	۴,۷۶۴۹۴	۴,۹۵۳۰۳	0.۱۸۸۰۹
	۲.۰۰	۷,۰۳۹۹۹	۷,۳۸۹۰۶	0.۳۴۹۰۷

تمام نقاطی که حساب شده، داده شوند، فقط در چند نقطه داده شده‌اند. توجه کنید که در هر نقطه، وقتی h نصف می‌شود، خطا نیز تقریباً نصف می‌شود.

۲. معادله

$$y' = \frac{1}{1+x^2} - 2y^2 \quad y(0) = 0 \quad (۸.۲.۶)$$

را در نظر می‌گیریم. جواب درست مسأله

$$Y(x) = \frac{x}{1+x^2}$$

است. نتایج در جدول ۲.۶ داده شده‌اند. باز هم به رفتار خطا وقتی h کاهش می‌یابد توجه کنید.

تحلیل همگرایی در هر مرحله از روش اویلر، یک خطای برشی اضافی (۵.۲.۶) وارد می‌شود. ما اثر جمعی این خطاها را تحلیل می‌کنیم. خطای $Y(x) - y(x)$ را خطای کلی نامند، و ستونهای آخر جدولهای ۱.۶ و ۲.۶ مثالهایی از خطای کلی‌اند.

جدول ۲.۶ روش اویلر برای مثال (۲)

	x	$y_h(x)$	$Y(x)$	$Y(x) - y_h(x)$
$h = 0.2$	0.00	0.0	0.0	0.0
	0.40	0.37631	0.34483	-0.03148
	0.80	0.52228	0.48780	-0.05448
	1.20	0.52709	0.49180	-0.03529
	1.60	0.46632	0.44944	-0.01689
	2.00	0.40682	0.40000	-0.00682
$h = 0.1$	0.40	0.36085	0.34483	-0.01603
	0.80	0.51371	0.48780	-0.02590
	1.20	0.50961	0.49180	-0.01781
	1.60	0.45872	0.44944	-0.00928
	2.00	0.40419	0.40000	-0.00419
$h = 0.05$	0.40	0.35287	0.34483	-0.00804
	0.80	0.50049	0.48780	-0.01268
	1.20	0.50073	0.49180	-0.00892
	1.60	0.45425	0.44944	-0.00481
	2.00	0.40227	0.40000	-0.00227

برای آنکه یک برداشت شهودی از رفتار خطای کلی روش اویلر پیدا کنیم، مسأله خیلی ساده زیر را در نظر می‌گیریم

$$y' = 2x \quad y(0) = 0 \quad (9.2.6)$$

$Y(x) = x^2$ جواب آن است. روش اویلر برای این مسأله چنین است

$$y_{n+1} = y_n + 2hx_n \quad y_0 = 0$$

سپس با استقراء می‌توان تحقیق کرد که

$$y_n = x_{n-1}x_n \quad n \geq 1$$

برای خطا

$$Y(x_n) - y_n = x_n^2 - x_nx_{n-1} = hx_n \quad (10.2.6)$$

پس خطای کلی در هر نقطه ثابت x متناسب با h است. و این نکته با رفتار مثالها در جدولهای ۱.۶ و ۲.۶ توافق دارد، که در آنها هرگاه h نصف شود خطا در حدود نصف تقلیل می‌یابد. برای تحلیل کامل خطا، مطلب را با لم زیر آغاز می‌کنیم، که در تحلیل روشهای تفاضلات متناهی بسیار سودمند است.

لم ۱ به‌ازای هر مقدار حقیقی x

$$1 + x \leq e^x$$

و به‌ازای هر $x \geq -1$

$$0 \leq (1 + x)^m \leq e^{mx} \quad (11.2.6)$$

برهان با استفاده از قضیه تیلر

$$e^x = 1 + x + \frac{x^2}{2} e^{\xi}$$

با ξ بین 0 و x . چون باقیمانده همیشه مثبت است، قسمت اول ثابت می‌شود. به دنبال آن فرمول (۱۱.۲.۶) به سادگی نتیجه می‌شود. ■

در بقیه این فصل، فرض می‌شود که تابع $f(x, y)$ در شرط قویتر لپیشیتس زیر صدق می‌کند:

$$|f(x, y_1) - f(x, y_2)| \leq k |y_1 - y_2| \quad -\infty < y_1, y_2 < \infty \quad x_0 \leq x \leq b \quad (12.2.6)$$

برای مقداری از $K \geq 0$. اگرچه این شرط از شرط لازم قویتر است، ولی اثبات را آسان می‌کند. و اگر تابع $f(x, y)$ که در شرط ضعیفتر (۱.۱.۶) صدق می‌کند، و یک جواب $Y(x)$ در مسأله مقدار اولیه (۱.۰.۶) داده شده باشد، تابع f را می‌توان به گونه‌ای تغییر داد که در شرط (۱۲.۲.۶) صدق نماید، بی‌آنکه جواب $Y(x)$ تغییر یابد یا ویژگی اساسی مسأله (۱.۰.۶) و جواب عددی آن عوض شود. [برای جزئیات شمایین^۱ و گوردون^۲ (۱۹۷۵) (ص ۲۴) را ببینید].

قضیه ۳.۶ فرض می‌کنیم که $Y(x)$ ، جواب (۱.۰.۶)، دارای مشتق مرتبه دوم کراندار در $[x_0, b]$ باشد. در این صورت جواب $\{y_h(x_n) \mid x_0 \leq x_n \leq b\}$ که از روش (۳.۲.۶)ی اوپلر به دست

آمده است در رابطه زیر صدق می‌کند

$$\text{Max}_{x_0 \leq x_n \leq b} |Y(x_n) - y_h(x_n)| \leq e^{(b-x_0)K} |e_0| + \left[\frac{e^{(b-x_0)K} - 1}{K} \right] \tau(h) \quad (۱۳.۲.۶)$$

که

$$\tau(h) = \frac{h}{4} \|Y''\|_{\infty} \quad (۱۴.۲.۶)$$

$$e_0 = Y_0 - y_h(x_0) \text{ و}$$

به علاوه اگر $h \rightarrow 0$ وقتی

$$|Y_0 - y_h(x_0)| \leq c_1 h \quad (۱۵.۲.۶)$$

برای مقداری از $c_1 \geq 0$ (مثلاً $c_1 = Y_0 - y_0$) به ازای جمیع مقادیر h برقرار باشد، آنگاه یک ثابت $B \geq 0$ وجود دارد که برای آن

$$\text{Max}_{x_0 \leq x_n \leq b} |Y(x_n) - y_h(x_n)| \leq Bh \quad (۱۶.۲.۶)$$

برهان گیریم $e_n = Y(x_n) - y(x_n)$ ، $n \geq 0$ و تعریف $N(h)$ در ابتدای بخش را به یاد می‌آوریم. براساس خطای برشی در (۵.۲.۶) تعریف می‌کنیم

$$\tau_n = \frac{h}{4} Y''(\xi_n) \quad 0 \leq n \leq N(h) - 1$$

با استفاده از (۱۴.۲.۶) به سادگی نتیجه می‌شود

$$\text{Max}_{0 \leq n \leq N-1} |\tau_n| \leq \tau(h)$$

از روابط (۴.۲.۶)، (۲.۲.۶) و (۳.۲.۶) داریم

$$Y_{n+1} = Y_n + hf(x_n, Y_n) + h\tau_n \quad (۱۷.۲.۶)$$

$$y_{n+1} = y_n + hf(x_n, y_n) \quad 0 \leq n \leq N(h) - 1 \quad (۱۸.۲.۶)$$

که از نماد معمولی $Y_n \equiv Y(x_n)$ استفاده کرده‌ایم. از کم کردن (۱۸.۲.۶) از (۱۷.۲.۶)،

$$e_{n+1} = e_n + h[f(x_n, Y_n) - f(x_n, y_n)] + h\tau_n \quad (۱۹.۲.۶)$$

به دست می آید. و کران آن را با استفاده از (۱۲.۲.۶) حساب می کنیم

$$|e_{n+1}| \leq |e_n| + hK |Y_n - y_n| + h |\tau_n|$$

$$|e_{n+1}| \leq (\lambda + hK) |e_n| + h\tau(h) \quad 0 \leq n \leq N(h) - 1 \quad (20.2.6)$$

رابطه (۲۰.۲.۶) را می توانیم به صورت بازگشتی به کار برده به دست آوریم،

$$|e_n| \leq (\lambda + hK)^n |e_0| + \{1 + (\lambda + hK) + \dots + (\lambda + hK)^{n-1}\} h\tau(h)$$

با استفاده از فرمول مجموع سری هندسی متناهی

$$1 + r + r^2 + \dots + r^{n-1} = \frac{r^n - 1}{r - 1} \quad r \neq 1 \quad (21.2.6)$$

خواهیم داشت

$$|e_n| \leq (\lambda + hK)^n |e_0| + \left[\frac{(\lambda + hK)^n - 1}{K} \right] \tau(h) \quad (22.2.6)$$

با استفاده از لم ۱،

$$(\lambda + hK)^n \leq e^{nhK} = e^{(x_n - x_0)K} \leq e^{(b - x_0)K}$$

و این رابطه با (۲۲.۲.۶) قسمت اصلی قضیه، (۱۳.۲.۶)، قضیه را به دست خواهد داد.

باقی قضیه، (۱۶.۲.۶)، یک نتیجه نمایان (۱۳.۲.۶) است که ثابت B با

$$B = e_1 e^{(b - x_0)K} + \left[\frac{e^{(b - x_0)K} - 1}{K} \right] \cdot \frac{\|Y''\|_\infty}{2}$$

داده شده است. و این اثبات را کامل می نماید.

نتیجه (۱۶.۲.۶) ایجاب می کند که هرگاه طول گام h نصف شود، خطا دست کم با ضریب یک دوّم کاهش یابد. این مطلب با مثالهای جدولهای ۱.۶ و ۲.۶ تأیید شده است. بعداً در این بخش نشان داده شده است که (۱۶.۲.۶) دقیقاً نرخ درست همگرایی را به دست می دهد (مسأله ۷ را نیز ببینید).

کران (۱۶.۲.۶) سرعت درست همگرایی در روش اویلر را به دست می دهد، ولی ضریب ثابت B در اغلب معادلات خیلی بیش از اندازه بزرگ است. مثلاً، با مثال قبلی (۸.۲.۶)، فرمول

(۱۳.۲.۶) پیشگویی می‌کند که خطا با b بزرگ می‌شود. ولی به روشنی در جدول ۲.۶ مشاهده می‌شود، وقتی x افزایش می‌یابد، خطا کاهش پیدا می‌کند. ما اصلاح زیر را برای (۱۳.۲.۶) ارائه می‌کنیم که در خیلی حالات مانند (۸.۲.۶) به‌کار خواهد رفت.

فرع همان فرضهای قضیه ۳.۶ را در نظر می‌گیریم؛ به‌علاوه فرض می‌کنیم

$$\frac{\partial f(x, y)}{\partial y} \leq 0 \quad (23.2.6)$$

برای $x_0 \leq x \leq b$ و $-\infty < y < \infty$. در این صورت به‌ازای جمیع مقادیر h به اندازه کافی کوچک

$$|Y(x_n) - y_h(x_n)| \leq |e_0| + \frac{h}{4}(x_n - x_0) \text{Max}_{x_0 \leq x_n \leq b} |Y''(x)| \quad (24.2.6)$$

به‌ازای $x_0 \leq x_n \leq b$

برهان قضیه مقدار میانگین را برای معادله خطای (۱۹.۲.۶) به‌کار می‌بریم:

$$e_{n+1} = \left[1 + h \frac{\partial f(x_n, \zeta_n)}{\partial y} \right] e_n + \frac{h^2}{2} Y''(\xi_n) \quad (25.2.6)$$

با ζ_n بین $y_h(x_n)$ و $Y(x_n)$. از همگرایی $y_h(x)$ به $Y(x)$ در $[x_0, b]$ ، می‌فهمیم که مشتقهای جزئی $\frac{\partial f(x_n, \zeta_n)}{\partial y}$ به $\frac{\partial f(x, Y(x))}{\partial y}$ میل می‌کنند، و بنابراین باید از نظر قدرمطلق در $[x_0, b]$ کراندار باشند. $h_0 > 0$ را طوری می‌گیریم که برای تمام مقادیر $h \leq h_0$

$$1 + h \frac{\partial f(x_n, \zeta_n)}{\partial y} \geq -1 \quad x_0 \leq x_n \leq b \quad (26.2.6)$$

از (۲۳.۲.۶) نتیجه می‌گیریم که سمت چپ نیز به‌ازای جمیع مقادیر h از بالا با مقدار ۱ کراندار شده است. این نتایج را در (۲۵.۲.۶) به‌کار می‌بریم تا به‌دست آید

$$|e_{n+1}| \leq |e_n| + \frac{h^2}{2} |Y''(\xi_n)| \quad (27.2.6)$$

با استقرا می‌توانیم نشان دهیم

$$|e_n| \leq |e_0| + \frac{h^2}{2} [|Y''(\xi_0)| + \dots + |Y''(\xi_{n-1})|]$$

که به‌سادگی به (۲۴.۲.۶) منجر می‌شود.

نتیجه (۲۴.۲.۶) خیلی بهتر از کران قبلی (۱۳.۲.۶) است، به جای $\exp(K(b-x_0))$ مقدار $b-x_0$ (کران $x_n - x_0$) گذاشته شده است، که با افزایش b با سرعت کمتری بزرگ می‌شود. این قضیه مستقیماً برای مثال قبلی (۸.۲.۶) کاربرد ندارد، ولی با یک بررسی دقیق برهان در این حالت، نشان داده می‌شود که برهان هنوز معتبر است.

تحلیل پایداری تحلیل پایداری مسأله مقدار اولیه را که در قضیه ۲.۶ داده شده بود به یاد آورید. برای در نظر گرفتن طرح مشابهی برای روش اویلر، روش عددی زیر را با $z_0 = y_0 + \varepsilon$ در نظر می‌گیریم

$$z_{n+1} = z_n + h[f(x_n, z_n) + \delta(x_n)] \quad 0 \leq n \leq N(h) - 1 \quad (28.2.6)$$

مقایسه این دو فرمول مانند مقایسه (۵.۱.۶) با (۱۰.۰.۶) است، که پایداری مسأله مقدار اولیه را نشان می‌دهد. ما دو جواب عددی $\{z_n\}$ و $\{y_n\}$ را وقتی $h \rightarrow 0$ مقایسه می‌کنیم. بگیریم $e_n = z_n - y_n$ ، $n \geq 0$ ، پس $\varepsilon = e_0$ ، و با کم کردن (۳.۲.۶) از (۲۸.۲.۶)،

$$e_{n+1} = e_n + h[f(x_n, z_n) - f(x_n, y_n)] + h\delta(x_n)$$

این فرمول درست شکل (۱۹.۲.۶) را دارد. اگر شیوه‌ای مشابه آنچه را که به دنبال (۱۹.۲.۶) آمده است به کار بریم، خواهیم داشت

$$\text{Max}_{0 \leq n \leq N(h)} |z_n - y_n| \leq e^{(b-x_0)K} |\varepsilon| + \left[\frac{e^{(b-x_0)K} - 1}{K} \right] \|\delta\|_\infty$$

نتیجتاً ثابتهایی مانند k_1 و k_2 مستقل از h ، وجود دارند به طوری که

$$\text{Max}_{0 \leq n \leq N(h)} |z_n - y_n| \leq k_1 |\varepsilon| + k_2 \|\delta\|_\infty \quad (29.2.6)$$

این نتیجه مشابه نتیجه (۴.۱.۶) برای مسأله اصلی (۱۰.۰.۶) است. این رابطه می‌گوید که روش اویلر یک روش عددی پایدار برای حل مسأله مقدار اولیه (۱۰.۰.۶) است. تأکید می‌کنیم که تمام روشهای عددی برای مسائل مقدار اولیه چنین شکل پایداری دارند و از پایداری مسأله اصلی (۱۰.۰.۶) پیروی می‌کنند. به علاوه، به پایداری به شکلهای دیگر نیز نیاز داریم، که بعداً معرفی می‌شوند. در آینده، می‌گیریم $\delta(x) \equiv 0$ و فقط اثر اختلال در مقدار اولیه Y_0 را در نظر می‌گیریم. این امر موجب سادگی تحلیل می‌شود و نتایج به همان اندازه سودمند خواهند بود.

تحلیل خطای گردکردن در هر مرحله از روش اویلر خطایی وارد می‌کنیم، خطایی که از خطاهای گردکردن اعمال انجام شده حاصل شده است. این عدد، که با نماد ρ_n نشان داده می‌شود، خطای گردکردن موضعی خوانده می‌شود. اگر نتیجه مقادیر عددی را \tilde{y}_n بنامیم، خواهیم داشت

$$y_{n+1} = \tilde{y}_n + hf(x_n, \tilde{y}_n) + \rho_n \quad n = 0, 1, \dots, N(h) - 1 \quad (30.2.6)$$

مقادیر \tilde{y}_n اعدادی با ارقام اعشاری متناهی هستند که عملاً در رایانه به دست آمده‌اند، و y_n مقداری است که اگر حساب دقیق به کار می‌رفت به دست می‌آمد. گیریم $\rho(h)$ کرانی برای خطاهای گردکردن باشد،

$$\rho(h) = \max_{0 \leq n \leq N(h)-1} |\rho_n| \quad (31.2.6)$$

در یک وضعیت عملی، با استفاده از رایانه‌یی که طول کلمه ثابت داشته باشد، کران $\rho(h)$ وقتی $h \rightarrow 0$ ، کاهش نمی‌یابد. بلکه تقریباً ثابت می‌ماند و $\|Y\|_\infty \rho(h)$ با واحد خطای گردکردن در رایانه، متناسب خواهد بود، یعنی u کوچکترین عدد رایانه‌یی است که برای آن $1 + u > 1$ [فرمول (۱۲.۲.۱) فصل ۱ را ببینید].

برای آنکه اثر خطاهای گردکردن در (۳۰.۲.۶) را ببینیم، این رابطه را از جواب درست در (۴.۲.۶) کم می‌کنیم تا حاصل شود:

$$\tilde{e}_{n+1} = \tilde{e}_n + h[f(x_n, Y_n) - f(x_n, \tilde{y}_n)] + h\tau_n - \rho_n$$

که در آن $\tilde{e}_n = Y(x_n) - \tilde{y}_n$. مانند اثبات قضیه ۳.۶ عمل می‌کنیم ولی $\tau_n - \rho_n/h$ را با τ_n که در برهان قبلی بود یکی می‌گیریم. در این صورت به دست می‌آوریم

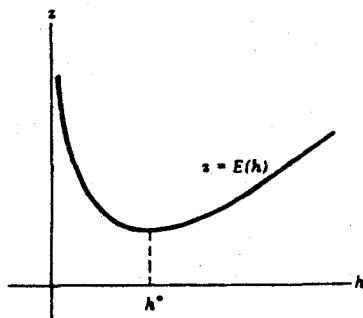
$$|\tilde{e}_n| \leq e^{(b-x_0)K} |Y_0 - \tilde{y}_0| + \left[\frac{e^{(b-x_0)K} - 1}{K} \right] \left[\tau(h) + \frac{\rho(h)}{h} \right] \quad (32.2.6)$$

برای بررسی بیشتر این کران، گیریم $\rho(h)/\|Y\|_\infty = u$ همان‌گونه که قبلاً بحث شد. پس

$$|\tilde{e}_n| \leq c \left[\frac{h}{2} \|Y''\|_\infty + \frac{u \|Y\|_\infty}{h} \right] \equiv E(h) \quad (33.2.6)$$

رفتار کیفی $E(h)$ در نمودار شکل ۵.۶ نشان داده شده است. در h^* مینیمم می‌شود، و هر کاهش دیگر، باعث افزایش خطای کران $E(h)$ خواهد شد.

این $E(h)$ ، بدترین حالت ممکن را برای اثر خطاهای گردکردن می‌دهد. در عمل، هم اندازه و هم علامت خطاهای گردکردن تغییر می‌کنند. برآیند حذف باعث می‌شود که \tilde{e}_n با سرعتی کمتر از



شکل ۵.۶ خم خطا برای (۳۳.۲.۶)

سرعتی که بر اثر جمله $1/h$ که در $E(h)$ و (۳۲.۲.۶) حاصل می‌شود، افزایش یابد. ولی هنوز هم یک مقدار بهینه h^* وجود دارد و برای h کمتر از آن، خطا باز افزایش می‌یابد. برای یک تحلیل کاملتر هنریچی (۱۹۶۲، صص ۳۵-۵۹) را ببینید.

معمولاً خطای گردکردن مسأله مهمی نیست. ولی، اگر دقت مطلوب به بهترین دقتی که به علت محدودیت طول کلمه در رایانه می‌توان دست یافت، نزدیک باشد، به اثرات گردکردن باید توجه بیشتری بشود. در یک رایانه بزرگ IBM (جدول ۱.۱ را ببینید) با حساب دقت مضاعف، اگر $h \geq 10^{-4}$ ، آنگاه ماکسیم u/h برابر 22×10^{-13} خواهد شد که u واحد گردکردن است. بنابراین خطای گردکردن معمولاً مسأله قابل توجهی به وجود نخواهد آورد مگر آنکه تحملهای خطای بسیار کوچکی خواسته شده باشد. ولی در دقت ساده، با همین محدودیت روی h ، ماکسیم u/h برابر $10^{-4} \times 10^5$ می‌شود، و با تحمل خطایی به این اندازه، (که یک مقدار غیرمنطقی نیست)، خطای گردکردن عامل مؤثرتری خواهد بود.

مثال مسأله

$$y' = -y + 2 \cos(x) \quad y(0) = 1$$

را که جواب درست آن $Y(x) = \sin(x) + \cos(x)$ است حل می‌کنیم. این مسأله را با روش اویلر و استفاده از سه حساب مختلف حل می‌کنیم: (۱) حساب با ممیز شناور با چهار رقم اعشار و قطع کردن؛ (۲) حساب ممیز شناور با چهار رقم اعشار و گردکردن؛ و (۳) حساب دقیق، یا حساب با دقت بسیار بالا. در دو حالت اول، واحد خطاهای گردکردن به ترتیب $u = 10^{-4}$ و $u = 10^{-5}$ هستند. کران (۳۲.۲.۶) را برای حالت‌های (۱) و (۲) به کار می‌بریم، حال آنکه

جدول ۳.۶ مثال اثرات خطای گردکردن در روش اویلر

		اعشاری قطع شده	اعشاری گرد شده	حساب دقیق
h	x			
۰.۰۴	۱	$-۱.۰۰E-۲$	$-۱.۷۰E-۲$	$-۱.۷۰E-۲$
	۲	$۱.۱۷E-۲$	$۱.۸۳E-۲$	$-۱.۸۳E-۲$
	۳	$-۱.۲۰E-۳$	$-۲.۸۰E-۳$	$-۲.۷۸E-۳$
	۴	$۱.۰۰E-۲$	$۱.۶۰E-۲$	$۱.۵۳E-۲$
	۵	$۱.۱۳E-۲$	$۱.۹۶E-۲$	$۱.۹۴E-۲$
۰.۰۲	۱	$۷.۰۰E-۳$	$-۹.۰۰E-۳$	$-۸.۴۶E-۳$
	۲	$۴.۰۰E-۳$	$-۹.۱۰E-۳$	$-۹.۱۳E-۳$
	۳	$۲.۳۰E-۳$	$-۱.۴۰E-۳$	$-۱.۴۰E-۳$
	۴	$-۶.۰۰E-۳$	$۸.۰۰E-۳$	$۷.۶۲E-۳$
	۵	$-۶.۰۰E-۳$	$۸.۵۰E-۳$	$۹.۶۳E-۳$
۰.۰۱	۱	$۲.۸۰E-۲$	$-۳.۰۰E-۳$	$-۴.۲۲E-۳$
	۲	$۲.۲۸E-۲$	$-۴.۳۰E-۳$	$-۴.۵۶E-۳$
	۳	$۷.۴۰E-۳$	$-۴.۰۰E-۴$	$-۷.۰۳E-۴$
	۴	$-۲.۳۰E-۲$	$۳.۰۰E-۳$	$۳.۸۰E-۳$
	۵	$-۲.۴۱E-۲$	$۴.۶۰E-۳$	$۴.۸۱E-۳$

حالت ۳ در کران نظری (۲۴.۲.۶) صدق می‌نماید. خطاها برای هر سه شکل روش اویلر در جدول ۳.۶ داده شده‌اند. خطا برای جوابهای حاصله با استفاده از حساب اعشاری، بر پایهٔ جواب درست $Y(x)$ که تا چهار رقم گرد شده است محاسبه شده‌اند.

برای حالت حساب اعشاری قطع شده، با $h = ۰.۰۲$ ، تأثیر روی خطاها آغاز شده است. با $h = ۰.۰۱$ ، خطای قطع کردن، اثر قابل ملاحظه‌ای بر خطای کلی دارد. برعکس، خطاها، با استفاده از حساب گرد شده به تدریج کاهش می‌یابند، اگرچه $h = ۰.۰۱$ قدری متأثر شده است. در این مسأله مانند بسیاری دیگر، استفاده از حساب گرد شده بسیار برتر از حساب قطع شده است.

تحلیل خطای مجانبی یک برآورد خطای مجانبی در روش اویلر، با چشم‌پوشی از اثرات خطای گردکردن به دست آمده است. قبل از هر چیز، برای آسان کردن عملیات جبری در این تحلیل یک نمادگذاری لازم است. اگر $B(x, h)$ تابعی باشد که برای $x \leq x \leq b$ و به ازای جمیع مقادیر به اندازهٔ کافی کوچک h تعریف شده باشد، آنگاه نماد

$$B(x, h) = O(h^p)$$

به ازای مقداری از $p > 0$ ، بدین معناست که به ازای جمیع مقادیر به اندازه کافی کوچک h ، مقدار ثابتی مانند c وجود دارد به طوری که

$$|B(x, h)| \leq ch^p \quad x_0 \leq x \leq b \quad (34.2.6)$$

اگر B فقط به h بستگی داشته باشد، همین نوع کران به دست می آید.

قضیه ۴.۶ فرض می کنیم که $Y(x)$ جواب مسأله مقدار اولیه (۱.۰.۶) و سه بار مشتق پذیر باشد. فرض می کنیم

$$f_y(x, y) \equiv \frac{\partial f(x, y)}{\partial y} \quad f_{yy}(x, y) \equiv \frac{\partial^2 f(x, y)}{\partial y^2}$$

به ازای $x_0 \leq x \leq b$ و $-\infty < y < \infty$ ، پیوسته و کراندار باشند. گیریم مقدار اولیه $y_h(x_0)$ در رابطه زیر صدق کند

$$Y_0 - y_h(x_0) = \delta \cdot h + O(h^2) \quad (35.2.6)$$

معمولاً این خطا صفر است و بنابراین $\delta = 0$.

در این صورت خطا در روش (۳.۲.۶) اویلر در رابطه زیر صدق می کند

$$Y(x_n) - y_h(x_n) = D(x_n)h + O(h^2) \quad (36.2.6)$$

که در آن $D(x)$ جواب مسأله مقدار اولیه خطی زیر است:

$$D'(x) = f_y(x, Y(x))D(x) + \frac{1}{2}Y''(x) \quad D(x_0) = \delta. \quad (37.2.6)$$

برهان قضیه تیلر را به کار می بریم،

$$Y(x_{n+1}) = Y(x_n) + hY'(x_n) + \frac{h^2}{2}Y''(x_n) + \frac{h^3}{6}Y^{(3)}(\xi_n)$$

به ازای یک $x_n \leq \xi_n \leq x_{n+1}$. رابطه (۳.۲.۶) را از رابطه فوق کم کرده از (۲.۲.۶) استفاده می کنیم تا به دست آوریم

$$e_{n+1} = e_n + h[f(x_n, Y_n) - f(x_n, y_n)] + \frac{h^2}{2}Y''(x_n) + \frac{h^3}{6}Y^{(3)}(\xi_n) \quad (38.2.6)$$

$f(x_n, y_n)$ را به عنوان تابعی از y_n در نظر گرفته قضیه تیلر را برای آن به کار می‌بریم،

$$f(x_n, y_n) = f(x_n, Y_n) + (y_n - Y_n)f_y(x_n, Y_n) + \frac{1}{2}(y_n - Y_n)^2 f_{yy}(x_n, \zeta_n)$$

به‌ازای مقداری از ζ_n بین y_n و Y_n . اگر از این رابطه در (۳۸.۲.۶) استفاده کنیم

$$e_{n+1} = [1 + hf_y(x_n, Y_n)]e_n + \frac{h^2}{2}Y''(x_n) + B_n$$

$$B_n = \frac{h^2}{6}Y^{(3)}(\xi_n) - \frac{1}{2}hf_{yy}(x_n, \zeta_n)e_n^2 \quad (39.2.6)$$

با استفاده از (۱۶.۲.۶)

$$B_n = O(h^2) \quad (40.2.6)$$

چون B_n نسبت به سایر جملات (۳۹.۲.۶) کوچک است، قسمت عمده خطا را با حذف

B_n به دست می‌آوریم. گیریم g_n معرف قسمت عمده خطا، به‌طور ضمنی با

$$g_{n+1} = [1 + hf_y(x_n, Y_n)]g_n + \frac{h^2}{2}Y''(x_n) \quad (41.2.6)$$

تعریف شده است، که

$$g_n = \delta_n h \quad (42.2.6)$$

قسمت عمده خطا در مقدار اولیه (۳۵.۲.۶) است. چون انتظار داریم که $g_n = e_n$ و چون

$e_n = O(h)$ ، دنباله $\{\delta_n\}$ را که به‌طور ضمنی این گونه معرفی می‌کنیم

$$g_n = h\delta_n \quad (43.2.6)$$

در (۴۱.۲.۶) می‌گذاریم و h را حذف کرده دوباره مرتب می‌نماییم، به دست می‌آید

$$\delta_{n+1} = \delta_n + h[f_y(x_n, Y_n)\delta_n + \frac{1}{2}Y''(x_n)] \quad x_n \leq x_{n+1} \leq b \quad (44.2.6)$$

مقدار اولیه δ_0 در (۳۵.۲.۶) مستقل از h تعریف شده است.

معادله (۴۴.۲.۶) روش اویلر برای حل مسأله مقدار اولیه (۳۷.۲.۶) است؛ بنابراین به موجب

قضیه ۳.۶،

$$D(x_n) - \delta_n = O(h) \quad x_n \leq x_{n+1} \leq b \quad (45.2.6)$$

این رابطه را با (۴۳.۲.۶) ترکیب می‌کنیم،

$$g_n = D(x_n)h + O(h^2) \quad (۴۶.۲.۶)$$

برای کامل کردن اثبات، باید نشان دهیم که g_n در واقع قسمت اصلی خطای e_n است. نماد زیر را معرفی می‌کنیم

$$k_n = e_n - g_n \quad (۴۷.۲.۶)$$

در این صورت $k_n = e_n - g_n = O(h^2)$ ، که از (۳۵.۲.۶) و (۴۲.۲.۶) نتیجه می‌شود. رابطه (۴۱.۲.۶) را از (۳۹.۲.۶) کم می‌کنیم و (۴۰.۲.۶) را به‌کار می‌بریم

$$k_{n+1} = [1 + hf_y(x_n, Y_n)]k_n + B_n$$

$$|k_{n+1}| \leq (1 + hK) |k_n| + O(h^2)$$

این رابطه به شکل (۲۰.۲.۶) در برهان قضیه ۳.۶ است که به‌جای جمله $h\tau(h)$ جمله $O(h^2)$ گذاشته شده است. با استفاده از همان محاسبه

$$|k_n| = O(h^2) \quad (۴۸.۲.۶)$$

از ترکیب (۴۶.۲.۶) - (۴۸.۲.۶) داریم

$$e_n = g_n + k_n = [hD(x_n) + O(h^2)] + O(h^2)$$

که رابطه (۳۶.۲.۶) را ثابت می‌کند.

تابع $D(x)$ به‌ندرت صریحاً به‌دست می‌آید، ولی شکل خطا در (۳۶.۲.۶) اطلاعات کیفی مفیدی در اختیار می‌گذارد. این اغلب به‌عنوان پایهٔ شیوه‌های برون‌یابی به‌کار می‌رود، که بعضی از آنها در بخشهای بعد مورد بحث قرار گرفته‌اند.

مثال مسأله

$$y' = -y \quad y(0) = 1$$

را با جواب $Y(x) = e^{-x}$ در نظر می‌گیریم. معادله برای $D(x)$ چنین است

$$D'(x) = -D(x) + \frac{1}{x}e^{-x} \quad D(0) = 0$$

و جواب این معادله

جدول ۴.۶ مثال (۴۹.۲.۶) از قضیه ۴.۶

x_n	$Y_n - y_n$	$hD(x_n)$
۰٫۴	۰٫۰۰۶۸۹	۰٫۰۰۶۷۰
۰٫۸	۰٫۰۰۹۲۰	۰٫۰۰۸۹۹
۱٫۲	۰٫۰۰۹۲۱	۰٫۰۰۹۰۴
۱٫۶	۰٫۰۰۸۱۹	۰٫۰۰۸۰۸
۲٫۰	۰٫۰۰۶۸۲	۰٫۰۰۶۷۷

$$D(x) = \frac{1}{4} x e^{-x}$$

است. این رابطه فرمول مجانبی زیر را برای خطا در روش اویلر به دست می دهد

$$Y(x_n) - y_h(x_n) \doteq \frac{h}{4} x_n e^{-x_n} \quad (49.2.6)$$

جدول ۴.۶ شامل خطاهای واقعی و خطاهایی است که با (۴۹.۲.۶) به ازای $h = ۰٫۵$ پیش بینی شده است. توجه نمایید که با افزایش x ، خطا کاهش می یابد، درست مثل خطای جواب $Y(x)$. ولی خطای نسبی به طور خطی نسبت به x ، افزایش می یابد،

$$\frac{Y(x_n) - y_h(x_n)}{Y(x_n)} \doteq \frac{h}{4} x_n$$

همچنین، برآورد (۴۹.۲.۶) به مراتب بهتر از کرانی است که توسط (۱۳.۲.۶) در قضیه ۳.۶ داده شده است. آن کران چنین است

$$|Y(x_n) - y_h(x_n)| \leq \frac{h}{4} (e^{x_n} - 1)$$

و به طور نمایی با x افزایش می یابد.

دستگاه معادلات برای ساده کردن بیان، فقط دستگاه مرتبه دوم زیر را در نظر می گیریم:

$$\begin{aligned} y'_1 &= f_1(x, y_1, y_2) & y_1(x_0) &= Y_{1,0} \\ y'_2 &= f_2(x, y_1, y_2) & y_2(x_0) &= Y_{2,0} \end{aligned} \quad (50.2.6)$$

تعمیم آن به دستگاههای مراتب بالاتر ساده است. روش اویلر برای حل (۵۰.۲.۶) چنین است

$$\begin{aligned} y_{1,n+1} &= y_{1,n} + h f_1(x_n, y_{1,n}, y_{2,n}) \\ y_{2,n+1} &= y_{2,n} + h f_2(x_n, y_{1,n}, y_{2,n}) \end{aligned} \quad (51.2.6)$$

که یک تعمیم روشن از (۳.۲.۶) است.

تمام نتایج پیشین این بخش، به صورت (۵۱.۲.۶) تعمیم پیدا می‌کنند، و شکل تعمیم بسیار روشن خواهد بود اگر نمادگذاری برداری (۱۳.۱.۶) را که در بخش قبل بیان شد برای (۵۰.۲.۶) و (۵۱.۲.۶) به کار بریم. می‌نویسیم

$$y' = f(x, y) \quad y(x_0) = Y.$$

به جای قدر مطلق، نرم (۱۶.۱.۱) از فصل ۱ را به کار می‌بریم

$$\|y\| = \text{Max}_i |y_i|$$

برای تعمیم شرط (۱۲.۲.۶) لیبیشیتس، قضیهٔ تیلر (قضیهٔ ۱.۵) برای توابع چندمتغیره را به کار می‌بریم تا به دست آوریم

$$\|f(x, z) - f(x, y)\| \leq K \|z - y\| \quad (۵۲.۲.۶)$$

$$K = \text{Max}_i \sum_j \text{Max}_{\substack{x, \leq x \leq b \\ -\infty < w_1, w_2 < \infty}} \left| \frac{\partial f_i(x, w_1, w_2)}{\partial w_j} \right| \quad (۵۳.۲.۶)$$

ماتریس ژاکوبی زیر نقش $\partial f(x, y) / \partial y$ را دارد

$$f_y(x, y) = \begin{bmatrix} \frac{\partial f_1}{\partial y_1} & \frac{\partial f_1}{\partial y_2} \\ \frac{\partial f_2}{\partial y_1} & \frac{\partial f_2}{\partial y_2} \end{bmatrix} \quad (۵۴.۲.۶)$$

به عنوان مثال، فرمول خطای مجانبی (۳۶.۲.۶) چنین می‌شود

$$Y(x_n) - y_h(x_n) = hD(x_n) + R_n \quad \|R_n\| = O(h^2) \quad (۵۵.۲.۶)$$

که $D(x)$ جواب دستگاه خطی زیر است

$$D'(x) = f_y(x, Y(x))D(x) + \frac{1}{h} Y''(x) \quad D(x_0) = \delta. \quad (۵۶.۲.۶)$$

که از ماتریس قبلی $f_y(x, y)$ استفاده شده است.

مثال مسألهٔ آونگ زیر را حل کنید

$$\theta''(t) = -\sin(\theta(t)) \quad \theta(0) = \frac{\pi}{4} \quad \theta'(0) = 0$$

جدول ۵.۶ روش اویلر برای مثال (۵۷.۲.۶)

h	x_n	$y_{1,n}$	$Y_1(x_n)$	خطاها	$y_{2,n}$	$Y_2(x_n)$	خطاها
۰.۲	۰.۲	۱.۵۷۰۸	۱.۵۵۰۸	-۰.۰۲۰۰	-۰.۲۰۰۰۰	-۰.۱۹۹۹۹۹	۰.۰۰۰۰۰۰۱
	۰.۶	۱.۴۵۰۸	۱.۳۹۱۰	-۰.۰۵۹۸	-۰.۵۹۹۸۴	-۰.۵۹۸۰۶	۰.۰۰۰۱۷۸
	۱.۰	۱.۱۷۱۱	۱.۰۷۴۹	-۰.۰۹۶۲	-۰.۹۹۲۶۷	-۰.۹۷۵۵۰	۰.۰۰۱۷۱۷
۰.۱	۰.۲	۱.۵۶۰۸	۱.۵۵۰۸	-۰.۰۱۰۰	-۰.۲۰۰۰۰	-۰.۱۹۹۹۹۹	۰.۰۰۰۰۰۰۱
	۰.۶	۱.۴۲۰۸	۱.۳۹۱۰	-۰.۰۲۹۸	-۰.۵۹۹۲۷	-۰.۵۹۸۰۶	۰.۰۰۰۱۲۱
	۱.۰	۱.۱۲۲۳	۱.۰۷۴۹	-۰.۰۴۷۴	-۰.۹۸۵۶۸	-۰.۹۷۵۵۰	۰.۰۰۱۰۱۸

با انتخاب $\theta = y_1$ و $\theta' = y_2$ و گذاشتن x به جای t ، معادله فوق به دستگاه زیر تبدیل می‌شود

$$\begin{aligned} y_1' &= y_2 & y_1(0) &= \frac{\pi}{4} \\ y_2' &= -\sin(y_1) & y_2(0) &= 0 \end{aligned} \quad (57.2.6)$$

نتایج عددی در جدول ۵.۶ داده شده‌اند. توجه کنید وقتی h نصف می‌شود خطا تقریباً نصف می‌شود.

۳.۶ روشهای چندگامی

این بخش شامل مقدمه‌ای بر روشهای چندگامی است. بعضی روشهای خاص با جزئیات بیشتر در بخشهای بعدی بررسی شده‌اند. یک نظریه کاملتر در بخش ۸.۶ داده شده است.

مانند قبل، گیریم $h > 0$ و گره‌ها را با $x_n = x_0 + nh$ ، $n \geq 0$ ، تعریف می‌کنیم. شکل کلی روشهای چندمرحله‌یی که در نظر گرفته خواهند شد به قرار زیر است

$$y_{n+1} = \sum_{j=0}^p a_j y_{n-j} + h \sum_{j=-1}^p b_j f(x_{n-j}, y_{n-j}) \quad n \geq p \quad (1.3.6)$$

ضرایب $a_0, \dots, a_p, b_{-1}, \dots, b_p$ ثابت‌اند و $p \geq 0$. اگر $a_p \neq 0$ یا $b_p \neq 0$ برقرار باشد روش را $p+1$ مرحله‌یی خوانند، زیرا $p+1$ مقدار قبلی، در محاسبه y_{n+1} به‌کار گرفته شده‌اند. مقادیر y_1, \dots, y_p باید به وسایل دیگری به‌دست آیند؛ این مطلب در بخشهای دیگر بحث شده است. روش اویلر یک مثال از روش یک گامی با شرطهای $p=0$ و

$$a_0 = 1 \quad b_0 = 1 \quad b_{-1} = 0$$

اگر $b_{-1} = 0$ ، آنگاه y_{n+1} فقط در سمت چپ معادله (۱.۳.۶) ظاهر می‌شود. چنین فرمولهایی را روشهای صریح نامند. اگر $b_{-1} \neq 0$ ، آنگاه y_{n+1} در دو طرف معادله (۱.۳.۶) وجود دارد و فرمول را روش ضمنی خوانند. وجود جواب y_{n+1} را برای h به اندازه کافی کوچک، می‌توان با استفاده از نظریه نقطه ثابت بخش ۵.۲ نشان داد. روشهای ضمنی معمولاً با روشهای بارستی حل می‌شوند که برای روش دوزنقه‌یی در بخش ۵.۶ بحث شده است.

مثال ۱. روش میانگاهی با رابطه

$$y_{n+1} = y_{n-1} + 2hf(x_n, y_n) \quad n \geq 1 \quad (2.3.6)$$

تعریف شده است که یک روش دوگامی صریح است. این روش در بخش ۴.۶ بررسی شده است.

۲. روش دوزنقه‌یی با رابطه زیر تعریف شده است،

$$y_{n+1} = y_n + \frac{h}{4}[f(x_n, y_n) + f(x_{n+1}, y_{n+1})] \quad n \geq 0 \quad (3.3.6)$$

این روش یک روش یک‌گامی ضمنی است و در بخشهای ۵.۶ و ۶.۶ مورد بحث قرار گرفته است. برای هر تابع مشتقپذیر $Y(x)$ ، خطای برشی برای انتگرالگیری $Y'(x)$ را با رابطه زیر تعریف می‌کنیم

$$T_n(Y) = Y(x_{n+1}) - \left[\sum_{j=0}^p a_j Y(x_{n-j}) + h \sum_{j=-1}^p b_j Y'(x_{n-j}) \right] \quad n \geq p \quad (4.3.6)$$

تابع $\tau_n(Y)$ را با رابطه زیر تعریف می‌کنیم

$$\tau_n(Y) = \frac{1}{h} T_n(Y) \quad (5.3.6)$$

برای اثبات همگرایی جواب تقریبی $\{y_n \mid x_0 \leq x_n \leq b\}$ از رابطه (۱.۳.۶) به جواب $Y(x)$ از مسأله مقدار اولیه (۱.۰.۶) لازم است که داشته باشیم

$$\tau(h) \equiv \max_{x_p \leq x_n \leq b} |\tau_n(Y)| \rightarrow 0 \quad \text{وقتی } h \rightarrow 0 \quad (6.3.6)$$

این رابطه اغلب شرط سازگاری برای روش (۱.۳.۶) خوانده می‌شود. سرعت همگرایی جواب $\{y_n\}$ به جواب درست $Y(x)$ وابسته به سرعت همگرایی در (۶.۳.۶) است، و بنابراین نیاز داریم

شرایطی را بدانیم که تحت آنها برای انتخاب مطلوبی از $m \geq 1$ داشته باشیم

$$\tau(h) = O(h^m) \quad (۷.۳.۶)$$

ما اکنون استلزامهای (۶.۳.۶) و (۷.۳.۶) برای ضرایب در (۱.۳.۶) را تحقیق می‌کنیم. نتیجه همگرایی برای (۱.۳.۶) بعداً به صورت قضیه ۶.۶ داده شده است.

قضیه ۵.۶ گیریم عدد صحیح $m \geq 1$ داده شده باشد. برای آنکه (۶.۳.۶) به‌ازای همه توابع پیوسته مشتقپذیر $Y(x)$ برقرار باشد، یعنی برای اینکه روش (۱.۳.۶) سازگار باشد، لازم و کافی است که

$$\sum_{j=0}^p a_j = 1 \quad - \sum_{j=0}^p j a_j + \sum_{j=-1}^p b_j = 1 \quad (۸.۳.۶)$$

و برای آنکه (۷.۳.۶) به‌ازای همه توابع $m + 1$ بار پیوسته مشتقپذیر $Y(x)$ معتبر باشد، لازم و کافی است که (۸.۳.۶) برقرار بوده و

$$\sum_{j=0}^p (-j)^i a_j + i \sum_{j=-1}^p (-j)^{i-1} b_j = 1 \quad i = 2, \dots, m \quad (۹.۳.۶)$$

برهان توجه نمایید که به‌ازای جمیع مقادیر ثابت α و β و همه توابع مشتقپذیر Y و W ، داریم

$$T_n(\alpha Y + \beta W) = \alpha T_n(Y) + \beta T_n(W) \quad (۱۰.۳.۶)$$

برای بررسی نتایج (۶.۳.۶) و (۷.۳.۶)، با استفاده از قضیه ۴.۱ تیلر، $Y(x)$ را حول x_n بسط می‌دهیم تا به‌دست آوریم

$$Y(x) = \sum_{i=0}^m \frac{1}{i!} (x - x_n)^i Y^{(i)}(x_n) + R_{m+1}(x) \quad (۱۱.۳.۶)$$

با فرض آنکه $Y(x)$ $m + 1$ بار پیوسته مشتقپذیر است. با گذاردن در (۴.۳.۶) و استفاده از (۱۰.۳.۶) داریم

$$T_n(Y) = \sum_{i=0}^m \frac{1}{i!} Y^{(i)}(x_n) T_n((x - x_n)^i) + T_n(R_{m+1})$$

لازم است که $T_n((x - x_n)^i)$ را به‌ازای $i \geq 0$ محاسبه کنیم.

برای $i = 0$

$$T_n(1) = c_0 \equiv 1 - \sum_{j=0}^p a_j \quad (12.3.6)$$

برای $i \geq 1$

$$\begin{aligned} T_n((x - x_n)^i) &= (x_{n+1} - x_n)^i - \left[\sum_{j=0}^p a_j (x_{n-j} - x_n)^i + h \sum_{j=-1}^p b_j i (x_{n-j} - x_n)^{i-1} \right] \\ &= c_i h^i \end{aligned} \quad (13.3.6)$$

$$c_i = 1 - \left[\sum_{j=0}^p (-j)^i a_j + i \sum_{j=-1}^p (-j)^{i-1} b_j \right] \quad i \geq 1$$

و از آنجا نتیجه می‌گیریم

$$T_n(Y) = \sum_{i=0}^m \frac{c_i}{i!} h^i Y^{(i)}(x_n) + T_n(R_{m+1}) \quad (14.3.6)$$

اگر باقیمانده $R_{m+1}(x)$ را به شکل زیر بنویسیم

$$R_{m+1}(x) = \frac{1}{(m+1)!} (x - x_n)^{m+1} Y^{(m+1)}(x_n) + \dots$$

آنگاه

$$T_n(R_{m+1}) = \frac{c_{m+1}}{(m+1)!} h^{m+1} Y^{(m+1)}(x_n) + O(h^{m+2}) \quad (15.3.6)$$

با فرض اینکه Y ، $m+2$ بار مشتقپذیر است.

برای پیدا کردن شرط سازگاری (۶.۳.۶)، لازم است که $\tau(h) = O(h)$ و این ایجاب می‌نماید $T_n(Y) = O(h^2)$ با استفاده از (۱۴.۳.۶) با $m = 1$ باید داشته باشیم $c_0, c_1 = 0$ که مجموعه معادلات (۸.۳.۶) را خواهد داد. در بعضی کتابها، این معادلات را شرایط سازگاری نامیده‌اند. برای به دست آوردن (۷.۳.۶) برای مقداری از $m \geq 1$ باید داشته باشیم $T_n(Y) = O(h^{m+1})$ که درستی آن از (۱۴.۳.۶) و (۱۳.۳.۶) به دست می‌آید اگر و فقط اگر $c_i = 0$ ، $i = 0, 1, \dots, m$ ، این شرایط (۹.۳.۶) را ثابت و برهان را کامل می‌کند. ■

بزرگترین مقدار m که برای آن (۷.۳.۶) برقرار باشد مرتبه یا مرتبه همگرایی روش (۱.۳.۶) خوانده می‌شود. در بخش ۷.۶ پیدا کردن روشهایی از هر مرتبه دلخواه را، بررسی خواهیم نمود.

اکنون یک قضیه همگرایی برای جواب (۱.۳.۶) خواهیم آورد. گرچه، این قضیه همه روشهای چندگامی را که همگرا هستند شامل نمی‌شود، بیشتر روشهای سودمند رایج را در برمی‌گیرد. علاوه بر آن، برهان آن بسیار ساده‌تر از برهان قضیه کلی ۸.۶ از بخش ۸.۶ است.

قضیه ۶.۶ حل مسأله مقدار اولیه زیر را با استفاده از روش چندگامی (۱.۳.۶) در نظر می‌گیریم

$$y' = f(x, y) \quad y(x_0) = Y. \quad x_0 \leq x \leq b$$

گیریم خطاهای اولیه در رابطه زیر صدق می‌کنند

$$\eta(h) \equiv \text{Max}_{0 \leq i \leq p} |Y(x_i) - y_h(x_i)| \rightarrow 0 \quad h \rightarrow 0 \quad \text{وقتی} \quad (16.3.6)$$

فرض کنید که این روش سازگار است، یعنی در (۶.۳.۶) صدق می‌کند. و بالاخره، فرض کنید ضرایب a_j همگی نامنفی هستند،

$$a_j \geq 0 \quad j = 0, 1, \dots, p \quad (17.3.6)$$

در این صورت روش (۱.۳.۶) همگراست و برای مقادیر مناسبی از c_1 و c_2

$$\text{Max}_{x_0 \leq x_n \leq b} |Y(x_n) - y_h(x_n)| \leq c_1 \eta(h) + c_2 \tau(h) \quad (18.3.6)$$

اگر روش (۱.۳.۶) از مرتبه m باشد و اگر خطاهای اولیه در $\eta(h) = O(h^m)$ صدق کنند، آنگاه سرعت همگرایی این روش برابر $O(h^m)$ است.

برهان رابطه (۴.۳.۶) را دوباره می‌نویسیم و از $Y'(x) = f(x, Y(x))$ استفاده می‌کنیم تا به دست آوریم

$$Y(x_{n+1}) = \sum_{j=0}^p a_j Y(x_{n-j}) + h \sum_{j=-1}^p b_j f(x_{n-j}, Y(x_{n-j})) + h \tau_n(Y)$$

با کم کردن (۱.۳.۶) از آن و استفاده از نماد $e_i = Y(x_i) - y_i$

$$e_{n+1} = \sum_{j=0}^p a_j e_{n-j} + h \sum_{j=-1}^p b_j [f(x_{n-j}, Y_{n-j}) - f(x_{n-j}, y_{n-j})] + h \tau_n(Y)$$

شرط لیبشیتس و فرض (۱۷.۳.۶) را به کار می‌بریم تا به دست آوریم

$$|e_{n+1}| \leq \sum_{j=0}^p a_j |e_{n-j}| + hK \sum_{j=-1}^p |b_j| |e_{n-j}| + h\tau(h)$$

تابع خطای کرانی زیر را وارد می‌کنیم

$$f_n = \max_{0 \leq i \leq n} |e_i| \quad n = 0, 1, \dots, N(h)$$

با به کار بردن این تابع،

$$|e_{n+1}| \leq \sum_{j=0}^p a_j f_n + hK \sum_{j=-1}^p |b_j| f_{n+1} + h\tau(h)$$

و با استفاده از (۸.۳.۶)

$$|e_{n+1}| \leq f_n + hc f_{n+1} + h\tau(h) \quad c = K \sum_{j=-1}^p |b_j|$$

روشن است که طرف راست یک کران برای f_n است، و بنابراین

$$f_{n+1} \leq f_n + hc f_{n+1} + h\tau(h)$$

برای $hc \leq 1/2$ که وقتی $h \rightarrow 0$ باید درست باشد،

$$\begin{aligned} f_{n+1} &\leq \frac{f_n}{1-hc} + \frac{h}{1-hc} \tau(h) \\ &\leq (1+2hc)f_n + 2h\tau(h) \end{aligned}$$

با توجه به اینکه $f_p = \eta(h)$ ، مانند (۲۰.۲.۶) که در برهان قضیه ۳.۶ آمده است، عمل می‌کنیم.

پس

$$f_n \leq e^{2c(b-x_n)} \eta(h) + \left[\frac{e^{2c(b-x_n)} - 1}{c} \right] \tau(h) \quad x_0 \leq x_n \leq b \quad (19.3.6)$$

برهان را کامل می‌کند.

نتیجه‌های این قضیه را می‌توان با مفروضات ضعیفتری اثبات نمود؛ به‌ویژه، به جای ((۱۷.۳.۶)) می‌توان شرط بسیار ضعیفتری گذاشت. این نتایج در بخش ۸.۶ داده شده‌اند. برای پیدا کردن یک نرخ همگرایی $O(h^m)$ برای روش (۱.۳.۶)، لازم است که هر مرحله دارای خطای

$$T_n(Y) = O(h^{m+1})$$

باشد. ولی مقادیر اولیه y_0, \dots, y_p فقط لازم است با دقت $O(h^m)$ محاسبه شوند، زیرا $\eta(h) = O(h^m)$ برای (۱۸.۳.۶) کافی است. مثالهایی برای نشان دادن استفاده از روش مرتبه پایبتر برای تولید مقادیر اولیه y_0, \dots, y_n در بخشهای آینده داده شده است.

نتیجه (۱۹.۳.۶) را می‌توان برای حالات خاص بهبود بخشید، ولی سرعت همگرایی به همان حالت خواهد ماند. مثالهایی از این قضیه در بخشهای بعد آمده است. همانند تحلیل در روش اویلر، یک تحلیل کامل پایداری می‌توان ارائه داد، از جمله نتیجه‌ای را به شکل (۲۹.۲.۶). این برهان تغییر یافته ساده برهان قضیه ۶.۶ است. همچنین یک تحلیل خطای مجانبی می‌توان ارائه داد؛ مثالهایی در دو بخش بعد داده شده‌اند.

۴.۶ روش میانگاهی

ما روش میانگاهی را تعریف و تحلیل می‌نماییم، و آن را برای نشان دادن بعضی مفاهیم که با روش اویلر ممکن نبود، به‌کار می‌بریم. همانند روش اویلر می‌توانیم روش میانگاهی را به چندین راه پیدا کنیم، و در اینجا از مشتقگیری عددی استفاده می‌کنیم. با توجه به (۱۱.۷.۵) فصل ۵، داریم

$$g'(a) = \frac{g(a+h) - g(a-h)}{2h} - \frac{h^2}{6} g^{(3)}(\xi)$$

برای مقداری از $a-h \leq \xi \leq a+h$. با به‌کار بردن این رابطه در

$$Y'(x_n) = f(x_n, Y(x_n))$$

داریم

$$\frac{Y(x_{n+1}) - Y(x_{n-1}))}{2h} - \frac{h^2}{6} Y^{(3)}(\xi_n) = f(x_n, Y(x_n))$$

با $x_{n-1} \leq \xi_n \leq x_{n+1}$. از حل این رابطه برحسب $Y(x_{n+1})$

$$Y(x_{n+1}) = Y(x_{n-1}) + 2hf(x_n, Y(x_n)) + \frac{1}{3}h^3 Y^{(3)}(\xi_n) \quad (۱.۴.۶)$$

روش میانگاهی با حذف آخرین جمله به‌دست می‌آید:

$$y_{n+1} = y_{n-1} + 2hf(x_n, y_n) \quad n \geq 1 \quad (۲.۴.۶)$$

این یک روش صریح دوگامی است و مرتبه همگرایی آن دو است. مقدار y_1 باید با روش دیگری محاسبه شود.

روش میانگاهی را می‌توانستیم از به‌کار بردن قاعده انتگرالگیری عددی میانگاهی (۱۷.۲.۵) در انتگرال زیر که صورت انتگرالی معادله دیفرانسیل (۱۰.۶) است به‌دست آوریم:

$$Y(x_{n+1}) = Y(x_{n-1}) + \int_{x_{n-1}}^{x_{n+1}} f(t, Y(t)) dt \quad (3.4.6)$$

ما از ذکر جزئیات صرف‌نظر می‌کنیم. این راه در بخش ۷.۶، برای پیدا کردن روشهای چندگامی دیگر به‌کار رفته است.

برای تحلیل همگرایی (۲.۴.۶)، از قضیه ۶.۶ استفاده می‌کنیم. برای خطای برشی، به راحتی از (۱.۴.۶) به‌دست می‌آوریم که

$$\tau_n(Y) = \frac{1}{3} h^2 Y^{(3)}(\xi_n) \quad x_{n-1} \leq \xi_n \leq x_n \quad (4.4.6)$$

یک شکل برهان بهتر قضیه ۶.۶ برای روش میانگاهی، نتیجه می‌دهد که

$$\max_{x_0 \leq x_n \leq b} |Y(x_n) - y_h(x_n)| \leq e^{\gamma K(b-x_0)} \eta(h) + \left[\frac{e^{\gamma K(b-x_0)} - 1}{\gamma K} \right] \left[\frac{1}{3} h^2 \|Y^{(3)}\|_{\infty} \right] \quad (5.4.6)$$

$$\eta(h) = \max\{|Y_0 - y_0|, |Y(x_1) - y_h(x_1)|\}$$

با فرض آنکه برای جميع مقادیر h ، $y_0 = Y_0$ ، نیاز داریم که $Y(x_1) - y_h(x_1) = O(h^2)$ تا مرتبه کلی همگرایی $O(h^2)$ را در (۵.۴.۶) داشته باشیم. با توجه به (۴.۲.۶)، یکی از روشهای تک‌گامی اویلر این ویژگی مورد نظر را داراست:

$$y_1 = y_0 + hf(x_0, y_0) \quad y_0 = Y_0 \quad (6.4.6)$$

$$Y(x_1) - y_1 = \frac{h^2}{2} Y''(\zeta) \quad x_0 \leq \zeta \leq x_1 \quad (7.4.6)$$

با این مقدار اولیه $y_1 = y_h(x_1)$ ، خطای نتیجه (۵.۴.۶) ایجاب می‌کند که

$$\max_{x_0 \leq x_n \leq b} |Y(x_n) - y_h(x_n)| = O(h^2) \quad (8.4.6)$$

می‌توانیم یک تحلیل پایداری کامل، مشابه با آنچه برای روش اویلر داده شد، برای روش میانگاهی بیاوریم. اگر برای سادگی فرض کنیم که $\eta(h) = O(h^2)$ ، آنگاه فرمول خطای مجانبی

زیر را خواهیم داشت

$$Y(x_n) - y_h(x_n) = D(x_n)h^r + O(h^r) \quad x_0 \leq x_n \leq b$$

$$D' = f_y(x, Y(x))D + \frac{1}{\epsilon} Y^{(r)}(x) \quad D(x_0) = 0 \quad (9.4.6)$$

تفاوت اثبات این نتایج با آنچه در روش اویلر دیدیم کم است و ما از ذکر اثبات صرف نظر می‌کنیم. پایداری ضعیف همان گونه که قبلاً اشاره شد، روش میانگاهی دارای همان نوع پایداری است که برای روش اویلر در (۲۸.۲.۶) و (۲۹.۲.۶) نشان داده شد. ولی به هر حال، این برای منظورهای عملی کافی نیست. نشان خواهیم داد که روش میانگاهی برای مفهوم دیگری از پایداری که باید تعریف شود، رضایتبخش نیست.

راه حل عددی مسأله زیر را در نظر می‌گیریم

$$y' = \lambda y \quad y(0) = 1 \quad (10.4.6)$$

که دارای جواب $Y(x) = e^{\lambda x}$ است. این مسأله به‌عنوان یک مسأله نمونه برای حالت کلیتر مسأله (۱۰.۰.۶) به‌کار خواهد رفت، موضوعی که در بخش ۸.۶ توضیح خواهیم داد. در اینجا کافی است توجه کنیم که اگر یک روش عددی برای مسأله ساده‌ای چون (۱۰.۴.۶)، بد عمل کند، آنگاه غیرمحتمل است که چنین روشی برای معادلات دیفرانسیل پیچیده‌تر، خوب عمل کند. روش میانگاهی برای معادله (۱۰.۴.۶) چنین است

$$y_{n+1} = y_{n-1} + 2h\lambda y_n \quad n \geq 1 \quad (11.4.6)$$

ما جواب درست این معادله را پیدا و با جواب $Y(x) = e^{\lambda x}$ مقایسه می‌کنیم. معادله (۱۱.۴.۶) یک مثال از معادله تفاضلی خطی مرتبه دو است. یک نظریه عمومی برای معادلات تفاضلی خطی مرتبه p ام وجود دارد که با نظریه معادلات دیفرانسیل خطی مرتبه p مشابه است. بیشتر روشهای حل معادلات دیفرانسیل، مشابهی برای حل معادلات تفاضلی دارند و این خود یک راهنما برای حل معادله (۱۱.۴.۶) خواهد بود. ما کار را با جستجوی جوابهای مستقل خطی برای معادله تفاضلی آغاز می‌کنیم. سپس این جوابها را برای یافتن جواب کلی ترکیب خواهیم کرد. برای نظریه کلی معادلات تفاضلی خطی هنریچی (۱۹۶۲، صص ۲۱۰-۲۱۵) را ببینید.

در قیاس با جوابهای نمایی برای معادلات دیفرانسیل خطی، جوابی برای معادله (۱۱.۴.۶) به شکل زیر به‌ارای مقدار مجهول r فرض می‌کنیم

$$y_n = r^n \quad n \geq 0 \quad (12.4.6)$$

و آن را در معادله (۱۱.۴.۶) قرار می‌دهیم تا شرایط لازم برای r را پیدا کنیم

$$r^{n+1} = r^{n-1} + 2h\lambda r^n$$

با حذف r^{n-1} داریم،

$$r^2 = 1 + 2h\lambda r \quad (۱۳.۴.۶)$$

این استدلال برگشت‌پذیر است. اگر r در معادله درجه دوم (۱۳.۴.۶) صدق کند، آنگاه (۱۲.۴.۶) در (۱۱.۴.۶) صدق خواهد کرد.

معادله (۱۳.۴.۶) را معادله مشخصه روش میانگامی نامند. ریشه‌های آن چنین‌اند

$$r_0 = h\lambda + \sqrt{1 + h^2\lambda^2} \quad r_1 = h\lambda - \sqrt{1 + h^2\lambda^2} \quad (۱۴.۴.۶)$$

جواب عمومی (۱۱.۴.۶) عبارتست از

$$y_n = \beta_0 r_0^n + \beta_1 r_1^n \quad n \geq 0 \quad (۱۵.۴.۶)$$

ضرایب β_0 و β_1 به نحوی انتخاب می‌شوند که مقادیر y_0 و y_1 که در ابتدا داده شده بودند با مقادیری که با (۱۵.۴.۶) محاسبه شده‌اند مطابقت نمایند:

$$\beta_0 + \beta_1 = y_0$$

$$\beta_0 r_0 + \beta_1 r_1 = y_1$$

جواب عمومی این دستگاه چنین است

$$\beta_0 = \frac{y_1 - r_1 y_0}{r_0 - r_1} \quad \beta_1 = \frac{y_0 r_0 - y_1}{r_0 - r_1}$$

برای آنکه درک مستقیمی از این فرمولها به دست آوریم مقادیر اولیه دقیق را در نظر می‌گیریم

$$y_0 = 1 \quad y_1 = e^{\lambda h}$$

سپس با استفاده از قضیه تیلر،

$$\begin{aligned} \beta_0 &= \frac{e^{\lambda h} - r_1}{2\sqrt{1 + h^2\lambda^2}} = 1 + O(h^2\lambda^2) \\ \beta_1 &= \frac{r_0 - e^{\lambda h}}{2\sqrt{1 + h^2\lambda^2}} = O(h^2\lambda^2) \end{aligned} \quad (۱۶.۴.۶)$$

از این مقادیر، وقتی $h \rightarrow 0$ داریم $\beta \rightarrow 1$ و $\beta_1 \rightarrow 0$. بنابراین در فرمول (۱۵.۴.۶)، جمله $\beta \cdot r^n$ باید متناظر جواب درست $e^{\lambda x_n}$ باشد، زیرا وقتی $h \rightarrow 0$ ، جمله $\beta_1 r^n \rightarrow 0$ در واقع

$$r^n = e^{\lambda x_n} [1 + O(h)] \quad (17.4.6)$$

که اثبات آن به عهده خواننده واگذار می‌شود.

برای آنکه دشواری حل عددی $y' = \lambda y$ را در استفاده از (۱۵.۴.۶) ببینیم، اندازه‌های نسبی r_1 و r_0 را به دقت بررسی می‌کنیم. ما فقط حالتی را که λ حقیقی است در نظر می‌گیریم. برای $0 < \lambda < \infty$ و همه مقادیر h ,

$$r_0 > |r_1| > 0$$

بنابراین مقدار r^n با سرعتی کمتر از r^n افزایش می‌یابد، و جمله درست در جواب عمومی (۱۵.۴.۶) یعنی $\beta \cdot r^n$ غالب می‌شود.

ولی، برای $-\infty < \lambda < 0$ ، خواهیم داشت

$$0 < r_0 < 1 \quad r_1 < -1 \quad h > 0$$

و در نتیجه، $\beta_1 r^n$ با افزایش n و مقدار ثابت h ، سرانجام بر $\beta \cdot r^n$ غالب خواهد شد، بدون توجه به اینکه h در ابتدا چقدر کوچک انتخاب شده‌باشد. وقتی $m \rightarrow \infty$ ، جمله $\beta \cdot r^n \rightarrow 0$ ؛ حال آنکه $\beta_1 r^n$ از نظر قدر مطلق افزایش می‌یابد و علامت آن با افزایش n یک در میان عوض می‌شود. جمله $\beta_1 r^n$ ، جواب مزاحم روش عددی (۱۱.۴.۶) نامیده می‌شود، زیرا متناظر با هیچ یک از جوابهای معادله دیفرانسیل اصلی $y' = \lambda y$ نیست. معادله دیفرانسیل اصلی یک خانواده جواب تک‌پارامتری دارد که به مقدار اولیه Y بستگی دارد، ولی جواب تقریبی (۱۱.۴.۶) یک خانواده جواب دو پارامتری (۱۵.۴.۶) دارد که وابسته به مقادیر y_0 و y_1 است. جواب جدید $\beta_1 r^n$ یک آفرینش روش عددی است؛ برای مسأله (۱۰.۴.۶) با $\lambda < 0$ ، این موجب می‌شود که جواب عددی وقتی $x_n \rightarrow \infty$ ، از جواب درست مسأله، دور شود. به علت این رفتار گوییم روش میانگاهی فقط ضعیف - پایدار است.

پس از آنکه قدری از نظریه را که لازم داریم ارائه کردیم، مجدداً در بخش ۸.۶ به بحث فوق باز می‌گردیم. عملی بودن قابلیت کاربرد مسأله نمونه (۱۰.۴.۶) را با در نظر گرفتن علامت $\partial f(x, Y(x)) / \partial Y$ تعمیم می‌دهیم. اگر منفی باشد، آنگاه ناپایداری ضعیف روش میانگاهی در حل مسأله مقدار اولیه متناظر، معمولاً ظاهر می‌شود. این موضوع در مثال دوم در زیر نشان داده شده است.

جدول ۶.۶ مثال (۱) ناپایداری روش میانگاهی

x_n	y_n	$Y(x_n)$	خطاها
۰٫۲۵	۰٫۷۵۰۰	۰٫۷۷۸۸	۰٫۰۲۸۸
۰٫۵۰	۰٫۶۲۵۰	۰٫۶۰۶۵	-۰٫۰۱۸۵
۰٫۷۵	۰٫۴۳۷۵	۰٫۴۷۲۴	۰٫۰۳۴۹
۱٫۰۰	۰٫۴۰۶۳	۰٫۳۶۷۹	-۰٫۰۳۸۴
۱٫۲۵	۰٫۲۳۴۴	۰٫۲۸۶۵	۰٫۰۵۲۱
۱٫۵۰	۰٫۲۸۹۱	۰٫۲۲۳۱	-۰٫۰۶۵۹
۱٫۷۵	۰٫۰۸۹۸	۰٫۱۷۳۸	۰٫۰۸۳۹
۲٫۰۰	۰٫۲۴۴۱	۰٫۱۳۵۳	-۰٫۱۰۸۸
۲٫۲۵	-۰٫۰۳۲۲	۰٫۱۰۵۴	۰٫۱۳۷۶

مثال ۱. مسأله نمونه (۱۰.۴.۶) را با $\lambda = -1$ در نظر می‌گیریم. نتایج عددی در جدول ۶.۶ به‌ازای $h = 0.25$ داده شده است. مقدار y_1 با استفاده از روش اویلر، مانند (۶.۴.۶)، به‌دست آمده است. از مقادیر جدول مشاهده می‌شود که قدم‌مطلق جواب مزاحم افزایش می‌یابد. برای $x_n = 2.25$ ، جواب عددی y_n منفی می‌شود، و علامت آن با هر مرحله متوالیاً تغییر می‌کند.

۲. مسأله زیر را در نظر می‌گیریم

$$y' = x - y^2 \quad y(0) = 0$$

جواب $Y(x)$ به‌ازای $x \geq 0$ اکیداً صعودی است؛ به‌ازای مقدار بزرگ x ، $Y(x) \doteq \sqrt{x}$ گرچه این جواب صعودی است

$$\frac{\partial f(x, y)}{\partial y} = -2y < 0 \quad y > 0 \text{ به‌ازای}$$

بنابراین انتظار داریم که روش میانگاهی نوعی ناپایداری نشان دهد. این مطلب توسط نتایج داده‌شده در جدول ۷.۶ با طول گام $h = 0.25$ ، تأیید شده است. جواب عددی در اطراف $x_n = 2.25$ روبه‌کاهش می‌گذارد و در $x_n = 3.25$ ، جواب y_n منفی می‌شود.

۵.۶ روش دوزنقه‌یی

روش دوزنقه‌یی را برای معرفی روشهای ضمنی و مفاهیم وابسته به آنها، به‌کار می‌بریم. به‌علاوه، روش دوزنقه‌یی، به‌دلیل ویژگی پایداری خاصی که دارد، خود سودمند است. برای معرفی قاعده دوزنقه‌یی از انتگرالگیری عددی استفاده می‌کنیم.

جدول ۷.۶ مثال (۲) ناپایداری میانگاهی

x_n	y_n	$Y(x_n)$	خطاها
۰٫۲۵	۰٫۰	۰٫۰۳۱۲	۰٫۰۳۱۲
۰٫۵۰	۰٫۱۲۵۰	۰٫۱۲۳۵	-۰٫۰۱۵
۰٫۷۵	۰٫۲۴۲۲	۰٫۲۷۰۰	۰٫۰۲۷۸
۱٫۰۰	۰٫۴۷۰۷	۰٫۴۵۵۵	-۰٫۰۱۵۱
۱٫۲۵	۰٫۶۳۱۴	۰٫۶۵۸۵	۰٫۰۲۷۱
۱٫۵۰	۰٫۸۹۶۳	۰٫۸۵۷۴	-۰٫۰۳۸۹
۱٫۷۵	۰٫۹۷۹۷	۱٫۰۳۷۶	۰٫۰۵۷۹
۲٫۰۰	۱٫۲۹۱۴	۱٫۱۹۳۶	-۰٫۰۹۷۸
۲٫۲۵	۱٫۱۴۵۹	۱٫۳۲۶۴	۰٫۱۸۱
۲٫۵۰	۱٫۷۵۹۹	۱٫۴۴۰۵	-۰٫۳۱۹
۲٫۷۵	۰٫۸۴۷۲	۱٫۵۴۰۴	۰٫۶۹۳
۳٫۰۰	۲٫۷۷۶۰	۱٫۶۳۰۲	-۱٫۱۵
۳٫۲۵	-۱٫۵۰۵۸	۱٫۷۱۲۵	۳٫۲۲

از معادلهٔ دیفرانسیل $Y'(t) = f(t, Y(t))$ روی $[x_n, x_{n+1}]$ انتگرال می‌گیریم تا به دست آوریم

$$Y(x_{n+1}) = Y(x_n) + \int_{x_n}^{x_{n+1}} f(t, Y(t)) dt$$

قاعدهٔ دوزنقه‌یی ساده (۱۲.۱.۵) و (۱۴.۱.۵) را به‌کار می‌بریم تا به دست آوریم

$$Y(x_{n+1}) = Y(x_n) + \frac{h}{4} [f(x_n, Y(x_n)) + f(x_{n+1}, Y(x_{n+1}))] - \frac{h^3}{12} Y^{(3)}(\xi_n) \quad (۱.۵.۶)$$

به‌ازای یک مقدار $x_n \leq \xi_n \leq x_{n+1}$. با حذف جملهٔ باقی‌مانده، روش دوزنقه‌یی

$$y_{n+1} = y_n + \frac{h}{4} [f(x_n, y_n) + f(x_{n+1}, y_{n+1})] \quad n \geq 0 \quad (۲.۵.۶)$$

را به دست می‌آوریم. این روش تک‌گامی با مرتبهٔ همگرایی $O(h^2)$ است. و نیز یک مثال ساده از روش ضمنی است، زیرا y_{n+1} در هر دو طرف (۲.۵.۶) ظاهر شده است. یک مثال عددی در آخر این بخش داده شده است.

راه حلّ بارستی فرمول (۲.۵.۶) یک معادله غیرخطی با ریشه y_{n+1} است، و هر یک از روشهای کلی فصل ۲ را می توان برای حل آن به کار برد. بارست خطی ساده (بخش ۵.۲ را ببینید) مناسبترین و معمولاً کافی است. گیریم $y_{n+1}^{(0)}$ یک حدس اولیه خوب برای جواب y_{n+1} باشد، و تعریف می کنیم

$$y_{n+1}^{(j+1)} = y_n + \frac{h}{4} [f(x_n, y_n) + f(x_{n+1}, y_{n+1}^{(j)})] \quad j = 0, 1, \dots \quad (3.5.6)$$

حدس اولیه معمولاً با استفاده از یک روش روشن به دست می آید.

برای تحلیل بارست و تعیین شرایطی که تحت آنها، بارست همگرا می شود، رابطه (۳.۵.۶) را از (۲.۵.۶) کم می کنیم تا به دست آوریم

$$y_{n+1} - y_{n+1}^{(j+1)} = \frac{h}{4} [f(x_{n+1}, y_{n+1}) - f(x_{n+1}, y_{n+1}^{(j)})] \quad (4.5.6)$$

شرط لپشیتس (۱۲.۲.۶) را به کار می بریم تا کران آن را با مقدار زیر به دست آوریم

$$|y_{n+1} - y_{n+1}^{(j+1)}| \leq \frac{hK}{4} |y_{n+1} - y_{n+1}^{(j)}| \quad j \geq 0 \quad (5.5.6)$$

اگر

$$\frac{hK}{4} < 1 \quad (6.5.6)$$

آنگاه بازستهای $y_{n+1}^{(j)}$ وقتی $j \rightarrow \infty$ به y_{n+1} همگرا می شوند. یک برآورد دقیقتر از نرخ همگرایی با کار بستن قضیه مقدار میانگین برای (۴.۵.۶) به دست می آید:

$$y_{n+1} - y_{n+1}^{(j+1)} = \frac{h}{4} f_y(x_{n+1}, y_{n+1}) [y_{n+1} - y_{n+1}^{(j)}] \quad (7.5.6)$$

اغلب در عمل طول گام h و حدس اولیه $y_{n+1}^{(0)}$ به گونه ای انتخاب می شوند که مسلم شود فقط به محاسبه یک بارست نیاز است، و در آن صورت $y_{n+1} = y_{n+1}^{(1)}$ انتخاب می شود.

محاسبه y_{n+1} از y_n شامل یک خطای برشی $O(h^3)$ است [۱.۵.۶] را ببینید]. برای نگهداری این مرتبه دقت، بارست نهایی $y_{n+1}^{(i)}$ که برای نشان دادن y_{n+1} انتخاب شده است، باید در رابطه $|y_{n+1} - y_{n+1}^{(i)}| = O(h^3)$ صدق کند. و اگر خواهیم خطای بارست اهمیت کمتری داشته باشد (همانطور که در بخش بعد عمل خواهیم کرد)، باید $y_{n+1}^{(i)}$ چنان انتخاب شود که در رابطه زیر صدق نماید:

$$|y_{n+1} - y_{n+1}^{(i)}| = O(h^4) \quad (8.5.6)$$

برای تحلیل خطا در انتخاب یک حدس اولیه y_{n+1} ، باید مفهوم حل موضعی را معرفی نماییم. این مفهوم برای روشن ساختن اینکه دقیقاً چه جوابی توسط اغلب برنامه های رایانه ای برای حل معادلات

دیفرانسیل به دست می‌آید نیز مهم خواهد بود. گیریم $u_n(x)$ معرفّ جواب معادلهٔ دیفرانسیلی باشد که از نقطهٔ (x_n, y_n) عبور می‌کند:

$$u'_n(x) = f(x, u_n(x)) \quad u_n(x_n) = y_n \quad (۹.۵.۶)$$

در مرحلهٔ x_n ، با دانستن y_n ، کمیت $u_{n+1}(x)$ است که سعی در محاسبهٔ آن داریم، نه $Y(x_{n+1})$. با به‌کار بردن روشی که به (۱.۵.۶) منجر شد به‌ازای مقداری چون $x_n \leq \xi_n \leq x_{n+1}$ ، داریم:

$$u_n(x_{n+1}) = y_n + \frac{h}{\gamma} [f(x_n, y_n) + f(x_{n+1}, u_n(x_{n+1}))] - \frac{h^{\gamma}}{\gamma^2} u_n^{(\gamma)}(\xi_n) \quad (۱۰.۵.۶)$$

گیریم $\tilde{e}_{n+1} = u_n(x_{n+1}) - y_{n+1}$ ، که آن را خطای موضعی در محاسبهٔ y_{n+1} از y_n می‌نامیم. رابطهٔ (۲.۵.۶) را از رابطهٔ بالا کم می‌کنیم تا به دست آوریم

$$\begin{aligned} \tilde{e}_{n+1} &= \frac{h}{\gamma} [f(x_{n+1}, u_n(x_{n+1})) - f(x_{n+1}, y_{n+1})] - \frac{h^{\gamma}}{\gamma^2} u_n^{(\gamma)}(\xi_n) \\ &= \frac{h}{\gamma} f_y(x_{n+1}, y_{n+1}) \tilde{e}_{n+1} + O(h \tilde{e}_{n+1}^{\gamma}) - \frac{h^{\gamma}}{\gamma^2} u_n^{(\gamma)}(x_n) + O(h^{\gamma}) \end{aligned}$$

که در آن دوبار قضیهٔ مقدار میانگین را به‌کار برده‌ایم. می‌توان نشان داد که به‌ازای جمیع مقادیر به اندازهٔ کافی کوچک h ،

$$\tilde{e}_{n+1} = O(h^{\gamma})$$

به‌طور دقیقتر

$$\begin{aligned} \tilde{e}_{n+1} &= \left[1 - \frac{h}{\gamma} f_y(x_{n+1}, y_{n+1}) \right]^{-1} \cdot \left[-\frac{h^{\gamma}}{\gamma^2} u_n^{(\gamma)}(x_n) + O(h^{\gamma}) \right] \\ u_n(x_{n+1}) - y_{n+1} &= -\frac{h^{\gamma}}{\gamma^2} u_n^{(\gamma)}(x_n) + O(h^{\gamma}) \quad (۱۱.۵.۶) \end{aligned}$$

این رابطه نشان می‌دهد که خطای موضعی، در اصل همان خطای برشی است. اگر روش اویلر برای محاسبهٔ $y_{n+1}^{(o)}$ به‌کار رفته باشد،

$$y_{n+1}^{(o)} = y_n + hf(x_n, y_n) \quad (۱۲.۵.۶)$$

در آن صورت $u_n(x_{n+1})$ را می‌توان بسط داده نشان داد که

$$u_n(x_{n+1}) - y_{n+1}^{(o)} = \frac{h^{\gamma}}{\gamma} u_n^{(\gamma)}(\xi_n) \quad x_n \leq \xi_n \leq x_{n+1} \quad (۱۳.۵.۶)$$

از ترکیب با رابطه (۱۱.۵.۶)،

$$y_{n+1} - y_{n+1}^{(0)} = O(h^2) \quad (۱۴.۵.۶)$$

برای برقراری (۸.۵.۶) از کران (۵.۵.۶) نتیجه می‌شود که باید دو بارست محاسبه و سپس از $y_{n+1}^{(2)}$ برای معرفی y_{n+1} استفاده کنیم.

با استفاده از روش میانگاهی می‌توانیم حدس اولیه دقیقتری به دست آوریم

$$y_{n+1}^{(0)} = y_{n-1} + 2hf(x_n, y_n) \quad (۱۵.۵.۶)$$

برای برآورد خطا، از همان راهی که به (۱.۴.۶) منجر شد استفاده می‌کنیم تا به دست آوریم

$$u_n(x_{n+1}) = u_n(x_{n-1}) + 2hf(x_n, u_n(x_n)) + \frac{h^2}{3} u_n^{(2)}(\eta_n)$$

برای مقداری از $x_{n-1} \leq \eta_n \leq x_{n+1}$ (۱۱.۵.۶) را از آن کم می‌کنیم،

$$u_n(x_{n+1}) - y_{n+1}^{(0)} = u_n(x_{n-1}) - y_{n-1} + \frac{h^2}{3} u_n^{(2)}(\eta_n)$$

کمیت $u_n(x_{n-1}) - y_{n-1}$ را می‌توان به طریقی مشابه راهی که در (۱۱.۵.۶) به کار برده شد محاسبه نمود و نتیجه مشابهی به دست آورد:

$$u_n(x_{n-1}) - y_{n-1} = \frac{h^2}{12} u_n^{(2)}(x_n) + O(h^4) \quad (۱۶.۵.۶)$$

پس

$$u_n(x_{n+1}) - y_{n+1}^{(0)} = \frac{5h^2}{12} u_n^{(2)}(x_n) + O(h^4)$$

و از ترکیب این رابطه با (۱۱.۵.۶)

$$y_{n+1} - y_{n+1}^{(0)} = \frac{h^2}{3} u_n^{(2)}(x_n) + O(h^4) \quad (۱۷.۵.۶)$$

با حدس اولیه (۱۵.۵.۶)، یک بارست از (۳.۵.۶) برای برقراری (۸.۵.۶) بر مبنای کران در (۵.۵.۶) کافی خواهد بود.

فرمولهای (۱۲.۵.۶) و (۱۵.۵.۶) را فرمولهای پیشگو نامند، و فرمول بارستی ذوزنقه‌یی (۳.۵.۶) فرمول تصحیح‌کننده خوانده می‌شود. این دو فرمول روی هم روش پیشگو-تصحیح‌کننده را تشکیل می‌دهند و پایه روشی هستند که برای کنترل اندازه خطای موضعی می‌توان به کار برد. این مطلب در بخش آینده نشان داده شده است.

همگرایی و نتایج پایداری همگرایی روش دوزنقه‌ی با قضیه ۶.۶ تضمین شده است. فرض کنیم $hk \leq 1$

$$\begin{aligned} \max_{x_* \leq x_n \leq b} |Y(x_n) - y_h(x_n)| \leq e^{\tau K(b-x_*)} |e_*| \\ + \left[\frac{e^{\tau K(b-x_*)} - 1}{K} \right] \left[\frac{h^2}{12} \|Y''\|_\infty \right] \quad (18.5.6) \end{aligned}$$

پیداکردن فرمول خطای مجانبی مشابه پیداکردن آن در روش اویلر است. با فرض $e_* = \delta h^2 + O(h^3)$ می‌توانیم نشان دهیم

$$\begin{aligned} Y(x_n) - y_h(x_n) &= D(x_n)h^2 + O(h^3) \\ D'(x) &= f_y(x, Y(x))D(x) - \frac{1}{12}Y''(x) \quad D(x_*) = \delta. \end{aligned} \quad (19.5.6)$$

برای روش دوزنقه‌ی نیز می‌توان یک فرمول پایداری داد، به شکل استاندارد، مانند آنچه در (۲۸.۲.۶) و (۲۹.۲.۶) برای روش اویلر داده شده است. اثبات را به عهده خواننده واگذار می‌کنیم. همان‌طور که در روش میانگاهی عمل کردیم، می‌توانیم اثر به‌کار بردن روش دوزنقه‌ی را در مورد معادله نمونه

$$y' = \lambda y \quad y(0) = 1 \quad (20.5.6)$$

که جواب آن $Y(x) = e^{\lambda x}$ است، بررسی کنیم. برای آنکه انگیزش انجام چنین کاری را بیشتر نشان دهیم، روش دوزنقه‌ی را که برای مسأله خطی زیر به‌کار رفته در نظر می‌گیریم

$$y' = \lambda y + g(x) \quad y(0) = Y. \quad (21.5.6)$$

یعنی

$$y_{n+1} = y_n + \frac{h}{4} [\lambda y_n + g(x_n) + \lambda y_{n+1} + g(x_{n+1})] \quad n \geq 0 \quad (22.5.6)$$

با $y_* = Y$. در این صورت روش عددی اختلال یافته را در نظر می‌گیریم

$$z_{n+1} = z_n + \frac{h}{4} [\lambda z_n + g(x_n) + \lambda z_{n+1} + g(x_{n+1})] \quad n \geq 0$$

با $z_* = Y_* + \varepsilon$. برای تحلیل اثر اختلال در مقدار اولیه گیریم $w_n = z_n - y_n$. با کم‌کردن داریم

$$w_{n+1} = w_n + \frac{h}{4} [\lambda w_n + \lambda w_{n+1}] \quad n \geq 0 \quad w_* = \varepsilon \quad (23.5.6)$$

این همان روش ذوزنقه‌بندی است که برای مسأله نمونه به‌کار رفته است، جز اینکه مقدار اولیه به جای یک، ε است. جواب عددی (۲۳.۵.۶) در واقع ε برابر آن جواب عددی است که برای (۲۰.۵.۶) به‌دست آمده است. بنابراین رفتار جواب عددی مسأله نمونه، رفتار پایداری روش ذوزنقه‌بندی را که برای (۲۱.۵.۶) به‌کار رفته بود به‌دست می‌دهد.

مسائل نمونه جالب برای ما آنهایی هستند که در آنها λ حقیقی و منفی یا λ مختلط با قسمت حقیقی منفی است. دلیل چنین انتخابی این است که، در چنین صورتی مسأله معادله دیفرانسیل (۲۱.۵.۶) خوش-وضع است، همان‌گونه که در (۸.۱.۶) ذکر شده است، و جالبترین حالتی که مستثنا شده‌اند $\lambda = 0$ و λ انگاری محض هستند. با به‌کار بردن قاعده ذوزنقه‌بندی برای (۲۰.۵.۶)،

$$y_{n+1} = y_n + \frac{h\lambda}{2} [y_n + y_{n+1}] \quad y_0 = 1$$

پس

$$y_{n+1} = \left[\frac{1 + (h\lambda/2)}{1 - (h\lambda/2)} \right] y_n \quad n \geq 0$$

به طریق استقراء

$$y_n = \left[\frac{1 + (h\lambda/2)}{1 - (h\lambda/2)} \right]^n \quad n \geq 0 \quad (24.5.6)$$

به شرطی که $h\lambda \neq 2$. برای حالت $\lambda < 0$ حقیقی، می‌نویسیم

$$r = \frac{1 + (h\lambda/2)}{1 - (h\lambda/2)} = 1 + \frac{h\lambda}{1 - (h\lambda/2)} = -1 + \frac{2}{1 - (h\lambda/2)}$$

این رابطه نشان می‌دهد که به‌ازای جمیع مقادیر $h > 0$ ، داریم $-1 < r < 1$. بنابراین

$$\lim_{n \rightarrow \infty} y_n = 0 \quad (25.5.6)$$

برای کراندار بودن $\{y_n\}$ هیچ محدودیتی برای h وجود ندارد، و بنابراین پایداری روش عددی (۲۲.۵.۶) برای همه مقادیر $h > 0$ و همه مقادیر $\lambda < 0$ ، تضمین می‌شود. این یک حکم قویتری است از آنچه که ممکن است برای بیشتر روشهای عددی گفته شود، معمولاً h باید به اندازه کافی کوچک باشد تا پایداری تضمین شود. در برخی کاربردها، مثل معادلات دیفرانسیل سرسخت، این یک نکته مهمی است. این ویژگی که (۲۵.۵.۶) به‌ازای جمیع مقادیر $h > 0$ و همه مقادیر مختلط λ با $\text{Real}(\lambda) < 0$ برقرار است، A -پایداری (یا پایداری مطلق-م.) خوانده می‌شود. ما آن را در بخش ۸.۶ و مسأله ۳۷ بیشتر باز خواهیم کرد.

برآورد خطای ریچاردسن این برآورد خطا در بخش ۴.۵ معرفی، و هم برای پیشگویی خطا [مانند برآورد خطا در (۴۲.۴.۵)] و هم برای به‌دست‌آوردن یک روش انتگرالگیری عددی همگرای سریعتر [مانند برونایی در (۴۰.۴.۵)] به‌کار گرفته شد. این برآورد را می‌توان به هر دو طریق فوق برای حل معادلات دیفرانسیل به‌کار برد، اگرچه ما بیشتر آن را برای پیشگویی خطا به‌کار می‌گیریم. گیریم $y_h(x)$ و $y_{2h}(x)$ معرف جوابهای عددی $y' = f(x, y)$ در $[x_0, b]$ باشند، که با استفاده از روش ذوزنقه‌یی (۲.۵.۶) به‌دست آمده‌اند. در این صورت با استفاده از (۱۹.۵.۶)،

$$Y(x_n) - y_h(x_n) = D(x_n)h^2 + O(h^2)$$

$$Y(x_n) - y_{2h}(x_n) = 4D(x_n)h^2 + O(h^2)$$

از ضرب اولین معادله در ۴ و کم کردن آن از دومین معادله، و حل نسبت به $Y(x_n)$ به‌دست می‌آوریم:

$$Y(x_n) = \frac{1}{3}[4y_h(x_n) - y_{2h}(x_n)] + O(h^2) \quad (26.5.6)$$

سمت راست فرمول یک مرتبه همگرایی بالاتری نسبت به روش ذوزنقه‌یی دارد ولی باید توجه کرد که برای آن نیاز به محاسبه $y_h(x_n)$ و $y_{2h}(x_n)$ در تمام نقاط گرهی x_n در $[x_0, b]$ است. کاربرد بیشتر فرمول (۲۶.۵.۶) در پیشگویی خطای فراگیر در $y_h(x)$ است. با استفاده از (۲۶.۵.۶)،

$$Y(x_n) - y_h(x_n) = \frac{1}{3}[y_h(x_n) - y_{2h}(x_n)] + O(h^2)$$

با توجه به رابطه (۱۹.۵.۶)، طرف چپ معادله فوق $O(h^2)$ است و بنابراین اولین جمله سمت راست نیز باید $O(h^2)$ باشد. بنابراین

$$Y(x_n) - y_h(x_n) \doteq \frac{1}{3}[y_h(x_n) - y_{2h}(x_n)] \quad (27.5.6)$$

یک برآورد مجانبی خطاست. این یک شیوه عملی در برآورد خطای فراگیر است، اگرچه روشی که ما برای پیدا کردن آن به‌کار بردیم برای گره‌های با طول گام متغیر، مجاز نخواهد بود.

مثال مسأله زیر را در نظر می‌گیریم

$$y' = -y^2 \quad y(0) = 1$$

جدول ۸.۶ روش دوزنقه‌ی و برآورد خطای ریچاردسن

x	$y_{2h}(x)$	$Y(x) - y_{2h}(x)$	$y_h(x)$	$Y(x) - y_h(x)$	$\frac{1}{3}[y_h(x) - y_{2h}(x)]$
۱.۰	۰.۴۸۳۱۴۴	۰.۰۱۶۸۵۶	۰.۴۹۶۰۲۱	۰.۰۰۳۹۷۹	۰.۰۰۴۲۹۲
۲.۰	۰.۳۲۳۶۱۰	۰.۰۰۹۷۲۳	۰.۳۳۰۹۹۱	۰.۰۰۲۳۴۲	۰.۰۰۲۴۶۰
۳.۰	۰.۲۴۳۸۹۰	۰.۰۰۶۱۱۰	۰.۲۴۸۵۲۱	۰.۰۰۱۴۷۹	۰.۰۰۱۵۴۳
۴.۰	۰.۱۹۴۸۳۸	۰.۰۰۴۱۶۲	۰.۱۹۸۹۹۱	۰.۰۰۱۰۰۹	۰.۰۰۱۰۵۱
۵.۰	۰.۱۶۳۶۵۸	۰.۰۰۳۰۰۸	۰.۱۶۵۹۳۷	۰.۰۰۰۷۳۰	۰.۰۰۰۷۵۹

که دارای جواب $Y(x) = 1/(1+x)$ است. نتایج جدول ۸.۶ برای طول گامهای $h = ۰.۲۵$ و ۰.۵۰ هستند. ستون آخر برآورد خطای (۲۷.۵.۶) است که یک برآوردگر دقیق از خطای واقعی $Y(x) - y_h(x)$ است.

۶.۶ الگوریتمی پیشگو-تصحیح‌کننده از مرتبه پایین

در این بخش، الگوریتم نسبتاً ساده‌ای برای حل مسأله مقدار اولیه (۱.۰.۶) آورده شده است. در این الگوریتم از روش دوزنقه‌ی (۲.۵.۶) استفاده می‌شود و اندازه خطای موضعی با تغییر طول گام h کنترل می‌شود. این روش به علت مرتبه پایین همگرایی‌اش مفید نیست، ولی برخی از مفاهیم و تکنیکهای موجود در ساختن الگوریتمی پیشگو-تصحیح‌کننده، با طول گام متغیر را نشان می‌دهد. هر مرحله از x_n به x_{n+1} عبارت از ساختن y_{n+1} از y_n و y_{n-1} است و y_{n+1} یک جواب تقریبی (۲.۵.۶) بر پایه استفاده از بارستی از (۳.۵.۶) خواهد بود. یک مرحله عادی دارای طول گام $h = x_n - x_{n-1} = x_{n+1} - x_n$ است، پیشگوی میانگاهی (۱۵.۵.۶) به‌کار گرفته شده و خطای موضعی از تفاضل فرمولهای پیشگو و تصحیح‌کننده، تخمین زده شده است. وقتی تغییر طول گام انجام شده باشد، پیشگوی (۱۲.۵.۶) ی اویلر به‌کار گرفته شده است.

استفاده‌کننده از این الگوریتم باید علاوه بر تعیین پارامترهایی که مسأله معادله دیفرانسیل (۱.۰.۶) را تعریف می‌کنند، چندین پارامتر را نیز معین نماید. طول گام h تغییر می‌کند و کاربر باید h_{Min} و h_{Max} را که طول گام را محدود می‌نمایند معین کند. کاربر همچنین باید یک مقدار اولیه برای h معین کند و این مقدار باید طوری باشد که به‌ازای آن $\frac{1}{3}h f_y(x_0, y_0)$ از لحاظ قدرمطلق، به اندازه کافی کوچکتر از یک، مثلاً $۱۰^{-۶}$ باشد. این کمیت سرعت همگرایی بارستی را در (۳.۵.۶) مشخص می‌کند که بعداً در همین بخش به دنبال یک مثال عددی مورد بحث قرار گرفته است. یک تحمل خطای ε باید داده شود، و طول گام h باید طوری انتخاب شده باشد که خطای برشی

موضعی در هر مرحله در رابطه زیر، صدق کند

$$\frac{1}{4}\varepsilon h \leq |\text{trunc}| \leq \varepsilon h \quad (۱.۶.۶)$$

این فرمول را کنترل‌کنندهٔ خطا در طول گام واحد گویند. اهمیت آن در نزدیک به پایان همین بخش مورد بحث قرار گرفته است.

نمادگذاری بخش قبل ادامه دارد. تابع $u_n(x)$ جواب $y' = f(x, y)$ است که از نقطهٔ (x_n, y_n) می‌گذرد. خطای موضعی که باید برآورد و کنترل شود، فرمول (۱۱.۵.۶) است که خطای پیدا کردن $u_n(x_{n+1})$ با استفاده از روش ذوزنقه‌یی است:

$$u_n(x_{n+1}) - y_{n+1} = -\frac{h^2}{12} u_n^{(2)}(x_n) + O(h^4) \quad h = x_{n+1} - x_n \quad (۲.۶.۶)$$

اگر y_n به اندازهٔ کافی به $Y(x_n)$ نزدیک باشد، آنگاه این عبارت یک تقریب خوب برای خطای برشی بسیار وابسته در (۱.۵.۶) است:

$$-\frac{h^2}{12} Y^{(2)}(\xi_n)$$

و (۲.۶.۶) تنها کمیتهی است که اطلاعات لازم برای کنترل آن را در دست داریم.

انتخاب طول گام اولیه مسأله عبارت است از پیدا کردن مقدار اولیهٔ h و نقطهٔ گرهی $x_1 = x_0 + h$ که برای آن $|y_1 - Y(x_1)|$ در کرانه‌های (۱.۶.۶) صدق کند. با مقدار اولیهٔ h که توسط کاربر داده می‌شود، مقدار $y_h(x_1)$ با استفاده از پیشگوی (۱۲.۵.۶) اوایلر و دوبار بارستن در (۳.۵.۶) به دست می‌آید. با همین روند، مقادیر $y_{h/2}(x_0 + h/2)$ و $y_{h/2}(x_1)$ نیز محاسبه می‌شوند. روند برونیایی ریچاردسن برای پیشگویی خطا در $y_h(x)$ به کار گرفته می‌شود،

$$Y(x_1) - y_h(x_1) \doteq \frac{4}{3} [y_{h/2}(x_1) - y_h(x_1)] \quad (۳.۶.۶)$$

اگر این خطا در کرانه‌های (۱.۶.۶) صدق نماید، مقدار h پذیرفته می‌شود و مرحلهٔ ذوزنقه‌یی معمولی با استفاده از پیشگوی میانگاهی (۱۱.۵.۶) آغاز می‌شود. ولی اگر (۳.۶.۶) در (۱.۶.۶) صدق نکند، آنگاه یک مقدار جدید برای h انتخاب می‌شود.

با استفاده از مقادیر

$$f_0 = f(x_0, y_0) \quad f_1 = f\left(x_0 + \frac{h}{4}, y_{h/2}\left(x_0 + \frac{h}{4}\right)\right) \quad f_2 = f(x_1, y_{h/2}(x_1))$$

تقریب زیر را پیدا می‌کنیم

$$Y^{(3)}(x_0) \doteq D_{\tau}y \equiv \frac{(f_2 - 2f_1 + f_0)}{(h^2/4)} \quad (4.6.6)$$

این تقریب، تقریبی است که در آن از تفاضل منقسم مرتبهٔ دوم $Y' = f(x, y)$ استفاده می‌شود؛ برای مثال لم ۲ بخش ۴.۳ را به‌کار می‌بریم. برای هر طول گام کوچک h ، خطای برشی در $x_0 + h$ به خوبی با فرمول زیر تقریب زده می‌شود:

$$-\frac{h^3}{12} Y^{(3)}(\xi_0) \doteq -\frac{h^3}{12} D_{\tau}y$$

طول گام جدید به گونه‌ای انتخاب می‌شود که

$$\begin{aligned} \left| \frac{h^3}{12} D_{\tau}y \right| &= \frac{1}{4} \varepsilon h \\ h &= \sqrt{\frac{6\varepsilon}{|D_{\tau}y|}} \end{aligned} \quad (5.6.6)$$

این باید خطای برشی اولیه را تقریباً در وسط ناحیهٔ (۱.۶.۶) به‌عنوان محک خطا در گام واحد، قرار دهد. با این مقدار جدید h ، آزمون (۳.۶.۶)، به‌عنوان کنترل درستی عمل، مجدداً تکرار می‌شود. با انتخاب h به گونه‌ای که وقتی خطای برشی نصف یا دو برابر می‌شود، در کران (۱.۶.۶) صدق کند، عدم نیاز به تغییر طول گام را برای چندین گام تأمین کرده‌ایم، به شرطی که مشتق $Y^{(3)}(x)$ به‌سرعت تغییر ننماید. تغییر طول گام پرهزینه‌تر از یک گام معمولی است، و ما می‌خواهیم نیاز به چنین تغییری را مینیمم نماییم.

گام پیشگو-تصحیح‌کنندهٔ معمولی طول گام h در $x_{n+1} - x_n = x_n - x_{n-1} = h$ صدق می‌کند. برای پیدا کردن مقدار y_{n+1} ، قاعدهٔ پیشگوی میانگاهی (۱۵.۵.۶) را به‌کار می‌بریم و در (۳.۵.۶) یک بارست اجرا می‌کنیم. خطای موضعی (۱۱.۵.۶) با استفاده از (۱۷.۵.۶) برآورد می‌شود:

$$\begin{aligned} -\frac{1}{6}(y_{n+1} - y_{n+1}^{(o)}) &= -\frac{h^3}{12} u_n^{(3)}(x_n) + O(h^4) \\ &= [u_n(x_{n+1}) - y_{n+1}] + O(h^4) \end{aligned} \quad (6.6.6)$$

بنابراین خطای موضعی را با استفاده از رابطهٔ

$$\text{trunc} \equiv -\frac{1}{6}(y_{n+1} - y_{n+1}^{(o)}) \quad (7.6.6)$$

اندازه می‌گیریم. اگر خطای برشی در (۱.۶.۶) صدق کند، h تغییر داده نمی‌شود و محاسبات با این روند گام منظم ادامه خواهد یافت. ولی وقتی (۱.۶.۶) برقرار نباشد، مقادیر x_{n+1} و y_{n+1} کنار گذاشته می‌شوند و طول گام جدیدی بر پایه مقدار خطای برشی، انتخاب می‌شود.

تغییر طول گام با استفاده از (۶.۶.۶)،

$$\frac{u_n^{(r)}(x_n) \cdot \text{trunc}}{12 h^r}$$

که h طول گام به‌کار برده شده در محاسبه trunc را نشان می‌دهد. برای یک طول گام دلخواه h ، خطای موضعی در به‌دست آوردن y_{n+1} با عبارت زیر برآورد می‌شود:

$$u_n(x_n + h) - y_n(x_{n+1}) \doteq - \frac{h^r}{12} u_n^{(r)}(x_n) \doteq \left[\frac{h}{h_0} \right]^r \text{trunc}$$

h را به گونه‌ای انتخاب می‌کنیم که

$$\left[\frac{h}{h_0} \right]^r | \text{trunc} | = \frac{1}{4} \varepsilon h$$

$$h = \sqrt{\frac{\varepsilon h_0^r}{2 | \text{trunc} |}} \quad (۸.۶.۶)$$

y_{n+1} را با استفاده از پیشگوی اوایلر و دو مرتبه بارستن در (۳.۵.۶) محاسبه می‌کنیم و سپس به گام پیشگو- تصحیح‌کننده عادی باز می‌گردیم. برای احتراز از تغییرات سریع در h که می‌تواند به خطاهای قابل‌ملاحظه‌ای منجر شود، مقدار جدید h هرگز نباید بیش از دو برابر مقدار قبلی افزایش یابد. اگر مقدار جدید h از h_{\min} کمتر باشد، آنگاه محاسبه متوقف می‌شود. ولی اگر h بزرگتر از h_{\max} باشد، فقط $h = h_{\max}$ را انتخاب کرده محاسبه را ادامه می‌دهیم. این امر ممکن است مشکلاتی در برداشته باشد، که به دنبال مثال عددی مورد بحث قرار گرفته است.

الگوریتم $\text{Detrap}(f, x_0, y_0, x_{\text{end}}, \varepsilon, h, h_{\min}, h_{\max}, \text{ier})$

۱. تبصره: مسأله‌ای که باید حل شود $Y' = f(x, y)$ و $Y(x_0) = y_0$ با ازای $x_0 \leq x \leq x_{\text{end}}$ است، با روشی که قبلاً در این بخش گفته شد. مقادیر جواب تقریبی در هر نقطه گرهی درج می‌شود. در مورد پارامتر خطای ε و پارامترهای طول گام قبلاً در این بخش بحث شد. متغیر ier یک شاخص خطا و یک خروجی در پایان الگوریتم است: $\text{ier} = 0$ یعنی یک بازگشت

معمولی؛ $ier = 1$ یعنی در بعضی نقاط گرهی $h = h_{Max}$ ؛ و $ier = 2$ یعنی چون $h < h_{Min}$ ،
 انتگرالگیری خاتمه یافته است زیرا $h \geq h_{Min}$ یک شرط لازم است.

۲. آغاز کنید: $ier := 0$ و $loop := 1$

۳. تبصره: یک مقدار آغازی برای h انتخاب کنید

۴. با استفاده از روش (۲.۵.۶) مقادیر $y_{h/2}(x_0 + h)$ ، $y_{h/2}(x_0 + h/2)$ ، $y_h(x_0 + h)$ را محاسبه کنید. در هر مورد، از پیشگویی (۱۲.۵.۶) اویلر استفاده کرده آن را با دو بارست (۳.۵.۶) دنبال کنید

۵. برای خطا در $y_h(x_0 + h)$ از

$$trunc := \frac{4}{3}[y_{h/2}(x_0 + h) - y_h(x_0 + h)]$$

استفاده کنید.

۶. اگر $|\text{trunc}| \leq \epsilon h$ یا اگر $\frac{1}{4}\epsilon h \leq \text{trunc} \leq \epsilon h$ ، $loop = 2$ ، آنگاه $x_1 := x_0 + h$ ، $y_1 := y_h(x_0 + h)$ و x_1 را درج کنید و به مرحله ۱۰ بروید.

۷. از (۴.۶.۶)، $D_2 y = y^{(2)}(x_0)$ را حساب کنید. اگر $D_2 y \neq 0$ ، آنگاه

$$h := \left[\frac{6\epsilon}{|D_2 y|} \right]^{1/2}$$

اگر $D_2 y = 0$ ، آنگاه $h := h_{Max}$ و $loop := 2$

۸. اگر $h < h_{Min}$ ، آنگاه $ier := 2$ و خارج شوید. اگر $h > h_{Max}$ ، آنگاه $h := h_{Max}$ و $loop := 2$ ، $ier := 1$

به مرحله ۴ بروید.

۱۰. تبصره: این بخش الگوریتم شامل مرحله پیشگو-تصحیح کننده عادی با کنترل خطاست.

۱۱. بگیرید $x_2 := x_1 + h$ و $y_2 := y_1 + 2hf(x_1, y_2)$ را با یک بارست (۵.۶.۳)

به دست آورید.

$$trunc := -\frac{1}{6}(y_2 - y_2^{(0)}) \quad 12$$

۱۳. اگر $|\text{trunc}| > \epsilon h$ یا $|\text{trunc}| < \frac{1}{4}\epsilon h$ به مرحله ۱۶ بروید.

۱۴. x_2 و y_2 را درج کنید.

۱۵. $x_0 := x_1$ ، $x_1 := x_2$ ، $y_0 := y_1$ و $y_1 := y_2$ ، اگر $x_1 < x_{end}$ ، به مرحله ۱۱ بروید.

وگرنه خارج شوید.

۱۶. تبصره: طول گام را تغییر دهید.

۱۷. $x_1 := x_0, y_1 := y_0, h_1 := h_0$ و با استفاده از (۸.۶.۶)، h را محاسبه کنید.

۱۸. $h := \text{Min}\{h, 2h_0\}$.

۱۹. اگر $h < h_{\text{Min}}$ ، آنگاه $ier := 2$ و خارج شوید. اگر $h > h_{\text{Max}}$ ، آنگاه $ier := 1$ و

$h := h_{\text{Max}}$.

۲۰. $y_1^{(0)} := y_0 + hf(x_0, y_0)$ و y_1 را با انجام دو بارست (۳.۵.۶) محاسبه کنید.

۲۱. x_1 و y_1 را درج کنید.

۲۲. اگر $x_1 < x_{\text{end}}$ ، به مرحله ۱۰ برگردید، وگرنه خارج شوید.

در مثال زیر یک اجرای Detrap به‌کار رفته که در آن trunc نیز درج شده است. جزیی به برنامه‌نویسی اضافه شده است تا خطای برشی y_1 در مرحله ۲۰ پیشگویی شود.

مثال مسأله زیر را در نظر می‌گیریم

$$y' = \frac{1}{1+x^2} - 2y^2 \quad y(0) = 0 \quad (9.6.6)$$

جواب درست آن چنین است:

$$Y(x) = \frac{x}{1+x^2}$$

این یک مسأله جالب برای آزمون Detrap است و به خوبی اجرا می‌شود. این معادله در بازه $[0, 10]$ با $h_{\text{Min}} = 0.001$ و $h_{\text{Max}} = 10$ ، $h = 0.1$ ، $\varepsilon = 0.0005$ حل شده است. جدول ۹.۶ شامل بعضی از نتایج است، از جمله خطای فراگیر درست، و خطای موضعی درست که با True le نشان داده شده است. این خطای آخری با استفاده از یک روش عددی دقیقتر محاسبه شده است. به علت کمبود جا، فقط قسمتهایی از خروجی را در این جدول آورده‌ایم.

چندین نکته را با این مثال نشان خواهیم داد. نخست، مرحله ۱۸ برای احتراز از طول گامهایی بسیار بزرگ لازم است. برای مسأله (۹.۶.۶) داریم

$$Y^{(2)}(x) = \frac{-6(x^4 - 6x^2 + 1)}{(1+x^2)^2}$$

که در $x = \pm 0.414$ و $x = \pm 2.414$ صفر می‌شود. بنابراین، بر پایه رابطه (۲.۶.۶) و نزدیک بودن $u_n(x)$ به $Y(x)$ ، خطای موضعی در حل مسأله (۹.۶.۶) در نزدیکی این نقاط بسیار کوچک

جدول ۹.۶ مثال از الگوریتم Detrap

x_n	h	y_n	$Y(x_n) - y_n$	trunc	True le
۰٫۰۲۲۷	۰٫۰۲۲۷	۰٫۰۲۲۶۸۹	۵٫۸۴E-۶	۵٫۸۴E-۶	۵٫۸۴E-۶
۰٫۰۴۵۴	۰٫۰۲۲۷	۰٫۰۴۵۳۰۸	۱٫۱۷E-۵	۵٫۸۳E-۶	۵٫۸۴E-۶
۰٫۰۶۸۱	۰٫۰۲۲۷	۰٫۰۶۷۷۸۷	۱٫۷۴E-۵	۵٫۷۶E-۶	۵٫۷۵E-۶
۰٫۰۹۰۸	۰٫۰۲۲۷	۰٫۰۹۰۰۶۰	۲٫۲۸E-۵	۵٫۶۲E-۶	۵٫۶۱E-۶
۰٫۱۱۳۵	۰٫۰۲۲۷	۰٫۱۱۲۵۹۴	۵٫۱۶E-۵	۲٫۹۶E-۶	۲٫۸۵E-۶
۰٫۱۳۶۵	۰٫۰۳۴۰	۰٫۱۳۵۱۲۵	۵٫۶۶E-۵	۶٫۷۴E-۶	۶٫۷۹E-۶
۰٫۱۵۹۵	۰٫۰۳۴۰	۰٫۱۵۸۰۸۴	۶٫۰۱E-۵	۶٫۲۱E-۶	۵٫۷۳E-۶
۰٫۱۸۲۵	۰٫۰۳۴۰	۰٫۱۸۱۴۱۱	۶٫۱۱E-۵	۴٫۲۸E-۶	۳٫۵۴E-۶
۰٫۲۰۵۵	۰٫۰۶۶۲	۰٫۲۰۴۹۰۱۹	۵٫۰۵E-۵	-۶٫۵۶E-۶	-۵٫۲۰E-۶
۰٫۲۲۸۵	۰٫۰۶۶۲	۰٫۲۲۸۲۹۷	۲٫۴۴E-۵	۱٫۰۴E-۵	-۲٫۱۲E-۵
۰٫۲۵۱۵	۰٫۰۶۶۲	۰٫۲۵۱۶۶۹	-۲٫۰۳E-۵	-۲٫۹۲E-۵	-۴٫۲۱E-۵
۰٫۲۷۴۵	۰٫۰۴۰۶	۰٫۲۷۴۵۸۷۹	-۲٫۹۹E-۵	-۱٫۱۲E-۵	-۱٫۱۰E-۵
۰٫۲۹۷۵	۰٫۱۳۵	۰٫۲۹۷۴۹۸۲	-۱٫۰۲E-۴	-۱٫۶۴E-۵	-۱٫۶۷E-۵
۰٫۳۲۰۵	۰٫۱۳۵	۰٫۳۲۰۸۹۴۴	-۱٫۰۳E-۴	-۱٫۷۹E-۵	-۲٫۱۱E-۵
۰٫۳۴۳۵	۰٫۲۲۳	۰٫۳۴۳۸۶۴	-۶٫۵۷E-۵	۱٫۲۷E-۵	۸٫۱۵E-۶
۰٫۳۶۶۵	۰٫۴۴۶	۰٫۳۶۶۴۴۹	۳٫۴۴E-۴	۴٫۴۱E-۴	۳٫۷۸E-۴
۰٫۳۸۹۵	۰٫۳۰۳	۰٫۳۸۹۴۴۷	۳٫۲۱E-۴	۹٫۳۹E-۵	۸٫۴۱E-۵
۰٫۴۱۲۵	۰٫۶۷۲	۰٫۴۱۲۷۳۹۶	۳٫۸۷E-۴	۸٫۷۷E-۵	۱٫۱۲E-۴
۰٫۴۳۵۵	۱٫۰۰۰	۰٫۴۳۵۱۰۰	۳٫۹۶E-۴	۱٫۷۳E-۴	۱٫۵۷E-۴
۰٫۴۵۸۵	۱٫۰۰۰	۰٫۴۵۸۱۶۲۵	۴٫۲۷E-۴	۱٫۱۸E-۴	۱٫۶۸E-۴
۰٫۴۸۱۵	۱٫۰۰۰	۰٫۴۸۱۲۲۷۳	۴٫۱۱E-۴	۹٫۴۵E-۵	۱٫۲۱E-۴

می‌شود. این امر موجب پیشگویی یک h بسیار بزرگ، در (۸.۶.۶) می‌شود، مقداری که برای نقاط بعدی x_n بسیار بزرگ است. در $x_n = ۲٫۷۶۳۲$ ، گام ۱۸ الگوریتم، برای احتراز از مقدار بزرگ گمراه‌کننده h لازم بود. همان‌گونه که می‌توان مشاهده کرد، خطای موضعی در $x_n = ۲٫۷۶۳۲$ به لحاظ مقدار بزرگتر h ، بسیار افزایش می‌یابد. کمی بعد از آن، طول گام h کوچک می‌شود تا اندازه خطای موضعی را کاهش دهد.

در تمام نتیجه‌گیریهای این بخش و بخش بعد، برآوردهایی انجام شده‌اند که اگر h به اندازه کافی کوچک می‌بود بسیار دقیق می‌بودند. در بیشتر حالات کمیت قاطع، در واقع $h.f_y(x_n, y_n)$ است،

همانند (۷.۵.۶)، هنگامی که نرخ همگرایی بارست (۳.۵.۶) تحلیل می‌شد. در حالت بارست دوزنقه‌یی (۳.۵.۶)، این نرخ همگرایی چنین است

$$\text{Rate} \doteq \frac{1}{4} h f_y(x_n, y_n) \quad (۱۰.۶.۶)$$

و برای مسأله (۹.۶.۶) این نرخ چنین است:

$$\text{Rate} \doteq -2hy_n$$

اگر این کمیت نزدیک ۱ باشد چندین بارست برای به‌دست آوردن یک مقدار دقیق y_{n+1} لازم خواهد بود. مطابق جدول، اندازه این نرخ، وقتی h افزایش می‌یابد، تقریباً بزرگ می‌شود. در $x = ۲۳۱۷۲$ این نرخ حدود ۱۶۲° است. این نرخ به اندازه کافی کوچک به‌نظر می‌آید، ولی خطای موضعی بیشتر از قبل نادقیق است، و این شاید بر اثر بارست نادقیق به‌دست آمده در (۳.۵.۶) باشد. می‌توان الگوریتم را پیچیده‌تر نمود تا مسائل با یک h بسیار بزرگ شناسایی شود، ولی انتخاب یک h_{Max} با اندازه‌های معقول، نیز مفید خواهد بود.

خطای کلی با دادن یک کران خطا، مشابه کران (۱۸.۵.۶) برای یک طول گام ثابت، مسأله را آغاز می‌کنیم. می‌نویسیم

$$Y(x_{n+1}) - y_{n+1} = [Y(x_{n+1}) - u_n(x_{n+1})] + [u_n(x_{n+1}) - y_{n+1}] \quad (۱۱.۶.۶)$$

برای آخرین جمله، فرض می‌کنیم که محک خطا در طول گام واحد (۱.۶.۶) برقرار باشد:

$$|u_n(x_{n+1}) - y_{n+1}| \leq \varepsilon(x_{n+1} - x_n) \quad (۱۲.۶.۶)$$

برای جمله دیگر در (۱۱.۶.۶)، شکلهای معادلات انتگرالی

$$Y(x) = Y(x_n) + \int_{x_n}^x f(t, Y(t)) dt$$

$$u_n(x) = y_n + \int_{x_n}^x f(t, u_n(x)) dt \quad x \geq x_n \quad (۱۳.۶.۶)$$

را معرفی می‌کنیم. این دو رابطه را از هم کم می‌کنیم و کرانها را با استفاده از شرط لیبشیتس می‌گیریم

$$|Y(x) - u_n(x)| \leq |e_n| + K(x - x_n) \text{Max}_{x_n \leq t \leq x} |Y(t) - u_n(t)| \quad x \geq x_n$$

که در آن $e_n = Y(x_n) - y_n$ با استفاده از این رابطه، به دست می‌آوریم

$$|Y(x_{n+1}) - u_n(x_{n+1})| \leq \frac{1}{1 - K(x_{n+1} - x_n)} |e_n| \quad (۱۴.۶.۶)$$

می‌نویسیم

$$H = \max_{x_n \leq x_{n+1} \leq b} (x_{n+1} - x_n) = \max h$$

و فرض می‌کنیم

$$HK < 1$$

از ترکیب روابط (۱۱.۶.۶)، (۱۲.۶.۶) و (۱۴.۶.۶)، به دست می‌آوریم

$$|e_{n+1}| \leq \frac{1}{1 - HK} |e_n| + \varepsilon H \quad x_n \leq x_{n+1} \leq b$$

این رابطه، مانند آنچه در قضیه ۳.۶، دیدیم به سادگی حل می‌شود و برای یک مقدار مناسب $c > 0$ به دست می‌آوریم

$$|Y(x_n) - y_n| \leq e^{c(b-x_n)K} |Y(x_0) - y_0| + \left[\frac{e^{c(b-x_0)} - 1}{c} \right] \varepsilon \quad (۱۵.۶.۶)$$

این یک نتیجه اساسی خطا در استفاده از طول گام متغیر، و یک توجیه جزئی از محک خطا در (۱.۶.۶) است.

در بعضی از مواقع، می‌توانیم یک کران واقع بینانه‌تری پیدا کنیم. برای سادگی فرض می‌کنیم برای جميع مقادیر (x, y) ، $f_y(x, y) \leq 0$ روابط (۱۳.۶.۶) را از هم کم می‌کنیم

$$\begin{aligned} Y(x) - u_n(x) &= e_n + \int_{x_n}^x [f(t, Y(t)) - f(t, u_n(t))] dt \\ &= e_n + \int_{x_n}^x \frac{\partial f(t, \zeta(t))}{\partial y} [Y(t) - u_n(t)] dt \end{aligned}$$

در مرحله آخر از قضیه مقدار میانگین ۲.۱ استفاده می‌شود، و می‌توان نشان داد که $v(x) \equiv Y(x) - u_n(x)$ این نشان می‌دهد که $\partial f(t, \zeta(t)) / \partial y$ یک تابع پیوسته از t است. این نشان می‌دهد که یک جواب مسأله خطی زیر است

$$v'(x) = \frac{\partial f(x, \zeta(x))}{\partial y} v(x) \quad v(x_n) = e_n$$

جواب این مسأله خطی همراه با فرض $f_y(x, y) \leq 0$ به دست می‌دهد

$$|Y(x_{n+1}) - u_n(x_{n+1})| \leq |e_n|$$

شرط $f_y(x, y) \leq 0$ مرتبط است با مسائل مقدار اولیه خوش-وضع که قبلاً در (۸.۱.۶) به آن اشاره شد. از ترکیب با (۱۱.۶.۶) و (۱۲.۶.۶) به دست می آوریم

$$|Y(x_{n+1}) - y_{n+1}| \leq |Y(x_n) - y_n| + \varepsilon(x_{n+1} - x_n)$$

این نامعادله را حل می کنیم، یک کران واقع بینانه تر زیر به دست می آید

$$|Y(x_n) - y_n| \leq |Y(x_0) - y_0| + \varepsilon(x_n - x_0) \quad (۱۶.۶.۶)$$

این رابطه تا اندازه ای رفتار خوب مثال در جدول ۹-۶ را توضیح می دهد؛ نتایج نظری بهتری هم وجود دارد. ولی نتایج (۱۵.۶.۶) و (۱۶.۶.۶)، برای توجیه استفاده از آزمون (۱.۶.۶) که خطا را در هر گام کنترل می کند، کافی خواهند بود. در دستگاه معادلات $y' = f(x, y)$ ، به جای شرط $f_y(x, y) \leq 0$ این شرط که، ویژه مقادیرهای ماتریس ژاکوبی $f_y(x, Y(x))$ دارای اجزای حقیقی صفر یا منفی باشند، گذاشته می شود.

الگوریتم *Detrap*، را به چند طریق می توان بهتر کرد. ولی این ساختمان، یک الگوریتم پیشگو-تصحیح کننده با طول گام متغیر را نشان می دهد. خروجیها در یک مجموعه نامناسب از نقاط گرهی x_n درج می شوند ولی با یک الگوریتم ساده درونیابی می توان این مشکل را حل کرد. پیشگوها را می توان بهتر کرد ولی آن نیز الگوریتم را پیچیده تر می کند. در بخش بعد به بحث در الگوریتمهای پیشگو-تصحیح کننده عملی فعلاً موجود که اغلب آنها از نظر مرتبه نیز متغیرند باز می گردیم.

۷.۶ پیدا کردن روشهای چند گامی از مراتب بالاتر

فرمول عمومی روش $p+1$ گامی (۱.۳.۶) از بخش ۳.۶ برای حل مسأله مقدار اولیه (۱.۹.۶) را یادآوری می کنیم:

$$y_{n+1} = \sum_{j=0}^p a_j y_{n-j} + h \sum_{j=-1}^p b_j f(x_{n-j}, y_{n-j}) \quad n \geq p \quad (۱.۷.۶)$$

یک نظریه برای این روشها در بخش ۳.۶ داده شده است. بعضی روشهای مراتب بالاتر خاص را اکنون به دست می آوریم. دو راه اساسی برای به دست آوردن چنین فرمولهایی وجود دارند: (۱) روش ضرایب نامعین و (۲) انتگرالگیری عددی. روشهای انتگرالگیری عددی امروزه از همه رایجترند ولی دیدگاه روش ضرایب نامعین در تحلیل و مطالعه روشهای عددی هنوز مهم است.

فرمولهای ضمنی را می‌توان از راه بارستی، کاملاً مشابه با (۲.۵.۶) و (۳.۵.۶) در روش دوزنقه‌یی، حل کرد. در (۱.۷.۶) اگر $b_{-1} \neq 0$ ، این بارست به شکل زیر تعریف می‌شود:

$$y_{n+1}^{(l+1)} = \sum_{j=0}^p [a_j y_{n-j} + hb_j f(x_{n-j}, y_{n-j})] + hb_{-1} f(x_{n+1}, y_{n+1}^{(l)}) \quad l \geq 0 \quad (2.7.6)$$

این بارست همگراست اگر $hb_{-1}K < 1$ ، که K ثابت لپشیتس برای $f(x, y)$ است که در (۱۲.۲.۶) آمده است. نرخ همگرایی خطی، مشابه با (۵.۵.۶) در روش دوزنقه‌یی، با $hb_{-1}K$ کراندار می‌شود.

ما دنبال یک زوج فرمول، یک تصحیح‌کننده و یک پیشگو، هستیم. فرض کنید که فرمول تصحیح‌کننده از مرتبه m باشد، یعنی در هر مرحله خطای برشی $O(h^{m+1})$ باشد. اغلب فقط یک بارست در (۲.۷.۶) محاسبه می‌شود، و این بدان معناست که پیشگو باید حداقل از درجه $m - 1$ باشد تا خطای برشی در $y_{n+1}^{(1)}$ نیز $O(h^{m+1})$ باشد. بحث خطای بارستی روش دوزنقه‌یی را در بخش ۵.۶، بین (۳.۵.۶) و (۱۷.۵.۶)، ببینید که از پیشگوهای اوایلر و میانگاهی استفاده می‌کند. مفاهیم اصلی بدون هیچ تغییر قابل ملاحظه‌ای به (۱.۷.۶) و (۲.۷.۶) انتقال می‌یابند.

روش ضرایب نامعین اگر قرار باشد فرمول (۱.۷.۶) از مرتبه $m \geq 1$ باشد، آنگاه با توجه به قضیه ۵.۶ لازم و کافی است که

$$\sum_{j=0}^p a_j = 1$$

$$\sum_{j=0}^p a_j (-j)^i + i \sum_{j=-1}^p b_j (-j)^{i-1} = 1 \quad i = 1, 2, \dots, m \quad (3.7.6)$$

برای یک روش صریح، شرط اضافی $b_{-1} = 0$ لازم است.

برای یک روش ضمنی کلی، $2p + 3$ پارامتر $\{a_j, b_j\}$ وجود دارند که باید تعیین شوند و $m + 1$ معادله وجود دارند. ممکن است چنین تصور شود که می‌توان تساوی $m + 1 = 2p + 3$ را نوشت، ولی این امر از نقطه نظر پایداری و همگرایی (۱.۷.۶) بسیار نامعقول است. این نکته در بخش بعد و در مسائل نشان داده شده است. معمولاً، بهترین کار آن است که برای یک روش ضمنی $m \leq p + 2$ انتخاب شود. برای یک روش صریح، ملاحظات پایداری اغلب زیاد مهم نیست زیرا این روش معمولاً یک پیشگو برای یک فرمول ضمنی است.

مثال همه روشهای دوگامی مرتبه دوم را پیدا کنید. فرمول (۱.۷.۶) چنین است

$$y_{n+1} = a_0 y_n + a_1 y_{n-1} + h[b_{-1} f(x_{n+1}, y_{n+1}) + b_0 f(x_n, y_n) + b_1 f(x_{n-1}, y_{n-1})] \quad n \geq 1 \quad (4.7.6)$$

ضرایب باید در رابطه (۳.۷.۶) با $m = 2$ صدق نمایند.

$$a_0 + a_1 = 1 \quad -a_1 + b_{-1} + b_0 + b_1 = 1 \quad a_1 + 2b_{-1} - 2b_1 = 1$$

از حل این دستگاه داریم:

$$a_1 = 1 - a_0 \quad b_{-1} = 1 - \frac{1}{4}a_0 - \frac{1}{4}b_0 \quad b_1 = 1 - \frac{3}{4}a_0 - \frac{1}{4}b_0 \quad (5.7.6)$$

که a_0 و b_0 متغیرند. روش میانگامی حالت ویژه‌ای است که در آن $a_0 = 0$ ، $b_0 = 2$. ضرایب a_0 و b_0 را می‌توان طوری انتخاب کرد که پایداری بهبود یابد، یا خطای برشی کوچک شود، یا یک فرمول صریح به دست آید، یا ترکیبی از این موارد حاصل شود. شرایطی که پایداری و همگرایی را تضمین نمایند (غیر از $0 \leq a_0 \leq 1$ و استفاده از قضیه ۶.۶) پیش از نظریه کلی برای (۱.۷.۶) در بخش بعد، نمی‌توان داد.

با برقراری (۳.۷.۶)، قضیه ۵.۶ ایجاب می‌کند که خطای برشی، برای همه مقادیر $Y(x)$ که $m+1$ بار بر $[x_0, b]$ پیوسته مشتقپذیر است، در

$$T_n(Y) = O(h^{m+1})$$

صدق کند. این شرط برای بیشتر منظوره‌های نظری و عملی کافی است. ولی گاهی در ساختن الگوریتمهای پیشگو-تصحیح‌کننده، ترجیح داده می‌شود که خطا دارای شکل زیر باشد

$$T_n(Y) = d_m h^{m+1} Y^{(m+1)}(\xi_n) \quad x_{n-p} \leq \xi_n \leq x_{n+1} \quad (6.7.6)$$

که d_m مقدار ثابتی است مستقل از n . مثالها عبارت‌اند از روش اویلر، روش میانگامی و روش ذوزنقه‌یی.

برای به دست آوردن فرمول (۶.۷.۶) با فرض برقراری (۳.۷.۶)، ابتدا خطای برشی $T_n(y)$ را به شکل یک انتگرال با استفاده از مفهوم هسته پتانو که در آخرین قسمت بخش ۱.۵ داده شده،

بیان می‌کنیم. $Y(x)$ را حول x_{n-p} بسط می‌دهیم

$$Y(x) = \sum_{i=0}^m \frac{(x - x_{n-p})^i}{i!} Y^{(i)}(x_{n-p}) + R_{m+1}(x)$$

$$R_{m+1}(x) = \frac{1}{m!} \int_{x_{n-p}}^{x_{n+1}} (x-t)_+^m Y^{(m+1)}(t) dt$$

$$(x-t)_t^m = \begin{cases} (x-t)^m & x \geq t \\ 0 & x < t \end{cases}$$

این بسط را در فرمول (۴.۳.۶) به جای $T_n(Y)$ قرار می‌دهیم. با استفاده از این فرض که این روش از مرتبه m است خواهیم داشت

$$T_n(Y) = T_n(R_{m+1})$$

$$= R_{m+1}(x_{n+1}) - \left[\sum_{j=0}^p a_j R_{m+1}(x_{n-j}) + h \sum_{j=-1}^p b_j R'_{m+1}(x_{n-j}) \right]$$

$$T_n(Y) = \int_{x_{n-p}}^{x_{n+1}} G(t - x_{n-p}) Y^{(m+1)}(t) dt \quad (۷.۷.۶)$$

با هسته پتانو

$$G(s) = \frac{1}{m!} \left\{ (x_{p+1} - s)_+^m - \left[\sum_{j=0}^p a_j (x_{p-j} - s)_+^m + hm \sum_{j=-1}^p b_j (x_{p-j} - s)_+^{m-1} \right] \right\}$$

$$= T_p((x-s)_+^m) \quad 0 \leq s \leq x_{p+1} \quad (۸.۷.۶)$$

تابع $G(s)$ اغلب تابع نفوذ نامیده می‌شود.

$m+1$ شرط (۱.۷.۶) برای ضرایب (۳.۷.۶) بدین معناست که $r = (2p+3) - (m+1)$ پارامتر آزاد وجود دارد [برای فرمول صریح $r = (2p+3) - (m+2)$]. بنابراین r پارامتر آزاد برای تعیین $G(s)$ در دست است. ما مقادیری برای این پارامترها انتخاب می‌کنیم که $G(s)$ در $[0, x_{p+1}]$ دارای یک علامت شود. و آنگاه با استفاده از قضیه مقدار میانگین، خواهیم داشت

$$T_{n+1}(Y) = Y^{(m+1)}(\xi_n) \int_{x_{n-p}}^{x_{n+1}} G(t - x_{n-p}) dt$$

با شرط $x_{n-p} \leq \xi_n \leq x_{n+1}$ با عملیات بیشتر، (۶.۷.۶) را با

$$d_m = \frac{1}{m!} \int_0^{p+1} \left[\nu^m - \sum_{j=0}^p a_j (\nu - j - 1)_+^m - m \sum_{j=1}^p b_j (\nu - j - 1)_+^{m-1} \right] d\nu \quad (9.7.6)$$

به دست می آوریم. باز هم این رابطه وابسته به $G(s)$ است که به ازای $0 \leq s \leq x_{p+1}$ دارای یک علامت است.

مثال فرمول (۴.۷.۶) را در نظر گرفته فرض کنید که فرمول صریح است ($b_{-1} = 0$). در این صورت

$$y_{n+1} = a_0 y_n + a_1 y_{n-1} + h[b_0 f(x_n, y_n) + b_1 f(x_{n-1}, y_{n-1})] \quad n \geq 1 \quad (10.7.6)$$

با

$$a_1 = 1 - a_0 \quad b_0 = 2 - \frac{1}{\gamma} a_0 \quad b_1 = -\frac{1}{\gamma} a_0$$

فرمول قبلی را به کار می بریم

$$G(s) = \frac{1}{\gamma} \left[(x_2 - s)_+^2 - a_0 (x_1 - s)_+^2 - a_1 (x_0 - s)_+^2 - 2hb_0 (x_1 - s)_+ - 2hb_1 (x_0 - s)_+ \right] \\ = \begin{cases} \frac{1}{\gamma} s[s(1 - a_0) + a_0 h] & 0 \leq s \leq h \\ \frac{1}{\gamma} (x_2 - s)_+^2 & h \leq s \leq 2h \end{cases}$$

شرط اینکه $G(s)$ در فاصله $[0, 2h]$ دارای یک علامت باشد برقرار است، اگر و تنها اگر $a_0 \geq 0$.

پس

$$T_{n+1}(Y) = \left(\frac{1}{\gamma} a_0 + \frac{1}{\gamma} \right) h^2 Y^{(2)}(\xi_n) \quad x_{n-1} \leq \xi_n \leq x_{n+1} \quad (11.7.6)$$

توجه کنید که خطای برشی مینیمم است وقتی که $a_0 = 0$ ؛ بنابراین در میان همه روشهای صریح دوگامی مرتبه دوم، روش میانگامی برای داشتن خطای برشی کوچک، از همه بهتر است. برای بحث بیشتر از این مثال، مسأله ۳۱ را ببینید.

روشهای مبتنی بر انتگرالگیری عددی موضوع اصلی مسألهٔ زیر است. معادلهٔ دیفرانسیل را با انتگرالگیری در بازه‌ای چون $[x_{n-r}, x_{n+1}]$ به صورت زیر به دست می‌آوریم

$$Y(x_{n+1}) = Y(x_{n-r}) + \int_{x_{n-r}}^{x_{n+1}} f(t, Y(t)) dt \quad (۱۲.۷.۶)$$

بهازای یک $r \geq 0$ و همهٔ مقادیر $n \geq r$. یک چندجمله‌یی $P(t)$ درست می‌کنیم که انتگرالده $Y'(t) = f(t, Y(t))$ را در یک مجموعه نقاط گرهی $\{x_i\}$ درونیابی کند و سپس از $P(t)$ در بازهٔ $[x_{n-r}, x_{n+1}]$ انتگرال می‌گیریم تا تقریبی برای (۱۲.۷.۶) به دست آید. هر سه روش قبلی، اویلر، میانگاهی، و ذوزنقه‌یی، را می‌توان از این طریق به دست آورد.

مثال قاعدهٔ انتگرالگیری سیمپسون، (۱۳.۱.۵) و (۱۵.۱.۵) را برای معادلهٔ زیر به کار می‌بریم

$$Y(x_{n+1}) = Y(x_{n-1}) + \int_{x_{n-1}}^{x_{n+1}} f(t, Y(t)) dt$$

این کار چنین نتیجه می‌دهد

$$Y_{n+1} = Y_{n-1} + \frac{h}{3} [Y'_{n-1} + 4Y'_n + Y'_{n+1}] - \frac{h^5}{90} Y^{(5)}(\xi_n)$$

برای مقداری چون $x_{n-1} \leq \xi_n \leq x_{n+1}$. این تقریب بر پایهٔ انتگرالگیری چندجمله‌یی درجهٔ دوم که $Y' = f(t, Y(t))$ را در نقاط گرهی x_{n-1}, x_n, x_{n+1} درونیابی می‌نماید، قرار دارد. برای سادگی نمادگذاری $Y_j = Y(x_j)$ و $Y'_j = Y'(x_j)$ را به کار می‌بریم. از حذف جملهٔ خطا، فرمول ضمنی مرتبهٔ چهار زیر به دست می‌آید

$$y_{n+1} = y_{n-1} + \frac{h}{3} [f(x_{n-1}, y_{n-1}) + 4f(x_n, y_n) + f(x_{n+1}, y_{n+1})] \quad n \geq 1 \quad (۱۳.۷.۶)$$

این فرمول یک فرمول مشهور و تصحیح‌کنندهٔ یک الگوریتم پیشگو-تصحیح‌کنندهٔ کلاسیک است که به روش میلن^۱ معروف است. بنابر قضیهٔ ۶.۶ این روش همگراست و

$$\max_{x_n \leq x_n \leq b} |Y(x_n) - y_n| = O(h^4)$$

ولی این روش، به همان شکلی که برای روش نقطهٔ میانی در بخش ۴.۶ بیان شد پایداری ضعیفی دارد.

روشهای آدامز این روشها، روشهای چند گامی هستند که بیشترین کاربرد را دارند، و برای تولید الگوریتمهای پیشگو-تصحیح کننده‌ای به کار می‌روند که در آنها خطا، با دو تغییر، یکی طول گام h و دیگری مرتبه روش، کنترل می‌شود. این موضوع مفصلاً در این بخش بحث شده و یک مثال عددی هم داده شده است.

برای به دست آوردن این روشها، از صورت انتگرالی زیر استفاده می‌کنیم

$$Y_{n+1} = Y_n + \int_{x_n}^{x_{n+1}} f(t, Y(t)) dt \quad (۱۴.۷.۶)$$

چند جمله‌بیهایی که $Y'(t) = f(t, Y(t))$ را درونیایی می‌کنند ساخته می‌شوند و سپس از آنها بر بازه $[x_n, x_{n+1}]$ انتگرالگیری می‌شود تا تقریبی برای Y_{n+1} به دست آید. مطلب را از فرمولهای صریح یا پیشگو آغاز می‌کنیم.

حالت ۱. روشهای آدامز-بشفورت^۱ گیریم $P_p(t)$ یک چند جمله‌یی از درجه نایزگتر از P باشد که $Y'(t)$ را در x_n, \dots, x_{n-p} درونیایی نماید. مناسبترین شکل برای $P_p(t)$ فرمول تفاضلی پسرو نیوتن، (۱۱.۳.۳) است که حول x_n بسط داده شده است:

$$P_p(t) = Y'_n + \frac{(t-x_n)}{h} \nabla Y'_n + \frac{(t-x_n)(t-x_{n-1})}{2!h^2} \nabla^2 Y'_n + \dots + \frac{(t-x_n)\dots(t-x_{n-p+1})}{p!h^p} \nabla^p Y'_n \quad (۱۵.۷.۶)$$

که در آن $\nabla Y'_n = Y'(x_n) - Y'(x_{n-1})$ ، و تفاضلات پسرو از مراتب بالاتر به طور مشابه تعریف می‌شوند [۱۰.۳.۳] را ببینید. خطای درونیایی با رابطه زیر داده می‌شود

$$E_p(t) = (t-x_{n-p})\dots(t-x_n) Y'[x_{n-p}, \dots, x_n, t] = \frac{(t-x_{n-p})\dots(t-x_n)}{(p+1)!} Y^{(p+2)}(\zeta_t) \quad x_{n-p} \leq \zeta_t \leq x_{n+1} \quad (۱۶.۷.۶)$$

به شرطی که $x_{n-p} \leq t \leq x_{n+1}$ ، و $Y(t)$ ، $(p+2)$ بار پیوسته مشتقپذیر باشد. فرمولهای (۸.۲.۳) و (۱۱.۱.۳) از فصل ۳ را برای اثبات درستی (۱۶.۷.۶) ببینید.

انتگرال $P_p(t)$ چنین داده شده است:

$$\int_{x_n}^{x_{n+1}} P_p(t) dt = h Y'_n + \sum_{j=1}^p \frac{1}{j! h^j} \nabla^j Y'_n \int_{x_n}^{x_{n+1}} (t-x_n)\dots(t-x_{n+1-j}) dt = h \sum_{j=0}^p \gamma_j \nabla^j Y'_n \quad (۱۷.۷.۶)$$

جدول ۱۰.۶ ضرایب آدامز-بشورت

γ_0	γ_1	γ_2	γ_3	γ_4	γ_5
۱	$\frac{1}{2}$	$\frac{5}{12}$	$\frac{3}{8}$	$\frac{251}{720}$	$\frac{95}{288}$

ضرایب γ_j با تبدیل متغیر $s = (t - x_n)/h$ به صورت زیر به دست می‌آیند

$$\gamma_j = \frac{1}{j!} \int_0^1 s(s+1)\dots(s+j-1)ds \quad j \geq 1 \quad (18.7.6)$$

که γ_0 برابر ۱ اختیار می‌شود.

جدول ۱۰.۶ شامل چند مقدار اول γ_j است. گی‌یر^۱ (۱۹۷۱، صص ۱۰۴-۱۱۱) شامل اطلاعات بیشتری، از جمله یک تابع مولد برای ضرایب است.

خطای برشی در به‌کار بردن چند جمله‌ی درونیاب $P_p(t)$ در (۱۴.۷.۶) با رابطه زیر داده می‌شود

$$T_n(Y) = \int_{x_n}^{x_{n+1}} E_p(t) dt = \int_{x_n}^{x_{n+1}} (t-x_n)\dots(t-x_{n-p}) Y' [x_{n-p}, \dots, x_n, t] dt$$

با فرض اینکه $Y(x)$ به ازای $x_n \leq x \leq x_{n+1}$ $p+2$ بار پیوسته مشتق‌پذیر است، از قضیه ۳.۳، فصل ۳ نتیجه می‌شود که تفاضل منقسم در انتگرال اخیر یک تابع پیوسته از t است. چون چند جمله‌ی $(t-x_n)\dots(t-x_{n-p})$ در $[x_n, x_{n+1}]$ نامنفی است، می‌توان قضیه مقدار میانگین در انتگرالها (قضیه ۳.۱) را به‌کار برد و برای مقداری چون $x_{n-p} \leq \zeta \leq x_{n+1}$ به دست آورد

$$T_n(Y) = Y' [x_{n-p}, \dots, x_n, \zeta] \int_{x_n}^{x_{n+1}} (t-x_n)\dots(t-x_{n-p}) dt$$

رابطه (۱۸.۷.۶) را برای محاسبه انتگرال و (۱۲.۲.۳) را برای تبدیل تفاضل منقسم به‌کار می‌بریم. پس

$$T_n(Y) = \gamma_{p+1} h^{p+2} Y^{(p+2)}(\xi_n) \quad x_{n-p} \leq \xi_n \leq x_{n+1} \quad (19.7.6)$$

یک شکل دیگری وجود دارد که برای برآورد خطای برشی بسیار مفید است. با به‌کار بردن

قضیه مقدار میانگین (قضیه ۲.۱) و لم مشابه تفاضلات پس‌رواز لم ۲ از بخش ۴.۳ داریم:

$$T_n(Y) = h \gamma_{p+1} \nabla^{p+1} Y'_n + O(h^{p+2}) \quad (20.7.6)$$

قسمت اصلی این خطا، $h \gamma_{p+1} \nabla^{p+1} Y'_n$ ، آخرین جمله می‌شد اگر به جای $P_p(t)$ چند جمله‌ی درونیاب $P_{p+1}(t)$ به‌کار برده می‌شد. این شکل خطای یک ابزار اصلی در الگوریتمهایی است که مرتبه روش را برای کنترل خطای برشی، تغییر می‌دهد.

جدول ۱۱.۶ فرمولهای آدامز-بشفورت

$p = 0$	$Y_{n+1} = Y_n + hY'_n + \frac{1}{2}h^2Y''(\xi_n)$
$p = 1$	$Y_{n+1} = Y_n + \frac{h}{2}[3Y'_n - Y'_{n-1}] + \frac{5}{12}h^2Y^{(2)}(\xi_n)$
$p = 2$	$Y_{n+1} = Y_n + \frac{h}{12}[23Y'_n - 16Y'_{n-1} + 5Y'_{n-2}] + \frac{3}{8}h^3Y^{(3)}(\xi_n)$
$p = 3$	$Y_{n+1} = Y_n + \frac{h}{24}[55Y'_n - 59Y'_{n-1} + 37Y'_{n-2} - 9Y'_{n-3}] + \frac{251}{720}h^4Y^{(4)}(\xi_n)$

با استفاده از معادلات (۱۷.۷.۶) و (۱۹.۷.۶)، معادله (۱۴.۷.۶) به ازای یک $x_{n-p} \leq$

$\xi_n \leq x_{n+1}$ چنین می شود

$$Y_{n+1} = Y_n + h \sum_{j=0}^p \gamma_j \nabla^j Y'_n + \gamma_{p+1} h^{p+2} Y^{(p+2)}(\xi_n) \quad (21.7.6)$$

روش عددی متناظر آن چنین است

$$y_{n+1} = y_n + h \sum_{j=0}^p \gamma_j \nabla^j y'_n \quad n \geq p \quad (22.7.6)$$

در فرمول فوق، $y'_j \equiv f(x_j, y_j)$ ، و به عنوان یک مثال از تفاضلات پسرو داریم $\nabla y'_j = y'_j - y'_{j-1}$. جدول ۱۱.۶ شامل این فرمولها برای $p = 0, 1, 2, 3$ است. این معادلات به شکل معمولتر (۱.۷.۶) نوشته شده اند، که در آنها بستگی به مقدار $f(x_{n-j}, y_{n-j})$ به طور صریح نشان داده شده است. توجه کنید که حالت $p = 0$ همان روش اوایلر است.

فرمولهای آدامز-بشفورت در فرضهای قضیه ۶.۶ صدق می نمایند، بنابراین روشهایی همگرا و پایدارند. گذشته از آن، ناپایداری از نوع موجود در روش میانگامی (۲.۴.۶) و روش (۱۳.۷.۶) سیمپسون را ندارند. اثبات این مطلب در بخش بعد داده شده است و بحث بیشتر تا نزدیک به انتهای آن بخش، هنگام معرفی مفهوم پایداری نسبی به تعویق انداخته شده است.

حالت ۲. روشهای آدامز-مولتن^۱ بازم فرمول انتگرالی (۱۴.۷.۶) را به کار می بریم ولی $Y'(t) = f(t, Y(t))$ را در $p+1$ نقطه $x_{n+1}, \dots, x_{n-p+1}$ ، به ازای $p \geq 0$ ، درونیایی می نماییم. نحوه عمل دقیقاً مانند روشهای آدامز-بشفورت است، و ما فقط نتایج نهایی را می آوریم. معادله (۱۴.۷.۶) به معادله زیر تبدیل می شود

$$Y_{n+1} = Y_n + h \sum_{j=0}^p \delta_j \nabla^j Y'_{n+1} + \delta_{p+1} h^{p+2} Y^{(p+2)}(\xi_n) \quad (23.7.6)$$

جدول ۱۲.۶ ضرایب آدامز-مولتن

δ_0	δ_1	δ_2	δ_3	δ_4	δ_5
۱	$-\frac{1}{2}$	$-\frac{1}{12}$	$-\frac{1}{24}$	$-\frac{19}{720}$	$-\frac{3}{160}$

جدول ۱۳.۶ فرمولهای آدامز-مولتن

$p = 0$	$Y_{n+1} = Y_n + hY'_{n+1} - \frac{1}{2}h^2Y''(\zeta_n)$
$p = 1$	$Y_{n+1} = Y_n + \frac{h}{2}[Y'_{n+1} + Y'_n] - \frac{1}{12}h^2Y^{(2)}(\zeta_n)$
$p = 2$	$Y_{n+1} = Y_n + \frac{h}{12}[\delta Y'_{n+1} + 8Y'_n - Y'_{n-1}] - \frac{1}{24}h^2Y^{(2)}(\zeta_n)$
$p = 3$	$Y_{n+1} = Y_n + \frac{h}{24}[9Y'_{n+1} + 19Y'_n - 5Y'_{n-1} + Y'_{n-2}] - \frac{19}{720}h^3Y^{(3)}(\zeta_n)$

با $x_{n-p+1} \leq \zeta_n \leq x_{n+1}$ ضرایب δ_j با رابطه زیر تعریف شده‌اند

$$\delta_j = \frac{1}{j!} \int_0^1 (s-1)s(s+1)\dots(s+j-2)ds \quad j \geq 1 \quad (24.7.6)$$

با $\delta_0 = 1$ ، و چند مقدار آنها در جدول ۱۲.۶ داده شده‌اند. خطای برشی را درست مانند (۲۰.۷.۶) می‌توان به شکل زیر نوشت

$$T_{n+1}(Y) = h\delta_{p+1} \nabla^{p+1} Y'_{n+1} + O(h^{p+2}) \quad (25.7.6)$$

روش عددی متناظر با (۲۳.۷.۶) چنین است

$$y_{n+1} = y_n + h \sum_{j=0}^p \delta_j \nabla^j y'_{n+1} \quad n \geq p-1 \quad (26.7.6)$$

که در آن، مانند قبل، $y'_j \equiv f(x_j, y_j)$ جدول ۱۳.۶ شامل چند فرمول مرتبه پایین برای $p = 0, 1, 2, 3$ همان روش ذوزنقه‌بی است. توجه کنید که حالت $p = 1$ روش ذوزنقه‌بی است.

فرمول (۲۶.۷.۶) یک روش ضمنی است، و بنابراین برای حل آن با روش بارستی یک پیشگو لازم است. مفاهیم اصلی در اینجا درست همان مفاهیم موجود در بخش ۵.۶ برای حل بارستی روش ذوزنقه‌بی است. اگر یک الگوریتم پیشگو-تصحیح‌کننده با مرتبه ثابت خواسته شده باشد، و اگر فقط یک بارست باید حساب شود، آنگاه در یک فرمول آدامز-مولتن از مرتبه $m \geq 2$ می‌توان یک پیشگو از مرتبه m یا $m-1$ به‌کار برد. مزیت به‌کارگیری پیشگوی مرتبه $m-1$ آن است که پیشگو و تصحیح‌کننده هر دو از مقادیر مشتق در نقاط گرهی $x_{n-m+2}, \dots, x_{n-1}, x_n$

استفاده می‌کنند. برای مثال، فرمول مرتبهٔ دوم آدامز-مولتن با فرمول مرتبهٔ یک آدامز-بشفورت به‌عنوان پیشگو، درست همان روش ذوزنقه‌یی با پیشگویی اوایلر است. این موضوع در بخش ۵.۶ مورد بحث واقع و نشان داده شده که مناسب است؛ در هر دو روش تنها از یک مقدار قبلی مشتق، $f(x_n, y_n)$ استفاده می‌شود.

یک مثال جالبتر، روش مرتبه چهار زیر است؛

$$y_{n+1}^{(0)} = y_n + \frac{h}{12} [23f(x_n, y_n) - 16f(x_{n-1}, y_{n-1}) + 5f(x_{n-2}, y_{n-2})]$$

$$y_{n+1}^{(j+1)} = y_n + \frac{h}{24} [9f(x_{n+1}, y_{n+1}^{(j)}) + 19f(x_n, y_n) - 5f(x_{n-1}, y_{n-1}) + f(x_{n-2}, y_{n-2})]$$

(۲۷.۷.۶)

معمولاً فقط یک بارست محاسبه می‌شود، گرچه این امر، شکل خطای برشی در (۲۳.۷.۶) را تغییر می‌دهد. گیریم $u_n(x)$ جواب $y' = f(x, y)$ گذرنده بر نقطهٔ (x_n, y_n) باشد. پس خطای برشی در به‌کاربردن تقریب $y_{n+1}^{(1)}$ عبارتست از،

$$u_n(x_{n+1}) - y_{n+1}^{(1)} = [u_n(x_{n+1}) - y_{n+1}] + [y_{n+1} - y_{n+1}^{(1)}]$$

با استفاده از (۲۳.۷.۶) و بسط خطای بارستی،

$$u_n(x_{n+1}) - y_{n+1}^{(1)} = \delta_4 h^5 u_n^{(5)}(x_n) + \frac{3h}{8} f_y(x_n, y_n)(y_{n+1} - y_{n+1}^{(0)}) + O(h^6)$$

(۲۸.۷.۶)

دو جملهٔ اول طرف راست، پس از علامت تساوی، از مرتبهٔ h^5 اند. اگر یا (۱) بارستهای بیشتری محاسبه شود، یا (۲) پیشگویی از مرتبهٔ بالاتری به‌کار گرفته شود، آنگاه قسمت اصلی خطای برشی فقط $\delta_4 h^5 u_n^{(5)}(x_n)$ خواهد شد. و این از نظر برآورد خطا، حالت مطلوبتری است.

برای روشهای از مرتبهٔ مشابه، خطای برشی فرمولهای آدامز-مولتن به مراتب کوچکتر از خطای آدامز-بشفورت است. برای مثال، خطای برشی فرمول مرتبهٔ چهار آدامز-مولتن 0.076×10^{-6} خطای برشی فرمول آدامز-بشفورت مرتبهٔ چهار است. این است دلیل عمدهٔ استفاده از روشهای ضمنی، اگرچه، ملاحظات دیگری نیز مطرح است. همچنین توجه کنید که خطای برشی فرمول مرتبهٔ چهار آدامز-مولتن بیش از دو برابر خطای برشی روش سیمپسون است. دلیل استفاده از فرمول آدامز-مولتن این است که ویژگیهای پایداری آن خیلی بهتر از ویژگیهای پایداری روش سیمپسون است.

جدول ۱۴.۶ مثال عددی برای روش آدامز-مولتن

x	خطا برای	خطا برای	نسبت
	$h = 0.125$	$h = 0.0625$	
۲.۰	$2.07E-5$	$1.21E-6$	۱۷.۱
۴.۰	$2.21E-6$	$1.20E-7$	۱۸.۳
۶.۰	$3.74E-7$	$2.00E-8$	۱۸.۷
۸.۰	$1.00E-7$	$5.24E-9$	۱۹.۱
۱۰.۰	$3.58E-8$	$1.83E-9$	۱۹.۶

مثال روش (۲۷.۷.۶) برای حل معادله

$$y' = \frac{1}{1+x^2} - 2y^2 \quad y(0) = 0 \quad (29.7.6)$$

که دارای جواب $Y(x) = x/(1+x^2)$ است، به کار رفته است. مقادیر اولیه y_1 ، y_2 و y_3 را برای ساده کردن مثال، مقادیر درست انتخاب کرده اند. مقادیر جوابها با دو مقدار h محاسبه شده اند، و خطاهای حاصل در چند نقطه گرهی در جدول ۱۴.۶ داده شده اند. ستونی که با نسبت مشخص شده است نسبت خطا با $h = 0.125$ به خطا در $h = 0.0625$ است. توجه نمایید که مقادیر نسبت نزدیک به ۱۶ هستند، که باید انتظار می داشتیم، زیرا روش (۲۷.۷.۶) از مرتبه چهار است.

روشهای متغیر-مرتبه در حال حاضر مشهورترین الگوریتم های پیشگو-اصلاحگر خطای برشی را هم با تغییر طول گام و هم با تغییر مرتبه روش کنترل می نمایند، و در تمام این الگوریتم ها از فرمولهای خانواده آدامز استفاده می کنند. اولین برنامه های رایانه ای از این نوع، که به طور گسترده مورد استفاده قرار گرفته بود DIFSUB از گی پر (۱۹۷۱، صص ۱۵۸-۱۶۶) و برنامه های کروا بودند. پس از آن، چنین برنامه های متغیر-مرتبه آدامز نوشته شدند که، با GEAR از هایندرمارش^۲ (۱۹۷۴) و با DE/STEP از شمپاین و گوردن (۱۹۷۵) از متداولترین آنها هستند. برنامه GEAR باز هم بهبود یافت و به برنامه LSODE تبدیل شد که قسمتی از یک بسته برنامه بزرگی است که ODEPACK خوانده می شود [هایندمارش (۱۹۸۳) را برای توضیح ببینید]. برنامه DE/STEP بهبود بیشتری یافته و به برنامه DDEABM تبدیل شد و قسمتی از یک بسته کلی دیگری به نام DEPAC است [برای توضیح به شمپاین و واتس^۳ (۱۹۸۰) مراجعه کنید]. در همه حالات، برنامه هایی که قبلاً گفتیم برآورد خطایی را به کار می برند که بر پایه فرمولهای (۲۳.۷.۶) و (۲۵.۷.۶) استوار است، اگرچه این

فرمولها [و فرمولهای (۲۲.۷.۶) و (۲۶.۷.۶)] برای طول گام متغیر، ممکن است نیاز به اصلاح داشته باشند. برنامه‌ها از بعضی جنبه‌های مهم ولی فنی تغییر می‌کنند، از جمله تعداد بارستهای اصلاحگر آدامز-مولتن و شکلی که در آن اطلاعات گذشته درباره مشتق f پیش کشیده می‌شود و شکل درونیابی جواب برای نقطه گرهی خروجی جاری. به علت کمبود جا، و پیچیدگی موضوعات وارده، بیش از این در موضوع تفاوت‌های این برنامه‌ها بحث نمی‌کنیم. برای بعضی ملاحظات دیگر در مورد این برنامه‌ها گویتا^۱ و همکاران (۱۹۸۵، صص ۱۶-۱۹) را ببینید.

با مجاز دانستن تغییر مرتبه، هیچ مشکلی برای پیدا کردن مقادیر اولیه روشهای مرتبه بالاتر آدامز نخواهد بود. برنامه‌ها با فرمول مرتبه دوم دوزنقه‌یی با پیشگوی اوپلر آغاز می‌شوند، آنگاه با فراهم آمدن مقادیر اولیه افزایش می‌یابند. اگر جواب به سرعت تغییر یابد، آنگاه برنامه معمولاً فرمولی از مرتبه پایین انتخاب می‌کند، حال آنکه برای یک جواب با تغییرات آهسته‌تر و هموارتر، مرتبه معمولاً بالاتر است. در برنامه DE شمپاین و گوردن (۱۹۸۵، صص ۱۸۶-۲۰۹)، خطای برشی در x_{n+1} که $trunc$ خوانده شده، باید در رابطه زیر صدق کند

$$|trunc_j| \leq ABSERR + RELERR * |y_{n,j}| \quad (۳۰.۷.۶)$$

این رابطه باید برای هر مؤلفه خطای برشی و هر $y_{n,j}$ ، مؤلفه متناظر جواب y_n دستگاه معادلات دینفرانسیل داده شده برقرار باشد. مقادیر ABSERR و RELERR را کاربر می‌دهد. مقدار $trunc$ با فرض یکنواخت بودن فاصله‌ها، تقریباً با

$$trunc \approx h \delta_{p+1} \nabla^{p+1} y'_{n+1}$$

داده می‌شود. این عبارت، خطای برشی فرمول p مرحله‌یی زیر است

$$y_{n+1}^{(p)} = y_n + h \sum_{j=0}^p \delta_j \nabla^j y'_{n+1}$$

وقتی $trunc$ در آزمون (۳۰.۷.۶) صدق کند، مقدار y_{n+1} چنین می‌شود

$$y_{n+1} \equiv y_{n+1}^{(p+1)} = y_{n+1}^{(p)} + trunc \quad (۳۱.۷.۶)$$

بنابراین خطای برشی واقعی $O(h^{p+3})$ است، و چنانچه با (۳۰.۷.۶) ترکیب شود، می‌توان نشان داد که خطای برشی y_{n+1} در یک محک خطا در گام واحد صدق می‌نماید، که مشابه خطای

برشی است که با (۱.۶.۶) برای الگوریتم Detrap در بخش ۶.۶ بیان شده است. برای یک بحث مفصل شمایین و گوردن (۱۹۷۵، ص ۱۰۰) را ببینید.

برنامه DE (و جانشین آن DDEABM) در کنترل خطا، از جمله انتخاب مرتبه و طول گام، خیلی پیچیده است. نمی توان از آن به قدر کافی در این جای محدود، در این کتاب بحث کرد ولی بهترین مرجع، کتاب شمایین و گوردن (۱۹۷۵) است که به الگوریتمهای متغیر-مرتبه آدامز اختصاص یافته است. برنامه های DE و DDEABM از لحاظ کنترل خطا و راحتی کاربر خوب طراحی شده اند. هر کدام به شکل قابل حمل هم نوشته شده است، و به طور کلی برنامه ای است که برای حل معادلات دیفرانسیل خیلی توصیه شده است.

مثال مسأله زیر را در نظر می گیریم

$$y' = \frac{y}{4} \left(1 - \frac{y}{20} \right) \quad y(0) = 1 \quad (32.7.6)$$

که دارای جواب زیر است

$$Y(x) = \frac{20}{(1 + 19e^{-(x/4)})}$$

برای حل این مسأله DDEABM با مقادیر خروجی در $x = 2, 4, 6, \dots, 20$ به کار برده شده است. سه مقدار ABSERR به کار برده شده و در هر حالت RELERR=0. خطاهای کلی واقعی در جدول ۱۵.۶ داده شده اند. ستونی که با NFE مشخص شده تعداد محاسبات $f(x, y)$ را که از x_0 شروع می شود و برای به دست آوردن $y_h(x)$ لازم است، نشان می دهد.

خطای کلی (فراگیر) برنامه های رایانه ای خودکار که قبلاً توضیح داده شدند خطای موضعی یا خطای برشی را کنترل می کنند. ولی خطای کلی جواب را کنترل نمی کنند. معمولاً خطای برشی

جدول ۱۵.۶ مثال برنامه خودکار DDEABM

x	ABSERR = 10^{-2}		ABSERR = 10^{-6}		ABSERR = 10^{-9}	
	خطا	NFE	خطا	NFE	خطا	NFE
4_r	$-3.26E-5$	۱۵	$۱.24E-6$	۲۸	$2.86E-10$	۵۲
8_r	$6.00E-4$	۲۱	$3.86E-6$	۴۲	$-1.98E-9$	۷۶
12_r	$۱.70E-3$	۲۵	$4.93E-6$	۵۴	$-2.41E-9$	۱۰۲
16_r	$9.13E-4$	۳۱	$3.73E-6$	۶۴	$-1.86E-9$	۱۲۴
20_r	$9.16E-4$	۳۷	$۱.79E-6$	۷۴	$-9.58E-10$	۱۳۸

در این برنامه‌ها آنقدر کوچک نگهداشته می‌شوند که خطای کلی در یک محدوده، قابل قبول است، گرچه این امر تضمین نمی‌شود. دلایل برای این خطای کلی کوچک خیلی شبیه دلایلی است که در بخش ۶.۶ تشریح شده است؛ به‌ویژه (۱۵.۶.۶) و (۱۶.۶.۶) را به یاد آورید.

خطای کلی را می‌توان کنترل کرد، ما مثالی در ذیل می‌آوریم. ولی حتی با یک برآورد خطای کلی، نمی‌توانیم آن را برای بسیاری از معادلات کنترل نماییم. علت آن این است که خطای کلی ترکیب شده است از اثرات همه خطاهای برشی گذشته، و کاهش طول گام در مرحله جاری خطاهای مراحل گذشته را تغییر نخواهد داد. معمولاً اگر خطای کلی بسیار بزرگ باشد، معادله باید دوباره با طول گام کوچکتری حل شود.

برای نمایش خطای کلی چندین روش ارائه شده است. یکی از آنها، برونیابی ریچاردسن، در بخش ۱۰.۶ برای روش رونگه-کوتا توضیح داده شده است. در زیر روش دیگری برای روش بخش ۶.۶ ارائه خواهیم داد. برای یک بازنگری کلی از موضوع به اسکیل^۱ (۱۹۸۶) مراجعه کنید. برای روش ذوزنقه‌یی، جواب درست $Y(x)$ در

$$Y(x_{n+1}) = Y(x_n) + \frac{h}{4} [f(x_n, Y(x_n)) + f(x_{n+1}, Y(x_{n+1}))] - \frac{h^3}{12} Y^{(3)}(\xi_n)$$

صدق می‌کند که $h = x_{n+1} - x_n$ و $x_n \leq \xi_n \leq x_{n+1}$. اگر قاعده ذوزنقه‌یی

$$y_{n+1} = y_n + \frac{h}{4} [f(x_n, y_n) + f(x_{n+1}, y_{n+1})]$$

را از آن کم کنیم خواهیم داشت

$$e_{n+1} = e_n + \frac{h}{4} \{ [f(x_n, y_n + e_n) - f(x_n, y_n)] + [f(x_{n+1}, y_{n+1} + e_{n+1}) - f(x_{n+1}, y_{n+1})] \} - \frac{h^3}{12} Y^{(3)}(\xi_n) \quad (33.7.6)$$

با فرض $e_n = Y(x_n) - y_n$ ، $n \geq 0$. این معادله، معادله خطا برای روش ذوزنقه‌یی است و ما سعی می‌کنیم آن را به‌طور تقریبی برای محاسبه e_{n+1} حل کنیم.

به الگوریتم Detrap از بخش ۶.۶ بازمی‌گردیم، به‌جای عبارت برشی در (۳۳.۷.۶) متغیر Trunc را که در Detrap محاسبه شده می‌گذاریم. سپس (۳۳.۷.۶) را نسبت به \hat{e}_{n+1} که تقریبی برای خطای کلی درست e_{n+1} است، حل می‌کنیم. می‌توانیم این معادله را نسبت به \hat{e}_{n+1} با روشهای

ریشه‌یابی گوناگون حل کنیم، ما بارست‌های نقطه - ثابت ساده را به‌ازای $z \geq 0$ به‌کار می‌بریم

$$\hat{e}_{n+1}^{(j+1)} = \hat{e}_n + \frac{h}{\gamma} \{ [f(x_n, y_n + \hat{e}_n) - f(x_n, y_n)] + [f(x_{n+1}, y_{n+1} + \hat{e}_{n+1}^{(j)}) - f(x_{n+1}, y_{n+1})] \} + \text{trunc} \quad (34.7.6)$$

$\hat{e}_{n+1}^{(0)} = \hat{e}_n$ را به‌کار می‌بریم و چون این فقط به‌منظورهای توضیحی است، در (34.7.6) چندین بارست را انجام می‌دهیم. این فکر ساده خیلی به روشهای تصحیح تفاضلی اسکیل (۱۹۸۶) وابسته است.

مثال محاسبات جدول ۹.۶ برای Detrap در حل معادله (۹.۶.۶) را تکرار می‌کنیم. ما همان پارامترها را برای Detrap به‌کار می‌بریم. نتایج در جدول ۱۶.۶ برای همان مقادیر x_n در جدول ۹.۶ داده شده‌اند. این نتایج نشان می‌دهند که e_n و \hat{e}_n تقریباً همیشه به اندازه معقولی، البته از نظر قدر مطلق، بهم نزدیک‌اند. تقریب $e_n \doteq \hat{e}_n$ در نزدیکی $x = 0.5$ ضعیف است و این به دلیل ضعف برآورد خطای برشی در Detrap است. حتی این نتایج ضعیف هم در این مسأله از بین می‌روند و برای مقادیر بزرگتر x_n ، تقریب $e_n \doteq \hat{e}_n$ هنوز سودمند است.

۸.۶ نظریه همگرایی و پایداری برای روشهای چندگامی

در این بخش یک نظریه کامل همگرایی و پایداری برای روش چندگامی زیر ارائه شده است

$$y_{n+1} = \sum_{j=0}^p a_j y_{n-j} + h \sum_{j=-1}^p b_j f(x_{n-j}, y_{n-j}) \quad x_{p+1} \leq x_{n+1} \leq b \quad (1.8.6)$$

در این بررسی، عملیات بخش ۳.۶ تعمیم داده خواهد شد، و ابزار ریاضی لازم به‌وجود خواهد آمد تا بررسی شود که آیا ضعیف - پایداری روش (۱.۸.۶)، فقط به دلیل ناپایداری از نوع مرتبط با روش میانگامی است یا به دلیل دیگر.

مطلب را با چند تعریف آغاز می‌کنیم. مفهوم پایداری که با روش اوایل معرفی شد [۲۸.۲.۶] و (۲۰.۲.۶) را ببینید]، اکنون تعمیم می‌یابد. گیریم $\{y_n \mid 0 \leq n \leq N(h)\}$ جواب (۱.۸.۶) برای معادله دیفرانسیلی چون $y' = f(x, y)$ برای همه مقادیر به اندازه کافی کوچک h ، $h \leq h_0$ باشد. یادآور می‌شویم که $N(h)$ بزرگترین زیرنمایه N است که برای آن $x_N \leq b$. به‌ازای هر مقدار $h \leq h_0$ اختلافی در مقادیر اولیه y_0, \dots, y_p به مقادیر جدید z_0, \dots, z_p به‌وجود می‌آوریم

جدول ۱۶.۶ محاسبه خطای کلی برای Detrap

x_n	h	e_n	\hat{e}_n	trunc
۰٫۰۲۲۷	۰٫۰۲۲۷	۵٫۸۴E-۶	۵٫۸۳E-۶	۵٫۸۴E-۶
۰٫۰۴۵۴	۰٫۰۲۲۷	۱٫۱۷E-۵	۱٫۱۶E-۵	۵٫۸۳E-۶
۰٫۰۶۸۱	۰٫۰۲۲۷	۱٫۷۴E-۵	۱٫۷۳E-۵	۵٫۷۶E-۶
۰٫۰۹۰۸	۰٫۰۲۲۷	۲٫۲۸E-۵	۲٫۲۸E-۵	۵٫۶۲E-۶
۰٫۱۲۲۵	۰٫۰۲۲۷	۵٫۱۶E-۵	۵٫۱۹E-۵	۲٫۹۶E-۶
۰٫۱۳۰۶۵	۰٫۰۳۴۰	۵٫۶۶E-۵	۵٫۶۶E-۵	۶٫۷۴E-۶
۰٫۱۳۴۰۵	۰٫۰۳۴۰	۶٫۰۱E-۵	۶٫۰۵E-۵	۶٫۲۱E-۶
۰٫۱۳۷۴۶	۰٫۰۳۴۰	۶٫۱۱E-۵	۶٫۲۱E-۵	۴٫۲۸E-۶
۰٫۱۴۴۰۸	۰٫۰۶۶۲	۵٫۰۵E-۵	۵٫۰۴E-۵	-۶٫۵۶E-۶
۰٫۱۵۰۷۰	۰٫۰۶۶۲	۲٫۴۴E-۵	۳٫۵۷E-۵	-۱٫۰۴E-۵
۰٫۱۵۷۳۲	۰٫۰۶۶۲	-۲٫۰۳E-۵	۴٫۳۴E-۶	-۲٫۹۲E-۵
۰٫۱۶۱۳۸	۰٫۰۴۰۶	-۲٫۹۹E-۵	-۶٫۷۶E-۶	-۱٫۱۲E-۵
۱٫۹۵۹۵	۰٫۱۳۵	-۱٫۰۲E-۴	-۸٫۷۶E-۵	-۱٫۶۴E-۵
۲٫۰۹۴۲	۰٫۱۳۵	-۱٫۰۳E-۴	-۸٫۶۸E-۵	-۱٫۷۹E-۵
۲٫۳۱۷۲	۰٫۲۲۳	-۶٫۵۷E-۵	-۵٫۰۸E-۵	۱٫۲۷E-۵
۲٫۷۶۳۲	۰٫۴۴۶	۳٫۴۴E-۴	۳٫۱۷E-۴	۴٫۴۱E-۴
۳٫۰۶۶۴	۰٫۳۰۳	۳٫۲۱E-۴	۲٫۹۶E-۴	۹٫۳۹E-۵
۷٫۶۹۵۹	۰٫۶۷۲	۳٫۸۷E-۴	۲٫۶۹E-۴	۸٫۷۷E-۵
۸٫۶۹۵۹	۱٫۰۰۰	۳٫۹۶E-۴	۳٫۰۵E-۴	۱٫۷۳E-۴
۹٫۶۹۵۹	۱٫۰۰۰	۴٫۲۷E-۴	۲٫۹۴E-۴	۱٫۱۸E-۴
۱۰٫۶۹۵۹	۱٫۰۰۰	۴٫۱۱E-۴	۲٫۷۷E-۴	۹٫۴۵E-۵

به طوری که شرط

$$\text{Max}_{0 \leq n \leq p} |y_n - z_n| \leq \varepsilon \quad 0 < h \leq h_0 \quad (۲.۸.۶)$$

برقرار باشد. توجه کنید که مقادیر آغازی احتمالاً به h بستگی دارند. گوییم یک خانواده جواب $\{y_n \mid 0 \leq n \leq N(h)\}$ پایدار است اگر مقدار ثابتی چون c مستقل از h وجود داشته باشد به طوری که به ازای جميع مقادیر به اندازه کافی کوچک ε داشته باشیم

$$\text{Max}_{0 \leq n \leq N(h)} |y_n - z_n| \leq c\varepsilon \quad 0 < h \leq h_0 \quad (۳.۸.۶)$$

همه مسائل معادله دیفرانسیل

$$y' = f(x, y) \quad y(x_0) = Y_0 \quad (۴.۸.۶)$$

را در نظر می‌گیریم که در آنها $f(x, y)$ مشتق پیوسته باشد و در شرط (۱۲.۲.۶) لیشیتس صدق کند، و فرض می‌کنیم که جوابهای تقریبی $\{y_n\}$ همگی پایدارند. در این صورت می‌گوییم (۱.۸.۶) یک روش عددی پایدار است.

به منظور تعریف همگرایی برای یک مسأله داده‌شده (۴.۸.۶)، فرض می‌کنیم که مقادیر اولیه y_0, \dots, y_p در رابطه زیر صدق می‌کنند

$$\eta(h) \equiv \text{Max}_{0 \leq n \leq p} |Y(x_n) - y_n| \rightarrow 0 \quad \text{وقتی } h \rightarrow 0 \quad (۵.۸.۶)$$

در این صورت گفته می‌شود جواب $\{y_n\}$ به $Y(x)$ همگراست اگر

$$\text{Max}_{x_0 \leq x_n \leq b} |Y(x_n) - y_n| \rightarrow 0 \quad \text{وقتی } h \rightarrow 0 \quad (۶.۸.۶)$$

اگر (۱.۸.۶) برای تمام مسائل (۴.۸.۶) همگرا باشد، آنگاه آن را یک روش عددی همگرا نامند. تعریف سازگاری را که در بخش ۳.۶ داده شده به یاد می‌آوریم. روش (۱.۸.۶) سازگار است اگر برای تمام توابع پیوسته مشتقپذیر $Y(x)$ در $[x_0, b]$ ،

$$\frac{1}{h} \text{Max}_{x_p \leq x_n \leq b} |T_n(Y)| \rightarrow 0 \quad \text{وقتی } h \rightarrow 0$$

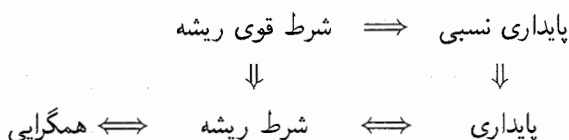
یا هم‌ارز با آن از قضیه ۵.۶، ضرایب $\{a_j\}$ و $\{b_j\}$ باید در روابط زیر صدق کنند

$$\sum_{j=0}^p a_j = 1 \quad - \sum_{j=0}^p j a_j + \sum_{j=-1}^p b_j = 1 \quad (۷.۸.۶)$$

می‌توان نشان داد که همگرایی، سازگاری را ایجاب می‌کند؛ در نتیجه، ما فقط روشهایی را در نظر می‌گیریم که در (۷.۸.۶) صدق می‌کنند. به عنوان یک مثال از برهان لزوم (۷.۸.۶)، فرض همگرایی (۱.۸.۶) برای مسأله

$$y' \equiv 0 \quad y(0) = 1$$

اولین شرط (۷.۸.۶) را ایجاب خواهد کرد. کافی است بگیریم $y_0 = \dots = y_p = 1$ و به نتایج همگرایی y_{p+1} به $Y(x) \equiv 1$ توجه کنیم.



شکل ۶.۶ شرح اجمالی نظریه روشهای چندگامی سازگار

همگرایی و پایداری (۱.۸.۶) به ریشه‌های چندجمله‌یی زیر وابسته‌اند

$$\rho(r) = r^{p+1} - \sum_{j=0}^p a_j r^{p-j} \quad (۸.۸.۶)$$

ملاحظه می‌کنید که با توجه به شرط سازگاری (۷.۸.۶) داریم $\rho(1) = 0$. گیریم r_0, \dots, r_p معرف ریشه‌های $\rho(r)$ باشند، که هر ریشه مطابق چندگانی خود تکرار شده است، و گیریم $r_0 = 1$. روش (۱.۸.۶) در شرط ریشه صدق می‌نماید اگر

۱.

$$|r_j| \leq 1 \quad j = 0, 1, \dots, p \quad (۹.۸.۶)$$

۲.

$$|r_j| = 1 \Rightarrow \rho'(r_j) \neq 0 \quad (۱۰.۸.۶)$$

شرط اول ایجاب می‌کند که تمام ریشه‌های $\rho(r)$ در دایره واحد $\{z : |z| \leq 1\}$ واقع در صفحه مختلط باشند. شرط (۱۰.۸.۶) بیان می‌کند که تمام ریشه‌های واقع بر مرز دایره باید ریشه‌های ساده $\rho(r)$ باشند.

نتایج عمده این بخش در شکل ۶.۶ نشان داده شده‌اند، گرچه بعضی از آنها اثبات نشده‌اند. شرط قوی ریشه و مفهوم پایداری نسبی بعداً در این بخش معرفی خواهد شد.

نظریه پایداری تمام روشهای عددی که در بخشهای قبل داده شده‌اند پایدار بودند. اکنون یک مثال از روش ناپایدار می‌آوریم. این امر به منظور ایجاد انگیزه برای مطالعه نظریه کلی پایداری صورت گرفته است.

مثال فرمول کلی (۱۰.۷.۶) را برای یک روش صریح مرتبه دوم دوگامی به‌یاد آورید و a_0 را ۳

انتخاب کنید، در این صورت روش

$$y_{n+1} = 3y_n - 2y_{n-1} + \frac{h}{4} [f(x_n, y_n) - 3f(x_{n-1}, y_{n-1})] \quad n \geq 1 \quad (11.8.6)$$

را با خطای برشی

$$T_n(Y) = \frac{Y}{12} h^3 Y^{(3)}(\xi_n) \quad x_{n-1} \leq \xi_n \leq x_{n+1}$$

را به دست می آوریم. حل مسئله $y' \equiv 0$ ، $y(0) = 0$ ، را در نظر می گیریم که $Y(x) \equiv 0$ جواب آن است. اگر $z_0 = y_0 = 0$ استفاده کنیم، واضح می شود که $z_n = y_n = 0$ ، $n \geq 0$ ، جواب عددی آن است. اختلالی در داده های آغازی به اندازه $z_0 = \varepsilon/2$ ، $z_1 = \varepsilon$ ، به ازای یک $\varepsilon \neq 0$ ایجاد می نماییم. در این صورت می توان نشان داد که جواب عددی متناظر، عبارت است از

$$z_n = \varepsilon \times 2^{n-1} \quad n \geq 0 \quad (12.8.6)$$

استدلالی که منجر به پیداشدن این جواب شده است بعداً در شرایط کلیتری داده خواهد شد. برای ملاحظه اثر این اختلال در جواب اصلی،

$$\text{Max}_{x_0 \leq x_n \leq b} |y_n - z_n| = \text{Max}_{0 \leq x_n \leq b} |\varepsilon| 2^{n-1} = |\varepsilon| 2^{N(h)-1}$$

از آنجا که وقتی $h \rightarrow 0$ ، $N(h) \rightarrow \infty$ ، انحراف $\{z_n\}$ از $\{y_n\}$ وقتی $h \rightarrow 0$ ، به طور افزایشی زیاد می شود. روش (11.8.6) ناپایدار است، و هرگز نباید مورد استفاده قرار گیرد. همچنین باید توجه کرد که شرط ریشه نیز نقض شده است، زیرا $\rho(r) = r^2 - 3r + 2$ دارای ریشه های $r_0 = 1$ و $r_1 = 2$ است.

برای بررسی پایداری (1.8.6)، فقط معادله خاص

$$y' = \lambda y \quad y(0) = 1 \quad (13.8.6)$$

را در نظر می گیریم، که $Y(x) = e^{\lambda x}$ جواب آن است. نتایج به دست آمده به مطالعه پایداری مسئله معادله دیفرانسیل کلی انتقال می یابد. یک دلیل شهودی، به سادگی به دست می آید. $Y'(x) = f(x, y(x))$ را حول (x_0, y_0) بسط می دهیم تا به دست آوریم

$$\begin{aligned} Y'(x) &\doteq f(x_0, Y_0) + f_x(x_0, Y_0)(x - x_0) + f_y(x_0, Y_0)(Y(x) - Y_0) \\ &= \lambda(Y(x) - Y_0) + g(x) \end{aligned} \quad (14.8.6)$$

که در آن $\lambda = f_y(x_0, Y_0)$ و $g(x) = f(x_0, Y_0) + f_x(x_0, Y_0)(x - x_0)$. اگر $(x - x_0)$ به اندازه کافی کوچک باشد، این تقریب یک تقریب معتبر است. عبارت $V(x) = Y(x) - Y_0$ را وارد می‌کنیم

$$V'(x) \doteq \lambda V(x) + g(x) \quad (15.8.6)$$

جمله ناهمگن $g(x)$ در به دست آوردن پایداری عددی همه جا حذف می‌شود، زیرا ما با تفاضلهای جوابیهای معادله سروکار داریم. حذف $g(x)$ در (15.8.6)، معادله نمونه (13.8.6) را به دست می‌دهد. برای توضیح بیشتر، به نتایج پایداری (5.1.6) تا (10.1.6) و به نتایج پایداری دوزنقه‌یی در (20.5.6) تا (23.5.6) رجوع کنید.

در حالتی که $y' = f(x, y)$ یک دستگاه m معادله دیفرانسیل مانند (13.1.6) را نشان می‌دهد، مشتق جزئی $f_y(x, y)$ ، مانند (54.2.6)، یک ماتریس ژاکوبی می‌شود:

$$[f_y(x, y)]_{ij} = \frac{\partial f_i}{\partial y_j} \quad 1 \leq i, j \leq m$$

پس مسأله نمونه، به یک دستگاه m معادله دیفرانسیل خطی

$$y' = \Lambda y + g(x) \quad (16.8.6)$$

بدل می‌شود که $\Lambda = f_y(x_0, Y_0)$. می‌توان نشان داد که در بیشتر حالات، این دستگاه به دستگاه هم‌ارز زیر بدل می‌شود

$$z'_i = \lambda_i z_i + \gamma_i(x) \quad 1 \leq i \leq m \quad (17.8.6)$$

که $\lambda_1, \dots, \lambda_m$ ویژه مقدرهای ماتریس Λ هستند (مسأله 24 را ببینید). با (17.8.6) ما به معادله نمونه ساده (13.8.6) بازگشته‌ایم، به شرطی که λ عدد مختلط باشد تا بتواند تمام ویژه مقدرهای ممکن Λ را شامل شود.

با به‌کار بردن (1.8.6) برای معادله نمونه (13.8.6)، به دست می‌آوریم

$$y_{n+1} = \sum_{j=0}^p a_j y_{n-j} + h\lambda \sum_{j=-1}^p b_j y_{n-j} \quad (18.8.6)$$

$$(1 - h\lambda b_{-1})y_{n+1} - \sum_{j=0}^p (a_j + h\lambda b_j)y_{n-j} = 0 \quad n \geq p \quad (19.8.6)$$

این معادله یک معادله تفاضلی خطی همگن از مرتبه $p + 1$ است و نظریه حلپذیری آن کاملاً مشابه نظریه معادلات دیفرانسیل خطی همگن مرتبه $p + 1$ است. به عنوان یک مرجع کلی، هنریچی (۱۹۶۲، صص ۲۱۰-۲۱۵)، آیزکسون و کالر (۱۹۶۶، صص ۴۰۵ تا ۴۱۷) را ببینید. برای پیدا کردن جواب عمومی، سعی می‌کنیم ابتدا به دنبال جوابهای خاصی به شکل

$$y_n = r^n \quad n \geq 0$$

بگردیم. اگر بتوانیم $p + 1$ جواب مستقل خطی پیدا کنیم، آنگاه هر ترکیب خطی یک جواب عمومی (۱۹.۸.۶) خواهد بود.

با گذاردن $y_n = r^n$ در (۱۹.۸.۶) و حذف r^{n-p} ، به دست می‌آوریم

$$(1 - h\lambda b_{-1})r^{p+1} - \sum_{j=0}^p (a_j + h\lambda b_j)r^{p-j} = 0 \quad (20.8.6)$$

این معادله، معادله مشخصه نامیده می‌شود، و سمت چپ آن چندجمله‌یی مشخصه است. ریشه‌ها را ریشه‌های مشخصه نامند. تعریف می‌کنیم

$$\sigma(r) = b_{-1}r^{p+1} + \sum_{j=0}^p b_j r^{p-j}$$

و تعریف (۸.۸.۶) را برای $\rho(r)$ را به یاد می‌آوریم. معادله (۲۰.۸.۶) به معادله زیر تبدیل می‌شود

$$\rho(r) - h\lambda\sigma(r) = 0 \quad (21.8.6)$$

ریشه‌های مشخصه را با

$$r_0(h\lambda), r_1(h\lambda), \dots, r_p(h\lambda)$$

نشان می‌دهیم که می‌توان نشان داد به‌طور پیوسته به مقدار $h\lambda$ بستگی دارند. وقتی $h\lambda = 0$ ، معادله (۲۱.۸.۶) به صورت $\rho(r) = 0$ در می‌آید، برای r_j ، ریشه‌های قبلی $\rho(r) = 0$ داریم $r_j(0) = r_j$ ، $j = 0, 1, \dots, p$. چون $r_0 = 1$ یک ریشه $\rho(r)$ است، گیریم که $r_0(h\lambda)$ ریشه (۲۱.۸.۶) باشد که برای آن $r_0(0) = 1$. به دلایلی که بعداً روشن خواهد شد، ریشه $r_0(h\lambda)$ ، ریشه اصلی خوانده می‌شود. اگر ریشه‌های $r_j(h\lambda)$ همگی متمایز باشند، جواب عمومی (۱۹.۸.۶) چنین است

$$y_n = \sum_{j=0}^p \gamma_j [r_j(h\lambda)]^n \quad n \geq 0 \quad (22.8.6)$$

ولی اگر $r_j(h\lambda)$ یک ریشه با چندگانگی $\nu > 1$ باشد، آنگاه ν جواب مستقل خطی (۱۹.۸.۶) به صورت زیرند:

$$\{[r_j(h\lambda)]^n\}, \{n[r_j(h\lambda)]^{n-1}\}, \dots, \{n^{\nu-1}[r_j(h\lambda)]^{n-\nu}\}$$

با این جوابها همراه با جوابهایی که از سایر ریشهها به دست می آیند، می توان جواب عمومی (۱۹.۸.۶)، قابل مقایسه با (۲۲.۸.۶) را تولید کرد.

قضیه ۷.۶ شرط سازگاری (۷.۸.۶) مفروض است. روش چندگامی (۱.۸.۶) پایدار است اگر و تنها اگر شرط ریشه (۹.۸.۶) و (۱۰.۸.۶) برقرار باشد.

برهان ۱. اثبات را با نشان دادن لزوم شرط ریشه برای پایداری آغاز می کنیم. بدین منظور، خلاف آن را با فرض اینکه برای مقداری از j

$$|r_j(0)| > 1$$

در نظر می گیریم. معادله دیفرانسیل $y' \equiv 0$ ، $y(0) = 0$ ، با جواب $Y(x) \equiv 0$ را در نظر می گیریم. پس (۱.۸.۶) چنین می شود

$$y_{n+1} = \sum_{j=0}^p a_j y_{n-j} \quad n \geq p \quad (23.8.6)$$

اگر $y_0 = y_1 = \dots = y_p = 0$ ، روشن است که جواب عددی برای جمیع مقادیر $n \geq 0$ ، $y_n = 0$ خواهد شد. برای مقادیر اختلال یافته اولیه مقادیر زیر را می گیریم

$$z_0 = \varepsilon, z_1 = \varepsilon r_j(0), \dots, z_p = \varepsilon r_j(0)^p \quad (24.8.6)$$

برای این مقادیر اولیه

$$\text{Max}_{0 \leq n \leq p} |y_n - z_n| = \varepsilon |r_j(0)|^p$$

که یک کران یکنواخت برای همه مقادیر کوچک h است، زیرا سمت راست مستقل از h است. این کران نیز وقتی $\varepsilon \rightarrow 0$ ، به صفر میل می کند.

جواب (۲۴.۸.۶) با شرایط اولیه (۲۴.۸.۶) به صورت زیر است

$$z_n = \varepsilon [r_j(0)]^n \quad n \geq 0$$

برای انحراف از $\{y_n\}$ داریم

$$\text{Max}_{x_n \leq x_n \leq b} |y_n - z_n| = \varepsilon |r_j(\circ)|^{N(h)}$$

وقتی $h \rightarrow \circ, N(h) \rightarrow \infty$ و کران نامتناهی می شود. و این ثابت می کند که وقتی $|r_j(\circ)| > 1$ این روش ناپایدار است. اگر شرط ریشه با فرض نادرست بودن (۱۰.۸.۶) نقض شود، برهان مشابهی می توان ارائه نمود. این برهان به مسأله ۲۹ واگذار شده است.

۲. فرض می کنیم که شرط ریشه برقرار باشد. برهان پایداری را به معادله ای به صورت (۱۳.۸.۶) محدود می کنیم. یک برهان برای معادله کلی $y' = f(x, y)$ می توان داد که با کمی تغییر برهان زیر است. برهان کلی شامل حل معادلات تفاضلی خطی نا همگن است [برای یک بحث کلی آیزکسون و کور (۱۹۶۶، صص ۴۰۵-۴۱۷) را ببینید]. برای آنکه برهان باز هم ساده تر شود، فرض می کنیم که ریشه های $r_j(\circ)$ ، $j = 0, 1, \dots, p$ ، همگی متمایز باشند. این امر عیناً برای $r_j(h\lambda)$ درست است، مشروط بر آنکه مقدار h به اندازه کافی کوچک، مثلاً $0 \leq h \leq h_0$ ، نگهداشته شود. گیریم $\{y_n\}$ و $\{z_n\}$ دو جواب (۱۹.۸.۶) در $[x_0, b]$ باشند، و فرض می کنیم،

$$\text{Max}_{0 \leq n \leq p} |y_n - z_n| \leq \varepsilon \quad 0 < h \leq h_0. \quad (25.8.6)$$

خطای $e_n = y_n - z_n$ را وارد می کنیم. با کم کردن و به کار بردن (۱۹.۸.۶) برای هر جواب،

$$(1 - h\lambda b_{-1})e_{n+1} - \sum_{j=0}^p (a_j + h\lambda b_j)e_{n-j} = 0 \quad (26.8.6)$$

برای $x_{p+1} \leq x_{n+1} \leq b$. جواب عمومی چنین است

$$e_n = \sum_{j=0}^p \gamma_j [r_j(h\lambda)]^n \quad n \geq 0 \quad (27.8.6)$$

ضرایب $\gamma_0(h), \dots, \gamma_p(h)$ باید طوری انتخاب شوند که

$$\gamma_0 + \gamma_1 + \dots + \gamma_p = e_0$$

$$\gamma_0 r_0(h\lambda) + \dots + \gamma_p r_p(h\lambda) = e_1$$

⋮

$$\gamma_0 [r_0(h\lambda)]^p + \dots + \gamma_p [r_p(h\lambda)]^p = e_p$$

در این صورت جواب (۲۷.۸.۶) با اختلالات اولیه داده شده e_0, \dots, e_p مطابقت دارد و در معادله تفاضلی (۲۶.۸.۶) صدق می نماید. با استفاده از کران (۲۵.۸.۶) و نظریه دستگاه معادلات

خطی، بهسادگی می‌توان نشان داد که برای مقدار ثابتی چون $c_1 > 0$

$$\text{Max}_{\substack{i \\ \leq i \leq p}} |\gamma_i| \leq c_1 \varepsilon \quad 0 < h \leq h_0. \quad (28.8.6)$$

ما برهان را نمی‌آوریم گرچه با مفاهیمی که در فصل ۷ و ۸ معرفی شده‌اند می‌توان بهسادگی آن را انجام داد. برای کراندار نمودن e_n در $[x_0, b]$ ، باید هر جمله $[r_j(h\lambda)]^n$ را کراندار کنیم. برای این کار بسط

$$r_j(u) = r_j(0) + ur'_j(\zeta) \quad (29.8.6)$$

را به‌ازای مقداری از ζ بین 0 و u در نظر می‌گیریم. برای محاسبه $r'_j(u)$ ، از معادله مشخصه زیر مشتق می‌گیریم

$$\rho(r_j(u)) - u\sigma(r_j(u)) = 0$$

پس

$$r'_j(u) = \frac{\sigma(r_j(u))}{\rho'(r_j(u)) - u\sigma'(r_j(u))} \quad (30.8.6)$$

با فرض اینکه $r_j(0)$ یک ریشه ساده $\rho(r) = 0$ ، $0 \leq j \leq p$ ، نتیجه می‌شود که $\rho'(r_j(0)) \neq 0$ و با توجه به پیوستگی، به‌ازای جمیع مقادیر به اندازه کافی کوچک u ، $\rho'(r_j(u)) \neq 0$. مخرج $(30.8.6)$ ناصفر است و می‌توانیم $r'_j(u)$ را کراندار کنیم

$$|r'_j(u)| \leq c_2 \quad |u| \leq u_0$$

برای مقداری از $u_0 > 0$.

با استفاده از این نتیجه از $(29.8.6)$ و از شرط ریشه $(9.8.6)$ ، به‌ازای جمیع مقادیر

$$0 < h \leq h_0 \text{ داریم}$$

$$|r_j(h\lambda)| \leq |r_j(0)| + c_2 |h\lambda| \leq 1 + c_2 |h\lambda|$$

$$|[r_j(h\lambda)]^n| \leq [1 + c_2 |h\lambda|]^n \leq e^{c_2 n |h\lambda|} \leq e^{c_2 (b-x_0) |n\lambda|} \quad (31.8.6)$$

از ترکیب با $(27.8.6)$ و $(28.8.6)$ به‌ازای یک مقدار مناسب ثابت c_2 ، داریم

$$\text{Max}_{x_0 \leq x_n \leq b} |e_n| \leq c_2 |\varepsilon| e^{c_2 (b-x_0) |n\lambda|} \quad 0 < h \leq h_0.$$

برای مقدار مناسب ثابت c_2 .

نظریه همگرایی قضیه زیر تعمیم قضیه ۶.۶ از بخش ۳.۶، با شرایط لازم و کافی است که برای همگرایی روشهای چندگامی داده شده است.

قضیه ۸.۶ شرط سازگاری (۷.۸.۶) مفروض است. در این صورت روش چندگامی (۱.۸.۶) همگراست اگر و تنها اگر شرط ریشه (۹.۸.۶)، (۱۰.۸.۶) برقرار باشد.

برهان ۱. با نشان دادن لزوم شرط ریشه برای همگرایی آغاز و باز هم از مسأله $y' \equiv 0$ و $y(0) = 0$ با جواب $Y(x) \equiv 0$ استفاده می‌کنیم. روش چندگامی (۱.۸.۶) چنین می‌شود

$$y_{n+1} = \sum_{j=0}^p a_j y_{n-j} \quad n \geq p \quad (32.8.6)$$

با انتخاب y_0, \dots, y_p به نحوی که در شرط زیر صدق کنند

$$\eta(h) \equiv \max_{0 \leq n \leq p} |y_n| \rightarrow 0 \quad \text{وقتی } h \rightarrow 0 \quad (33.8.6)$$

فرض می‌کنیم که شرط ریشه نقض شده است. نشان می‌دهیم که (۳۲.۸.۶) به $Y(x) \equiv 0$ همگرا نیست.

فرض می‌کنیم که یک $|r_j(0)| > 1$. یک جواب قابل قبول (۳۲.۸.۶) چنین است

$$y_n = h[r_j(0)]^n \quad x_0 \leq x_n \leq b \quad (34.8.6)$$

شرط (۳۳.۸.۶) برقرار است زیرا

$$\eta(h) = h |r_j(0)|^p \rightarrow 0 \quad \text{وقتی } h \rightarrow 0$$

ولی جواب $\{y_n\}$ همگرا نیست. نخست،

$$\max_{x_0 \leq x_n \leq b} |Y(x_n) - y_n| = h |r_j(0)|^{N(h)}$$

آن مقادیر از h را که $h = b/N(h)$ در نظر می‌گیریم. می‌توان از قاعده هوییتال استفاده نموده نشان داد که

$$\lim_{N \rightarrow \infty} \frac{b}{N} |r_j(0)|^N = \infty$$

که نشان می‌دهد (۳۲.۸.۶) به جواب $Y(x) \equiv 0$ همگرا نیست.

فرض می‌کنیم از شرط ریشه، (۹.۸.۶) برقرار باشد، ولی یک $r_j(\circ)$ ریشه چندگانه $\rho(r)$ باشد و $|r_j(\circ)| = 1$. آنگاه اثباتی مشابه آنچه گفته شد کماکان برقرار است ولی باید از جواب زیر استفاده کنیم

$$y_n = hn[r_j(\circ)]^n \quad \circ \leq n \leq N(h)$$

و این برهان لزوم شرط ریشه را کامل می‌کند.

۲. فرض کنید که شرط ریشه برقرار است. مانند قضیه قبل، اثبات همگرایی برای یک معادله دیفرانسیل دلخواه خیلی مشکل است. می‌توانید به آیزکسون و کلر (۱۹۶۶، صص ۴۰۵-۴۱۷) مراجعه کنید. برهان ارائه شده به معادله‌ای به صورت (۱۳.۸.۶) محدود شده است، و باز هم برای آسان بودن اثبات فرض می‌کنیم که ریشه‌های $r_j(\circ)$ متمایزند.

روش چندگامی (۱.۸.۶) برای معادله نمونه $y' = \lambda y$ و $y(\circ) = 1$ به (۱۸.۸.۶) تبدیل می‌شود. نشان می‌دهیم که جمله $\gamma_\circ[r_\circ(h\lambda)]^n$ در جواب عمومی

$$y_n = \sum_{j=\circ}^p \gamma_j[r_j(h\lambda)]^n \quad (۳۵.۸.۶)$$

به جواب $Y(x) = e^{\lambda x}$ در $[0, b]$ ، همگراست. جملات دیگر $\gamma_j[r_j(h\lambda)]^n$ ، $j = 1, 2, \dots, p$ ، جوابهای مزاحم هستند و می‌توان نشان داد که وقتی $h \rightarrow \circ$ این جوابها به صفر می‌گریند (مسئله ۳۰ را ببینید).

با استفاده از قضیه تیلر، $r_\circ(h\lambda)$ را بسط می‌دهیم

$$r_\circ(h\lambda) = r_\circ(\circ) + h\lambda r'_\circ(\circ) + O(h^2)$$

از (۳۰.۸.۶) داریم

$$r'_\circ(\circ) = \frac{\sigma(1)}{\rho'(1)}$$

که با استفاده از شرط سازگاری (۷.۸.۶)، این تساوی به $r'_\circ(\circ) = 1$ منجر می‌شود. پس روی هر بازه متناهی $0 \leq x_n \leq b$

$$r_\circ(h\lambda) = 1 + h\lambda + O(h^2) = e^{\lambda h} + O(h^2)$$

$$[r_\circ(h\lambda)]^n = e^{\lambda nh} [1 + O(h^2)]^n = e^{\lambda x_n} [1 + O(h)] \quad (۳۶.۸.۶)$$

بنابراین

$$\text{Max}_{\circ \leq x_n \leq b} | [r_0(h\lambda)]^n - e^{\lambda x_n} | \rightarrow \circ \quad h \rightarrow \circ \quad \text{وقتی که} \quad (۳۷.۸.۶)$$

اکنون باید نشان دهیم که وقتی $h \rightarrow \circ$ ، $\gamma_0 \rightarrow ۱$.

ضرایب $\gamma_0(h), \dots, \gamma_p(h)$ در دستگاه خطی زیر صدق می نمایند

$$\gamma_0 + \gamma_1 + \dots + \gamma_p = y_0$$

$$\gamma_0 [r_0(h\lambda)] + \dots + \gamma_p [r_p(h\lambda)] = y_1 \quad (۳۸.۸.۶)$$

⋮

$$\gamma_0 [r_0(h\lambda)]^p + \dots + \gamma_p [r_p(h\lambda)]^p = y_p$$

مقادیر آغازی y_0, \dots, y_p بستگی به h دارند، و بنا به فرض باید در رابطه زیر صدق کنند

$$\eta(h) \equiv \text{Max}_{\circ \leq n \leq p} | e^{\lambda x_n} - y_n | \rightarrow \circ \quad h \rightarrow \circ \quad \text{وقتی}$$

ولی این امر ایجاب می کند که

$$\lim_{h \rightarrow \circ} y_n = ۱ \quad \circ \leq n \leq p \quad (۳۹.۸.۶)$$

ضریب γ_0 را می توان با استفاده از قاعده کرامر در حل (۳۸.۸.۶) به دست آورد:

$$\gamma_0 = \frac{\begin{vmatrix} y_0 & ۱ & \dots & ۱ \\ y_1 & r_1 & \dots & r_p \\ \vdots & \vdots & \dots & \vdots \\ y_p & r_1^p & \dots & r_p^p \end{vmatrix}}{\begin{vmatrix} ۱ & ۱ & \dots & ۱ \\ r_0 & r_1 & \dots & r_p \\ \vdots & \vdots & \dots & \vdots \\ r_0^p & r_1^p & \dots & r_p^p \end{vmatrix}} \quad (۴۰.۸.۶)$$

مخرج کسر برای $r_0(\circ), \dots, r_1(\circ), r_p(\circ) = ۱$ به دترمینان واندرموند میل می کند که ناصفر است زیرا ریشه ها متمایزند (مسئله ۱، فصل ۳ را ببینید). با استفاده از (۳۹.۸.۶)، صورت کسر نیز وقتی $h \rightarrow \circ$ به همان مقدار میل خواهد کرد. بنابراین وقتی $h \rightarrow \circ$ آنگاه $\gamma_0 \rightarrow ۱$. با استفاده

از این ویژگی همراه با (۳۷.۸.۶) و مسأله ۳۰، جواب $\{y_n\}$ به سمت $Y(x) = e^{\lambda x}$ در $[0, b]$ میل خواهد کرد.

آنچه ذیلاً می‌آید نتیجه مشهودی است و نتیجه نمایان قضایای ۷.۶ و ۸.۶ است.

فرع گیریم (۱.۸.۶) یک روش چندگامی سازگار باشد. در این صورت این روش همگراست اگر و تنها اگر پایدار باشد.

پایداری نسبی و پایداری ضعیف بازهم معادله نمونه (۱.۸.۶) و جواب عددی آن (۳۲.۸.۶) را در نظر می‌گیریم. قضیه قبل می‌گوید که جوابهای انگلی $\gamma_j[r_j(h\lambda)]^n$ وقتی $h \rightarrow 0$ ، به صفر می‌گیرند. ولی برای یک h ثابت و x_n صعودی، می‌خواهیم این جوابها نسبت به قسمت اصلی جواب $\gamma_0[r_0(h\lambda)]^n$ کوچک بمانند. این امر وقتی میسر است که ریشه‌های مشخصه، به‌ازای جمیع مقادیر به اندازه کافی کوچک h در رابطه زیر صدق کنند

$$|r_j(h\lambda)| \leq r_0(h\lambda) \quad j = 1, 2, \dots, p \quad (41.8.6)$$

این نابرابری ما را به تعریف پایداری نسبی زیر هدایت می‌کند.

گوییم روش (۱.۸.۶) نسبی - پایدار است اگر ریشه‌های مشخصه $r_j(h\lambda)$ به‌ازای مقادیر به اندازه کافی کوچک $|h\lambda|$ در رابطه (۴۱.۸.۶) صدق کنند. و می‌گویند روش در شرط قوی ریشه صدق می‌کند اگر

$$|r_j(0)| < 1 \quad j = 1, 2, \dots, p \quad (42.8.6)$$

این شرط، به‌سادگی کنترل می‌شود، و پایداری نسبی را ایجاب می‌کند. فقط کافی است از پیوستگی ریشه‌های $r_j(h\lambda)$ نسبت به $h\lambda$ استفاده کنیم تا ببینیم که شرط (۴۲.۸.۶) شرط (۴۱.۸.۶) را ایجاب می‌کند. پایداری نسبی شرط قوی ریشه را تضمین نمی‌کند، اگرچه برای بیشتر روشها، آنها هم‌ارزند (مسأله ۳۶ (ب) را ببینید). اگر یک روش چندگامی پایدار باشد ولی نسبی - پایدار نباشد آن را ضعیف - پایدار خوانند.

مثال ۱. برای روش میانگاهی

$$r_0(h\lambda) = 1 + h\lambda + O(h^2) \quad r_1(h\lambda) = -1 + h\lambda + O(h^2) \quad (43.8.6)$$

بنابر (۴۱.۸.۶) وقتی $\lambda < 0$ ، ضعیف- پایدار است که با آنچه قبلاً در بخش ۴.۶ بیان شد مطابقت می‌کند.

۲. روشهای آدامز- بشفورت و آدامز- مولتن، (۲۲.۷.۶) و (۲۶.۷.۶)، وقتی $h = 0$ ، دارای یک چندجمله‌یی مشخصه هستند،

$$\rho(r) = r^{p+1} - r^p \quad (44.8.6)$$

ریشه‌ها عبارتند از $r_0 = 1$ و $r_j = 0$ ، $j = 1, 2, \dots, p$ بنابراین شرط قوی ریشه برقرار است و روشهای آدامز نسبی- پایدارند.

ناحیه‌های پایداری در بحثهای قبلی پایداری، لازم بود مقادیر h به اندازه کافی کوچک باشند تا بتوان عملیات را ادامه داد. اصلاً اشاره نشده بود که h چه اندازه باید کوچک باشد. روشن است که اگر لازم شود h بسیار کوچک باشد، آنگاه این روش برای بسیاری از مسائل غیرعملی خواهد بود. بنابراین لازم است مقادیر مجاز h را بررسی کنیم. چون پایداری به ریشه‌های مشخصه بستگی دارد، و چون خود این ریشه‌ها نیز به $h\lambda$ وابسته‌اند، ما علاقه‌مندیم مقادیری از $h\lambda$ را پیدا کنیم که برای آنها روش چندگامی (۱.۸.۶) به تعبیری پایدار باشد. برای اینکه بتوانیم در همه وضعیتهایی که در حل دستگاه معادلات دیفرانسیل پیش می‌آیند، بحث کنیم لازم است، همان‌گونه که به دنبال (۱۷.۸.۶) اشاره شد، λ بتواند مقادیر مختلط اختیار کند.

برای دیدن انگیزه بحث اخیر، پایداری روش اویلر را در نظر می‌گیریم. روش اویلر را برای معادله

$$y' = \lambda y + g(x) \quad y(0) = Y. \quad (45.8.6)$$

به‌کار می‌بریم، نتیجه زیر به دست می‌آید

$$y_{n+1} = y_n + h[\lambda y_n + g(x_n)] \quad n \geq 0 \quad y_0 = Y. \quad (46.8.6)$$

سپس مسأله اختلال یافته را در نظر می‌گیریم

$$z_{n+1} = z_n + h[\lambda z_n + g(x_n)] \quad n \geq 0 \quad z_0 = Y + \varepsilon \quad (47.8.6)$$

برای معادله اصلی (۴۵.۸.۶)، این اختلال Y به جوابهای $Y(x)$ و $Z(x)$ می‌انجامد که در رابطه زیر صدق می‌کنند

$$Y(x) - Z(x) = \varepsilon e^{\lambda x} \quad x \geq 0$$

در این مسأله اصلی، معمولاً به حالتی توجه داریم که $\text{Real}(\lambda) \leq 0$ ، زیرا در این صورت $|Y(x) - Z(x)|$ وقتی $x \rightarrow 0$ ، کراندار باقی می ماند. توجه خود را به حالت $\text{Real}(\lambda) < 0$ معطوف می داریم. در این حالت وقتی $x \rightarrow \infty$ ، $Y(x) - Z(x) \rightarrow 0$ ، برای چنین λ ای می خواهیم مقادیری از h را پیدا کنیم که جوابهای عددی (۴۶.۸.۶) و (۴۷.۸.۶) رفتار وابسته به $Y(x)$ و $Z(x)$ را حفظ کنند.

گیریم $e_n = z_n - y_n$. با کم کردن (۴۶.۸.۶) از (۴۷.۸.۶)،

$$e_{n+1} = e_n + h\lambda e_n = (1 + h\lambda)e_n \quad e_0 = \varepsilon$$

و با استقراء

$$e_n = (1 + h\lambda)^n \varepsilon \quad (48.8.6)$$

پس وقتی $x_n \rightarrow \infty$ ، $e_n \rightarrow 0$ ، اگر و تنها اگر

$$|1 + h\lambda| < 1 \quad (49.8.6)$$

این رابطه یک مجموعه از مقادیر مختلط $h\lambda$ به دست می دهد که از نقاط درون دایره به شعاع یک به مرکز -1 در صفحه مختلط تشکیل شده است. اگر $h\lambda$ به این مجموعه تعلق داشته باشد، آنگاه وقتی $x_n \rightarrow \infty$ ، $y_n - z_n \rightarrow 0$ ، ولی نه در غیر این صورت.

برای آنکه ببینیم این بحث برای همگرایی نیز اهمیت دارد، توجه می کنیم که معادله دیفرانسیل اصلی را می توان به عنوان یک حالت اختلال یافته معادله تقریب زنده (۴۶.۸.۶) منظور کرد. از به کار بردن (۱۷.۲.۶) در (۴۵.۸.۶)،

$$Y_{n+1} = Y_n + h[\lambda Y_n + g(x_n)] + \frac{h^2}{2} Y''(\xi_n) \quad (50.8.6)$$

در اینجا در هر مرحله یک اختلال معادله (۴۶.۸.۶) را داریم، نه فقط یک اختلال در نقطه اولیه $x_0 = 0$. با این حال، می توان نشان داد که تحلیل پایداری قبلی را برای این اختلال (۴۶.۸.۶) می توان به کار برد. فرمول خطای (۴۸.۸.۶) باید به نحو مناسبی اصلاح شود، ولی باز هم فرمول به کران (۴۹.۸.۶) زیاد وابسته خواهد بود. (مسأله ۴۰ را ببینید).

مثال روش اویلر را برای مسأله

$$y' = \lambda y + (1 - \lambda) \cos(x) - (1 + \lambda) \sin(x) \quad y(0) = 1 \quad (51.8.6)$$

جدول ۱۷.۶ روش اویلر برای (۵۱.۸.۶)

λ	x	خطا: $h = 0.5$	خطا: $h = 0.1$	خطا: $h = 0.01$
-۱	۱	$-2.46E-1$	$-4.32E-2$	$-4.22E-3$
	۲	$-2.55E-1$	$-4.64E-2$	$-4.55E-3$
	۳	$-2.66E-2$	$-6.78E-3$	$-7.22E-4$
	۴	$2.27E-1$	$3.91E-2$	$3.78E-3$
	۵	$2.72E-1$	$4.91E-2$	$4.81E-3$
-۱°	۱	$3.98E-1$	$-6.99E-3$	$-6.99E-4$
	۲	$6.90E+0$	$-2.90E-3$	$-3.08E-4$
	۳	$1.11E+2$	$3.86E-3$	$3.64E-4$
	۴	$1.77E+3$	$7.07E-3$	$7.04E-4$
	۵	$2.83E+4$	$3.78E-3$	$3.97E-4$
-۵°	۱	$3.26E+0$	$1.06E+3$	$-1.39E-4$
	۲	$1.88E+3$	$1.11E+9$	$-5.16E-5$
	۳	$1.08E+6$	$1.17E+15$	$8.25E-5$
	۴	$6.24E+8$	$1.23E+21$	$1.41E-4$
	۵	$3.59E+11$	$1.28E+27$	$7.00E-5$

که $Y(x) = \sin(x) + \cos(x)$ جواب صحیح آن است، به کار برید. ما نتایج را برای چندین مقدار λ و h داده ایم. برای $\lambda = -1, -1^\circ, -5^\circ$ کران (۴۹.۸.۶) برای h به ترتیب کرانهایی به شکل زیر به دست می دهد.

$$0 < h < 2 \quad 0 < h < \frac{1}{5} = 0.2 \quad 0 < h < \frac{1}{25} = 0.04$$

همان گونه که از جدول ۱۷.۶ می توان دید، استفاده از مقادیر بزرگتر h نتایج عددی نامطلوبی به بار می آورد.

آنچه در بالا با روش اویلر به دست آوردیم انگیزه روش کلی، برای پیدا کردن مجموعه همه $h\lambda$ هایی است که برای آنها روش (۱.۸.۶) پایدار است. چون ما فقط حالت های $\text{Real}(\lambda) < 0$ را در نظر می گیریم، $\{y_n\}$ ، جواب عددی (۱.۸.۶) در حل معادله نمونه $y' = \lambda y$ را می خواهیم که با $x_n \rightarrow \infty$ برای همه مقادیر اولیه y_0, \dots, y_p به صفر میل کند. مجموعه همه $h\lambda$ هایی که برای آنها این موضوع درست باشد ناحیه پایداری مطلق روش (۱.۸.۶) خوانده می شود. هرچه

این ناحیه وسیعتر باشد، محدودیت h برای داشتن جواب عددی پایدار، کمتر است.

وقتی (۱.۸.۶) برای معادله نمونه به کار رود، معادله قبلی (۱۸.۸.۶) را به دست می آوریم، و به شرطی که ریشه های مشخصه $r_0(h\lambda), \dots, r_p(h\lambda)$ متمایز باشند جواب آن با (۲۲.۸.۶)، یعنی

$$y_n = \sum_{j=0}^p \gamma_j [r_j(h\lambda)]^n \quad n \geq 0$$

داده می شود. برای آنکه وقتی $n \rightarrow \infty$ ، این جواب به ازای همه مقادیر اولیه y_0, \dots, y_p به صفر میل کند، لازم و کافی است که داشته باشیم

$$|r_j(h\lambda)| < 1 \quad j = 0, 1, \dots, p \quad (52.8.6)$$

مجموعه همه مقادیر $h\lambda$ که در این مجموعه نامساویها صدق کنند نیز ناحیه پایداری مطلق خوانده می شود. این ناحیه در مجموعه ای که در پاراگراف قبلی تعریف شده قرار دارد، و معمولاً با آن مجموعه برابر است. در پیدا کردن ناحیه پایداری مطلق، فقط با (52.8.6) کار می کنیم.

مثال روش مرتبه دو آدامز - بشفورت

$$y_{n+1} = y_n + \frac{h}{4} [3y'_n - y'_{n-1}] \quad n \geq 1 \quad (53.8.6)$$

را در نظر می گیریم. معادله مشخصه چنین است

$$r^2 - (1 + \frac{3}{4}h\lambda)r + \frac{1}{4}h\lambda = 0$$

و ریشه ها عبارت اند از

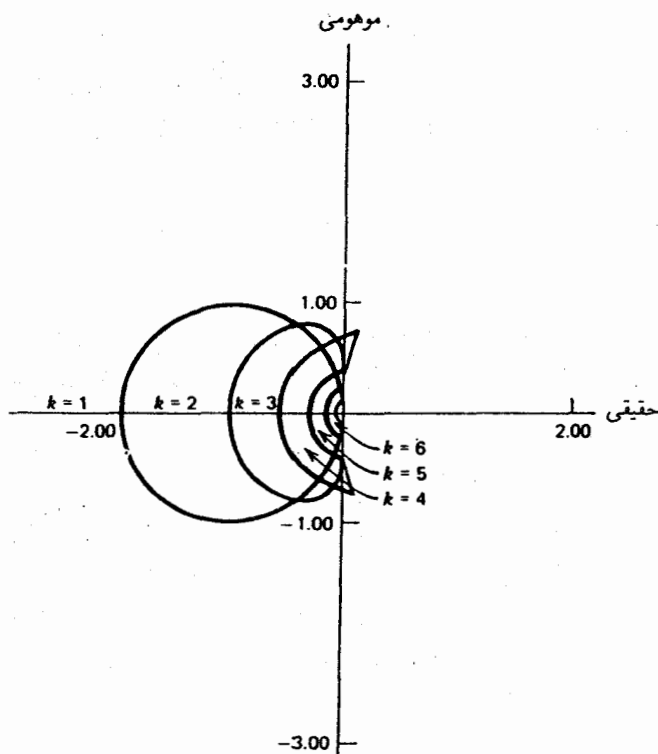
$$r_0 = \frac{1}{4} \left\{ 1 + \frac{3}{4}h\lambda + \sqrt{1 + h\lambda + \frac{1}{4}h^2\lambda^2} \right\}$$

$$r_1 = \frac{1}{4} \left\{ 1 + \frac{3}{4}h\lambda - \sqrt{1 + h\lambda + \frac{1}{4}h^2\lambda^2} \right\}$$

ناحیه پایداری مطلق، مجموعه $h\lambda$ هاست که برای آنها

$$|r_0(h\lambda)| < 1 \quad |r_1(h\lambda)| < 1$$

به ازای مقادیر حقیقی λ ، مقادیر قابل قبول $h\lambda$ ، عبارت اند از $-1 < h\lambda < 0$.

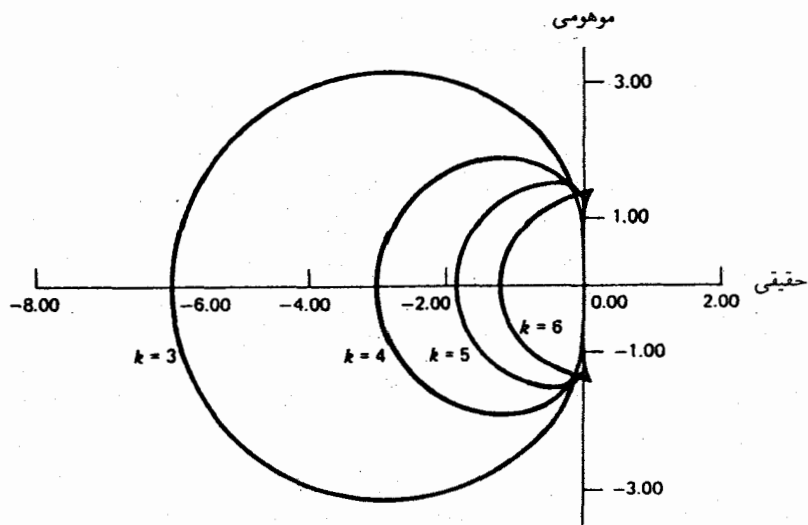


شکل ۷.۶ نواحی پایداری برای روشهای آدامز-بشفورت. روش مرتبه k در داخل ناحیه‌ای که در سمت چپ مرکز مشخص شده پایدار است [با اجازه، از صفحه ۱۳۱ گی‌پر (۱۹۷۱) گرفته شده است.]

کرانه‌های نواحی پایداری مطلق روشهای آدامز-بشفورت و آدامز-مولتن به ترتیب در شکل‌های ۷.۶ و ۸.۶ داده شده‌اند. در مورد فرمولهای آدامز-مولتن با یک بارست از پیشگویی آدامز-بشفورت، نواحی پایداری مطلق در شمایین و گوردون (۱۹۷۵، صص ۱۳۵-۱۴۰) داده شده‌اند.

از این نمودارها، دیده می‌شود که وقتی مرتبه روش افزایش یابد، ناحیه همگرایی مطلق کوچکتر می‌شود. و برای فرمولهای هم‌مرتبه، در دو فرمول آدامز-مولتن و آدامز-بشفورت، ناحیه همگرایی مطلق در اولی به مراتب بزرگتر است. اندازه این ناحیه‌ها، معمولاً از نقطه نظر عملی به خوبی قابل قبول‌اند. برای مثال، مقادیر حقیقی $h\lambda$ در ناحیه پایداری مطلق برای فرمول مرتبه چهار آدامز-مولتن با $0 < h\lambda < 3$ داده شده است. در اغلب حالات، این محدودیت زیادی برای h نخواهد بود.

خانواده فرمولهای آدامز، برای ایجاد الگوریتمهای متغیر-مرتبه، بسیار مناسب‌اند، و ناحیه‌های پایداری کاملاً قابل قبول‌اند. این روشها برای مسائلی که λ منفی و قدرمطلق آن بزرگ باشد، مشکلاتی



شکل ۸.۶ نواحی پایداری برای روشهای آدامز-بسفورت. روش مرتبه k در داخل ناحیه‌ای که مشخص شده پایدار است. [از صفحه ۱۳۱ گی‌پر (۱۹۷۱) با اجازه برداشته شده است.]

دارند، و این نوع مسائل با روشهای دیگر بهتر حل می‌شوند، که در بخش بعد ملاحظه خواهیم نمود. روشهای ویژه‌ای وجود دارند که برای آنها ناحیه پایداری مطلق شامل تمام مقادیر مختلط $h\lambda$ با $\text{Real}(\lambda) < 0$ است. این روشها A -پایدار خوانده می‌شوند، و با آنها برای داشتن پایداری از نوعی که ملاحظه کرده‌ایم، هیچ محدودیتی برای h وجود نخواهد داشت. قاعده دوزنقه‌یی یک مثال از چنین روش است. [(۲۴.۵.۶) - (۲۵.۵.۶) را ببینید].

مثال روش پسرو اویلر را در نظر می‌گیریم:

$$y_{n+1} = y_n + hf(x_{n+1}, y_{n+1}) \quad n \geq 0 \quad (۵۴.۸.۶)$$

با به‌کار بردن این روش در معادله نمونه $y' = \lambda y$ و حل آن، داریم

$$y_n = \left[\frac{1}{1 - h\lambda} \right]^n y_0 \quad n \geq 0 \quad (۵۵.۸.۶)$$

بنابراین وقتی $x_n \rightarrow \infty$ ، $y_n \rightarrow 0$ اگر و تنها اگر

$$\left| \frac{1}{1 - h\lambda} \right| < 1$$

جدول ۱۸.۶ مثال قاعدهٔ دوزنقه: $h = 0.5$

x	خطا: $\lambda = -1$	خطا: $\lambda = -10$	خطا: $\lambda = -50$
۲	$-1.13E - 2$	$-2.78E - 3$	$-7.91E - 4$
۴	$-1.43E - 2$	$-8.91E - 5$	$-8.91E - 5$
۶	$2.02E - 2$	$2.77E - 3$	$4.72E - 4$
۸	$-2.86E - 3$	$-2.22E - 3$	$-5.11E - 4$
۱۰	$-1.79E - 2$	$-9.23E - 4$	$-1.56E - 4$

این رابطه به‌ازای تمام مقادیر $h\lambda$ با $\text{Real}(\lambda) < 0$ درست است و روش پسر اوایلر یک روش A -پایدار است.

مثال روش دوزنقه‌یی را برای مسئله (۵۱.۸.۶)، که قبلاً با روش اوایلر حل شده به‌کار می‌بریم. از طول گام $h = 0.5$ برای $\lambda = -1, -10, -50$ استفاده می‌کنیم. نتایج در جدول ۱۸.۶ داده شده‌اند. این نتایج نشان می‌دهند که وقتی $|\lambda|$ افزایش می‌یابد $\text{Real}(\lambda) < 0$ ، قاعدهٔ دوزنقه‌یی ناپایدار نمی‌شود.

داشتن روشهای A -پایدار از مراتب بزرگتر از ۲ مفید خواهد بود. ولی یک قضیه از دال‌کویست^۱ (۱۹۶۳) نشان می‌دهد که چنین روشهایی وجود ندارند. ما در بخش آینده بعضی روشهای مراتب بالاتر را که بیشترین ویژگیهای لازم پایداری را دارا هستند بررسی خواهیم کرد.

۹.۶ معادلات دیفرانسیل سرسخت و روش خطوط

حل عددی معادلات دیفرانسیل سرسخت، در طی ده پانزده سال اخیر بسیار مورد مطالعه قرار گرفته است. این‌گونه معادلات (از جمله دستگاههای معادلات دیفرانسیل) در کاربردهای روزافزونی در موضوعات بسیار دور از هم مثل جنبش‌شناسی شیمیایی و حل عددی معادلات با مشتقات جزئی ظاهر شده‌اند. در این بخش ما بعضی از نکات اصلی این موضوع را مطرح می‌نماییم و رابطهٔ آن را با حل عددی یک معادلهٔ سادهٔ گرما نشان می‌دهیم.

تعاریف زیادی از مفهوم معادلهٔ دیفرانسیل سرسخت وجود دارد. مهمترین جنبهٔ مشترک این تعاریف آن است که چنین معادلاتی وقتی با روشهای معمولی (مثلاً روشهای آدامز، بخش ۷.۶) حل شوند، طول گام h اجباراً باید بسیار کوچک گرفته شود تا پایداری تأمین گردد. بسیار کوچکتر از آنچه ممکن است در مورد خطای برشی، لازم به‌نظر آید. یک گواه این، حل معادله (۵۱.۸.۶)

به روش اویلر است، که در جدول ۱۷.۶ داده شده است. در آن حالت، مجهول $Y(x)$ با λ تغییر نمی‌کرد و بنابراین خطای برشی مستقل از λ بود. ولی خطای واقعی قویاً تحت تأثیر قدرمطلق λ قرار گرفته بود، زیرا لازم بود که λ در شرط پایداری $|1 + h\lambda| < 1$ صدق نماید تا همگرایی به دست آید. وقتی $|\lambda|$ افزایش می‌یافت، اندازه h باید مطابق با آن کوچک می‌شد. این رفتار وقتی از روشهای عددی استاندارد برای معادلات دیفرانسیل سرسخت استفاده می‌شود امری عادی است، با این تفاوت عمده که مقادیر واقعی $|\lambda|$ بسیار بزرگترند، برای مثال $\lambda = -10^6$.

اکنون به متداولترین ردهٔ این معادلات دیفرانسیل می‌نگریم، بر این پایه که بررسی خود را به خطی نمودن دستگاه $y' = f(x, y)$ ، که در (۱۴.۸.۶) تا (۱۷.۸.۶) شرح داده شده معطوف می‌داریم:

$$y' = \Lambda y + g(x) \quad (1.9.6)$$

با $\Lambda = f_y(x_0, Y_0)$ ، ماتریس ژاکوبی f . معادلهٔ دیفرانسیل $y' = f(x, y)$ را سرسخت گوئیم اگر λ ، یکی از ویژه مقدرهای Λ ، یا به‌طور کلیتر $f_y(x, y)$ ، دارای قسمت حقیقی منفی با اندازهٔ بسیار بزرگ باشد. ما روشهای عددی برای معادلات دشوار را با در نظر گرفتن اثر آنها روی معادلهٔ نمونهٔ

$$y' = \lambda y + g(x) \quad (2.9.6)$$

مطالعه می‌کنیم که $\text{Real}(\lambda)$ منفی و از لحاظ اندازه بسیار بزرگ است. این روش، محدودیتهای خاص خود را دارد که بعضی از آنها را بعداً نشان خواهیم داد، ولی به ما امکاناتی می‌دهد که روشهای غیر رضایتبخش را کنار بگذاریم، و همچنین ما را به فکر بعضی روشهای رضایتبخش ممکن می‌رساند. مفهوم ناحیهٔ پایداری مطلق، که در بخش اخیر وارد شد، ابزار اولیه در مطالعهٔ پایداری یک روش عددی در حل معادلات دیفرانسیل سرسخت است. ما روشهایی را جستجو می‌کنیم که ناحیهٔ پایداری آنها شامل تمام محور حقیقی منفی، و همچنین تا آنجا که ممکن است، شامل ناحیهٔ چپ صفحهٔ مختلط باشد. برای پیدا کردن این گونه روشها چند راه وجود دارد، ولی ما فقط یکی از آنها را، که فرمولهای مشتقگیری پسرو^۱ ($BDFs$) را به دست می‌دهد، مورد بحث قرار می‌دهیم.

گیریم $P_p(x)$ یک چندجمله‌یی از درجهٔ نابزرگتر از p باشد که $Y(x)$ را در نقاط $x_n, x_{n+1}, \dots, x_{n-p+1}$ برای مقداری از $p \geq 1$ ، درونیابی نماید:

$$P_p(x) = \sum_{j=-1}^{p-1} Y(x_{n-j}) l_{j,n}(x) \quad (3.9.6)$$

جدول ۱۹.۶ ضرایب روش BDF (۶.۹.۶)

p	β	α_0	α_1	α_2	α_3	α_4	α_5
۱	۱	۱					
۲	۲	۴	-۱				
	۳	۳	-۳				
۳	۶	۱۸	۹	۲			
	۱۱	۱۱	-۱۱	۱۱			
۴	۱۲	۴۸	۳۶	۱۶	۳		
	۲۵	۲۵	-۲۵	۲۵	-۲۵		
۵	۶۰	۳۰۰	۳۰۰	۲۰۰	۷۵	۱۲	
	۱۳۷	۱۳۷	-۱۳۷	۱۳۷	-۱۳۷	۱۳۷	
۶	۶۰	۳۶۰	۴۵۰	۴۰۰	۲۲۵	۷۲	۱۰
	۱۴۷	۱۴۷	-۱۴۷	۱۴۷	-۱۴۷	۱۴۷	-۱۴۷

که در آن $\{l_{j,n}(x)\}$ توابع پایه درونیایی لاگرانژی برای نقاط گرهی $x_{n+1}, \dots, x_{n-p+1}$ هستند (۵.۱.۳) را ببینید. از رابطه زیر استفاده می‌کنیم

$$P'_p(x_{n+1}) \doteq Y'(x_{n+1}) = f(x_{n+1}, Y(x_{n+1})) \quad (۴.۹.۶)$$

از ترکیب با (۳.۹.۶) و حل آن نسبت به $Y(x_{n+1})$ داریم

$$Y(x_{n+1}) \doteq \sum_{j=0}^{p-1} \alpha_j Y(x_{n-j}) + h\beta f(x_{n+1}, Y(x_{n+1})) \quad (۵.۹.۶)$$

روش p گامی فرمول مشتقگیری پسرو با رابطه زیر داده می‌شود

$$y_{n+1} = \sum_{j=0}^{p-1} \alpha_j y_{n-j} + h\beta f(x_{n+1}, y_{n+1}) \quad (۶.۹.۶)$$

ضرایب در حالت‌های $p = 1, \dots, 6$ در جدول ۱۹.۶ داده شده‌اند. حالت $p = 1$ همان روش پسرو اویلر است، که به دنبال (۵۴.۸.۶) در بخش اخیر بحث شد. خطای برشی برای (۶.۹.۶) را از فرمول خطا برای مشتقگیری عددی، که در (۵.۷.۵) داده شده بود، می‌توان به دست آورد:

$$T_n(Y) = -\frac{\beta}{p+1} h^{p+1} Y^{(p+1)}(\xi_n) \quad (۷.۹.۶)$$

برای مقداری از $x_{n-p+1} \leq \xi_n \leq x_{n+1}$

ناحیه‌های پایداری مطلق برای فرمولهای جدول ۱۹.۶ در گی‌پر (۱۹۷۱، صص ۲۱۵-۲۱۶) داده شده‌اند. برای ساختن این نواحی، باید تمام مقادیر $h\lambda$ را که برای آنها

$$|r_j(h\lambda)| < 1 \quad j = 0, \dots, p \quad (۸.۹.۶)$$

پیدا کنیم، که ریشه‌های مشخصه $r_j(h\lambda)$ جوابهای معادله

$$r^p = \sum_{j=0}^{p-1} \alpha_j r^{p-1-j} + h\lambda \beta r^p \quad (۹.۹.۶)$$

هستند. می‌توان نشان داد که برای $p = 1$ و $p = 2$ ، فرمولهای مشتفگیری پسر A - پایداری، و برای $3 \leq p \leq 6$ ، ناحیه پایداری مطلق با افزایش p ، کوچک می‌شود، گرچه در هر حالت تمام محور حقیقی منفی را در بر می‌گیرد. برای $p \geq 7$ ، ناحیه‌های پایداری مطلق، برای حل مسائل سرسخت، قابل قبول نیستند. برای بحث بیشتر در مورد این ناحیه‌های پایداری، گی‌پر (۱۹۷۱، فصل ۱۱) و لامبرت (۱۹۷۳، فصل ۸) را ببینید.

باز هم با روشهای BDF و با سایر روشهایی که فقط بر پایه ناحیه پایداری مطلق انتخاب شده‌اند، اشکالاتی وجود دارند. ابتدا، با معادله نمونه $y' = \lambda y$ آغاز می‌کنیم، اگر $\text{Real}(\lambda)$ منفی و قدرمطلق آن بزرگ باشد، آنگاه جواب $Y(x)$ بسیار سریع به صفر می‌گراید، و چنانچه $\text{Real}(\lambda)$ افزایش یابد، $Y(x)$ سریعتر به صفر همگرا می‌شود، ما می‌خواهیم که همین رفتار برای جواب عددی معادله نمونه $\{y_n\}$ برقرار باشد. ولی با قاعده دوزنقه‌یی A - پایدار، جواب [از (۲۴.۵.۶)] چنین است

$$y_n = \left[\frac{1 + \frac{h\lambda}{2}}{1 - \frac{h\lambda}{2}} \right]^n y_0 \quad n \geq 0$$

اگر $|\text{Real}(\lambda)|$ بزرگ باشد، آنگاه کسر داخل کروشه نزدیک -1 می‌شود، و y_n خیلی آهسته به صفر کاهش می‌یابد. با استفاده از نوع استدلالی که برای روش اوایلر در (۴۵.۸.۶) - (۵۰.۸.۶) به‌کار رفت، اثر اختلالات، برای مقادیر بزرگ λ ، به سرعت به صفر کاهش نمی‌یابد. بنابراین ممکن است روش دوزنقه‌یی برای مسائل سرسخت، یک انتخاب کاملاً رضایتبخش نباشد. در مقایسه، روش پسر A - پایدار اوایلر رفتار مطلوب را داراست. با توجه به (۵۵.۸.۶)، جواب مسأله نمونه چنین است

$$y_n = \left[\frac{1}{1 - h\lambda} \right]^n y_0 \quad n \geq 0$$

وقتی $|\lambda|$ افزایش می‌یابد، دنباله $\{y_n\}$ سریعتر به صفر می‌گراید. بنابراین جواب پسر و اوایل رفتار جواب درست معادله نمونه را بهتر منعکس می‌نماید.

اشکال دوم در مورد ناحیه‌های پایداری این است که این ناحیه‌ها بر مبنای مقدار ثابت λ و مسائل خطی پایه‌گذاری شده‌اند. خطی‌کردن (۱۰.۹.۶) اغلب معتبر است ولی نه همیشه. برای مثال، مسأله خطی مرتبه دوم زیر را که در آن یکی از ضرایب ثابت نیست در نظر می‌گیریم:

$$y'' + ay' + (1 + b \cdot \cos(2\pi x))y = g(x) \quad x \geq 0 \quad (10.9.6)$$

آن را به دستگاه هم‌ارز زیر برمی‌گردانیم

$$\begin{aligned} y_1' &= y_2 \\ y_2' &= -(1 + b \cos(2\pi x))y_1 - ay_2 + g(x) \end{aligned} \quad (11.9.6)$$

فرض می‌کنیم که $a > 0$ و $|b| < 1$. ویژه مقادیرهای معادله همگن $[g(x) \equiv 0]$ عبارت‌اند از

$$\lambda = \frac{-a \pm \sqrt{a^2 - 4[1 + b \cdot \cos(2\pi x)]}}{2} \quad (12.9.6)$$

این ویژه مقادیرها، اعداد حقیقی منفی یا اعداد مختلط با قسمت حقیقی منفی هستند. براساس نظریه پایداری برای حالت ضرایب ثابت (یا Λ ثابت)، فرض خواهیم کرد که اثر تمام اختلالات در داده‌های اولیه وقتی $x \rightarrow \infty$ ، از بین می‌روند. ولی در واقع، قسمت همگن (۱۰.۹.۶) جوابهای بیکران خواهد داشت. بنابراین اختلالاتی در مقادیر اولیه وجود دارند که موجب جوابهای اختلال یافته بیکران در (۱۰.۹.۶) می‌شوند. این موضوع، اعتبار استفاده از معادله نمونه $y' = \lambda y + g(x)$ را زیر سؤال می‌برد. کاربرد آن روشهایی را می‌طلبد که ممکن است خواهیم بیشتر درباره آنها مطالعه کنیم، ولی این روش بخودی خود، برای دربرگرفتن انواع فراوان مسائل خطی و غیرخطی کافی نخواهد بود. مثال (۱۰.۹.۶) از ایکن^۱، (۱۹۸۵، ص ۲۶۹) گرفته شده است.

حل روش تفاضلات متناهی ما این مسأله را با در نظر گرفتن روش پسر و اوایل نشان می‌دهیم:

$$y_{n+1} = y_n + hf(x_{n+1}, y_{n+1}) \quad n \geq 0 \quad (13.9.6)$$

اگر فرمول معمولی بارستی

$$y_{n+1}^{(j+1)} = y_n + hf(x_{n+1}, y_{n+1}^{(j)}) \quad j \geq 0 \quad (14.9.6)$$

به‌کار رود، آنگاه

$$y_{n+1} - y_{n+1}^{(j+1)} \doteq h \frac{\partial f(x_{n+1}, y_{n+1})}{\partial y} [y_{n+1} - y_{n+1}^{(j)}]$$

برای همگرایی نیاز داریم که داشته باشیم

$$\left| h \frac{\partial f(x_{n+1}, y_{n+1})}{\partial y} \right| < 1 \quad (۱۵.۹.۶)$$

ولی در معادله‌های سرسخت، لازمهٔ این رابطه بسیار کوچک بودن h است، چیزی که ما می‌خواهیم از آن پرهیز نماییم. بنابراین روش ریشه‌یابی دیگری باید برای پیدا کردن y_{n+1} در (۱۳.۹.۶) به کار بریم.

متداولترین روشها برای حل (۱۳.۹.۶) روشهایی هستند که بر پایهٔ روش نیوتن استوارند. برای یک معادلهٔ دیفرانسیل، روش نیوتن برای پیدا کردن y_{n+1} چنین است

$$y_{n+1}^{(j+1)} = y_{n+1}^{(j)} - [1 - hf_y(x_{n+1}, y_{n+1}^{(j)})]^{-1} [y_{n+1}^{(j)} - y_n - hf(x_{n+1}, y_{n+1}^{(j)})] \quad (۱۶.۹.۶)$$

برای $z \geq 0$ یک حدس اولیهٔ خام $y_{n+1}^{(0)} = y_n$ است، اگرچه معمولاً این جواب را می‌توان بهتر کرد. برای دستگاه معادلات دیفرانسیل، روش نیوتن خیلی گران می‌شود. برای کاهش هزینه، ماتریس

$$I - hf_y(x_{n+1}, z) \quad z \doteq y_n \quad \text{برای مقداری از} \quad (۱۷.۹.۶)$$

را برای جمیع مقادیر z و برای چندین مقدار پی‌درپی n به کار می‌بریم. بنابراین روش نیوتن [بخش ۱۱.۲ را ببینید] برای حل دستگاه نمونهٔ (۱۳.۹.۶) با رابطهٔ زیر تقریب زده می‌شود

$$\begin{aligned} [I - hf_y(x_{n+1}, z)] \delta^{(j)} &= y_{n+1}^{(j)} - y_n - hf(x_{n+1}, y_{n+1}^{(j)}) \\ y_{n+1}^{(j+1)} &= y_{n+1}^{(j)} + \delta^{(j)} \end{aligned} \quad (۱۸.۹.۶)$$

برای $z \geq 0$ این مسأله به حل چند دستگاه خطی که یک ماتریس ضرایب دارند، می‌انجامد. این کار را در مقایسه با وقتی که ماتریس اصلاح شده است، بسیار آسانتر می‌توان انجام داد (مطالب بخش ۱.۸ را ببینید). ماتریس رابطهٔ (۱۷.۹.۶) باید متناوباً روزآمد شود، ولی در مقایسه با روش دقیق نیوتن، صرفه‌جویی هنوز هم خیلی مهم است. برای بحث بیشتر این موضوع ایکن (۱۹۸۵) و گوپتا و همکاران (۱۹۸۵، صص ۲۲-۲۵) را ببینید. برای یک بازنگری برنامه‌های رایانه‌ای در حل معادلات دیفرانسیل سرسخت، ایکن (۱۹۸۵، فصل ۴) را ببینید.

روش خطوط معادله دیفرانسیل جزئی سهموی زیر را در نظر می‌گیریم

$$U_t = U_{xx} + G(x, t) \quad 0 < x < 1 \quad t > 0 \quad (۱۹.۹.۶)$$

$$U(0, t) = d_0(t) \quad U(1, t) = d_1(t) \quad t \geq 0 \quad (۲۰.۹.۶)$$

$$U(x, 0) = f(x) \quad 0 \leq x \leq 1 \quad (۲۱.۹.۶)$$

نمادهای U_t و U_{xx} به ترتیب مشتقات جزئی نسبت به t و x را نشان می‌دهند. تابع مجهول $U(x, t)$ به زمان t و متغیر مکانی x بستگی دارد. شرایط (۲۰.۹.۶)، شرایط مرزی خوانده می‌شوند، و (۲۱.۹.۶) شرط اولیه نامیده می‌شود. جواب U را می‌توان دمای یک میله عایق‌دار به طول ۱ تعبیر نمود، که $U(x, t)$ دما در مکان x و زمان t است؛ بنابراین (۱۹.۹.۶) اغلب معادله گرما نامیده می‌شود. توابع G ، d_0 ، d_1 و f داده شده و هموار فرض می‌شوند. برای مطالعه مبانی نظری (۱۹.۹.۶) - (۲۱.۹.۶)، ویدر^۱ (۱۹۷۵) یا هر کتاب مقدماتی معمولی در معادلات دیفرانسیل جزئی را ببینید. ما در اینجا روش خطوط را برای حل U می‌آوریم، روشی عددی که در ده پانزده سال اخیر بسیار متداول شده است. این روش هم به حل یک دستگاه معادلات دیفرانسیل معمولی سرسخت می‌انجامد. گیریم $m > 0$ یک عدد صحیح باشد، $1/m$ را با δ نشان می‌دهیم و تعریف می‌کنیم،

$$x_j = j\delta \quad j = 0, 1, \dots, m$$

معادله (۱۹.۹.۶) را با تقریب‌زدن مشتق مکانی، گسسته می‌سازیم. فرمولهای (۱۷.۷.۵) و (۱۸.۷.۵) را برای تقریب‌زدن مشتقات مرتبه دوم به یاد می‌آوریم. با استفاده از این فرمولها،

$$U_{xx}(x_j, t) = \frac{U(x_{j+1}, t) - 2U(x_j, t) + U(x_{j-1}, t))}{\delta^2} - \frac{\delta^2}{12} \frac{\partial^2 U(\xi_j, t)}{\partial x^2}$$

برای $j = 1, 2, \dots, m-1$ با گذاردن در (۱۹.۹.۶) خواهیم داشت:

$$U_t(x_j, t) = \frac{U(x_{j+1}, t) - 2U(x_j, t) + U(x_{j-1}, t))}{\delta^2} + G(x_j, t) - \frac{\delta^2}{12} \frac{\partial^2 U(\xi_j, t)}{\partial x^2} \quad 1 \leq j \leq m-1 \quad (۲۲.۹.۶)$$

معادله (۱۹.۹.۶) را باید در هر نقطه گرهی داخلی x_j تقریب بزیم. مجهول ξ_j به $[x_{j-1}, x_{j+1}]$ تعلق دارد: $\xi_j \in [x_{j-1}, x_{j+1}]$.

جمله آخر در (۲۲.۹.۶)، خطای برشی در مشتگیری عددی، را حذف می‌کنیم. با تحمیل تساوی در معادله تقریبی حاصل شده، به دست می‌آوریم

$$u'_j(t) = \frac{1}{\delta^2} [u_{j+1}(t) - 2u_j(t) + u_{j-1}(t)] + G(x_j, t) \quad (23.9.6)$$

برای $j = 1, 2, \dots, m-1$. توابع $u_j(t)$ تقریبهای $U(x_j, t)$ ، $1 \leq j \leq m-1$ ، در نظر گرفته شده‌اند. این فرمول تقریب (۱۹.۹.۶) با روش خطوط است، و یک دستگاه $m-1$ معادله دیفرانسیل معمولی است. توجه می‌کنیم که $u_m(t)$ و $u_0(t)$ که به‌ازای $j = m-1$ و $j = 1$ در (۲۳.۹.۶) مورد نیازند، با استفاده از (۲۰.۹.۶) به دست آمده‌اند:

$$u_0(t) = d_0(t) \quad u_m(t) = d_1(t) \quad (24.9.6)$$

شرط اولیه برای (۲۳.۹.۶) به توسط (۲۱.۹.۶) داده شده است:

$$u_j(0) = f(x_j) \quad 1 \leq j \leq m-1 \quad (25.9.6)$$

نامگذاری روش خطوط از حل $U(x, t)$ در طول خطوط (x_j, t) ، $1 \leq j \leq m-1$ ، $t \geq 0$ ، در صفحه (x, t) گرفته شده است.

تحت مفروضات همواری مناسب در توابع d_0 ، d_1 ، G و f ، می‌توان نشان داد که

$$\text{Max}_{\substack{0 \leq j \leq m \\ 0 \leq t \leq T}} |U(x_j, t) - u_j(t)| \leq C_T \delta^2 \quad (26.9.6)$$

بنابراین برای تکمیل فرایند حل، فقط لازم است دستگاه (۲۳.۹.۶) را حل کنیم.

راحت‌تر است که (۲۳.۹.۶) را به شکل ماتریسی بنویسیم. نمادهای زیر را وارد می‌کنیم

$$\mathbf{u}(t) = [u_1(t), \dots, u_{m-1}(t)]^T \quad \mathbf{u}_0 = [f(x_1), \dots, f_{m-1}(x)]^T$$

$$\mathbf{g}(t) = \left[\frac{1}{\delta^2} d_0(t) + G(x_1, t), G(x_2, t), \dots, G(x_{m-2}, t), \frac{1}{\delta^2} d_1(t) + G(x_{m-1}, t) \right]^T$$

$$\Lambda = \frac{1}{\delta^2} \begin{bmatrix} -2 & 1 & 0 & 0 & \dots & 0 \\ 1 & -2 & 1 & 0 & & \vdots \\ 0 & 1 & -2 & 1 & & \\ \vdots & & & \ddots & & 0 \\ & & & & 1 & -2 & 1 \\ 0 & 0 & \dots & 0 & 1 & -2 \end{bmatrix} \quad (27.9.6)$$

ماتریس Λ از مرتبه $m - 1$ است. در تعریفهای \mathbf{u} و \mathbf{g} ، اندیس فوقانی T ترانزاده ماتریس را نشان می‌دهد، لذا \mathbf{u} و \mathbf{g} بردارهای ستونی به طول $m - 1$ هستند. با استفاده از این ماتریسها، معادله (۲۳.۹.۶) - (۲۵.۹.۶) را می‌توان دوباره به شکل زیر نوشت

$$\mathbf{u}'(t) = \Lambda \mathbf{u}(t) + \mathbf{g}(t) \quad \mathbf{u}(0) = \mathbf{u}. \quad (28.9.6)$$

اگر از روش اولر استفاده کنیم، روش عددی زیر را خواهیم داشت:

$$\mathbf{V}_{n+1} = \mathbf{V}_n + h[\Lambda \mathbf{V}_n + \mathbf{g}(t_n)] \quad \mathbf{V}_0 = \mathbf{u}. \quad (29.9.6)$$

که در آن $t_n = nh$ و $\mathbf{V}_n \doteq \mathbf{u}(t_n)$. این روش یک روش عددی معروفی برای معادله گرماست، که روش صریح ساده نامیده می‌شود. پایداری (۲۹.۹.۶) و چند روش دیگر برای حل (۲۸.۹.۶) را بررسی می‌نماییم.

معادله (۲۸.۹.۶) به شکل معادله نمونه (۱.۹.۶) است، و بنابراین به‌ویژه مقادیرهای λ برای امتحان سرسختی دستگاه نیاز داریم. می‌توان نشان داد که این ویژه مقادیرها همگی حقیقی‌اند و با رابطه زیر داده می‌شوند:

$$\lambda_j = -\frac{4}{\delta^2} \sin^2 \left(\frac{j\pi}{2m} \right) \quad 1 \leq j \leq m - 1 \quad (30.9.6)$$

اثبات این را به مسأله ۶ فصل ۷ احاله می‌نماییم. اگر مستقیماً این فرمول را امتحان کنیم، داریم

$$\lambda_{m-1} \leq \lambda_j \leq \lambda_1 \quad (31.9.6)$$

$$\lambda_{m-1} = -\frac{4}{\delta^2} \sin^2 \left(\frac{(m-1)\pi}{2m} \right) \doteq -\frac{4}{\delta^2}$$

$$\lambda_1 = -\frac{4}{\delta^2} \sin^2 \left(\frac{\pi}{2m} \right) \doteq -\pi^2$$

این تقریبها برای مقادیر بزرگ m هم معتبرند. می‌توان دید که اگر δ کوچک باشد، (۲۸.۹.۶) یک دستگاه سرسخت خواهد بود.

با استفاده از (۳۱.۹.۶) و (۴۹.۸.۶) برای تحلیل پایداری (۲۹.۹.۶)، داریم

$$|1 + h\lambda_j| < 1 \quad j = 1, \dots, m - 1$$

با به‌کار بردن (۳۰.۹.۶)، این رابطه ما را به حکم هم‌ارز زیر می‌رساند

$$0 < \frac{4h}{\delta^2} \sin^2 \left(\frac{j\pi}{2m} \right) < 2 \quad 1 \leq j \leq m - 1$$

این رابطه وقتی برقرار می‌شود که $2 \leq 4 \frac{h}{\delta^2}$ یا

$$h \leq \frac{1}{4} \delta^2 \quad (۳۲.۹.۶)$$

اگر δ کوچک، مثلاً $\delta = 0.1$ ، باشد، آنگاه مرحله زمانی h باید خیلی کوچک انتخاب شود تا پایداری به دست آید.

برخلاف این محدودیت (۳۲.۹.۶) با روش اویلر، روش پسرو اویلر هیچ محدودیتی ندارد زیرا A - پایدار است. این روش چنین می‌شود

$$\mathbf{V}_{n+1} = \mathbf{V}_n + h[\Delta \mathbf{V}_{n+1} + g(t_{n+1})] \quad \mathbf{V}_0 = \mathbf{u}. \quad (۳۳.۹.۶)$$

برای حل این دستگاه خطی بر حسب \mathbf{V}_{n+1} داریم،

$$(\mathbf{I} - h\Delta)\mathbf{V}_{n+1} = \mathbf{V}_n + hg(t_{n+1}) \quad (۳۴.۹.۶)$$

این یک دستگاه سه قطری معادلات خطی است (بخش ۳.۸ را ببینید). این دستگاه را می‌توان خیلی سریع حل کرد، با تقریباً $5m$ عمل حساب در هر مرحله زمانی، صرفنظر از بهای محاسبه طرف راست (۳۴.۸.۶). هزینه حل روش (۲۹.۹.۶) اویلر تقریباً به همین اندازه است، و بنابراین حل (۳۴.۹.۶)، زیاد وقت‌گیر نیست.

مثال مسأله معادله دیفرانسیل جزئی (۱۹.۹.۶) - (۲۱.۹.۶) را با توابع G ، d_1 ، d_0 و f که از جواب معلوم مسأله

$$U = e^{-\alpha t} \sin(\pi x) \quad 0 \leq x \leq 1 \quad t \geq 0 \quad (۳۵.۹.۶)$$

به دست می‌آیند، حل کنید. نتایج روش (۲۹.۹.۶) اویلر در جدول ۲۰.۶ و نتایج روش (۳۳.۹.۶) پسرو اویلر در جدول ۲۱.۶ داده شده‌اند.

برای روش اویلر، m را برابر ۴، ۸، ۱۶ می‌گیریم و برای حفظ پایداری، h را با توجه به (۳۲.۹.۶) مساوی $\delta^2/2$ اختیار می‌کنیم. باید توجه کرد که این انتخاب به ترتیب به طول مرحله‌های زمانی $h=0.31$ ، $h=0.078$ و $h=0.020$ می‌انجامد. از (۲۶.۹.۶) و فرمول خطا برای روش اویلر، انتظار داریم که خطا متناسب با δ^2 باشد، زیرا $h = \delta^2/2$. این امر ایجاب می‌کند که وقتی m دو برابر می‌شود، خطا بر ۴ تقسیم شود، و نتایج جدول ۲۰.۶، این امر را تأیید می‌کند. در این جدول، ستون خطا معرف خطای ماکسیمم در نقاط گرهی (x_j, t) ، $0 \leq j \leq n$ ، برای زمان داده شده t است.

جدول ۲۰.۶ روش خطوط: روش اویلر

t	خطا		خطا		خطا	
	m = 4	نسبت	m = 8	نسبت	m = 16	نسبت
۱.۰	۳.۸۹E-۲	۴.۰۹	۹.۵۲E-۳	۴.۰۲	۲.۳۷E-۳	۴.۰۲
۲.۰	۳.۱۹E-۲	۴.۰۹	۷.۷۹E-۳	۴.۰۲	۱.۹۴E-۳	۴.۰۲
۳.۰	۲.۶۱E-۲	۴.۰۹	۶.۳۸E-۳	۴.۰۱	۱.۵۹E-۳	۴.۰۱
۴.۰	۲.۱۴E-۲	۴.۱۰	۵.۲۲E-۳	۴.۰۲	۱.۳۰E-۳	۴.۰۲
۵.۰	۱.۷۵E-۲	۴.۰۹	۴.۲۸E-۳	۴.۰۴	۱.۰۶E-۳	۴.۰۴

جدول ۲۱.۶ روش خطوط: روش پسر اویلر

t	خطا		خطا		خطا	
	m = 4	نسبت	m = 8	نسبت	m = 16	نسبت
۱.۰	۴.۴۵E-۲	۴.۰۹	۱.۱۰E-۲	۴.۰۲	۲.۸۶E-۳	۴.۰۲
۲.۰	۳.۶۵E-۲	۴.۰۹	۹.۰۱E-۳	۴.۰۲	۲.۳۴E-۳	۴.۰۲
۳.۰	۲.۹۹E-۲	۴.۰۹	۷.۳۷E-۳	۴.۰۱	۱.۹۲E-۳	۴.۰۱
۴.۰	۲.۴۵E-۲	۴.۱۰	۶.۰۴E-۳	۴.۰۲	۱.۵۷E-۳	۴.۰۲
۵.۰	۲.۰۰E-۲	۴.۰۹	۴.۹۴E-۳	۴.۰۴	۱.۲۹E-۳	۴.۰۴

برای حل (۲۸.۹.۶) با روش پسر اویلر، دیگر هیچ نیازی به ارتباط بین مرحله مکانی δ و مرحله زمانی h ، نخواهد بود. با بررسی فرمول خطای (۲۶.۹.۶) در روش خطوط و فرمول خطای برشی (۷.۹.۶) ($p = 1$) می‌گیریم) برای روش پسر اویلر، ملاحظه می‌کنیم که خطا در حل مسأله (۱۹.۹.۶) - (۲۱.۹.۶) متناسب با $h + \delta^2$ است. برای تابع مجهول U ی (۳۴.۹.۶)، یک تغییر آهسته با تغییر t وجود دارد. بنابراین برای خطای برشی مربوط به انتگرالگیری زمانی، باید بتوانیم مرحله زمانی نسبتاً بزرگی برای h در مقایسه با مرحله مکانی δ ، انتخاب نماییم تا دو منبع خطا از نظر اندازه نسبتاً مساوی باشند. در جدول ۲۱.۶، از $h = 0.1$ و $m = 4, 8, 16$ استفاده کرده‌ایم. توجه نمایید که این مرحله زمانی بسیار بزرگتر از آن است که در جدول ۲۰.۶ برای روش اویلر به کار برده شده است و بنابراین روش پسر اویلر، برای این مسأله خاص بسیار کارا تر است.

برای بحث بیشتر درباره روش خطوط، ایکن (۱۹۸۵)، صص ۱۲۴-۱۴۸) را ببینید. برای بعضی برنامه‌های روش خطوط در حل دستگاه معادلات دیفرانسیل جزئی سهموی غیرخطی، با یک و دو متغیر مکانی، سینکاوک و مدسین^۱ (۱۹۷۵) و ملگارت^۲ و سینکاوک (۱۹۸۱) را ببینید.)

1. Sincovec & Madsen

2. Malgaard

۱۰.۶ روشهای تک گامی و روشهای رونگه - کوتا

روشهای تک گامی برای حل $y' = f(x, y)$ فقط نیاز به دانستن جواب عددی y_n دارند تا مقدار بعدی y_{n+1} را محاسبه نمایند. این امر برای روشهای تک گامی مزایای روشنی نسبت به روشهای $-p$ گامی که از چندین مقدار قبلی $\{y_n, \dots, y_{n-p+1}\}$ استفاده می کنند به وجود می آورد، زیرا مقادیر اولیه $\{y_1, \dots, y_{p-1}\}$ در روشهای چندگامی باید با روش عددی دیگری محاسبه شوند. یکی از معروفترین روشهای تک گامی، روشهای رونگه - کوتاست. برنامه نویسی آنها نسبتاً آسان و کنترل خطای برشی آنها در مقایسه با روشهای چند گامی، بسیار ساده است. برای روشهای چندگامی با مرتبه ثابت، که غالباً در گذشته به کار می رفت، روشهای رونگه - کوتا وسیله عادی برای محاسبه مقادیر اولیه مورد نیاز این روشهای چند گامی، بودند. بزرگترین عیب روشهای رونگه - کوتا آن است که، در مقایسه با روشهای چندگامی، برای رسیدن به یک دقت مفروض، به محاسبه خیلی بیشتر $f(x, y)$ نیاز دارد. بعداً بعضی نتایج مقایسه برنامه های متغیر - مرتبه آدامز و برنامه های ثابت - مرتبه رونگه - کوتا را یادآوری می نمایم.

ساده ترین روش تک گامی بر پایه استفاده از سری تیلر میتنی است. فرض کنید $Y(x)$ ، $(r+1)$ بار پیوسته مشتق پذیر باشد که $Y(x)$ جواب مسأله مقدار اولیه

$$y' = f(x, y) \quad y(x_0) = Y_0 \quad (۱۰.۱۰.۶)$$

است. با استفاده از قضیه تیلر، $Y(x_1)$ را حول x_0 بسط می دهیم:

$$Y(x_1) = Y(x_0) + hY'(x_0) + \dots + \frac{h^r}{r!} Y^{(r)}(x_0) + \frac{h^{r+1}}{(r+1)!} Y^{(r+1)}(\xi) \quad (۲.۱۰.۶)$$

به ازای مقداری چون $x_0 \leq \xi \leq x_1$. با حذف جمله باقیمانده، یک تقریب برای $Y(x_1)$ خواهیم داشت به شرطی که بتوانیم $Y''(x_0), \dots, Y^{(r)}(x_0)$ را حساب کنیم. از $Y'(x) = f(x, Y(x))$ مشتق می گیریم تا به دست آوریم

$$Y''(x) = f_x(x, Y(x)) + f_y(x, Y(x))Y'(x),$$

$$Y''(x) = f_x + f_y f$$

به همین ترتیب ادامه می دهیم تا مشتقات مراتب بالاتر $Y(x)$ را به دست آوریم.

جدول ۲۲.۶ مثال روش (۳.۱۰.۶) سری تیلر

		$h = 0.0625$	$h = 0.125$	
x	$y_h(x)$	$Y(x) - y_h(x)$	$Y(x) - y_h(x)$	نسبت
۲.۰	۰.۳۳۳۶۴۹	-۳.۲E-۴	-۱.۴E-۳	۴.۴
۴.۰	۰.۲۰۰۱۳۵	-۱.۴E-۴	-۵.۹E-۴	۴.۳
۶.۰	۰.۱۴۲۹۳۱	-۷.۴E-۵	-۳.۲E-۴	۴.۳
۸.۰	۰.۱۱۱۱۵۷	-۴.۶E-۵	-۲.۰E-۴	۴.۳
۱۰.۰	۰.۰۹۰۹۴۱	-۳.۱E-۵	-۱.۴E-۴	۴.۳

مثال مسأله زیر را با جواب $Y(x) = 1/(1+x)$ در نظر می‌گیریم،

$$y' = -y^2 \quad y(0) = 1$$

پس $Y'' = -2YY' = 2Y^3$ و (۲.۱۰.۶) با $r = 1$ به دست می‌دهد

$$Y(x_1) = Y_0 - hY_0^2 + h^2Y_0^3 + \frac{h^3}{6}Y^{(3)}(\xi_0) \quad x_0 \leq \xi_0 \leq x_1$$

جمله باقیمانده را حذف کرده‌ایم تا تقریبی برای $Y(x_1)$ به دست آوریم. به همین ترتیب، از این مقدار می‌توان برای محاسبه تقریب $Y(x_2)$ و غیره نیز استفاده کرد. روش عددی چنین است:

$$y_{n+1} = y_n - hy_n^2 + h^2y_n^3 \quad n \geq 0 \quad (3.10.6)$$

جدول ۲۲.۶ شامل خطاهای این جواب عددی در یک مجموعه انتخابی از نقاط گرهی است. برای فاصله نقاط در هر شبکه، $h = 0.125$ و $h = 0.0625$ به کار برده شده است و نسبت خطاهای به دست آمده نیز داده شده است. توجه نمایید وقتی h نصف می‌شود، نسبت تقریباً ۴ می‌شود. این را می‌توان به طور نظری توجیه کرد زیرا می‌توان نشان داد که نرخ همگرایی $O(h^2)$ است، برهان مشابه همان برهانی است که در قضیه ۳.۶ داده شده است یا در قضیه ۹.۶ بعداً در این بخش داده خواهد شد.

روش سری تیلر می‌تواند نتایج عالی به دست دهد. ولی استفاده از آن، به دلیل مشتقگیری تحلیلی $f(x, y)$ ، پردردسراست. محاسبه مشتقات ممکن است خیلی مشکل و محاسبه آنها به ویژه در دستگاه معادلات و تغییر باشد. این مشتقگیرها را می‌توان با استفاده از زبان عملیات نمادی با رایانه انجام داد، سپس به سادگی روش سری تیلر را تولید کرد. مع هذا، احتمالاً باز هم محاسبه مشتقات وقتگیر است، و به نظر می‌آید که روشهایی که بر پایه محاسبه فقط $f(x, y)$ هستند کارا تر باشند. برای پیروی از روش سری تیلر، وقتی فقط $f(x, y)$ محاسبه می‌شود، به فرمولهای رونگه - کوتا برمی‌گردیم.

روشهای رونگه- کوتا روشهای رونگه- کوتا دقیقاً به بسط سری تیلر $Y(x)$ در (۲.۱۰.۶) وابسته اند ولی مشتقگیری f در استفاده از این روشها هرگز لازم نیست. برای آسانی نمادگذاری، روش رونگه- کوتا را مختصراً با RK نمایش می دهیم. تمام روشهای RK بدین شکل نوشته می شوند:

$$y_{n+1} = y_n + hF(x_n, y_n, h; f) \quad n \geq 0 \quad (۴.۱۰.۶)$$

با مثالهایی از تابع F آغاز می کنیم، و بعداً فرضهای مربوط به آن را مورد بحث قرار می دهیم. ولی در اینجا، باید متوجه باشیم برای تمام مقادیر کوچک h می خواهیم داشته باشیم

$$F(x, Y(x), h; f) \doteq Y'(x) = f(x, Y(x)) \quad (۵.۱۰.۶)$$

خطای برشی برای (۴.۱۰.۶) را با

$$T_n(Y) = Y(x_{n+1}) - Y(x_n) - hF(x_n, Y(x_n), h; f) \quad n \geq 0 \quad (۶.۱۰.۶)$$

تعریف و $\tau_n(Y)$ را به طور ضمنی با رابطه زیر تعریف می کنیم

$$T_n(Y) = h\tau_n(Y)$$

رابطه (۶.۱۰.۶) را دوباره مرتب می کنیم و به دست می آوریم

$$Y(x_{n+1}) = Y(x_n) + hF(x_n, Y(x_n), h; f) + h\tau_n(Y) \quad n \geq 0 \quad (۷.۱۰.۶)$$

در قضیه ۹.۶، این رابطه با (۴.۱۰.۶) مقایسه می شود تا همگرایی $\{y_n\}$ به Y ثابت شود.

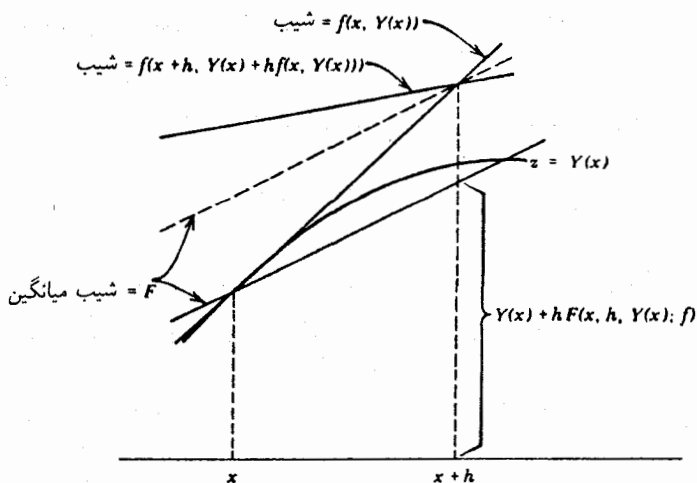
مثال ۱. روش ذوزنقهی را در نظر می گیریم و آن را با یک بارست با استفاده از روش اویلر به عنوان پیشگو حل می کنیم.

$$y_{n+1} = y_n + \frac{h}{2} [f(x_n, y_n) + f(x_{n+1}, y_n + hf(x_n, y_n))] \quad n \geq 0 \quad (۸.۱۰.۶)$$

در نمادگذاری (۴.۱۰.۶)،

$$F(x, y, h; f) = \frac{1}{2} [f(x, y) + f(x + h, y + hf(x, y))]$$

همان طور که در شکل ۹.۶ می توان دید، F یک شیب میانگین $Y(x)$ در $[x, x + h]$ است.



شکل ۹.۶ نمایش روش (۸.۱۰.۶) رونگه-کوتا

۲. روش زیر نیز بر پایهٔ به‌دست آوردن یک شیب میانگین برای جواب در $[x_n, x_{n+1}]$ است:

$$y_{n+1} = y_n + hf(x_n + \frac{1}{4}h, y_n + \frac{1}{4}hf(x_n, y_n)) \quad n \geq 0 \quad (9.10.6)$$

در این حالت

$$F(x, y, h; f) = f(x + \frac{1}{4}h, y + \frac{1}{4}hf(x, y))$$

پیدا کردن فرمولی برای خطای برشی به پیدا کردن این روشها مربوط است، و این مطلب برای روشهای RK از مراتب بالاتر نیز درست است. پیدا کردن روشهای RK با به‌دست آوردن یک خانواده از فرمولهای مرتبهٔ دوم شامل (۸.۱۰.۶) و (۹.۱۰.۶) روشن می‌شود. فرض می‌کنیم که F شکل کلی زیر را دارد

$$F(x, y, h; f) = \gamma_1 f(x, y) + \gamma_2 f(x + \alpha h, y + \beta hf(x, y)) \quad (10.10.6)$$

که در آن چهار ثابت $\gamma_1, \gamma_2, \alpha, \beta$ باید معین شوند.

از قضیهٔ تیلر (قضیهٔ ۵.۱) برای توابع دو متغیره استفاده کرده عبارت دوم طرف راست (۱۰.۱۰.۶) را تا جملات شامل مشتقات مرتبهٔ دوم بسط می‌دهیم. این بسط به ما خواهد داد

$$F(x, y, h; f) = \gamma_1 f(x, y) + \gamma_2 \{ f(x, y) + h[\alpha f_x + \beta f f_y] + h^2 [\frac{1}{2} \alpha^2 f_{xx} + \alpha \beta f_{xy} f + \frac{1}{2} \beta^2 f^2 f_{yy}] \} + O(h^3) \quad (11.10.6)$$

همچنین به بعضی از مشتقات $Y'(x) = f(x, Y(x))$ نیاز داریم، یعنی به

$$Y'' = f_x + f_y f$$

$$Y^{(r)} = f_{xx} + 2f_{xy}f + f_{yy}f^2 + f_y f_x + f_y^2 f \quad (۱۲.۱۰.۶)$$

برای خطای برشی،

$$\begin{aligned} T_n(Y) &= Y(x_{n+1}) - Y(x_n) - hF(x_n, Y(x_n), h; f) \\ &= hY'_n + \frac{h^2}{2}Y''_n + \frac{h^3}{6}Y^{(r)}_n + O(h^4) - hF(x_n, Y_n, h; f) \end{aligned}$$

اگر مقادیر (۱۱.۱۰.۶) و (۱۲.۱۰.۶) را در این عبارت قرار دهیم و برحسب توانهای h مرتب کنیم به دست می آوریم

$$\begin{aligned} T_n(Y) &= h[1 - \gamma_1 - \gamma_2]f + h^2[(\frac{1}{2} - \gamma_2\alpha)f_x + (\frac{1}{2} - \gamma_2\beta)f_y f] \\ &+ h^3[(\frac{1}{6} - \frac{1}{2}\gamma_2\alpha^2)f_{xx} + (\frac{1}{3} - \gamma_2\alpha\beta)f_{xy}f + (\frac{1}{6} - \frac{1}{2}\gamma_2\beta^2)f_{yy}f^2 \\ &+ \frac{1}{6}f_y f_x + \frac{1}{6}f_y^2 f] + O(h^4) \quad (۱۳.۱۰.۶) \end{aligned}$$

تمام مشتقات در (x_n, Y_n) محاسبه شده اند.

مایلم که خطای برشی در حد ممکن، با سرعت بیشتری به صفر همگرا شود. اگر f مجاز باشد به طور دلخواه تغییر نماید، ضریب h^3 نمی تواند در حالت کلی صفر شود. لازمه اینکه ضرایب h^2 و h صفر شوند چنین خواهد شد

$$\gamma_1 + \gamma_2 = 1 \quad \gamma_2\alpha = \frac{1}{2} \quad \gamma_2\beta = \frac{1}{2} \quad (۱۴.۱۰.۶)$$

و از آنجا نتیجه می شود

$$T_n(Y) = O(h^3)$$

دستگاه (۱۴.۱۰.۶) نامعین است و جواب عمومی آن

$$\gamma_1 = 1 - \gamma_2 \quad \alpha = \beta = \frac{1}{2\gamma_2} \quad (۱۵.۱۰.۶)$$

است، که γ_2 دلخواه است. هر دو حالت (۸.۱۰.۶) (با $\gamma_2 = 1/2$) و (۹.۱۰.۶) (با $\gamma_2 = 1$) حالت های خاص این جواب اند.

با گذاردن این مقادیر در (۱۳.۱۰.۶)، می‌توانیم جمله اصلی خطای برشی را که به γ_2 بستگی دارد پیدا کنیم. در بعضی حالات، در حالی که f می‌تواند به طور دلخواه تغییر کند، مقدار γ_2 طوری انتخاب می‌شود که ضریب h^3 در حد ممکن کوچک باشد. برای مثال، اگر (۱۳.۱۰.۶) را به شکل

$$T_n(Y) = c(f, \gamma_2)h^2 + O(h^4) \quad (16.10.6)$$

بنویسیم، نابرابری کوشی - شوارتس [(۸.۱.۷) را در فصل ۷ ببینید] را می‌توان به‌کار برده و نشان داد

$$|c(f, \gamma_2)| \leq c_1(f)c_2(\gamma_2) \quad (17.10.6)$$

که در آن

$$c_1(f) = [f_{xx}^2 + f_{xy}^2 f_y^2 + f_{yy}^2 f_x^2 + f_{yx}^2 f_x^2 + f_{yy}^2 f_x^2]^{1/2}$$

$$c_2(\gamma_2) = \left[\left(\frac{1}{6} - \frac{1}{4} \gamma_2 \alpha^2 \right)^2 + \left(\frac{1}{3} - \gamma_2 \alpha \beta \right)^2 + \left(\frac{1}{6} - \frac{1}{4} \gamma_2 \beta^2 \right)^2 + \frac{1}{18} \right]^{1/2}$$

که β و α در (۱۵.۱۰.۶) داده شده‌اند. مینیمم مقدار $c_2(\gamma_2)$ به‌ازای $\gamma_2 = 0.75$ حاصل می‌شود و $c_2(0.75) = 1/\sqrt{18}$ و روش عددی مرتبهٔ دوم به‌دست آمده چنین خواهد شد

$$y_{n+1} = y_n + \frac{h}{4} \left[f(x_n, y_n) + 3f\left(x_n + \frac{2}{3}h, y_n + \frac{2}{3}hf(x_n, y_n)\right) \right] \quad n \geq 0 \quad (18.10.6)$$

این روش بهینه است به این معنی که $c_2(\gamma_2)$ ، ضریب جملهٔ $h^3 c_1(f)h^3$ در خطای برشی را مینیمم می‌نماید. برای بحث مفصل از راه تحلیل خطای برشی در روشهای RK، شمایین (۱۹۸۶) را ببینید. فرمولهای مراتب بالاتر را می‌توان ساخت و به روشی مشابه تحلیل نمود، اگرچه محاسبات جبری بسیار پیچیده خواهند شد. فرض کنید فرمولی برای $F(x, y, h; f)$ به شکل زیر داده شده است

$$F(x, y, h; f) = \sum_{j=1}^p \gamma_j V_j \quad (19.10.6)$$

$$V_1 = f(x, y)$$

$$V_j = f\left(x + \alpha_j h, y + h \sum_{i=1}^{j-1} \beta_{ji} V_i\right) \quad j = 2, \dots, p \quad (20.10.6)$$

جدول ۲۳.۶ مرتبهٔ ماکسیمم روشهای رونگه - کوتا

تعداد ارزیابی‌های تابع	۱	۲	۳	۴	۵	۶	۷	۸
ماکسیمم مرتبهٔ روش	۱	۲	۳	۴	۴	۵	۶	۶

این ضرایب را می‌توان به گونه‌ای انتخاب نمود که جملات پیشرو در خطای برشی را صفر نمایند، همان‌گونه که برای (۱۰.۱۰.۶) و (۱۴.۱۰.۶) انجام شد. واضح است که بین تعداد محاسبات $f(x, y)$ که آن را p می‌نامیم، و مرتبهٔ ماکسیمم ممکن که می‌توان برای خطای برشی به دست آورد، ارتباطی وجود دارد. این ارتباط در جدول ۲۳.۶ داده شده است که بخشی از آن منسوب به بوچرا (۱۹۶۵) است. تا تقریباً ۱۹۷۰، متداولترین روش RK، احتمالاً فرمول کلاسیک اصلی بوده، که یک تعمیم روش سیمپسون است. این روش چنین است

$$y_{n+1} = y_n + \frac{h}{6} [V_1 + 2V_2 + 2V_3 + V_4] \quad (21.10.6)$$

$$V_1 = f(x_n, y_n) \quad V_2 = f\left(x_n + \frac{1}{4}h, y_n + \frac{1}{4}hV_1\right)$$

$$V_3 = f\left(x_n + \frac{1}{4}h, y_n + \frac{1}{4}hV_2\right) \quad V_4 = f(x_n + h, y_n + hV_3)$$

می‌توان ثابت کرد که این فرمول از مرتبهٔ چهار با $T_n(Y) = O(h^5)$ است. اگر $f(x, y)$ به y بستگی نداشته باشد، این فرمول به قاعدهٔ انتگرالگیری سیمپسون بدل می‌شود.

مثال مسأله

$$y' = \frac{1}{1+x^2} - 2y^2 \quad y(0) = 0 \quad (22.10.6)$$

با جواب $Y = x/(1+x^2)$ را در نظر می‌گیریم. روش (۲۱.۱۰.۶) با طول گام ثابت به کار برده شده و نتایج در جدول ۲۴.۶ داده شده‌اند. طول گامها $h = 0.25$ و $2h = 0.5$ بوده‌اند. ستون نسبت، نسبت خطا در نقاط گرهی مربوطه را، وقتی h نصف شده، نشان می‌دهد. ستون آخر مثالی است از فرمول (۲۴.۱۰.۶) از مطالب حاصل از برونمایی ریچاردسن. چون برای روش (۲۱.۱۰.۶) تساوی $T_n(Y) = O(h^5)$ برقرار است، قضیهٔ ۹.۶ ایجاب می‌کند که نرخ همگرایی $y_h(x)$ به $Y(x)$ باشد. مقدار نظری نسبت، ۱۶ است و وقتی h باز هم کاهش یابد، این نسبت به ۱۶ نزدیکتر می‌شود.

جدول ۲۴.۶ مثالی برای روش (۲۱.۱۰.۶) رونگه-کوتا

x	$y_h(x)$	$Y(x) - y_h(x)$	$Y(x) - y_{2h}(x)$	نسبت	$\frac{1}{15}[y_h(x) - y_{2h}(x)]$
۲.۰	۰.۳۹۹۹۵۶۹۹	۴.۳E-۵	۱.۰E-۳	۲۴	۶.۷E-۵
۴.۰	۰.۲۳۵۲۹۱۵۹	۲.۵E-۶	۷.۰E-۵	۲۸	۴.۵E-۶
۶.۰	۰.۱۶۲۱۶۱۷۹	۳.۷E-۷	۱.۲E-۵	۳۲	۷.۷E-۷
۸.۰	۰.۱۲۳۰۷۶۸۳	۹.۲E-۸	۳.۴E-۶	۳۶	۲.۲E-۷
۱۰.۰	۰.۰۹۹۰۰۹۸۷	۳.۱E-۸	۱.۳E-۶	۴۱	۸.۲E-۸

روشهای رونگه-کوتا، همانند روشهای چندگامی، دارای فرمولهای خطاهای مجانبی هستند. برای (۲۱.۱۰.۶)،

$$Y(x) - y_h(x) = D(x)h^4 + O(h^5) \quad (۲۳.۱۰.۶)$$

که $D(x)$ در مسأله مقدار اولیه معینی صدق می‌کند. برهان این قضیه، تعمیم برهان قضیه ۹.۶ و شبیه است به اشتقاق قضیه ۴.۶ برای روش اویلر، در بخش ۲.۶. نتیجه (۲۳.۱۰.۶) را می‌توان برای به‌دست آوردن یک برآورد خطا به‌کاربرد درست همان‌گونه که در روش ذوزنقه‌یی در فرمول (۲۶.۵.۶) از بخش ۵.۶ عمل شده بود. برای یک طول گام $2h$ ،

$$Y(x) - y_{2h}(x) = 16D(x)h^4 + O(h^5)$$

اگر مانند عمل قبلی برای (۲۷.۵.۶) عمل کنیم به‌دست می‌آوریم

$$Y(x) - y_h(x) = \frac{1}{15}[y_h(x) - y_{2h}(x)] + O(h^5) \quad (۲۴.۱۰.۶)$$

و اولین جمله سمت راست برآورد برای خطای سمت چپ است. این برآورد در آخرین ستون جدول ۲۴.۶ نشان داده شده است.

تحلیل همگرایی برای به‌دست آوردن همگرایی کلی (۴.۱۰.۶)، لازم است $\tau_n(Y) \rightarrow 0$ وقتی $h \rightarrow 0$ چون

$$\tau_n(Y) = \frac{Y(x_{n+1}) - Y(x_n)}{h} - F(x_n, Y(x_n), h; f)$$

می‌خواهیم

$$F(x, Y(x), h; f) \rightarrow Y'(x) = f(x, Y(x))$$

وقتی $h \rightarrow 0$ به عبارت دقیق‌تر، تعریف می‌کنیم

$$\delta(h) = \text{Max}_{\substack{x_0 \leq x \leq b \\ -\infty < y < \infty}} |f(x, y) - F(x, y, h; f)| \quad (25.10.6)$$

و فرض می‌کنیم

$$h \rightarrow 0 \quad \text{وقتی} \quad \delta(h) \rightarrow 0 \quad (26.10.6)$$

این شرط را گاهی اوقات شرط سازگاری روش (۴.۱۰.۶) رونگه - کوتا خوانند.

ما همچنین به یک شرط لیشیتس برای F نیاز داریم: به‌ازای تمام مقادیر $x_0 \leq x \leq b$ و $-\infty < y, z < \infty$ ، و برای جميع مقادیر کوچک $h > 0$

$$|F(x, y, h; f) - F(x, z, h; f)| \leq L |y - z| \quad (27.10.6)$$

این شرط معمولاً با استفاده از شرط (۱۲.۲.۶) لیشیتس برای $f(x, y)$ ، ثابت می‌شود. برای مثال، با روش (۹.۱۰.۶).

$$\begin{aligned} & |F(x, y, h; f) - f(x, z, h; f)| \\ &= \left| f\left(x + \frac{h}{\varphi}, y + \frac{h}{\varphi} f(x, y)\right) - f\left(x + \frac{h}{\varphi}, z + \frac{h}{\varphi} f(x, z)\right) \right| \\ &\leq K \left| y - z + \frac{h}{\varphi} [f(x, y) - f(x, z)] \right| \\ &\leq K \left(1 + \frac{h}{\varphi} K\right) |y - z| \end{aligned}$$

L را برابر $K(1 + \frac{1}{\varphi} K)$ برای $h \leq 1$ انتخاب می‌کنیم.

قضیه ۹.۶ فرض می‌کنیم روش (۴.۱۰.۶) رونگه - کوتا در شرط (۲۷.۱۰.۶) لیشیتس صدق می‌نماید. در این صورت برای مسأله مقدار آغازی (۱.۱۰.۶)، جواب $\{y_n\}$ در رابطه زیر صدق می‌کند

$$\text{Max}_{x_0 \leq x_n \leq b} |Y(x_n) - y_n| \leq e^{(b-x_0)L} |Y_0 - y_0| + \left[\frac{e^{(b-x_0)L} - 1}{L} \right] \tau(h) \quad (28.10.6)$$

که در آن

$$\tau(h) \equiv \text{Max}_{x_0 \leq x_n \leq b} |\tau_n(Y)| \quad (29.10.6)$$

اگر شرط سازگاری (۲۶.۱۰.۶) نیز برقرار باشد، جواب عددی $\{y_n\}$ به $Y(x)$ همگراست.

برهان رابطه (۴.۱۰.۶) را از (۷.۱۰.۶) کم می‌کنیم تا به دست آید

$$e_{n+1} = e_n + h[F(x_n, Y_n, h; f) - F(x_n, y_n, h; f)] + h\tau_n(Y) \quad (۳۰.۱۰.۶)$$

که در آن $e_n = Y(x_n) - y_n$. شرط لیبشیتس (۲۷.۱۰.۶) را به کار می‌بریم و از (۲۹.۱۰.۶) استفاده می‌کنیم تا رابطه زیر پیدا شود

$$|e_{n+1}| \leq (1 + hL) |e_n| + h\tau_n(h) \quad x_0 \leq x_N \leq b$$

با برهانی همانند برهان همگرایی برای روش اویلر، از رابطه فوق به آسانی به نتیجه (۲۸.۱۰.۶) می‌رسیم. [(۲۰.۲.۶) و (۲۲.۲.۶) را ببینید].

در بیشتر حالات، با محاسبه مستقیم معلوم می‌شود که وقتی $h \rightarrow 0$ ، $\tau(h) \rightarrow 0$ و در چنین حالتی، همگرایی $\{y_n\}$ به $Y(x)$ فوری ثابت می‌شود. ولی آنچه که لازم است بدانیم برقرار بودن (۲۶.۱۰.۶) است. برای ملاحظه این امر، می‌نویسیم

$$\begin{aligned} h\tau_n(Y) &= Y(x_{n+1}) - Y(x_n) - hF(x_n, Y(x_n), h; f) \\ &= hY'(x_n) + \frac{h^2}{2}Y''(\xi_n) - hF(x_n, Y(x_n), h; f) \\ h|\tau_n(Y)| &\leq h\delta(h) + \frac{h^2}{2}\|Y''\|_\infty \\ \tau(h) &\leq \delta(h) + \frac{1}{2}h\|Y''\|_\infty \end{aligned}$$

بنابراین وقتی $h \rightarrow 0$ ، $\tau(h) \rightarrow 0$ که برهان را کامل می‌کند.

فرع بعدی یک نتیجه بلافاصله (۲۸.۱۰.۶) است.

فرع اگر روش (۴.۱۰.۶) رونگه-کوتا دارای خطای برشی $T_n(Y) = O(h^{m+1})$ باشد، آنگاه نرخ همگرایی $\{y_n\}$ به $Y(x)$ برابر $O(h^m)$ است.

به دست آوردن یک فرمول خطای مجانبی برای روش (۴.۱۰.۶) رونگه-کوتا دشوار نیست مشروط بر آنکه یک چنین فرمولی برای خطای برشی معلوم باشد. فرض می‌کنیم

$$T_n(Y) = \varphi(x_n)h^{m+1} + O(h^{m+2}) \quad (۳۱.۱۰.۶)$$

که $\varphi(x)$ از روی $Y(x)$ و $f(x, Y(x))$ معلوم شده است. به عنوان یک مثال، نتیجه (۱۳.۱۰.۶) را در به دست آوردن این بسط برای روشهای مرتبه دوم RK ببینید. شکل‌های قویتر (۲۶.۱۰.۶) و (۲۷.۱۰.۶) نیز لازم‌اند. فرض می‌کنیم

$$F(x, y, h; f) - F(x, z, h; f) = \frac{\partial F(x, y, h; f)}{\partial y} (y - z) + O((y - z)^2) \quad (۳۲.۱۰.۶)$$

و همچنین فرض می‌کنیم

$$\delta_1(h) \equiv \max_{\substack{x \leq x \leq b \\ -\infty < y < \infty}} \left| \frac{\partial f(x, y)}{\partial y} - \frac{\partial F(x, y, h; f)}{\partial y} \right| \rightarrow 0 \quad h \rightarrow 0 \quad \text{وقتی} \quad (۳۳.۱۰.۶)$$

در عمل، هر دو این نتایج به سادگی اثبات می‌شوند. با این مفروضات، می‌توانیم فرمول زیر را به دست آوریم

$$Y(x) - y_h(x) = D(x)h^m + O(h^{m+1}) \quad (۳۴.۱۰.۶)$$

که $D(x)$ در مسأله مقدار آغازی زیر صدق می‌کند

$$D' = f_y(x, Y(x))D(x) + \varphi(x) \quad D(x_0) = 0 \quad (۳۵.۱۰.۶)$$

قضایای پایداری را، مشابه با قضایایی که برای روشهای چندگامی بیان شد، می‌توان برای روشهای رونگه - کوتا بیان کرد. نوع اساسی پایداری، که تقریباً در ابتدای بخش ۸.۶ تعریف شده، به سادگی ثابت می‌شود. برهان، یک تغییر ساده برهان پایداری روش (۲۹.۲.۶) اوایلر است، و ما آن را به مسأله ۴۹ واگذار می‌کنیم. یک تفاوت اساسی با نظریه چندگامی این است که با روشهای RK هیچ جواب انگلی ایجاد نمی‌شود؛ بنابراین مفهوم پایداری نسبی برای روشهای RK کاربرد ندارد. ناحیه‌های پایداری مطلق را مانند نظریه چندگامی می‌توان مطالعه کرد، ولی ما آن را در اینجا حذف و در ضمن مسائل مطرح می‌کنیم.

برآورد خطای برشی برای کنترل اندازه خطای برشی، ابتدا باید بتوانیم آن را برآورد کنیم. گیریم $u_n(x)$ معرف جواب $y' = f(x, y)$ گذرنده بر (x_n, y_n) باشد [۹.۵.۶] را ببینید]. می‌خواهیم خطا در $y_h(x_n + 2h)$ را نسبت به $u_n(x_n + 2h)$ برآورد نماییم.

با به کار بردن (۳۰.۱۰.۶) برای خطای $y_h(x)$ در مقایسه با $u_n(x)$ و استفاده از فرمول
مجانبی (۳۱.۱۰.۶)

$$\begin{aligned} u_n(x_{j+1}) - y_h(x_{j+1}) &= u_n(x_j) - y_h(x_j) \\ &+ h\{F(x_j, u_n(x_j), h; f) - F(x_j, y_h(x_j), h; f)\} \\ &+ \varphi_n(x_j)h^{m+1} + O(h^{m+2}) \quad j \geq n \end{aligned}$$

با توجه به این رابطه به سادگی می توان ثابت کرد که

$$\begin{aligned} u_n(x_n + h) - y_h(x_n + h) &= \varphi_n(x_n)h^{m+1} + O(h^{m+2}) \\ u_n(x_n + 2h) - y_h(x_n + 2h) &= 2\varphi_n(x_n)h^{m+1} + O(h^{m+2}) \end{aligned}$$

اگر همین شیوه را برای $y_{2h}(x_j)$ به کار ببریم،

$$u_n(x_n + 2h) - y_{2h}(x_n + 2h) = 2^{m+1}\varphi_n(x_n)h^{m+1} + O(h^{m+2})$$

از دو معادله اخیر،

$$\begin{aligned} u_n(x_n + 2h) - y_h(x_n + 2h) \\ = \frac{1}{2^m - 1} [y_h(x_n + 2h) - y_{2h}(x_n + 2h)] + O(h^{m+2}) \quad (۳۶.۱۰.۶) \end{aligned}$$

و اولین جمله سمت راست یک برآورد مجانبی برای خطای سمت چپ است.

محاسبه $y_h(x_n + 2h)$ از روی $y_n(x_n) \equiv y_n$ را یک مرحله الگوریتم در نظر می گیریم. فرض می کنیم که کاربرد یک تحمل خطای ε انتخاب کرده است و

$$\begin{aligned} \text{trunc} &\equiv u_n(x_n + 2h) - y_h(x_n + 2h) \\ &\doteq \frac{1}{2^m - 1} [y_h(x_n + 2h) - y_{2h}(x_n + 2h)] \quad (۳۷.۱۰.۶) \end{aligned}$$

باید در رابطه زیر صدق کند

$$|\text{trunc}| \leq 2\varepsilon h \quad (۳۸.۱۰.۶)$$

این فرمول، خطا در گام واحد را کنترل می کند و برطبق این رابطه خطا باید نه خیلی بزرگ باشد نه خیلی کوچک. مفهوم خطا در طول گام واحد را در (۱.۶.۶) بخش ۶.۶ به یاد آورید.

اگر (۳۸.۱۰.۶) برقرار باشد، محاسبات با همان h ادامه می‌یابد. ولی اگر برقرار نباشد، یک مقدار جدید \hat{h} باید انتخاب شود. گیریم \hat{h} چنین انتخاب شود

$$2 | \varphi_n(x_n) | \hat{h}^{m+1} \leq \varepsilon h$$

که در آن $\varphi_n(x_n)$ از رابطه زیر تعیین می‌شود

$$\varphi_n(x_n) \doteq \frac{\text{trunc}}{2h^{m+1}}$$

با این مقدار جدید \hat{h} ، خطای برشی جدید باید نزدیک میانگاهی (۳۸.۱۰.۶) باشد. الگوریتم به این صورت با چندین روش به‌کار گرفته شده است. برای مثال گی‌یر (۱۹۷۱، صص ۸۳-۸۴) را برای الگوریتم مشابهی برای یک روش مرتبه چهار RK ببینید.

در بسیاری از برنامه‌ها، برآورد خطای (۳۷.۱۰.۶) به مقادیر جاری محاسبات $y_h(x_n + 2h)$ اضافه می‌شود و یک نتیجه دقیقتری به‌دست می‌دهد. این روش را برونیابی موضعی نامند. در چنین صورتی، محک خطا در گام واحد جایگزین محک (۳۸.۱۰.۶) خطا در گام می‌شود:

$$0.5\varepsilon \leq |\text{trunc}| \leq 2\varepsilon \quad (39.10.6)$$

در چنین مواردی، می‌توان نشان داد که خطای موضعی مقادیر برونیابی شده $y_h(x_n + 2h)$ در یک محک اصلاح شده خطا در گام واحد صدق می‌نماید. [شمپاین و گوردون (۱۹۷۵، ص ۱۰۰) را ببینید]. در اجرای این روش، به‌نظر می‌رسد برنامه‌هایی که از برونیابی موضعی و محک خطا در گام استفاده می‌کنند کاراتر از برنامه‌هایی هستند که از (۳۸.۱۰.۶) استفاده می‌کنند و از برونیابی موضعی استفاده نمی‌کنند.

برای بهتری‌بودن به هزینه روشهای RK با برآورد خطایی که قبلاً داده شده‌است، فقط روشهای مرتبه چهار RK با چهار محاسبه $f(x, y)$ در هر گام RK را در نظر می‌گیریم. در رفتن از x_n به $x_n + 2h$ برای به‌دست آوردن $y_h(x_n + 2h)$ به هشت محاسبه نیاز است و سه محاسبه اضافی برای به‌دست آوردن $y_{2h}(x_n + 2h)$ لازم است. بنابراین یک گام الگوریتمهای متغیر گام به یازده محاسبه f نیاز دارد. اگرچه در مقایسه با روش چندگامی، روش متغیر گام RK نسبتاً گران است، ولی بسیار پایدار و قابل اعتماد است، و برنامه‌نویسی رایانه‌یی آن تا حدی آسان است.

روشهای رونگه- کوتا- فلبرگ^۱ روشهای رونگه- کوتا- فلبرگ روشهای RK هستند که در آنها خطای برشی از مقایسه جواب محاسبه شده y_{n+1} با نتیجه متناظری که از یک فرمول مرتبه بالاتر

RK به دست آمده، محاسبه می‌شود. رایجترین این روشها از فلبرک است [مثلاً فلبرک (۱۹۷۰) را ببینید]؛ در حال حاضر این روشها رایجترین روشهای RK هستند. برای روشن شدن مطلب، ما فقط رایجترین زوج فرمولهای رونگه - کوتا - فلبرک (RKF) از مرتبه ۴ و ۵ را در نظر می‌گیریم. این فرمولها همزمان محاسبه می‌شوند و اختلاف آنها را برآورد خطای برشی روش مرتبه چهار می‌گیرند. از جدول ۲۳.۶ دیده می‌شود که روش مرتبه پنج RK به شش محاسبه f در هر گام نیاز دارد. از این رو فلبرک در فرمول مرتبه چهار، به جای چهار محاسبه معمول پنج محاسبه f را انتخاب کرد. این آزادی انتخاب ۴ یا ۵ محاسبه f در فرمول مرتبه چهار امکان خطای برشی کوچکتر را به وجود آورده، و این موضوع بعداً نشان داده شده است.

مانند قبل، تعریف می‌کنیم

$$V_1 = f(x_n, y_n),$$

$$V_j = f\left(x_n + a_j h, y_n + h \sum_{i=1}^{j-1} \beta_{ji} V_i\right) \quad j = 2, \dots, 6 \quad (40.10.6)$$

فرمولهای مراتب چهار و پنج به ترتیب چنین‌اند

$$y_{n+1} = y_n + h \sum_{i=1}^5 \gamma_i V_i \quad (41.10.6)$$

$$\hat{y}_{n+1} = y_n + h \sum_{i=1}^6 \hat{\gamma}_i V_i \quad (42.10.6)$$

خطای برشی در y_{n+1} تقریباً برابر است با

$$\text{trunc} = \hat{y}_{n+1} - y_{n+1} = h \sum_{j=1}^6 c_j V_j \quad (43.10.6)$$

ضرایب در جدولهای ۲۵.۶ و ۲۶.۶ داده شده‌اند.

برای مقایسه روش سنتی (۲۱.۱۰.۶) RK و روش قبلی RKF، خطاهای برشی در حل

مسئله ساده

$$y' = x^r \quad y(0) = 0$$

را در نظر می‌گیریم. خطاهای برشی برای (۴۱.۱۰.۶) و (۲۱.۱۰.۶) به ترتیب چنین‌اند

$$\text{RKF}(41.10.6) : T_{n+1}(Y) \doteq \circ \cdot \circ \circ \circ 48 h^5$$

$$\text{RK}(21.10.6) : T_{n+1}(Y) \doteq - \circ \cdot \circ \circ 83 h^5 \quad n \geq 0$$

جدول ۲۵.۶ ضرایب β_{ji} و α_j برای روش RKF

j	α_j	β_{ji}				
		i = 1	۲	۳	۴	۵
۲	$\frac{1}{4}$	$\frac{1}{4}$				
۳	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{9}{8}$			
۴	$\frac{12}{13}$	$\frac{1932}{2197}$	$\frac{7200}{2197}$	$\frac{7296}{2197}$		
۵	۱	$\frac{439}{216}$	-۸	$\frac{3610}{513}$	$-\frac{145}{4104}$	
۶	$\frac{1}{2}$	$-\frac{8}{27}$	۲	$-\frac{3544}{2565}$	$\frac{1859}{4104}$	$-\frac{11}{40}$

جدول ۲۶.۶ ضرایب γ_j , $\hat{\gamma}_j$, c_j برای روش RKF

j	۱	۲	۳	۴	۵	۶
γ_j	$\frac{25}{216}$	۰	$\frac{1408}{2565}$	$\frac{2197}{4104}$	$-\frac{1}{5}$	
$\hat{\gamma}_j$	$\frac{16}{135}$	۰	$\frac{6656}{12825}$	$\frac{28561}{56430}$	$-\frac{9}{50}$	$\frac{2}{55}$
c_j	$\frac{1}{360}$	۰	$-\frac{128}{4275}$	$-\frac{2197}{75240}$	$\frac{1}{50}$	$\frac{2}{55}$

این نتایج نشان می‌دهند که روش RKF معمولاً دارای خطای برشی کوچکتری است، گرچه در عمل، تفاوت به این بزرگی نخواهد بود. باید توجه کرد که در روش سنتی (۲۱.۱۰.۶) با طول گام h و برآورد خطای (۳۷.۱۰.۶) به یازده محاسبه f برای رفتن از $y_h(x_n)$ به $y_h(x_n + 2h)$ نیاز است. و روش RKF (۴۱.۱۰.۶) برای رفتن از $y_h(x_n)$ به $y_h(x_n + 2h)$ دوازده محاسبه f را لازم دارد. در نتیجه تلاش محاسباتی در رفتن از x_n به $x_n + 2h$ قابل مقایسه‌اند، و درست است که خطای آنها را برای یک مقدار h ، مقایسه نماییم.

مثال روش (۴۱.۱۰.۶) RKF را برای حل مسئله (۲۲.۱۰.۶) به‌کار می‌بریم. این مسئله قبلاً یک مثال برای روش کلاسیک (۲۱.۱۰.۶) RK بود. مانند قبل، $h = 0.25$ ؛ و نتایج در جدول ۲۷.۶ داده شده‌اند. مقدار نظری نسبت بازهم ۱۶ است و واضح است که هنوز به آن مقدار نرسیده است. وقتی h کاهش یابد، این نسبت بیشتر به ۱۶ نزدیک می‌شود. استفاده

جدول ۲۷.۶ مثالی از روش RKF

x	$y_h(x)$	$Y(x) - y_h(x)$	$Y(x) - y_{rh}(x)$	نسبت	$\frac{1}{18}[y_h(x) - y_{rh}(x)]$
۲.۰	۰.۴۰۰۰۰۰۸۸۱	-۸.۸E-۶	-۵.۰E-۴	۵۷	-۳.۳E-۵
۴.۰	۰.۲۳۵۲۹۴۶۹	-۵.۸E-۷	-۴.۰E-۵	۶۹	-۲.۶E-۶
۶.۰	۰.۱۶۲۱۶۲۲۶	-۹.۵E-۸	-۷.۹E-۶	۸۳	-۵.۲E-۷
۸.۰	۰.۱۲۳۰۷۶۹۵	-۲.۶E-۸	-۲.۵E-۶	۹۵	-۱.۶E-۷
۱۰.۰	۰.۰۹۹۰۰۹۹۱	-۹.۴E-۹	-۱.۰E-۶	۱۰۶	-۶.۶E-۸

از فرمول درونیابی (۲۴.۱۰.۶) ریچاردسن در آخرین ستون داده شده است، و به روشنی، خطا بیش از اندازه واقعی برآورد شده است. مع هذا، باز هم یک برآورد خطای مفید است، به این دلیل که تصویری از اندازه خطای کلی به دست می دهد.

روش (۴۱.۱۰.۶) و (۴۲.۱۰.۶) معمولاً با برونیابی موضعی به کار برده می شود، همان گونه که بعداً نشان داده شده است. این روش خیلی مورد مطالعه قرار گرفته است تا معلوم شود آیا بهبودهایی امکان پذیر هست یا نه. اخیراً، شمپاین (۱۹۸۶) تحلیلی ارائه داده است که بعضی فرمولهای بهبودیافته RKF را که بر پایه چندین محک برای مقایسه فرمولهای رونگه-کوتا استوار است، پیشنهاد می نماید. تاکنون، این فرمولها در برنامه های رایانه ای با کیفیت بالا به کار برده نشده اند، گرچه انتظار می رود که مورد استفاده قرار گیرند.

برنامه های خودکار رونگه-کوتا-فلیبرگ با استفاده از (۴۳.۱۰.۶) می توان یک برنامه RKF با طول گام متغیر نوشت که خطای برشی در فرمول مرتبه چهار (۴۱.۱۰.۶) را برآورد و کنترل کند. یک چنین روشی را ل. شمپاین و ه. واتس نوشته اند و در شمپاین و واتس (۱۹۷۶a) توضیح داده شده است. ویژگیهای کلی آن به صورت زیر است. برنامه RKF45 نامگذاری شده و کاربر برنامه باید دو پارامتر خطای ABSERR و RELERR را مشخص کند. خطای برشی y_{n+1} در (۴۳.۱۰.۶) باید در محک خطا در گام

$$|\text{trunc}_j| \leq \text{ABSERR} + \text{RELERR} * |y_{n,j}| \quad (44.10.6)$$

برای هر مؤلفه y_n ، مانند $y_{n,j}$ صدق کند، که y_n جواب محاسبه شده دستگاه معادلات دیفرانسیلی است که باید حل شود. ولی بعد، نتیجه نهایی محاسبه، برای استفاده در محاسبات بعدی، به جای فرمول مرتبه چهارم y_{n+1} فرمول مرتبه پنجم \hat{y}_{n+1} از (۴۲.۱۰.۶) به کار می رود. بنابراین به جای جملات y_n در طرف راست (۴۱.۱۰.۶) و (۴۲.۱۰.۶) باید \hat{y}_n گذاشته شود. بنابراین RKF45

در واقع یک روش مرتبه پنج است، که طول گام h را، با کنترل کردن خطای برشی در فرمول مرتبه چهار (۴۱.۱۰.۶)، انتخاب می‌نماید. مانند قبل، این عمل برونیابی موضعی نامیده شده، و می‌توان نشان داد که کران (۴۴.۱۰.۶) برای y_{n+1} ، ایجاب می‌کند که در یک کران اصلاح شده خطا در طول گام واحد در خطای برشی صدق نماید. استدلال شبیه استدلالی است که در شمپاین و گوردن (۱۹۷۵، ص ۱۰۰) برای روش آدامز متغیر-مرتبه داده شده است.

آزمونهای شمپاین و همکاران (۱۹۷۹) نشان می‌دهند که RKF45 یک برنامه برتر RK است، برنامه‌ای که برای گنجاندن در کتابخانه برنامه‌های حل معادلات دیفرانسیل معمولی انتخابی بسیار عالی است. این برنامه خیلی متداول شده و در چندین کتاب درسی [مثلاً فورسایت و همکاران (۱۹۷۷)] آمده است.

به‌طور کلی، مقایسه‌هایی که در انزایت^۱ و هال^۲ (۱۹۷۶) داده شده، نشان می‌دهند که روشهای RKF برتر از روشهای دیگر RK هستند. مقایسه با روشهای چندگامی مشکلتر است. روشهای چندگامی، کمتر به محاسبه مشتق $f(x, y)$ نیاز دارند، ولی هزینه‌های بالاسری در گام برای روشهای چندگامی خیلی بیشتر از روشهای RK است. یک راه دآوری برای آنکه ببینیم کدام روش بهتر است به‌کار برده شود، این است که هزینه محاسبه $f(x, y)$ را با هزینه بالاسری در روشهای چندگامی، مقایسه کنیم. ملاحظات دیگری هم وجود دارند، برای مثال، اندازه دستگاه معادلات دیفرانسیلی است که باید حل شود. یک بحث کلی از این عوامل و اثرات آنها در انزایت و هال (۱۹۷۶) و شمپاین و همکاران (۱۹۷۶) داده شده است.

مثال مسأله

$$y' = \frac{y}{4} \left(1 - \frac{y}{20} \right) \quad y(0) = 1$$

را، که جواب آن

$$y(x) = \frac{20}{(1 + 19e^{-x/4})}$$

است، حل کنید.

مسأله با RKF45 حل شده و جوابها در نقاط $x = 2, 4, 6, \dots, 20$ داده شده‌اند. سه مقدار ABSERR به‌کار برده شده و در تمام حالات 10^{-12} RELERR = . خطاهای کلی نتیجه در جدول ۲۸.۶ داده شده است. ستونی که با NFE مشخص شده تعداد محاسبه‌های $f(x, y)$ را برای به‌دست آوردن جواب داده شده مشخص می‌کند. نتیجه را با آنچه که در جدول (۱۵.۶) که

جدول ۲۸.۶ مثالی از برنامه رونگه-کوتا-فلبرک RKF45

x	ABSERR = 10^{-3}		ABSERR = 10^{-6}		ABSERR = 10^{-9}	
	خطا	NFE	خطا	NFE	خطا	NFE
4_r°	$-1.4E-5$	۱۹	$-2.2E-7$	۴۳	$-5.5E-10$	۱۲۱
8_r°	$-3.4E-5$	۳۱	$-5.6E-7$	۷۹	$-1.2E-9$	۲۲۹
12_r°	$-3.7E-5$	۴۳	$-5.5E-7$	۱۰۳	$-8.6E-10$	۳۱۲
16_r°	$-3.3E-5$	۵۵	$-5.4E-7$	۱۲۷	$-1.2E-9$	۳۹۵
20_r°	$-1.8E-6$	۶۷	$-1.6E-7$	۱۶۳	$-4.3E-10$	۵۰۳

از برنامه چندگامی متغیر-مرتبه DDEABM به دست آمده بود مقایسه نمایید.

صورتی از RKF45 در دست است که خطای کلی در جواب محاسبه شده را، با استفاده از برآورد خطای ریچاردسن مانند (۲۴.۱۰.۶)، برآورد می نماید. برای بحث در مورد این برنامه که GERK خوانده می شود، شمپاین و واتس (۱۹۷۶b) را ببینید، که شامل یک بحث در مورد مسأله عمومی برآورد خطای کلی نیز هست.

روشهای ضمنی رونگه-کوتا ما فقط روشهای صریح RK را بررسی کرده ایم، زیرا اینها، روشهای اساسی در کاربردهای کنونی هستند. ولی فرمولهای ضمنی RK هم وجود دارند. با تعمیم فرمولهای صریح (۱۹.۱۰.۶) و (۲۰.۱۰.۶)، فرمولهای زیر را در نظر می گیریم

$$y_{n+1} = y_n + hF(x_n, y_n, h; f) \quad (45.10.6)$$

$$F(x, y, h; f) = \sum_{j=1}^p \gamma_j V_j$$

$$V_j = f \left(x + \alpha_j h, y + h \sum_{i=1}^p \beta_{ji} V_i \right) \quad j = 1, \dots, p \quad (46.10.6)$$

ضرایب $\{\alpha_j, \beta_{ji}, \gamma_j\}$ روش مورد نظر را مشخص می نمایند. در یک روش صریح برای $j \geq i$ داریم $\beta_{ji} = 0$.

در سالهای اخیر، روشهای ضمنی به طور دامنهداری مورد مطالعه قرار گرفته اند، زیرا بعضی از آنها دارای ویژگیهای پایداری مساعدی برای حل معادلات دیفرانسیل سرسخت هستند. برای هر مرتبه ای از همگرایی، روشهای A-پایداری از همان مرتبه وجود دارند. از این لحاظ، روشهای ضمنی RK بر روشهای چندگامی برتری دارند، زیرا در روشهای چندگامی هیچ روش A-پایدار بیش از مرتبه دو وجود ندارد. گسترده ترین برنامه های حل معادلات دیفرانسیل سرسخت، در حال

حاضر، بر پایه فرمولهای مشتقگیری پسرو قرار دارند که در بخش ۹.۶ توضیح داده شده است. ولی کار بسیار زیادی برای ایجاد برنامه‌های مشابه بر پایه روشهای ضمنی RK وجود دارد. برای یک مطالعه مقدماتی در این زمینه، ایکن (۱۹۸۵، بخش ۱.۳) را ببینید.

۱۱.۶ مسائل مقدار مرزی

تا اینجا، ما فقط روشهای عددی برای حل مسائل مقدار اولیه در معادلات دیفرانسیل را مطالعه کردیم. برای چنین مسائلی، شرایط جواب معادله دیفرانسیل در یک نقطه معین شده است، که نقطه اولیه خوانده می‌شود. اکنون مسائلی را بررسی می‌کنیم که در آنها، شرایط در بیش از یک نقطه بر جواب مجهول تحمیل شده است. چنین مسائلی در معادلات دیفرانسیل مسائل مقدار مرزی یا BVP خوانده می‌شوند. یک مسأله نمونه‌یی که مطالعه می‌کنیم BVPی دونقطه‌یی است:

$$y'' = f(x, y, y') \quad a < x < b$$

$$A \begin{bmatrix} u(a) \\ u'(a) \end{bmatrix} + B \begin{bmatrix} u(b) \\ u'(b) \end{bmatrix} = \begin{bmatrix} \gamma_1 \\ \gamma_2 \end{bmatrix} \quad (1.11.6)$$

جملات A و B معرّف ماتریسهای 2×2 هستند و γ_1 و γ_2 ثابتهای داده شده‌اند. نظریه BVP از این نوع بسیار پیچیده‌تر از نظریه مسأله مقدار اولیه است (قضیه‌های ۱.۶ و ۲.۶ در بخش ۱.۶ را ببینید). برای نمایاندن مشکلات ممکن، مثالهای زیر را می‌آوریم.

مثال ۱. BVPی خطی دونقطه‌یی زیر را در نظر می‌گیریم

$$y'' = -\lambda y \quad 0 < x < 1$$

$$y(0) = y(1) = 0 \quad (2.11.6)$$

اگر λ یکی از اعداد

$$\pi^2, 4\pi^2, 9\pi^2, \dots, n^2\pi^2, \dots \quad (3.11.6)$$

نباشد، این BVP دارای جواب یکتای $Y(x) \equiv 0$ است. در غیر این صورت، تعداد نامتناهی جواب

$$Y(x) = C \sin(\sqrt{\lambda}x) \quad (4.11.6)$$

وجود دارد، که λ از $(6-11-3)$ انتخاب می‌شود و C ثابت اختیاری است.

۲. همچنین مسألهٔ دونقطه‌یی زیر را در نظر می‌گیریم

$$y'' = -\lambda y + g(x) \quad 0 < x < 1$$

$$y(0) = y(1) = 0 \quad (5.11.6)$$

اگر λ از (۳.۱۱.۶) انتخاب نشده باشد و اگر $g(x)$ در $a \leq x \leq b$ پیوسته باشد، آنگاه این مسألهٔ یک جواب یکتای $Y(x)$ دارد که در $[0, 1]$ دو بار پیوسته مشتقپذیر است. به عکس اگر $\lambda = \pi^2$ ، آنگاه مسألهٔ (۵.۱۱.۶) دارای جواب است اگر و تنها اگر $g(x)$ در رابطهٔ زیر صدق کند

$$\int_0^1 g(x) \sin(\pi x) dx = 0$$

در حالتی که این رابطه برقرار باشد، جواب با رابطهٔ

$$Y(x) = C \sin(\pi x) + \frac{1}{\pi} \int_0^x g(t) \sin(\pi(x-t)) dt \quad (6.11.6)$$

داده می‌شود که C یک ثابت اختیاری است. نتیجهٔ مشابهی برای λ دیگری که از (۳.۱۱.۶) انتخاب شده باشد برقرار است.

چند قضیهٔ از نظریهٔ BVPی دو نقطه‌یی اکنون داده می‌شود، تا به فهم بیشتر مطالب و عرضهٔ روشهای عددی برای حل آنها، ما را یاری نمایند. ما کار را با مسألهٔ دونقطه‌یی برای معادلهٔ دیفرانسیل خطی مرتبهٔ دوم آغاز می‌کنیم:

$$y'' = p(x)y' + q(x)y + g(x) \quad a < x < b$$

$$A \begin{bmatrix} u(a) \\ u'(a) \end{bmatrix} + B \begin{bmatrix} u(b) \\ u'(b) \end{bmatrix} = \begin{bmatrix} \gamma_1 \\ \gamma_2 \end{bmatrix} \quad (7.11.6)$$

مسألهٔ همگن حالتی است که در آن $g(x) \equiv 0$ و $\gamma_1, \gamma_2 = 0$.

قضیهٔ ۱۰.۶ مسألهٔ ناهمگن (۷.۱۱.۶) دارای جواب یکتای $Y(x)$ بر $[a, b]$ ، برای هر مجموعه از داده‌های $\{g(x), \gamma_1, \gamma_2\}$ است اگر و تنها اگر مسألهٔ همگن مربوطهٔ آن فقط جواب بدیهی $Y(x) \equiv 0$ داشته باشد.

برهان استک‌گولد^۱ (۱۹۷۹، ص ۱۹۷) را ببینید. مثالهای قبل مثالهایی از این قضیه‌اند. ■

برای شرایطی که در آنها، مسأله همگن (۷.۱۱.۶)، فقط جواب صفر داشته باشد، مسأله خطی خاصتر زیر را در نظر می‌گیریم:

$$\begin{aligned} y'' &= p(x)y' + q(x)y + g(x) & a < x < b \\ a_0 y(a) - a_1 y'(a) &= \gamma_1 & b_0 y(b) + b_1 y'(b) &= \gamma_2 \end{aligned} \quad (۸.۱۱.۶)$$

در این مسأله می‌گوییم که شرایط مرزی جدا شده‌اند. فرض می‌کنیم که شرایط زیر برقرار باشند

$$\begin{aligned} q(x) &> 0 & a \leq x \leq b \\ a_0 a_1 &\geq 0 & b_0 b_1 &\geq 0 \end{aligned} \quad (۹.۱۱.۶)$$

$$|a_0| + |a_1| \neq 0 \quad |b_0| + |b_1| \neq 0 \quad |a_0| + |b_0| \neq 0$$

در این صورت مسأله همگن (۸.۱۱.۶) فقط دارای جواب صفر است، قضیه ۱۰.۶ را می‌توان به‌کار برد، و مسأله ناهمگن برای هر مجموعه از داده‌های $\{g(x), \gamma_1, \gamma_2\}$ دارای جواب یکتاست. برای اثبات این قضیه کِلر (۱۹۶۸، ص ۱۱) را ببینید. مثال (۵.۱۱.۶) توضیحی برای این قضیه است. این مثال همچنین نشان می‌دهد که شرایط (۹.۱۱.۶) لازم نیستند. مسأله (۵.۱۱.۶) برای بیشتر انتخابهای منفی $q(x) = \lambda$ یکتا حلپذیر است.

نظریه مسأله غیرخطی (۱۰.۱۱.۶) به مراتب پیچیده‌تر از نظریه مسأله خطی (۷.۱۱.۶) است. ما یک مقدمه مختصری از این نظریه برای مسأله محدودتر زیر می‌آوریم.

$$\begin{aligned} y'' &= f(x, y, y') & a < x < b \\ a_0 y(a) - a_1 y'(a) &= \gamma_1 & b_0 y(b) + b_1 y'(b) &= \gamma_2 \end{aligned} \quad (۱۰.۱۱.۶)$$

فرض بر این است که تابع f در همه نقاط (x, u, v) ، (x, u_i, v) ، در ناحیه

$$R = \{(x, u, v) \mid a \leq x \leq b, -\infty < u, v < \infty\}$$

در شرط لیبیشیس زیر صدق می‌کند

$$\begin{aligned} |f(x, u_1, v) - f(x, u_2, v)| &\leq K |u_1 - u_2| \\ |f(x, u, v_1) - f(x, u, v_2)| &\leq K |v_1 - v_2| \end{aligned} \quad (۱۱.۱۱.۶)$$

این شرط بسیار قویتر از شرط مورد نیاز است، ولی بیان قضیه زیر و تحلیل روشهای عددی را که بعداً داده خواهد شد آسان می‌کند.

قضیه ۱۱.۶ برای مسأله (۱۰.۱۱.۶)، فرض می‌کنیم $f(x, u, v)$ در ناحیه R پیوسته باشد و در شرط (۱۱.۱۱.۶) لپشیتس صدق کند. به علاوه فرض می‌کنیم که در R ، f در روابط زیر صدق می‌کند:

$$\frac{\partial f(x, u, v)}{\partial u} > 0 \quad \left| \frac{\partial f(x, u, v)}{\partial v} \right| \leq M \quad (12.11.6)$$

برای مقداری از ثابت $M > 0$. برای شرایط مرزی (۱۰.۱۱.۶)، فرض می‌کنیم

$$a_0 a_1 \geq 0 \quad b_0 b_1 \geq 0$$

$$|a_0| + |a_1| \neq 0 \quad |b_0| + |b_1| \neq 0 \quad |a_0| + |b_0| \neq 0 \quad (13.11.6)$$

در این صورت BVP (۱۰.۱۱.۶) دارای جواب یکتاست.

برهان کلر (۱۹۶۸، ص ۹) را ببینید. قضیه یکتایی قبلی برای مسأله خطی (۸.۱۱.۶) یک حالت خاص این قضیه است.

BVPهای غیرخطی ممکن است نایکتا حلپذیر باشند و فقط تعدادی متناهی جواب، داشته باشند. این وضعیت در مسائل خطی پیش نمی‌آید، همان‌طور که در (۲.۱۱.۶) و (۶.۱۱.۶) نشان داده شد نایکتایی زیادی است و همیشه به معنای بی‌نهایت جواب است. یک مثال نایکتا حلپذیری، مسأله مرتبه دوم زیر است:

$$\frac{d}{dx} \left[I(x) \frac{dy}{dx} \right] + \lambda \sin(y) = 0 \quad 0 < x < 1$$

$$y'(0) = y'(1) = 0 \quad (14.11.6)$$

که در مطالعه کمانش ستون پیش می‌آید. عامل λ متناسب با باری است که بر ستون گذاشته شده است؛ وقتی λ از حد معینی تجاوز نماید، مسأله (۱۴.۱۱.۶) یک جواب غیر از جواب صفر دارد. برای تفصیل بیشتر این مسأله، کلر و آنتمان (۱۹۶۹، ص ۴۳) را ببینید.

همانند مطالب قبلی درباره مسائل مقدار اولیه [(۱۱.۱.۶) تا (۱۵.۱.۶) بخش ۱.۶ را ببینید]، تمام مسائل مقدار مرزی از مراتب بالا را می‌توان به صورت مسائلی برای یک دستگاه معادلات مرتبه اول نوشت. این شکل کلی برای یک BVP دونقطه‌یی در یک دستگاه معادلات مرتبه اول چنین است:

$$y' = f(x, y) \quad a < x < b$$

$$Ay(a) + By(b) = \gamma \quad (۱۵.۱۱.۶)$$

که یک دستگاه n معادله مرتبه اول را نشان می‌دهد. کمیت‌های $y(x)$ ، $f(x, y)$ و γ بردارهایی با n مؤلفه و A و B ماتریسهایی از مرتبه $n \times n$ هستند. یک نظریه برای این‌گونه BVPها وجود دارد، مشابه با نظریه‌ای که برای مسأله دونقطه‌یی (۱۰.۱۱.۶) ذکر کردیم، ولی به دلایل کمبود جا، از ذکر آن صرف‌نظر می‌کنیم.

در بقیه این بخش، روشهای عددی اصلی برای حل BVPی دونقطه‌یی (۱۰.۱۱.۶) را به اختصار توضیح می‌دهیم. این روشها، برای دستگاه معادلات مرتبه اول از قبیل (۱۵.۱۱.۶) تعمیم می‌یابند، ولی به لحاظ کمبود جا، آن قضایا را ذکر نمی‌کنیم. در بیشتر مطالب ارائه شده از کلر (۱۹۶۸) پیروی کرده‌ایم، و نظریه دستگاههای مرتبه اول را در آنجا داده‌ایم. برخلاف وضعیت مسائل مقدار اولیه، اغلب بهتر است که با BVPی مرتبه بالاتر مستقیماً برخورد نماییم تا اینکه مسأله را به شکل یک دستگاه مرتبه اول به‌طور عددی حل کنیم. ارائه روشهای عددی مسأله دونقطه‌ای، کمتر پیچیده است، و بنابراین ما بحث درباره مسأله دو نقطه را به جای دستگاه (۱۵.۱۱.۶)، انتخاب نموده‌ایم.

روشهای پرتابی یکی از راههای متداول حل BVPی دونقطه‌یی، تبدیل آن است به مسأله‌ای که در آن از برنامه‌ای برای حل مسائل مقدار اولیه بتوان استفاده کرد. اکنون یک چنین روشی را برای BVPی (۱۰.۱۱.۶) به‌دست می‌دهیم.

مسأله مقدار اولیه زیر را در نظر می‌گیریم

$$y'' = f(x, y, y') \quad a < x < b$$

$$y(a) = a_1 s - c_1 \gamma_1 \quad y'(a) = a_0 s - c_0 \gamma_1 \quad (۱۶.۱۱.۶)$$

که به پارامتر s بستگی دارد، و c_0 و c_1 ثابتهای اختیاری هستند که در شرط زیر صدق می‌کنند

$$a_1 c_0 - a_0 c_1 = 1$$

جواب (۱۶.۱۱.۶) را با $Y(x; s)$ نشان می‌دهیم. در این صورت به‌سادگی دیده می‌شود که برای همه مقادیر s که Y وجود داشته باشد

$$a_0 Y(a; s) - a_1 Y'(a; s) = \gamma_1$$

چون Y یک جواب (۱۶.۱۱.۶) است، تنها چیزی که لازم است تا این Y جواب (۱۰.۱۱.۶) باشد این است که در شرط مرزی در b که باقیمانده صدق کند. این بدان معناست که $Y(x; s)$ باید در شرط زیر صدق کند

$$\varphi(s) \equiv b_0 Y(b; s) + b_1 Y'(b; s) - \gamma_2 = 0 \quad (17.11.6)$$

این معادله یک معادله غیرخطی برحسب s است. اگر s^* یک ریشه $\varphi(s)$ باشد، آنگاه $Y(x; s)$ در BVP (۱۰.۱۱.۶) صدق می‌کند. می‌توان نشان داد که با مفروضات مناسبی روی f و شرایط مرزی آن، معادله (۱۷.۱۱.۶) دارای یک جواب یکتای s^* خواهد بود [کلر (۱۹۶۸)، ص ۹] را ببینید. می‌توانیم یک روش ریشه‌یابی برای معادله غیرخطی به‌کار ببریم و s^* را پیدا کنیم. این طریق پیدا کردن جواب برای BVP روش پرتابی خوانده می‌شود. این نام از پرتاب‌شناسی گرفته شده که در آنجا سعی می‌کنند شرایط اولیه لازم در $x = a$ را معلوم نمایند تا مقدار معینی در $x = b$ به دست آید. هر یک از روشهای ریشه‌یابی فصل ۲ را می‌توان برای حل $\varphi(s) = 0$ به‌کار برد. هر محاسبه $\varphi(s)$ شامل حل (۱۶.۱۱.۶) بر $[a, b]$ است، و در نتیجه می‌خواهیم تعداد چنین محاسباتی را مینیمم سازیم. به‌عنوان یک مثال خاص از یک روش همگرایی سریع و مهم، به روش نیوتن می‌نگریم:

$$s_{m+1} = s_m - \frac{\varphi(s_m)}{\varphi'(s_m)} \quad m = 0, 1, \dots \quad (18.11.6)$$

برای محاسبه $\varphi'(s)$ ، از شرط (۱۷.۱۱.۶) مشتق می‌گیریم تا به دست آوریم

$$\varphi'(s) = b_0 \xi_s(b) + b_1 \xi'_s(b) \quad (19.11.6)$$

که

$$\xi_s(x) = \frac{\partial Y(x; s)}{\partial s} \quad (20.11.6)$$

برای پیدا کردن $\xi_s(x)$ از معادله زیر نسبت به s مشتق می‌گیریم

$$Y''(x; s) = f(x, Y(x; s), Y'(x; s))$$

در این صورت ξ_s در مسأله مقدار اولیه زیر صدق می‌کند

$$\begin{aligned} \xi_s''(x) &= f_2(x, Y(x; s), Y'(x; s)) \xi_s(x) \\ &+ f_3(x, Y(x; s), Y'(x; s)) \xi'_s(x) \end{aligned} \quad (21.11.6)$$

$$\xi_s(a) = a_1 \quad \xi'_s(a) = a_0$$

توابع f_1 و f_2 معرّف مشتقات جزئی $f(x, u, v)$ به ترتیب نسبت به u و v هستند. مقادیر اولیه همان مقادیری هستند که در (۱۶.۱۱.۶) و از تعریف ξ_s به دست آمده‌اند.

در عمل مسائل (۱۶.۱۱.۶) و (۲۱.۱۱.۶) را به یک دستگاه چهار معادله مرتبه اول با مجهولات Y, Y', ξ_s و ξ'_s برمی‌گردانیم. این دستگاه به‌طور عددی مثلاً با یک روش مرتبه p طول گام h حل می‌شود. با نمادگذاری مشابهی برای سایر مجهولات، فرض می‌کنیم $y_h(x; s)$ تقریب $Y(x; s)$ را نشان دهد. از قضایای پیشین برای حل مسائل مقدار اولیه، می‌توان نشان داد که این جوابهای تقریبی از لحاظ خطا از مرتبه $O(h^p)$ خواهند بود. با فرضهای مناسب در مسأله اصلی (۱۰.۱۱.۶)، می‌توان نشان داد که ریشه به دست آمده s_h^* نیز دارای خطای $O(h^p)$ خواهد بود، و همچنین است برای جواب تقریبی $y_h(x; s_h^*)$ وقتی با جواب $Y(x; s^*)$ مسأله مقدار مرزی مقایسه شود. برای جزئیات این تحلیل کُر (۱۹۶۸، صص ۴۷-۵۴) را ببینید.

مثال BVPی دو نقطه‌یی زیر را در نظر می‌گیریم

$$y'' = -y + \frac{2(y')^2}{y} \quad -1 < x < 1 \quad (22.11.6)$$

$$y(-1) = y(1) = (e + e^{-1})^{-1} \approx 0.324027137$$

جواب درست آن $Y(x) = (e^x + e^{-x})^{-1}$ است. مسأله مقدار اولیه (۱۵.۱۱.۶) برای روش پرتابی چنین است

$$y'' = -y + \frac{2(y')^2}{y} \quad -1 < x < 1 \quad (23.11.6)$$

$$y(-1) = (e + e^{-1})^{-1} \quad y'(-1) = s$$

مسأله مربوط به (۲۱.۱۱.۶) برای $\xi_s(x)$ عبارت است از:

$$\xi_s'' = \left[-1 - 2 \left(\frac{y'}{y} \right)^2 \right] \xi_s + 4 \frac{y'}{y} \xi_s' \quad (24.11.6)$$

$$\xi_s(-1) = 0 \quad \xi_s'(-1) = 1$$

در معادله ξ_s'' از $Y(x; s)$ ، جواب (۲۳.۱۱.۶)، استفاده می‌شود. تابع $\varphi(s)$ برای محاسبه s^* با رابطه زیر داده می‌شود

$$\varphi(s) \equiv Y(1; s) - (e + e^{-1})^{-1}$$

جدول ۲۹.۶ روش پرتابی برای حل (۲۲.۱۱.۶)

$n = \frac{2}{h}$	$s^* - s_h^*$	نسبت	E_h	نسبت
۴	$۴,۰۱E - ۳$		$۲,۸۳E - ۲$	
۸	$۱,۵۲E - ۳$	۲,۶۴	$۷,۳۰E - ۳$	۳,۸۸
۱۶	$۴,۶۴E - ۴$	۳,۲۸	$۱,۸۲E - ۳$	۴,۰۱
۳۲	$۱,۲۷E - ۴$	۳,۶۴	$۴,۵۴E - ۴$	۴,۰۱
۶۴	$۳,۳۴E - ۵$	۳,۸۲	$۱,۱۴E - ۴$	۴,۰۰

برای استفاده از تعریف روش نیوتن، داریم

$$\varphi'(s) = \xi_s(1)$$

با استفاده از Y ، جواب درست (۲۲.۱۱.۶) و شرط $s = y'(-1)$ در (۲۳.۱۱.۶)، به سادگی s^* ، ریشه مطلوب $\varphi(s)$ به دست می آید

$$s^* = Y'(-1) = \frac{e - e^{-1}}{(e + e^{-1})^2} \doteq ۰,۲۴۵۷۷۷۱۷۴$$

برای حل مسأله مقدار اولیه (۲۳.۱۱.۶) - (۲۴.۱۱.۶)، از روش مرتبه دوم (۹.۱۰.۶) رونگه-کوتا با طول گام $h = 2/n$ استفاده نموده ایم. نتایج برای چندین مقدار n در جدول ۲۹.۶ داده شده است. جواب (۲۴.۱۱.۶) با $y_h(x; s)$ نشان داده شده است و ریشه حاصل

$$\varphi_h(s) \equiv y_h(1; s) - (e + e^{-1})^{-1} = 0$$

با s_h^* نشان داده شده است. برای خطا در $y_h(x; s_h^*)$ ، گیریم

$$E_h = \max_{0 \leq i \leq n} |Y(x_i) - y_h(x_i; s_h^*)|$$

که $\{x_i\}$ ها نقاط گرهی اند که از آنها در حل مسأله مقدار اولیه استفاده شده است. ستونهایی که با نسبت مشخص شده اند ضرایبی را مشخص می کنند که با آنها، وقتی n دو برابر (یا h نصف) می شود، خطا کاهش می یابد. به طور نظری این ضرایب باید به ۴ نزدیک شوند زیرا روش رونگه-کوتا دارای خطایی از مرتبه $O(h^2)$ است. و همان گونه که انتظار می رود، در تجربه این ضرایب به $۴,۰$ نزدیک می شوند. برای بارست (۱۸.۱۱.۶) نیوتن، در هر حالت $۲,۰ = s$. به کار برده شده است. بارست وقتی پایان یافته است که آزمون

$$|s_{m+1} - s_m| \leq 10^{-1}$$

برآورده شده است. با این گزینشها، روش نیوتن در هر حالت به ۶ بارست نیاز داشته است جز حالت $n = 4$ (که هفت بارست لازم بوده است). ولی اگر $\epsilon = 0$ به کار برده می‌شد، برای حالت $n = 4$ ، ۲۵ بارست لازم می‌شد، که اهمیت یک گزینش خوب حدس اولیه s را نشان می‌دهد.

مسائل چندی در روش پرتابی پیش می‌آید. اول، یک حدس کلی s برای بارست نیوتن وجود ندارد، و با یک گزینش ضعیف، این بارست ممکن است واگرا شود. بدین دلیل، روش اصلاح‌شده (۱۱.۱۱.۲) - (۱۲.۱۱.۲) نیوتن در بخش ۱۱.۲ ممکن است لازم باشد تا همگرایی به وجود آید. مسأله دوم اینکه انتخاب $y_h(x; s)$ ممکن است نسبت به h و s و سایر ویژگیهای مسأله مقدار مرزی، خیلی حساس باشد. برای مثال، اگر خطی‌سازی مسأله مقدار آغازی (۱۶.۱۱.۶) ویژه مقدارهای مثبت بزرگی داشته باشد، انتخاب $Y(x; s)$ احتمالاً نسبت به تغییرات s حساس می‌شود. برای بحث در مورد تمام این مسائل، کلر (۱۹۶۸، فصل ۲)، اشتور^۱ و بورلیرش^۲ (۱۹۸۰، بخش ۳۰۷)، و فاکس^۳ (۱۹۸۰، صص ۱۸۰-۱۸۶) را ببینید. بررسی بعضی از این مسائل BVPهای خطی آسانتر است، مانند (۸.۱۱.۶)، آن‌گونه که در کلر (۱۹۶۸، فصل ۲) صورت گرفته است.

روشهای تفاضل متناهی BVP دونقطه‌یی

$$\begin{aligned} y'' &= f(x; y, y') & a < x < b \\ y(a) &= \gamma_1 & y(b) = \gamma_2 \end{aligned} \quad (25.11.6)$$

را در نظر می‌گیریم که جواب درست آن با $Y(x)$ نشان داده شده است. گیریم $n > 1$ به‌ازای $x_j = a + jh$ ، $h = (b - a)/n$ ، $j = 0, 1, \dots, n$ در هر نقطه گرهی داخلی x_i ، $0 < i < n$ و $Y''(x_i)$ و $Y'(x_i)$ را تقریب می‌زنیم:

$$\begin{aligned} Y''(x_i) &= \frac{Y_{i+1} - 2Y_i + Y_{i-1}}{h^2} - \frac{h^2}{12} Y^{(4)}(\xi_i) \\ Y'(x_i) &= \frac{Y_{i+1} - Y_{i-1}}{2h} - \frac{h^2}{6} Y^{(3)}(\eta_i) \end{aligned} \quad (26.11.6)$$

به‌ازای یک مقدار ξ_i و $x_{i-1} \leq \xi_i \leq x_{i+1}$ ، $i = 1, 2, \dots, n-1$. با حذف آخرین جمله‌های خطا، و استفاده از این تقریبات در معادله دیفرانسیل، دستگاه غیرخطی تقریبی زیر را خواهیم داشت:

$$\frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} = f\left(x_i, y_i, \frac{y_{i+1} - y_{i-1}}{2h}\right) \quad i = 1, \dots, n-1 \quad (27.11.6)$$

این دستگاه، یک دستگاه $n-1$ معادله غیرخطی برحسب $n-1$ مجهول y_2, y_1, \dots, y_{n-1} است. مقادیر $y_0 = \gamma_1$ و $y_n = \gamma_2$ با توجه به شرایط مرزی معلوم اند.

تحلیل خطا در $\{y_i\}$ در مقایسه با $\{Y(x_i)\}$ خیلی دشوارتر از آن است که در اینجا بتوان بیان کرد، زیرا نیاز به تحلیل حلپذیری دستگاه معادلات غیرخطی دارد. اساساً، اگر $Y(x)$ چهار بار مشتقپذیر باشد، اگر مسأله (۲۵.۱۱.۶) در ناحیه‌ای پیرامون نمودار $Y(x)$ بر $[a, b]$ ، به طور یکتا حلپذیر باشد، و اگر $f(x, u, v)$ به اندازه کافی مشتقپذیر باشد، آنگاه یک جواب (۲۷.۱۱.۶) وجود دارد و در رابطه زیر صدق می‌کند

$$\max_{\substack{i \leq n \\ i \geq 1}} |Y(x_i) - y_i| = O(h^2) \quad (28.11.6)$$

برای تحلیلی از این موضوع به کالر (۱۹۷۶، بخش ۲.۳) یا کالر (۱۹۶۸، بخش ۲.۳) مراجعه نمایید. وانگهی با فرضهای بیشتر روی f و همواری Y ، می‌توان نشان داد که

$$Y(x_i) - y_i = \tau(x_i)h^2 + O(h^4) \quad (29.11.6)$$

که $\tau(x)$ مستقل از h است. این را می‌توان برای توجیه برونمایی ریچاردسن در به دست آوردن نتایجی که سریعتر همگرا باشند، به کار برد. [در سایر روشها، بهبود همگرایی بر پایه تصحیح خطا در تقریبهای تفاضلی مرکزی (۲۷.۱۱.۶) قرار دارند.]

دستگاه (۲۷.۱۱.۶) به راههای گوناگونی حل می‌شود، که بعضی از آنها اصلاح ساده روشهایی

است که در بخش ۱۰.۲ و ۱۱.۲ داده شدند. در شکل ماتریسی، داریم

$$\frac{1}{h^2} \begin{bmatrix} -2 & 1 & 0 & \dots & 0 \\ 1 & -2 & 1 & & \vdots \\ \vdots & & \ddots & & \\ & & & 1 & -2 & 1 \\ 0 & \dots & & 0 & 1 & -2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{n-1} \end{bmatrix} = \begin{bmatrix} f\left(x_1, y_1, \frac{y_2 - y_0}{2h}\right) \\ f\left(x_2, y_2, \frac{y_3 - y_1}{2h}\right) \\ \vdots \\ f\left(x_{n-1}, y_{n-1}, \frac{y_n - y_{n-2}}{2h}\right) \end{bmatrix} - \begin{bmatrix} \frac{\gamma_1}{h^2} \\ 0 \\ \vdots \\ \frac{\gamma_2}{h^2} \end{bmatrix}$$

که آن را به شکل زیر نشان می‌دهیم.

$$\frac{1}{h^2} Ay = \hat{f}(y) + g \quad (۳۰.۱۱.۶)$$

ماتریس A ناتکین است [قضیه ۲.۸ فصل ۸ را ببینید]؛ و دستگاه خطی $Az = b$ به سادگی حلپذیر است، همان طور که قبل از قضیه ۲.۸ فصل ۸ آمده است. روش نیوتن (بخش ۱۱.۲ را ببینید) برای حل (۳۰.۱۱.۶) به صورت زیر داده می‌شود

$$y^{(m+1)} = y^{(m)} - \left[\frac{1}{h^2} A - F(y^{(m)}) \right]^{-1} \left[\frac{1}{h^2} Ay^{(m)} - \hat{f}(y^{(m)}) - g \right] \quad (۳۱.۱۱.۶)$$

$$F(y) = \left[\frac{\partial \hat{f}_i}{\partial y_j} \right] \quad 1 \leq i, j \leq n-1$$

ماتریس ژاکوبی به علت شکل خاص $\hat{f}(y)$ به طور قابل ملاحظه‌ای ساده می‌شود:

$$[F(y)]_{ij} = \frac{\partial f \left(x_i, y_i, \frac{y_{i+1} - y_{i-1}}{2h} \right)}{\partial y_i}$$

این عبارت صفر است مگر برای $j = i + 1$ یا $j = i - 1$:

$$[F(y)]_{ii} = f_{\tau} \left(x_i, y_i, \frac{y_{i+1} - y_{i-1}}{2h} \right) \quad 1 \leq i \leq n-1$$

$$[F(y)]_{i,i-1} = \frac{-1}{2h} f_{\tau} \left(x_i, y_i, \frac{y_{i+1} - y_{i-1}}{2h} \right) \quad 2 \leq i \leq n-1$$

$$[F(y)]_{i,i+1} = \frac{1}{2h} f_{\tau} \left(x_i, y_i, \frac{y_{i+1} - y_{i-1}}{2h} \right) \quad 1 \leq i \leq n-2$$

که $f_{\tau}(x, u, v)$ و $f_{\tau}(x, u, v)$ به ترتیب مشتقات جزئی نسبت به u و v را نشان می‌دهند. بنابراین ماتریسی که باید در (۳۱.۱۱.۶) معکوس شود به شکل خاصی است که سه قطری خوانده می‌شود. گیریم

$$B_m = \frac{1}{h^2} A - F(y^{(m)}) \quad (۳۲.۱۱.۶)$$

می‌توانیم (۳۱.۱۱.۶) را دوباره به شکل زیر بنویسیم

$$y^{(m+1)} = y^{(m)} - \delta^{(m)} \quad (۳۳.۱۱.۶)$$

$$B_m \delta^{(m)} = \frac{1}{h^2} Ay^{(m)} - f(y^{(m)}) - g$$

جدول ۳۰.۶ روش تفاضل متناهی برای حل (۲۲.۱۱.۶)

$n = \frac{2}{h}$	E_h	نسبت
۴	$۲,۶۳E - ۲$	
۸	$۵,۸۷E - ۳$	۴٫۴۸
۱۶	$۱,۴۳E - ۳$	۴٫۱۱
۳۲	$۳,۵۵E - ۴$	۴٫۰۳
۶۴	$۸,۸۶E - ۵$	۴٫۰۱

همانگونه که در بخش ۲.۸ فصل ۸ نشان داده شده است، این دستگاه خطی به سادگی و به سرعت حلپذیر است. می توان نشان داد که تعداد ضربها و تقسیمها در حدود $5n$ است، که در حل یک دستگاه $n - 1$ معادله خطی، زیاد نیست. صرفه جویی بیشتر را می توان، با تغییر ندادن B_m یا تغییر آن فقط پس از چند بارست (۳۳.۱۱.۶)، به وجود آورد. برای یک بازنگری وسیع و بحث در حل دستگاههای غیرخطی که در ارتباط با حل BVPها پیش می آید به دوینفل هارت^۱ (۱۹۷۹) مراجعه نمایید.

مثال روش تفاضل متناهی (۲۷.۱۱.۶) را قبلاً برای حل BVPی (۲۲.۱۱.۶)، که پیش از این برای مثال روش پرتابی استفاده شده بود، به کار برده ایم. نتایج برای دوبرابر کردنهای متوالی $n = 2/h$ در جدول ۳۰.۶ داده شده است. دستگاه غیرخطی (۲۷.۱۱.۶) همانگونه که در (۳۳.۱۱.۶) توضیح داده شد، با روش نیوتن، حل شده است. حدس آغازی چنین بوده است:

$$y_h^{(0)}(x_i) = (e + e^{-1})^{-1} \quad i = 0, 1, \dots, n$$

که بر پایه وصل کردن مقادیر مرزی به یکدیگر با یک خط مستقیم، حاصل شده است. کمیت

$$d_h = \max_{0 \leq i \leq n} |y_i^{(m+1)} - y_i^{(m)}|$$

برای هر بارست محاسبه شده و هرگاه شرط زیر برقرار شده، بارست خاتمه یافته است. در همه حالات،

$$d_h \leq 10^{-10}$$

تعداد بارستهای محاسبه شده ۵ یا ۶ بوده است. برای خطا، گیریم

$$E_h = \max_{0 \leq i \leq n} |Y(x_i) - y_h(x_i)|$$

که در آن y_h جواب (۲۷.۱۱.۶) است که از روش نیوتن به دست آمده است. به موجب (۲۸.۱۱.۶) و (۲۹.۱۱.۶) باید انتظار داشته باشیم که مقادیر E_h ، وقتی h نصف می‌شود، با ضربی حدود ۴ کاهش یابد و این چیزی است که در جدول مشاهده می‌نماییم.

روشهای مراتب بالاتر را می‌توان به چند راه به دست آورد. (۱) تقریبات از مرتبه بالاتر را برای مشتقات به کار برد تا (۲۶.۱۱.۶) بهبود یابد؛ (۲) برونیابی ریچاردسن را بر پایه (۲۹.۱۱.۶) به کار برد و مانند انتگرالگیری رامبرگ، برونیابی را تکرار نمود تا روشهای از مراتب به دلخواه بالاتر به دست آیند، (۳) خطاهای برشی در (۲۶.۱۱.۶) با تفاضلات مراتب بالاتر، با استفاده از مقادیر محاسبه شده y_h ، تقریب زدن با استفاده از این مقادیر برای اصلاحات (۲۷.۱۱.۶)، و به دست آوردن یک تقریب دقیقتر برای معادله دیفرانسیل در (۲۵.۱۱.۶)، که به جواب دقیقتر می‌انجامد. تمام این فنون به کار گرفته شده‌اند، و بعضی به صورت برنامه‌های رایانه‌یی کاملاً پیچیده در آمده‌اند. برای بحث بیشتر و برای مثالهایی از برنامه‌های رایانه‌یی، فاکس (۱۹۸۰، ص ۱۹۱)، جین^۱ (۱۹۸۴، فصل ۴)، و پیریرا^۲ (۱۹۷۹) را ببینید.

روشها و مسائل دیگر چند روش دیگر در حل مسائل مقدار مرزی به کار رفته‌اند. بهترین آنها احتمالاً روش هم‌مکانی است. برای بحث درباره روشهای هم‌مکانی ردین^۳ (۱۹۷۹)، دوپل هارت (۱۹۷۹) و آشر و راسل (۱۹۸۵) را ببینید. برای یک برنامه مهم رایانه‌یی هم‌مکانی آشر و همکاران (۱۹۸۱a، ۱۹۸۱b) را ببینید.

یک روش دیگر برای حل مسأله مقدار مرزی، این است که یک معادله انتگرالی صورت هم‌ارز آن را حل کنیم. دستاوردهای بسیار کمتری برای چنین روشهای عددی وجود دارند، ولی در بعضی موارد این روشها می‌توانند خیلی موثر باشند. برای آشنایی با این روش کلر (۱۹۶۸، فصل ۴) را ببینید.

همچنین انواع دیگری از مسائل مقدار مرزی وجود دارند، که بعضی شامل نوعی رفتارهای تکین هستند، که در اینجا بحث نکرده‌ایم. برای همه این مطالب به مقالات گزارشی آشر و راسل (۱۹۸۵)، عزیز^۴ (۱۹۸۵)، چایلدز و همکاران (۱۹۷۹) و گلاودول و سه‌ریز (۱۹۸۰) را ببینید؛ همچنین (فصل ۴) کلر (۱۹۷۶) را برای مسائل تکین ملاحظه نمایید. برای بحث نرم‌افزارها، چایلدز و همکاران (۱۹۷۹)، گلاودول^۵ و سه‌ریز^۶ (۱۹۸۰) و انزیت^۷ را ببینید.

گفتگو درباره منابع

معادلات دیفرانسیل معمولی و جزئی شکل اصلی مدل ریاضی در علوم و مهندسی هستند، و در نتیجه، حل عددی معادلات دیفرانسیل یک زمینه وسیع تحقیق است. دو کتاب درسی کلاسیک که معلومات در این زمینه را قبل از کاربرد فراگیر رایانه‌های رقمی، نشان می‌دهند کولاتس^۱ (۱۹۶۶) و میلن^۲ (۱۹۵۳) هستند. بعضی از کتابهای عمومی مهم که از ۱۹۶۰ به بعد در حل عددی معادلات دیفرانسیل معمولی نوشته شده‌اند عبارت‌اند از هنریچی (۱۹۶۲)، گی‌یر (۱۹۷۱)، لاپیدوس^۳ و زاینفلت^۴ (۱۹۷۱)، لامبرت^۵ (۱۹۷۳)، اشتتر^۶ (۱۹۷۳)، هال و وات (۱۹۷۶)، شمپاین و گوردن (۱۹۷۵)، واندرهوون^۷ (۱۹۷۷)، اورتگا و پول (۱۹۸۱) و بوچر^۸ (۱۹۸۷). یک ارزیابی مفید در گوپتا و همکاران (۱۹۸۵) داده شده است.

نظریه جدید همگرایی و پایداری روشهای چندگامی، که در بخش ۸.۶ معرفی شدند از دال‌کوئیست (۱۹۵۶) آغاز شده است. یک بیان تاریخی در دال‌کوئیست (۱۹۸۵) داده شده است. کتاب درسی هنریچی (۱۹۶۲) به صورت یک گزارش کلاسیک این نظریه، شامل بسط و کاربرد آن درآمده است. کتاب گی‌یر (۱۹۷۱) یک گزارش جدیدتر از همه این روشهاست، به‌ویژه روشهای متغیر-مرتب. اشتتر (۱۹۷۳) یک تحلیل مجرد کامل و عمومی برای نظریه عددی حل مسائل مقدار اولیه می‌دهد. یک گزارش کامل از روشهای رونگه-کوتا، بسط و تحلیل خطای آن روشها، تا سال ۱۹۷۰ توسط لاپیدوس و زاینفلت (۱۹۷۱) داده شده است. هال و وات (۱۹۷۶) یک بررسی از تمام جنبه‌های حل معادلات دیفرانسیل معمولی، شامل بسیاری از مباحث خاص کرده‌اند، که در ده سال اخیر مورد توجه بیشتر قرار گرفته‌اند.

اولین کاربرد قابل ملاحظه مفهوم روش متغیر-مرتب از آن گی‌یر (۱۹۷۱) و کرو^۹ (۱۹۶۹) است. چنین روشهایی نسبت به روش چندگامی ثابت مرتبه از نظر کارایی برتری دارند، و هیچ روش اضافی برای شروع انتگرالگیری یا برای تغییر طول گام نمی‌طلبند. یک بیان بسیار خوب از روش متغیر-مرتب آدامز در شمپاین و گوردن (۱۹۷۵)، شامل برنامه عالی DE/STEP داده شده است. برنامه‌های مهم دیگر پیشین که بر پایه فرمولهای خانواده آدامز، نهاده شده‌اند، برنامه‌هایی بوده‌اند در کرو (۱۹۶۹)، در DIFSUB از گی‌یر (۱۹۶۹) و در GEAR از هایندمارش^{۱۰} (۱۹۷۴) است. برنامه اخیر گی‌یر توسعه بیشتری یافته و به یک بسته بزرگ چندتابعی تبدیل شده است که ODEPACK نامیده می‌شود که در هایندمارش (۱۹۸۳) توضیح داده شده است. انواع مختلف

- | | | | |
|------------|---------------|------------------|-------------|
| 1. Collatz | 2. Milne | 3. Lapidus | 4. Seinfeld |
| 5. Lambert | 6. Stetter | 7. Vander Houwen | 8. Butcher |
| 9. Krogh | 10. Hindmarsh | | |

این برنامه‌ها و حل‌کننده‌های دیگر معادلات دیفرانسیل در کتابخانه‌های IMSL و NAG موجودند. روشهای رونگه-کوتا یک زمینه فعال پژوهشی نظری و توسعه برنامه‌ی هستند و یک توسعه بسیار کلی در پوچر (۱۹۸۷) داده شده است. روشهای جدید برای مسائل غیرسرخست در حال پدید آمدن هستند؛ برای مثال شمپاین (۱۹۸۶) و شمپاین و باکا^۱ (۱۹۸۶) را ببینید. همچنین علاقه زیادی به روشهای ضمنی رونگه-کوتا برای حل معادلات دیفرانسیل سرخست وجود دارد. برای مروری به این روشها ایکن (۱۹۸۵، صص ۷۰-۹۲) (۱۹۸۶) را ببینید. یک رقیب بسیار مهم برای برنامه RKF45، برنامه DVERK است که در هال و همکاران (۱۹۷۶) توضیح داده شده است. این برنامه براساس طرحی است از نوع فلبرک^۲ با یک جفت فرمولهای مرتبه ۵ و ۶.

دسته سوم از روشها در این کتاب نادیده گرفته شده است، روشهایی که بر پایه روش برونابی استوارند. تحقیقات فعلی در این زمینه توسط گرو^۳ (۱۹۶۸) و بولیرش و اشتور (۱۹۶۶) آغاز شده است. هدف اصلی، انجام برونابی مکرر از یک روش ساده است تا رفته رفته روشهایی از مراتب بالاتر و بالاتر به دست آید. در واقع، این عمل یک راه دیگری برای ایجاد روشهای متغیر-مرتبه به دست خواهد داد. این روشها در آزمونهای انزایت و هال (۱۹۷۶)، شمپاین و همکاران (۱۹۷۶)، نسبتاً خوب اجرا شده‌اند، ولی داوری درباره آنها این بوده که از جنبه نظری و عملی به اندازه روشهای چندگامی و رونگه-کوتا پیشرفته نیستند. برای یک بازنگری تازه در این زمینه، دویفل هارت (۱۹۸۵) را ببینید. همچنین شمپاین و باکا (۱۹۸۶) را ببینید که در آن روشهای برونابی به عنوان یک نمونه از روشهای متغیر-مرتبه رونگه-کوتا مورد بحث قرار گرفته است.

برآورد خطای کلی زمینه‌ای است که در آن مقالات نسبتاً کمی انتشار یافته‌اند. برای یک بررسی کلی اسکیل^۴ (۱۹۸۶) را ببینید. تا آنجا که اطلاع داریم، تنها برنامه رایانه‌ی متداول، برنامه GERK از شمپاین و واتس (۱۹۷۶a) است. کارهای بیشتری در این زمینه لازم است. بسیاری از کاربرهای بسته‌های خودکار تحت تاثیر این گمان اشتباه قرار دارند که برنامه خودکاری که آنها به کار می‌برند خطای کلی را کنترل می‌نماید، ولی این کنترل خطای کلی از لحاظ عملی امکان ندارد. در بسیاری از موارد داشتن تصویری از اندازه واقعی خطای کلی موجود در جواب عددی قاعدتاً مهم است.

مسائل مقدار مرزی برای معادلات دیفرانسیل معمولی مبحث مهم دیگری است، ولی هم نظریه عمومی آنها و هم تحلیل عددی آنها بسیار پیچیده‌تر از مسأله مقدار اولیه‌اند. کتابهای درسی مهم کلر (۱۹۷۸) و (۱۹۶۹) هستند، و گزارشهای چایلد و همکاران (۱۹۷۹) چندین مقاله مهم برای تهیه برنامه‌های رایانه‌ی ارائه می‌نماید. برای مقالات دیگر مجموعه عزیز (۱۹۷۵)، هال^۵ و وات (۱۹۷۶) و آشر^۶ و راسل^۷ (۱۹۸۵) را ببینید. این زمینه در مقایسه با مسأله مقدار اولیه برای معادلات غیرسرخست،

1. Baca

2. Fehlberg

3. Gragg

4. Skeel

5. Hall

6. Ascher

7. Russell

هنوز نسبتاً جدید است. مطالعه برنامه‌های رایانه‌ی در چندین جهت ادامه دارد و بعضی برنامه‌های نسبتاً خوب در سالهای اخیر تهیه شده‌اند. کارهای زیادتری در برنامه‌نویسی با استفاده از روش پرتابی انجام گرفته است ولی برنامه‌های عملی برای روشهای هم‌مکانی و تفاضل متناهی در حال تهیه هستند. برای بحثی از این برنامه‌ها انزایت (۱۹۸۵) و گلاادل و سه‌یژز (۱۹۸۰، صص ۲۷۳-۳۰۳) را ببینید. چند برنامه مقدار مرزی در مجموعه انتشارات IMSL و NAG داده شده‌اند.

معادلات دیفرانسیل سرسخت یکی از زمینه‌های خاصی است که در ده سال اخیر اهمیت خیلی بیشتری یافته‌اند. بهترین بررسی کلی در این زمینه در ایکن (۱۹۸۵) داده شده است. در این کتاب مثالهایی از چگونگی پیدایش این مسائل، نظریه روشهای عددی برای حل مسائل سرسخت، و یک بررسی از برنامه‌های رایانه‌ی موجود برای حل این مسائل داده شده است. بسیاری از کتابهای درسی دیگری که در فهرست منابع آمده‌اند، همچنین به مسأله معادلات دیفرانسیل سرسخت اشاره دارند. ما همچنین مطالعه کتاب شمپاین و گی‌یر (۱۹۷۹) را توصیه می‌کنیم.

معادلات با جواب خیلی نوسانی، در برخی از کاربردها ظاهر می‌شوند. برای بحث در مورد این مسائل، ایکن (۱۹۸۵، صص ۱۱۱-۱۲۳) را ببینید. روش خطوط، در حل معادلات دیفرانسیل با مشتقات جزئی وابسته به زمان، یک شیوه کلاسیک است که در سالهای اخیر خیلی متداول گشته است. این روش در ایکن (۱۹۸۵، صص ۱۲۴-۱۲۸)، سینکاوک و مدسین (۱۹۷۵) و ملخارت و سینکاوک (۱۹۸۱) مورد بحث قرار گرفته است.

باز یک زمینه دیگر قابل توجه، حل دستگاههای آمیخته معادلات دیفرانسیل و جبری (DAE) است. این نام به دستگاههایی داده می‌شود که در آنها، n مجهول، $m < n$ معادله دیفرانسیل، و $n - m$ معادله جبری، شامل همان n تابع مجهول، وجود دارند. چنین مسائلی در زمینه‌های کاربردی بسیاری ظاهر می‌شوند. یک چنین زمینه بسیار مورد توجه در سالهای اخیر طرح کمکی رایانه‌ی (CAD) است. برای مطالعه مقالاتی که در مورد این مسائل و مسائل دیگر DAE کاربرد دارند به راینبولت^۱ (۱۹۸۴) و (۱۹۸۶) مراجعه نمایید.

به دلیل ایجاد تعدادی برنامه‌های خودکار برای حل معادلات دیفرانسیل، مطالعات تجربی چندی انجام گرفته است تا اجرای این برنامه‌ها را برآورد کنند و آنها را با هم مقایسه نمایند. بعضی از مقایسه‌های مهم در انزایت و هال (۱۹۷۶)، انزایت و همکاران (۱۹۷۵) و شمپاین و همکاران (۱۹۷۶) داده شده‌اند. از کار آنها روشن می‌شود که برنامه‌ها و همچنین روشها باید مقایسه شوند. کاربردهای مختلف یک برنامه برای یک روش، ممکن است تفاوت‌های عمده‌ای در اجرا داشته باشد. هیچ نتایج مشابهی در مقایسه برنامه‌های مقدار مرزی در دست نیست.

- Aiken, R., Ed. (1985). *Stiff Computation*. Oxford Univ. Press, Oxford, England.
- Ascher, U. (1986). Collocation for two-point boundary value problems revisited, *SIAM J. Numer. Anal.* 23, 596-609.
- Ascher, U., J. Christiansen, and R. Russell (1981a). Collocation software for boundary-value ODEs, *ACM Trans. Math. Softw.* 7, 209-222.
- Ascher, U., J. Christiansen, and R. Russell (1981b). COLSYS: collocation software for boundary-value ODEs, *ACM Trans. Math. Softw.* 7, 223-229.
- Ascher, U., and R. Russell, Eds. (1985). *Numerical Boundary Value ODEs*. Birkhäuser, Basel.
- Aziz, A. K., Ed. (1975). *Numerical Solutions of Boundary Value Problems for Ordinary Differential Equations*. Academic Press, New York.
- Boyce, W., and R. DiPrima (1986). *Elementary Differential Equations and Boundary Value Problems*, 4th ed. Wiley, New York.
- Bulirsch, R., and J. Stoer (1966). Numerical treatment of ordinary differential equations by extrapolation, *Numer. Math.* 8, 1-13.
- Butcher, J. (1965). On the attainable order of Runge-Kutta methods, *Math. Comput.* 19, 408-417.
- Butcher, J. (1987). *The Numerical Analysis of Ordinary Differential Equations*. Wiley, New York.
- Childs, B., M. Scott, J. Daniel, E. Denman, and P. Nelson, Eds. (1979). *Codes for Boundary-Value Problems in Ordinary Differential Equations*. Lecture Notes in Computer Science 76, Springer-Verlag, New York.
- Coddington, E., and N. Levinson (1955). *Theory of Ordinary Differential Equations*. McGraw-Hill, New York.
- Collatz, L. (1966). *The Numerical Treatment of Differential Equations*, 3rd ed. Springer-Verlag, New York.
- Dahlquist, G. (1956). Numerical integration of ordinary differential equations, *Math. Scandinavica* 4, 33-50.
- Dahlquist, G. (1963). A special stability property for linear multistep methods, *BIT* 3, 27-43.
- Dahlquist, G. (1985). 33 years of numerical instability, part 1, *BIT* 25, 188-204.
- Deuffhard, P. (1979). Nonlinear equation solvers in boundary value problem codes. In B. Childs, M. Scott, J. Daniel, E. Denman, and P. Nelson (Eds.), *Codes for Boundary-Value Problems in Ordinary Differential Equations*, pp. 40-66. Lecture Notes in Computer Science 76, Springer-Verlag, New York.
- Deuffhard, P. (1985). Recent progress in extrapolation methods for ordinary differential equations, *SIAM Rev.* 27, 505-536.
- Enright, W. (1985). Improving the performance of numerical methods for two-point boundary value problems. In U. Ascher and R. Russell (Eds.), *Numerical Boundary Value ODEs*, pp. 107-120. Birkhäuser, Basel.

- Enright, W., and T. Hull (1976). Test results on initial value methods for non-stiff ordinary differential equations, *SIAM J. Numer. Anal.* 13, 944-961.
- Enright, W., T. Hull, and B. Lindberg (1975). Comparing numerical methods for stiff systems of O.D.E.'s, *BIT* 15, 10-48.
- Fehlberg, E. (1970). Klassische Runge-Kutta-Formeln vierter und niedrigerer Ordnung mit Schrittweiten-Kontrolle und ihre Anwendung auf Wärmeleitungsprobleme, *Computing* 6, 61-71.
- Forsythe, G., M. Malcolm, and C. Moler (1977). *Computer Methods for Mathematical Computations*. Prentice-Hall, Englewood Cliffs, N.J.
- Fox, L. (1980). Numerical methods for boundary-value problems. In I. Gladwell and D. Sayers (Eds.), *Computational Techniques for Ordinary Differential Equations*, pp. 175-217. Academic Press, New York.
- Gear, C. W. (1971). *Numerical Initial Value Problems in Ordinary Differential Equations*. Prentice-Hall, Englewood Cliffs, N.J.
- Gladwell, I., and D. Sayers, Eds. (1980). *Computational Techniques for Ordinary Differential Equations*. Academic Press, New York.
- Gragg, W. (1965). On extrapolation algorithms for ordinary initial value problems, *SIAM J. Numer. Anal.* 2, 384-403.
- Gupta, G., R. Sacks-Davis, and P. Tischer (1985). A review of recent developments in solving ODEs, *Comput. Surv.* 17, 5-47.
- Hall, G., and J. Watt, Eds. (1976). *Modern Numerical Methods for Ordinary Differential Equations*. Oxford Univ. Press, Oxford, England.
- Henrici, P. (1962). *Discrete Variable Methods in Ordinary Differential Equations*. Wiley, New York.
- Hindmarsh, A. (1974). GEAR: Ordinary differential equation solver. Lawrence Livermore Rep. UCID-30001, Rev. 3, Livermore, Calif.
- Hindmarsh, A. (1983). ODEPACK: A systematized collection of ODE solvers. In *Numerical Methods for Scientific Computation*, R. Stepleman, Ed. North-Holland, Amsterdam.
- Hull, T., W. Enright, and K. Jackson (1976). User's guide for DVERK: A subroutine for solving non-stiff ODEs. Dept. Computer Sci. Tech. Rep. 100, Univ. of Toronto, Toronto, Ont., Canada.
- Isaacson, E., and H. Keller (1966). *Analysis of Numerical Methods*. Wiley, New York.
- Jain, M. K. (1984). *Numerical Solution of Differential Equations*, 2nd ed. Halstead Press, New York.
- Keller, H. (1968). *Numerical Methods for Two-Point Boundary Value Problems*. Ginn (Blaisdell), Boston.
- Keller, H. (1976). *Numerical Solution of Two-Point Boundary Value Problems*. Regional Conference Series in Applied Mathematics 24, Society for Industrial and Applied Mathematics, Philadelphia.
- Keller, J., and S. Antman, Eds. (1969). *Bifurcation Theory and Nonlinear Eigenvalue Problems*. Benjamin, New York.

- Krogh, F. (1969), VODQ/SVDQ/DVDQ—Variable order integrators for the numerical solution of ordinary differential equations. Section 314 Sub-routine Writeup, Jet Propulsion Lab., Pasadena, Calif.
- Lambert, J. (1973). *Computational Methods in Ordinary Differential Equations*. Wiley, New York.
- Lapidus, L., and W. Schiesser, Eds. (1976). *Numerical Methods for Differential Equations: Recent Developments in Algorithms, Software, New Applications*, Academic Press, New York.
- Lapidus, L., and J. Seinfeld (1971). *Numerical Solution of Ordinary Differential Equations*. Academic Press, New York.
- Lentini, M., M. Osborne, and R. Russell (1985). The close relationship between methods for solving two-point boundary value problems, *SIAM J. Numer. Anal.* **22**, 280–309.
- Melgaard, D., and R. Sincovec (1981). Algorithm 565: PDETWO/PSETM/GEARB: Solution of systems of two-dimensional nonlinear partial differential equations, *ACM Trans. Math. Softw.* **7**, 126–135.
- Miller, R. K. and A. Michel (1982). *Ordinary Differential Equations*. Academic Press, New York.
- Milne, W. (1953). *Numerical Solution of Differential Equations*. Wiley, New York.
- Ortega, J., and W. Poole (1981). *An Introduction to Numerical Methods for Differential Equations*, Pitman, New York.
- Pereyra, V. (1979). PASVA3: An adaptive finite difference Fortran program for first order nonlinear, ordinary differential equation problems. In B. Childs, M. Scott, J. Daniel, E. Denman, and P. Nelson (Eds.), *Codes for Boundary-Value Problems in Ordinary Differential Equations*, pp. 67–88. Lecture Notes in Computer Science 76, Springer-Verlag, New York.
- Reddien, G. (1979). Projection methods. In B. Childs, M. Scott, J. Daniel, E. Denman, and P. Nelson (Eds.), *Codes for Boundary-Value Problems in Ordinary Differential Equations*, pp. 206–227. Lecture Notes in Computer Science 76, Springer-Verlag, New York.
- Rheinboldt, W. (1984). Differential-algebraic systems as differential equations on manifolds, *Math. Comput.* **43**, 473–482.
- Rheinboldt, W. (1986). *Numerical Analysis of Parametrized Nonlinear Equations*. Wiley, New York.
- Shampine, L. (1985). Local error estimation by doubling, *Computing* **34**, 179–190.
- Shampine, L. (1986). Some practical Runge–Kutta formulas, *Math. Comput.* **46**, 135–150.
- Shampine, L., and L. Baca (1986). Fixed versus variable order Runge–Kutta, *ACM Trans. Math. Softw.* **12**, 1–23.
- Shampine, L., and C. W. Gear (1979). A user's view of solving stiff ordinary differential equations, *SIAM Rev.* **21**, 1–17.
- Shampine, L., and M. Gordon (1975). *Computer Solution of Ordinary Differential Equations*. Freeman, San Francisco.

- Shampine, L., and H. Watts (1976a). Global error estimation for ordinary differential equations, *ACM Trans. Math. Softw.* **2**, 172-186.
- Shampine, L., and H. Watts (1976b). Practical solution of ordinary differential equations by Runge-Kutta methods. Sandia Labs. Tech. Rep. SAND 76-0585, Albuquerque, N.Mex.
- Shampine, L., and H. Watts (1980). DEPACK—Design of a user oriented package of ODE solvers. Sandia National Labs. Rep. SAND79-2374, Albuquerque, N.Mex.
- Shampine, L., H. Watts, and S. Davenport (1976). Solving nonstiff ordinary differential equations—The state of the art, *SIAM Rev.* **18**, 376-411.
- Sincovec, R., and N. Madsen (1975). Software for nonlinear partial differential equations, *ACM Trans. Math. Softw.* **1**, 232-260.
- Skeel, R. (1986). Thirteen ways to estimate global error, *Numer. Math.* **48**, 1-20.
- Stakgold, I. (1979). *Green's Functions and Boundary Value Problems*. Wiley, New York.
- Stetter, H. (1973). *Analysis of Discretization Methods for Ordinary Differential Equations*. Springer-Verlag, New York.
- Stoer, J., and R. Bulirsch (1980). *Introduction to Numerical Analysis*. Springer-Verlag, New York.
- Van der Houwen, P. J. (1977). *Construction of Integration Formulas for Initial Value Problems*. North-Holland, Amsterdam.
- Widder, D. (1975). *The Heat Equation*. Academic Press, New York.

مسائل

۱. میدان سودار $y' = x - y^2$ را رسم کنید، سپس یک نمونه از منحنیهای جواب را در آن رسم، و سعی کنید که رفتار جواب $Y(x)$ را وقتی $x \rightarrow \infty$ ، حدس بزنید.
۲. ثابتهای لیبیشیتس را برای توابع زیر مانند (۲.۱.۶) معین کنید.
- (الف) $x \geq 1$ ، $f(x, y) = 2y/x$
- (ب) $f(x, y) = \tan^{-1}(y)$
- (ج) $f(x, y) = (x^2 - 2)^{2y} / (17x^2 + 4)$
- (د) $|y| \leq 1$ ، $f(x, y) = x - y^2$
۳. مسائل زیر را به دستگاه معادلات مرتبه اول تبدیل کنید.
- (الف) $y'(\circ) = 1$ ، $y(\circ) = 1$ ، $y'' - 3y' + 2y = 0$
- (ب) $y'(\circ) = 0$ ، $y(\circ) = 1$ ، $y'' - 0.1(1 - y^2)y' + y = 0$ (معادله واندرپول)
- (ج) $r = \sqrt{x^2 + y^2}$ ، $y'' = -y/r^2$ ، $x'' = -x/r^2$ ، $x(\circ) = 4$ ، $x'(\circ) = 0$ (معادله مدار)، $y'(\circ) = 2$ ، $y(\circ) = 0$ ، $x'(\circ) = 0$

۴. گیریم $Y(x)$ ، جواب مسأله مقدار اولیه (۱.۰.۶)، در صورت وجود، باشد. با انتگرالگیری نشان دهید که Y در رابطه زیر صدق می‌کند

$$Y(x) = Y_0 + \int_{x_0}^x f(t, Y(t)) dt$$

به عکس، نشان دهید که اگر این معادله یک جواب پیوسته بر بازه $x_0 \leq x \leq b$ داشته باشد، مسأله مقدار اولیه (۱.۰.۶) نیز همان جواب را دارد.

۵. معادله انتگرالی مسأله ۴، حداقل به طور نظری، با استفاده از بارست

$$Y_{m+1}(x) = Y_0 + \int_{x_0}^x f(t, Y_m(t)) dt \quad x_0 \leq x \leq b$$

برای $m \geq 0$ با $Y_0(x) \equiv Y_0$ حل شده است. این بارست را بارست پیکار (Picard) گویند، و با مفروضات مناسبی، می‌توان نشان داد که بارستهای $\{Y_m(x)\}$ به طور یکتا به $Y(x)$ همگرا هستند. با محاسبه بارستهای Y_1, Y_2, Y_3 پیکار این حکم را برای مسائل زیر روشن کنید و آنها را با جواب درست Y مقایسه کنید.

(الف) $y(0) = 1, y' = -y$

(ب) $y(1) = 2, y' = -xy$

(ج) $y(0) = 1, y' = y + 2 \cos(x)$

۶. یک برنامه رایانه‌ی بنویسید که معادله $y' = f(x, y), y(x_0) = y_0$ را با روش اویلر حل کند. آن را به گونه‌ای بنویسید که برای تابع دلخواه f ، طول گام h ، و بازه $[x_0, b]$ بتوان به کار برد. با استفاده از این برنامه معادله $y' = x^2 - y, y(0) = 1$ را برای $0 \leq x \leq 4$ ، با طول گامهای $0.625, 0.125, 0.25, 0.5$ ، هر مقدار h ، جواب درست، جواب تقریبی، خطا، خطای نسبی را در نقاط گرهی $0, 1, 2, 3, 4, \dots, 0.75, 0.5, 0.25, 0$ را چاپ کنید. جواب درست چنین است: $Y(x) = x^2 - 2x + 2 - e^{-x}$. خروجی برنامه خود را تحلیل کنید و اظهارنظرهای کتبی برای آن بنویسید. تحلیل خروجی به اندازه بدست آوردن آن مهم است.

۷. برای مسأله $y' = y, y(0) = 1$ ، جواب صریح $\{y_n\}$ را برای تقریب اویلر این معادله بنویسید. با استفاده از آن نشان دهید که برای $x_n = 1$ ، وقتی $h \rightarrow 0$ ، داریم $Y(1) - y_n \approx (h/2)e$.

۸. حل مسأله زیر را با روش اویلر در نظر می‌گیریم

$$y' = -y^2 \quad x \geq 0 \quad y(0) = 1$$

کران (۱۳.۲.۶) را با به کار بردن $K = 2$ به عنوان ثابت لیشیتس، با برآورد مجانبی حاصل از (۳۶.۲.۶) مقایسه کنید. جواب درست $Y(x) = 1/(1+x)$ است.

۹. نشان دهید که روش اویلر در تقریب زدن جواب $Y(x) = (\frac{1}{2}x)^{3/2}$ ، $x \geq 0$ ، در مسأله $y' = y^{1/2}$ ، $y(0) = 0$ به شکست می انجامد. توضیح دهید چرا.

۱۰. برای معادلات مسأله ۳ معادلات تفاضلی تقریب زنده را که با استفاده از روش اویلر به دست می آیند بنویسید.

۱۱. مثال خطای گرد کردن به روش اویلر را با نتایجی که در جدول ۳.۶ نشان داده شده است به یاد آورید. سعی کنید رفتار مشابهی را در رایانه خود با کاهش تدریجی h ایجاد کنید تا اینکه خطا روبه افزایش گذارد.

۱۲. مسأله زیر را به یک دستگاه معادلات مرتبه اول تبدیل کنید.

$$y^{(3)} + 4y'' + 5y' + 2y = -4 \sin(x) - 2 \cos(x)$$

$$y(0) = 1 \quad y'(0) = 0 \quad y''(0) = -1$$

با استفاده از روش اویلر این دستگاه را حل و از راه تجربی، خطا را مطالعه کنید. جواب درست $Y = \cos x$ است.

۱۳. (الف) شرط (۵۲.۲.۶) لیشیتس را برای یک دستگاه دو معادله دیفرانسیل پیدا کنید.

(ب) ثابت کنید که روش (۵۱.۲.۶) به جواب (۵۰.۲.۶) همگراست.

۱۴. روش دوگامی زیر را در نظر می گیریم

$$y_{n+1} = \frac{1}{4}(y_n + y_{n-1}) + \frac{h}{4}[4y'_{n+1} - y'_n + 3y'_{n-1}] \quad n \geq 1$$

که در آن $y'_n \equiv f(x_n, y_n)$. نشان دهید که این روش مرتبه دو است و جمله پیشرو در خطای برشی را که به صورت (۱۵.۳.۶) نوشته شده پیدا کنید.

۱۵. فرض می کنیم که روش چندگامی (۱.۳.۶) سازگار است و همانند قضیه ۶.۶ در $a_j \geq 0$ ، $j = 0, 1, \dots, p$ ، $z = 0$ صدق می کند. پایداری (۱.۳.۶) را مشابه با (۲۸.۲.۶)، برای روش اویلر ثابت کنید، ولی δ را متحد با صفر بگیرید.

۱۶. برنامه ای برای حل $y' = f(x, y)$ ، $y(x_0) = y_0$ با استفاده از قاعده میانگامی (۲.۴.۶) بنویسید. از طول گام ثابت h استفاده کنید. برای مقدار اولیه y_1 ، روش اویلر را به کار برید:

$$y_1 = y_0 + hf(x_0, y_0)$$

با این برنامه مسائل زیر را حل کنید:

(الف) $Y(x) = 1/(1+x)$; $y(0) = 1$, $y' = -y^2$

(ب) $Y(x) = \frac{20}{[1 + 19e^{-x/2}]}$; $y(0) = 1$, $y' = \frac{y}{4} \left[1 - \frac{y}{20} \right]$

(ج) $Y(x) = \cos(x) + \sin(x)$; $y(0) = 1$, $y' = -y + 2 \cos(x)$

(د) $Y(x) = \cos(x) + \sin(x)$, $y(0) = 1$, $y' = y - 2 \sin(x)$

برای بازه $[0, 5]$ با $[x_0, b] = [0, 25]$ ، $h = 0.5$ ، 0.25 حل کنید. جواب عددی $y_n = y_h(x_n)$ را در هر نقطه گرهی همراه با خطای درست چاپ کنید. درباره نتایج خود بحث کنید.

۱۷. برنامه‌ای برای حل $y(x_0) = y_0$, $y' = f(x, y)$ با استفاده از قاعده ذوزنقه‌یی (۲.۵.۶) با طول گام ثابت h بنویسید. در بارست (۳.۵.۶)، روش میانگاهی را به عنوان پیشگوی $y_{n+1}^{(0)}$ به کار برید. تعداد بارستهای J را یک متغیر ورودی بگیرید، $J \geq 1$. معادلات مسأله ۱۶ را با $h = 0.5$ و $h = 0.25$ برای $[x_0, b] = [0, 10]$ حل کنید. ابتدا با $J = 1$ برای تمامی بازه حل کنید سپس فرایند را با $J = 2$ و $J = 3$ تکرار کنید. درباره نتایج بحث کنید. چگونه خطای کلی با J تغییر می‌کند؟
 ۱۸. فرمول خطای مجانبی (۱۹.۵.۶) را برای روش ذوزنقه‌یی به دست آورید.

۱۹. یک برنامه بنویسید که الگوریتم *Detrap* از بخش ۶.۶ را اجرا نماید. با استفاده از این برنامه معادلات داده شده در مسأله ۱۶ را با $\varepsilon = 0.001$ حل کنید.

۲۰. درونیاب درجه دوم را برای $Y'(x) = f(x, Y(x))$ در x_{n-2}, x_{n-1}, x_n به کار برید تا فرمول زیر به دست آید

$$Y(x_{n+1}) = Y_{n-2} + \frac{4h}{3} [2Y'_n - Y'_{n-1} + 2Y'_{n-2}] + \frac{28}{90} h^5 Y^{(5)}(\xi_n)$$

وقتی خطای برشی حذف شود روش

$$y_{n+1} = y_{n-2} + \frac{4h}{3} [2y'_n - y'_{n-1} + 2y'_{n-2}] \quad n \geq 3$$

را به دست می‌آوریم. این فرمول پیشگو برای روش میلن، (۱۳.۷.۶)، است.

۲۱. نشان دهید که روش (۱۳.۷.۶) سیمپسون تنها ضعیف-پایدار است، به همان تعبیری که در روش میانگاهی در بخش ۴.۶ آمده است.

۲۲. برنامه‌ای برای حل $x_0 \leq x \leq b$, $y(x_0) = y_0$, $y' = f(x, y)$ با استفاده از فرمول مرتبه چهار آدامز-مولتن با طول گام ثابت h بنویسید. از فرمول مرتبه چهار آدامز-بشفورت به عنوان پیشگو استفاده کنید. مقادیر اولیه y_1, y_2, y_3 را با استفاده از جواب درست مسأله تولید کنید.

معادلات مسأله ۱۶ را با $h = 0.5$ و $h = 0.25$ حل کنید. جوابهای به دست آمده و خطاهای درست را چاپ کنید. برای مقایسه با روش (۲۷.۷.۶)، مثال جدول ۱۴.۶ را نیز حل کنید. دقت برونیابی ریچاردسن را برای برآورد خطای

$$Y(x) - y_h(x) \doteq \frac{1}{15} [y_h(x) - y_{2h}(x)]$$

که بر پایه خطای کلی از مرتبه چهار استوار است کنترل نمایید. درباره تمام نتایج خود بحث کنید. ۲۳. (الف) برای ضرایب γ_i و δ_i از فرمولهای آدامز-بشفورت و آدامز-مولتن نشان دهید که $i \geq 1, \delta_i = \gamma_i - \gamma_{i-1}$.

(ب) برای فرمول p گامی (۲۶.۷.۶) آدامز-مولتن، ثابت کنید

$$y_{n+1} = y_{n+1}^{(0)} + h\gamma_p \nabla^{p+1} y'_{n+1} \quad (1)$$

که $y_{n+1}^{(0)}$ از فرمول $p+1$ گامی (۲۲.۷.۶) آدامز-بشفورت

$$y_{n+1}^{(0)} = y_n + h \sum_{j=0}^p \gamma_j \nabla^j y'_n \quad (2)$$

است. این فرمولها هر دو از مرتبه $p+1$ هستند. فرمول (۱) در به دست آوردن تصحیح کننده از پیشگو به کار می آید، و بر پایه اجرای تفاضلات پسرو معینی از یک مرحله به مرحله بعدی مبتنی است. یک نتیجه خیلی زیاد وابسته وجود دارد وقتی که یک پیشگوی p -گامی (مرتبه p) برای حل یک تصحیح کننده p گامی (مرتبه $p+1$) به کار برده می شود. [شمپاین و گوردن (۱۹۷۵، ص ۵۱) را ببینید].

۲۴. معادله نمونه (۱۶.۸.۶) را که در آن Λ یک ماتریس مربعی مرتبه m است در نظر می گیریم. فرض کنید که $\Lambda = P^{-1}DP$ ، که D یک ماتریس قطری با درایه های $\lambda_1, \dots, \lambda_m$ است. تابع برداری مجهول جدید $z = P^{-1}y(x)$ را وارد می کنیم. نشان دهید (۱۶.۸.۶) به شکلی که در (۱۷.۸.۶) داده شده تبدیل می شود، که نشان دهنده تبدیل به یک معادله نمونه یک بعدی است.

۲۵. با پیروی از مفاهیم بخشهای ۴.۶ و ۸.۶، یک شیوه کلی برای حل معادله تفاضلی خطی زیر ارائه دهید

$$y_{n+1} = a_0 y_n + a_1 y_{n-1}$$

این روش را برای حل معادلات زیر به کار برید

$$y_{n+1} = -\frac{1}{4}y_n + \frac{1}{4}y_{n-1} \quad (\text{الف})$$

$$y_{n+1} = y_n - \frac{1}{4}y_{n-1} \quad (\text{ب})$$

راهنمایی: فرمولی را که به دنبال (۲۲.۸.۶) آمده است ببینید.

۲۶. معادلهٔ تفاضلی خطی مرتبهٔ سه

$$u_{n+1} = u_n + c(u_{n-1} - u_{n-2}) \quad n \geq 2 \quad 0 < c < 1$$

را که u_2, u_1, u_0 داده شده‌اند، حل کنید. دربارهٔ $\lim_{n \rightarrow \infty} u_n$ چه می‌توان گفت؟

۲۷. روش عددی زیر را در نظر می‌گیریم

$$y_{n+1} = 4y_n - 3y_{n-1} - 2hf(x_{n-1}, y_{n-1}) \quad n \geq 1$$

مرتبهٔ آن را معین کنید. با یک مثال نشان دهید که این روش ناپایدار است.

۲۸. نشان دهید که روش دوگامی

$$y_{n+1} = 2y_{n-1} - y_n + h\left[\frac{5}{4}y'_n + \frac{1}{4}y'_{n-1}\right]$$

از مرتبهٔ ۲ و ناپایدار است. همچنین مستقیماً نشان دهید که در حل $y' = f(x, y)$ ممکن است همگرا نباشد.

۲۹. قسمت (۱) برهان قضیهٔ ۷.۶ را، که در آن شرط ریشه با فرض $|r_j| = 1$ و $\rho'(r_j) = 0$ برای مقداری از j ، نقض شده است کامل کنید.

۳۰. برای قسمت (۲) از برهان قضیهٔ ۸.۶ نشان دهید که هرگاه $h \rightarrow 0$ و $1 \leq j \leq p$ ، آنگاه $\gamma_j[r_j(h\lambda)]^n \rightarrow 0$.

۳۱. (الف) مقادیر a_0 را در روش صریح مرتبهٔ دوم (۱۰.۷.۶) که برای آن این روش پایدار است معین کنید.

(ب) اگر فقط خطای برشی (۱۱.۷.۶) در نظر گرفته شود، با رعایت محدودیت پایداری در قسمت (الف)، چگونه a_0 باید انتخاب شود؟

(ج) برای تامین یک ناحیهٔ بزرگ پایداری، با رعایت قسمت (الف)، a_0 چگونه باید انتخاب شود؟

۳۲. (الف) فرمول کلی را برای تمام روشهای دوگامی مرتبهٔ ۳ به دست آورید. این فرمولها یک خانواده از روشهای تک-پارامتری، مثلاً وابسته به ضریب a_1 ، خواهند بود.

(ب) برای پایدار بودن روش قسمت (الف) چه محدودیتهایی برای a_1 باید وجود داشته باشد؟

(ج) اگر خطای برشی به شکل

$$T_{n+1}(Y) = \beta h^r Y^{(r)}(x_n) + O(h^5)$$

نوشته شود فرمولی برای β برحسب a_1 به دست آورید. (لازم نیست که هسته پتانویا تابع اثر این روش را بسازید). a_1 چگونه انتخاب شود تا با رعایت محدودیت پایداری قسمت (ب)، خطای برشی مینیمم شود؟

(د) ناحیه پایداری مطلق برای روشهای قسمت (الف) را در نظر می‌گیریم. این ناحیه برای روش قسمت (ج) چه باشد تا ضریب β ی خطای برشی به حداقل برسد؟ مقدار دیگری از a_1 بدهید که روش پایدار به دست دهد و ناحیه پایداری مطلق بزرگتری داشته باشد. درباره انتخاب تقریبی a_1 برای یک ناحیه بهینه بحث کنید.

۳۳. (الف) تمام فرمولهای صریح مرتبه چهار به شکل زیر را پیدا کنید

$$y_{n+1} = a_0 y_n + a_1 y_{n-1} + a_2 y_{n-2} + h[b_0 y'_n + b_1 y'_{n-1} + b_2 y'_{n-2}] \quad n \geq 2$$

(ب) نشان دهید که هر چنین روشی ناپایدار است.

۳۴. یک روش چندگامی ضمنی مرتبه چهار غیر از آنچه در کتاب داده شده است به دست آورید. آن را نسبی - پایدار سازید.

۳۵. برای چندجمله‌یی $\rho(r) = r^{p+1} - \sum_{j=0}^p a_j r^{p-j}$ فرض کنید $a_j \geq 0$ ، $0 \leq j \leq p$ ، $\sum_{j=0}^p a_j = 1$ و نشان دهید که ریشه‌های $\rho(r)$ در شرط ریشه (۹.۸.۶) و (۱۰.۸.۶) صدق می‌نماید. این مستقیماً نشان می‌دهد که قضیه ۶.۶ یک فرع قضیه ۸.۶ است.

۳۶. (الف) روشهایی به شکل زیر را در نظر می‌گیریم

$$y_{n+1} = y_{n-q} + h \sum_{j=-1}^p b_j f(x_{n-j}, y_{n-j})$$

با $q \geq 1$. نشان دهید که این روشها در شرط قوی ریشه صدق نمی‌نمایند. در نتیجه، اغلب چنین روشهایی فقط ضعیف - پایدارند.

(ب) یک مثال با $q = 1$ پیدا کنید که نسبتاً پایدار باشد.

۳۷. نشان دهید که ناحیه پایداری مطلق روش دوزنقه‌یی مجموعه تمام کمیتهای مختلط $h\lambda$ با $\text{Real}(\lambda) < 0$ است.

۳۸. روش پرسرو اوپلر را برای حل مسأله (۵۱.۸.۶) به کار برید. چون این معادله خطی است، معادله ضمنی برای y_{n+1} را می‌توان دقیقاً حل کرد. نتایج خود را با آنچه در جدول ۱۷.۶ برای روش اوپلر داده شده مقایسه کنید.

۳۹. مسأله ۳۸ را با استفاده از فرمول مرتبه ۲ (۶.۹.۶) BDF تکرار کنید. برای پیدا کردن y جهت استفاده در (۶.۸.۶) از روش پسرو اوایلر استفاده کنید.

۴۰. معادله نمونه (۵۰.۸.۶) را که به عنوان یک اختلال از روش (۴۶.۸.۶) اوایلر در نظر گرفته شده است به یاد آورید. برای حالت خاص ثابت $Y''(x) \equiv Y''(x) - y_n$ رفتار خطای $Y(x_n) - y_n$ را که به $h\lambda$ وابسته است تحلیل نمایید. نشان دهید باز هم برای خوشرفتار بودن خطا شرط (۴۹.۸.۶) مورد نیاز است.

۴۱. فرمول خطای برشی (۷.۹.۶) را برای فرمولهای مشتقگیری پسرو به دست آورید.

۴۲. برای حل $y' = f(x, y)$ روش عددی زیر را در نظر می‌گیریم

$$y_{n+1} = y_n + \frac{h}{4}[y'_n + y'_{n+1}] + \frac{h^2}{12}[y''_n - y''_{n+1}] \quad n \geq 0$$

در اینجا $y'_n = f(x_n, y_n)$

$$y_n'' = \frac{\partial f(x_n, y_n)}{\partial x} + f(x_n, y_n) \frac{\partial f(x_n, y_n)}{\partial y}$$

که این فرمول براساس مشتقگیری از $Y'(x) = f(x, Y(x))$ به دست آمده است.

(الف) نشان دهید که این روش یک روش مرتبه چهار است: $T_n(Y) = O(h^4)$

(ب) نشان دهید که ناحیه پایداری مطلق شامل تمام محور حقیقی منفی از صفحه مختلط $h\lambda$ است.

۴۳. روش خطوط را که در (۲۳.۹.۶) - (۲۵.۹.۶) داده شده برای مسأله زیر تعمیم دهید:

$$U_t = a(x, t)U_{xx} + G(x, t, U(x, t)) \quad 0 < x < 1 \quad t > 0$$

$$U(0, t) = d_0(t) \quad U(1, t) = d_1(t) \quad t \geq 0$$

$$U(x, 0) = f(x) \quad 0 \leq x \leq 1$$

برای آنکه مسأله خوش-تعریف باشد فرض می‌کنیم $a(x, t) > 0$ و $0 \leq x \leq 1$ و $t \geq 0$.

۴۴. (الف) اگر یک حلال دستگاه جبری خطی سه قطری در اختیار داشته باشید، برنامه‌ای بنویسید که روش خطوط را برای مسأله (۱۹.۹.۶) - (۲۱.۹.۶) اجرا کند. مثال کتاب با مجهول

(۳۵.۹.۶) با روش پسرو اوایلر حل شده است. اکنون روش خطوط را با استفاده از قاعده ذوزنقه‌یی

اجرا کنید و نتایج خود را با آنچه در جدول ۲۰.۶ برای روش پسرو اوایلر داده شده مقایسه نمایید.

(ب) با روش BDF مرتبه دوم تکرار کنید.

۴۵. یک روش مرتبه سه سری تیلر برای حل $y' = -y^2$ به دست آورید. نتایج عددی را با نتایج

جدول ۲۲.۶ مقایسه نمایید.

۴۶. با استفاده از روش سری تیلر بخش ۱۰.۶، یک روش مرتبه چهار برای حل $y' = x - y^2$

$y(0) = 0$ ، تولید کنید. طول گامهای ثابت $h = 0.5$ ، $h = 0.25$ و $h = 0.125$ را پشت سرهم

به‌کار برده مسئله را برای $0 \leq x \leq 1$ حل کنید. خطای کلی را با برآورد خطای (۲۴.۱۰.۶) که بر پایهٔ برونمایی ریچاردسن استوار است برآورد کنید.

۴۷. برنامه‌ای برای حل $y(x_0) = y_0$, $y' = f(x, y)$ با استفاده از روش کلاسیک (۲۱.۱۰.۶) رونگه-کوتا بنویسید و طول گام h را ثابت بگیرید.

(الف) با استفاده از این برنامه، معادلات مسئلهٔ ۱۶ را حل کنید.

(ب) معادلهٔ $y' = x - y^2$, $y(0) = 0$ را برای $h = 0.5, 0.25, 0.125$ حل کنید. نتایج را با نتایج مسئلهٔ ۴۶ مقایسه کنید.

۴۸. فرمول سه‌گامی رونگه-کوتا را که ذیلاً داده شده در نظر می‌گیریم

$$y_{n+1} = y_n + h[\gamma_1 V_1 + \gamma_2 V_2 + \gamma_3 V_3]$$

$$V_1 = f(x_n, y_n), \quad V_2 = f(x_n + \alpha_2 h, y_n + h\beta_{21} V_1)$$

$$V_3 = f(x_n + \alpha_3 h, y_n + h(\beta_{31} V_1 + \beta_{32} V_2))$$

مجموعهٔ معادلاتی را که ضرایب $\{\gamma_j, \alpha_j, \beta_j\}$ باید در آنها صدق کنند تا این فرمول از مرتبهٔ ۳ باشد به‌دست آورید. یک جواب خاص این معادلات را پیدا کنید.

۴۹. ثابت کنید که اگر روش (۴.۱۰.۶) رونگه-کوتا در رابطهٔ (۲۷.۱۰.۶) صدق کند، روش پایدار است.

۵۰. روش کلاسیک (۲۱.۱۰.۶) رونگه-کوتا را برای مسئلهٔ آزمونی (۵۱.۸.۶) برای مقادیر چندی از λ و h به‌کار برید. برای مثال، $\lambda = -1, -10, -50$ و $h = 0.5, 0.1, 0.05$ را مانند جدول ۱۷.۶ امتحان کنید.

۵۱. قسمت حقیقی ناحیهٔ پایداری مطلق را برای روش (الف) (۸.۱۰.۶)، (ب) (۹.۱۰.۶) (ج) (۲۱.۱۰.۶) رونگه-کوتا به‌دست آورید. ما به رفتار جواب عددی معادلهٔ دیفرانسیل $y' = \lambda y$ با $\text{Real}(\lambda) < 0$ به‌ویژه به مقادیری از $h\lambda$ علاقه‌مندیم که برای آنها جواب عددی با $x_n \rightarrow \infty$ به صفر میل می‌کند.

۵۲. (الف) با به‌کار بردن روش (۸.۱۰.۶) رونگه-کوتا، معادلهٔ

$$y' = -y + x^{0.1} [1 + x] \quad y(0) = 0$$

را که جواب آن $Y(x) = x^{0.1}$ است حل کنید. این معادله را در $[0, 5]$ حل کنید و خطاها را در $x = 1, 2, 3, 4, 5$ طول گامهای $h = 0.1, 0.05, 0.025, 0.0125, 0.00625$ طول گامهای

به‌کار برید. خطاهایی را معین کنید که بر اثر آنها وقتی h نصف می‌شود خطاها کاهش یابند. چگونه این عمل با نرخ نظری معمولی همگرایی ($O(h^2)$) مقایسه می‌شود؟ نتایج خود را توضیح دهید.

(ب) چه مشکلی پیش می‌آید وقتی بخواهیم معادله قسمت (الف) را با روش ناکوچکتر از مرتبه ۲ تیلر حل کنیم؟ این نتیجه چه چیزی را دربارهٔ جواب به ما می‌گوید؟

۵۳. مسأله مقدار مرزی (۱.۱۱.۶) را به یک مسأله مقدار مرزی هم‌ارز برای یک دستگاه معادلات مرتبه اول، مانند (۱۵.۱۱.۶) تبدیل کنید.

۵۴. (الف) مسأله مقدار مرزی (۲۵.۱۱.۶) را در نظر می‌گیریم. برای تبدیل آن به یک مسأله هم‌ارز با شرایط مرزی صفر، می‌نویسیم $y(x) = z(x) + w(x)$ ، که $w(x)$ خط راستی است که در شرایط مرزی زیر صدق می‌نماید: $w(a) = \gamma_1$ ، $w(b) = \gamma_2$. یک مسأله جدید مقدار مرزی برای $z(x)$ به‌دست آورید.

(ب) این روند را برای مسأله (۱۰.۱۱.۶) تعمیم دهید. یک مسأله جدید با شرایط مرزی صفر به‌دست آورید. چه مفروضاتی، اگر فرضی لازم باشد، برای ضرایب a_0, a_1, b_0, b_1 لازم است؟

۵۵. با استفاده از روش پرتابی بخش ۱۱.۶، مسائل مقدار مرزی زیر را حل کنید. نرخ همگرایی را وقتی h تغییر می‌کند مطالعه کنید.

$$(الف) \quad 1 < x < 2, y'' = \frac{-2}{x} yy', \quad y(2) = \frac{2}{3}, \quad y(1) = \frac{1}{4}$$

جواب درست است. $Y(x) = x/(1+x)$

(ب) $0 < x < \frac{\pi}{4}, y(0) = 0, y(\pi/4) = 1$ و $y'' = 2yy'$ جواب $Y(x) = \tan(x)$ درست است.

۵۶. برنامه‌های معادله دیفرانسیل را که مرکز محاسبات خودتان تهیه دیده بررسی نمایید. آنهایی را در نظر بگیرید که با تغییر طول گام و احتمالاً مرتبه، خطای برشی را به‌طور خودکار کنترل می‌نمایند.

برنامه‌ها را برحسب اینکه چندگامی (ثابت - مرتبه یا متغیر - مرتبه)، رونگه - کوتا یا برونایی رده‌بندی کنید. یکی از این برنامه‌ها را با DDEABM [از بخش ۷.۶ و شمپاین و گوردن (۱۹۷۵)] و

RKF45 [از بخش ۹.۶ و شمپاین و واتس (۱۹۷۶b)] با حل کردن مسأله

$$y' = \frac{y}{4} \left[1 - \frac{y}{20} \right] \quad y(0) = 1$$

با خطاهای مطلق مطلوب 10^{-3} ، 10^{-6} و 10^{-9} مقایسه نمایید. نتایج را با نتایج جدولهای ۱۵.۶ و ۲۸.۶ مقایسه نمایید.

۵۷. مسأله

$$y' = \frac{1}{t+1} + c \cdot \tan^{-1}(y(t)) - \frac{1}{4} \quad y(0) = 0$$

را که c یک ثابت داده شده است، در نظر می‌گیریم. چون $y'(0) = \frac{1}{4}$ ، جواب $y(t)$ وقتی t افزایش پیدا کند بدون توجه به مقدار c در آغاز زیاد می‌شود. تا آنجا که می‌توانید، نشان دهید که مقداری برای c وجود دارد که آن را c^* می‌نامیم و برای آن (۱) اگر $c > c^*$ ، جواب $y(t)$ به بینهایت افزایش می‌یابد و (۲) اگر $c < c^*$ ، $y(t)$ ابتداً زیاد می‌شود و به حداکثر می‌رسد و بعد کاهش می‌یابد. c^* را با دقت 0.00005 تعیین و سپس جواب مربوط $y(t)$ را برای $0 \leq t \leq 5$ پیدا کنید.

۵۸. دستگاه زیر را در نظر بگیرید

$$x'(t) = Ax - Bxy \quad y'(t) = Cxy - Dy$$

این دستگاه به مدل شکار و شکارچی لوتکا-ولترا معروف است، که $x(t)$ تعداد شکار و $y(t)$ تعداد شکارچی در زمان t را مشخص می‌کند،

(الف) گیریم $A = 4$ ، $B = 2$ ، $C = 1$ ، $D = 3$. مدل را حداقل تا سه رقم معنی‌دار برای $0 \leq t \leq 5$ حل کنید. مقادیر اولیه $x(0) = 3$ و $y(0) = 5$ هستند. x و y را به صورت توابعی از t رسم کنید و x را نسبت به y رسم نمایید.

(ب) همان مدل را با $x(0) = 3$ و به ترتیب 2 ، 1.5 ، 1 حل کنید. x را نسبت به y در هر حالت رسم کنید. چه چیزی مشاهده می‌کنید؟ چرا نقطه $(3, 2)$ نقطه تعادل خوانده می‌شود؟

جبر خطی

حل دستگاه معادلات خطی و به دست آوردن ویژه مقدارها و ویژه بردارهای ماتریس دو مسأله بسیار مهم‌اند که در موقعیتهای گوناگون پیش می‌آیند. به عنوان مقدمه بحث در این گونه مسائل در فصلهای آینده، بعضی از قضایای جبرخطی را ارائه می‌دهیم. بخش نخست شامل مروری بر مطالب فضاهای برداری، ماتریسها، و دستگاههای خطی است که در بسیاری از دروس جبر خطی در دوره کارشناسی تدریس می‌شوند. این قضایا فقط خلاصه شده‌اند، و هیچ‌گونه راه به دست آوردن آنها ذکر نشده است. در بقیه بخشها، از ویژه مقدارها، شکلهای متعارف ماتریسها، نرمهای بردار و ماتریس و قضایای اختلال در وارون ماتریسها بحث می‌شود. اگر لازم باشد می‌توان از خواندن این فصل صرف نظر کرد و هر زمان در فصلهای ۸ و ۹ لازم شود به آن مراجعه کرد. نمادگذاری، بخش ۱.۷ و نماد نرم در بخش ۳.۷ باید به اجمال مورد توجه قرار گیرند.

۱.۷ فضاهای برداری، ماتریسها و دستگاههای خطی

قطع نظر از جزئیات، فضای برداری V مجموعه‌ای است از اشیائی به نام بردار که برای آنها جمع برداری و ضرب اسکالر تعریف شده‌اند. هر فضای برداری V یک مجموعه اسکالر وابسته دارد و در این کتاب این مجموعه ممکن است اعداد حقیقی \mathbf{R} یا اعداد مختلط \mathbf{C} باشد. عملیات برداری باید در قواعد

استاندهٔ شرکتپذیری، تعویضپذیری و توزیعپذیری معینی صدق نمایند، که در اینجا فهرست نمی‌شوند. W ، زیرمجموعهٔ فضای برداری V زیر فضای V خوانده می‌شود، اگر W خود یک فضای برداری باشد و از همان عملیات در V استفاده شود. برای یک بحث کامل در نظریهٔ فضاهای برداری، به هر یک از کتابهای جبر خطی دورهٔ کارشناسی می‌توان مراجعه نمود [برای مثال، آنتون^۱ (۱۹۸۴)، فصل ۳)، هالموس^۲ (۱۹۸۵)، فصل ۱)، نوبل^۳ (۱۹۶۹)، فصلهای ۴ و ۱۴)، استرانگ^۴ (۱۹۸۰)، فصل ۲).

مثال ۱. $V = \mathbb{R}^n$ ، مجموعهٔ تمام n -تاییهای (x_1, \dots, x_n) ، با درایه‌های حقیقی x_i و \mathbb{R} مجموعهٔ اسکالرهایی وابسته به آن است.

۲. $V = \mathbb{C}^n$ ، مجموعه تمام n -تاییها با درایه‌های مختلط و \mathbb{C} مجموعهٔ اسکالرهاست.

۳. $V =$ مجموعهٔ چندجمله‌بیهایی از درجهٔ نابزرگتر از n ، برای یک n مفروض، یک فضای برداری است. اسکالرها می‌توانند \mathbb{R} یا \mathbb{C} ، بسته به کاربرد مطلوب، باشند.

۴. $V = C[a, b]$ ، مجموعهٔ همهٔ توابع پیوستهٔ حقیقی مقدار [یا مختلط مقدار] بر بازهٔ $[a, b]$ با مجموعهٔ اسکالر \mathbb{R} [یا \mathbb{C}]، یک فضای برداری است. مثال (۳) یک زیر فضای $C[a, b]$ است.

تعریف گیریم V یک فضای برداری باشد و $v_1, v_2, \dots, v_m \in V$.

۱. گوئیم v_1, v_2, \dots, v_m خطی-وابسته‌اند اگر یک مجموعه از اسکالرهایی $\alpha_1, \dots, \alpha_m$ که حداقل یکی از آنها مخالف صفر است، وجود داشته باشد به طوری که

$$\alpha_1 v_1 + \dots + \alpha_m v_m = 0$$

چون حداقل یکی از اسکالرها، مثلاً α_i ، مخالف صفر است می‌توانیم بنویسیم

$$v_i = -\frac{\alpha_1}{\alpha_i} v_1 - \dots - \frac{\alpha_{i-1}}{\alpha_i} v_{i-1} - \frac{\alpha_{i+1}}{\alpha_i} v_{i+1} - \dots - \frac{\alpha_m}{\alpha_i} v_m$$

گوئیم v_i یک ترکیب خطی از بردارهای $v_1, \dots, v_{i-1}, v_{i+1}, \dots, v_m$ است. برای اینکه یک مجموعه از بردارها خطی-وابسته باشد، باید یکی از آنها یک ترکیب خطی از بقیه باشد.

۲. بردارهای v_1, \dots, v_m را خطی-مستقل گوئیم اگر وابسته نباشند. به عبارت دیگر تنها انتخاب

اسکالرهایی $\alpha_1, \dots, \alpha_m$ برای آنکه رابطه

$$\alpha_1 v_1 + \dots + \alpha_m v_m = 0$$

برقرار باشد، انتخاب نمایان $\alpha_1 = \dots = \alpha_m = 0$ باشد. هیچ یک از v_i ها را نتوان به صورت ترکیبی از بقیه نوشت.

۳. $\{v_1, \dots, v_m\}$ یک پایه برای V است اگر برای هر $v \in V$ ، یک انتخاب یکتا از اسکالرهایی $\alpha_1, \dots, \alpha_m$ موجود باشد به طوری که

$$v = \alpha_1 v_1 + \dots + \alpha_m v_m$$

توجه نمایید که از این رابطه نتیجه می شود که v_1, \dots, v_m مستقل اند. اگر چنین پایه متناهی وجود داشته باشد، گوئیم V متناهی - بُعد است. در غیر این صورت، آن را نامتناهی - بُعد خوانیم.

قضیه ۱.۷. اگر V یک فضای برداری با پایه $\{v_1, \dots, v_m\}$ باشد، هر پایه V دقیقاً m بردار دارد. عدد m بُعد V خوانده می شود.

مثال ۱. $\{1, x, x^2, \dots, x^n\}$ یک پایه برای فضای V از چند جمله‌بیهای از درجه نایزتر از n است. بنابراین بُعد V $n + 1$ است.

۲. R^n و C^n دارای پایه $\{e_1, \dots, e_n\}$ هستند، که در آن

$$e_i = (0, 0, \dots, 0, 1, 0, \dots, 0) \quad (1.1.7)$$

که ۱ در موضع i ام قرار دارد. بعد R^n و C^n برابر n است. این پایه را پایه استاندارد برای R^n و C^n خوانند و بردارهای آنها را بردارهای یکه نامند.

۳. $C[a, b]$ نامتناهی - بُعد است.

ماتریسها و دستگاههای خطی ماتریسها آرایه‌های مستطیلی از اعداد حقیقی یا مختلط هستند، و در حالت کلی، ماتریس از مرتبه $m \times n$ دارای شکل زیر است:

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & & & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix} \quad (2.1.7)$$

ماتریس مرتبه n ، نامگذاری مختصر ماتریس مربعی از مرتبه $n \times n$ است. ماتریسها با حروف بزرگ و درایه‌های آنها معمولاً با حروف کوچک نشان داده می‌شوند، که معمولاً نظایر نام ماتریس‌اند، همان‌طور که هم‌اکنون نشان داده شد. تعاریف زیر عملیات معمولی روی ماتریسها را به‌دست می‌دهند.

تعریف ۱. گیریم A و B از مرتبه $m \times n$ باشند. مجموع A و B ماتریس $C = A + B$ از مرتبه $m \times n$ با درایه‌های زیر است

$$c_{ij} = a_{ij} + b_{ij}$$

۲. گیریم A از مرتبه $m \times n$ و α یک اسکالر باشد. در این صورت $C = \alpha A$ ، مضرب اسکالر A از مرتبه $m \times n$ است و با

$$c_{ij} = \alpha a_{ij}$$

داده می‌شود.

۳. گیریم A از مرتبه $m \times n$ و B از مرتبه $n \times p$ باشد. در این صورت حاصلضرب $C = AB$ از مرتبه $m \times p$ است که به شکل زیر داده می‌شود

$$c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$$

۴. گیریم A از مرتبه $m \times n$ باشد. ترانزاده $C = A^T$ از مرتبه $n \times m$ است و با رابطه زیر داده می‌شود

$$c_{ij} = a_{ji}$$

ترانزاده مزدوج $C = A^*$ نیز از مرتبه $n \times m$ است و

$$c_{ij} = \bar{a}_{ji}$$

نماد \bar{z} معرف مزدوج مختلط عدد مختلط z است و z حقیقی است اگر و تنها اگر $\bar{z} = z$. ترانزاده مزدوج A^* ماتریس الحاقی A نیز خوانده می‌شود.

ویژگیهای حسابی ماتریسها که در زیر آمده‌اند، بدون اشکال زیادی ثابت می‌شوند و اثبات آنها به عهده خواننده واگذار می‌شود

$$\begin{aligned} (A+B)+C &= A+(B+C) & \text{(ب)} & \quad A+B=B+A & \text{(الف)} \\ A(BC) &= (AB)C & \text{(د)} & \quad A(B+C)=AB+AC & \text{(ج)} & \quad (۳.۱.۷) \\ (AB)^T &= B^T A^T & \text{(و)} & \quad (A+B)^T = A^T + B^T & \text{(ه)} \end{aligned}$$

توجه به این امر مهم است که بدانیم در بسیاری از کاربردها برای برقراری ویژگیهای فوق لازم نیست ماتریسها مربعی باشند.

فضاهای برداری R^n و C^n معمولاً با مجموعه بردارهای ستونی از مرتبه $1 \times n$ به ترتیب با درایه‌های حقیقی و مختلط تعریف می‌شوند. دستگاه خطی

$$\begin{aligned} a_{11}x_1 + \dots + a_{1n}x_n &= b_1 \\ &\vdots \\ a_{m1}x_1 + \dots + a_{mn}x_n &= b_m \end{aligned} \quad (۴.۱.۷)$$

را می‌توان به صورت $Ax = b$ نوشت که A ماتریس $(۲.۱.۷)$ است، و

$$x = [x_1, \dots, x_n]^T \quad b = [b_1, \dots, b_m]^T$$

بردار b یک بردار داده‌شده در R^m است و جواب x بردار مجهولی در R^n . استفاده از ضرب ماتریسی، دستگاه خطی $(۴.۱.۷)$ را به شکل ساده‌تر و شهودی‌تر $Ax = b$ در می‌آورد.

ما اکنون چند تعریف اضافی دیگر برای ماتریسها، از جمله بعضی ماتریسهای خاص می‌دهیم.

تعریف ۱. ماتریس صفر از مرتبه $m \times n$ همه درایه‌هایش برابر صفرند. آن را با نماد $0_{m \times n}$ یا به‌طور ساده‌تر با 0 نشان می‌دهند. برای هر ماتریس A از مرتبه $m \times n$

$$A + 0 = 0 + A = A$$

۲. ماتریس همانی از مرتبه n با $I = [\delta_{ij}]$ تعریف می‌شود،

$$\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases} \quad (۵.۱.۷)$$

برای جمیع مقادیر $1 \leq i, j \leq n$. برای هر ماتریس A از مرتبه $m \times n$ و هر ماتریس B از مرتبه $n \times p$

$$AI = A \quad IB = B$$

نماد δ_{ij} معرف تابع دلتای کرونکر است.

۳. گیریم A یک ماتریس مربعی از مرتبه n باشد. اگر یک ماتریس مربعی B از مرتبه n وجود داشته باشد به طوری که $AB = BA = I$ ، آنگاه A را وارونپذیر یا وارون B خوانند. می توان نشان داد که ماتریس B یکتاست و وارون A را با نماد A^{-1} نشان می دهیم.

۴. ماتریس A متقارن نامیده می شود اگر $A^T = A$ ، و ارمیتی خوانده می شود اگر $A = A^*$. اصطلاح متقارن معمولاً فقط برای ماتریسهای حقیقی به کار برده می شود. ماتریس A پادمتقارن نامیده می شود اگر $A^T = -A$. الزاماً همه ماتریسهای متقارن، ارمیتی، یا پادمتقارن باید مربعی نیز باشند.

۵. گیریم A یک ماتریس $m \times n$ باشد. رتبه سطری A تعداد سطرهای خطی مستقل A است که عناصری از R^n یا C^n تلقی می شوند، و رتبه ستونی تعداد ستونهای خطی مستقل است. می توان نشان داد (مسئله ۴) که این دو عدد همیشه مساوی اند و این را رتبه A خوانند.

برای تعریف و ویژگیهای دترمینان یک ماتریس مربعی A ، می توان به هر کتاب درسی جبر خطی مراجعه نمود [برای مثال، آنتون (۱۹۸۴)، فصل ۲؛ نوبل (۱۹۶۹)، فصل ۷؛ استرانگ (۱۹۸۰)، فصل ۴] را ببینید. بسیاری از قضایای وارون ماتریسها و حلپذیری دستگاههای خطی را در قضیه زیر خلاصه می کنیم.

قضیه ۲.۷. گیریم A یک ماتریس مربعی با عناصری از R (یا C) باشد، و فضای برداری (یا $V = C^n$) باشد. در این صورت عبارات زیر هم ارزند:

۱. $Ax = b$ برای هر $b \in V$ یک جواب یکتای $x \in V$ دارد.

۲. $Ax = b$ برای هر $b \in V$ یک جواب $x \in V$ دارد.

۳. $Ax = 0$ ایجاب می کند $x = 0$.

۴. A^{-1} وجود دارد.

۵. دترمینان $(A) \neq 0$.

۶. رتبه $(A) = n$.

برهانی در اینجا آورده نشده است، یک تمرین بسیار خوب، اثبات هم‌ارزی بعضی از این احکام است. مفاهیم استقلال خطی و پایه را همراه با قضیه ۱.۷ به کار برید. همچنین از تجزیه زیر استفاده کنید

$$Ax = x_1 A_{*1} + \dots + x_n A_{*n} \quad x \in \mathbf{R}^n \text{ یا } \mathbf{C}^n \quad (۶.۱.۷)$$

که A_{*j} معرّف ستون j ام در A است. این عبارت می‌گوید که فضای همه بردارهایی به شکل Ax توسط ستونهای A پدید می‌آیند، ولو اینکه خطی وابسته باشند.

فضاهای برداری با حاصلضرب داخلی یکی از دلایل مهم که مسائل را به صورت مسائل هم‌ارز در جبر خطی بیان می‌کنند عرضه یک دیدگاه هندسی مسأله است. آنچه در این فرایند مهم است، حاصلضرب داخلی و تعامد است.

تعریف ۱. حاصلضرب داخلی دو بردار $x, y \in \mathbf{R}^n$ چنین تعریف می‌شود

$$(x, y) = \sum_{i=1}^n x_i y_i = x^T y = y^T x$$

و برای بردارهای $x, y \in \mathbf{C}^n$ ، حاصلضرب داخلی با رابطه زیر تعریف می‌شود

$$(x, y) = \sum_{i=1}^n x_i \bar{y}_i = y^* x$$

۲. نرم اقلیدسی x در \mathbf{C}^n یا \mathbf{R}^n چنین تعریف می‌شود

$$\|x\|_r = \sqrt{(x, x)} = \sqrt{|x_1|^2 + \dots + |x_n|^2} \quad (۷.۱.۷)$$

اثبات قضایای زیر نسبتاً ساده است و به عهده خواننده واگذار شده است. گیریم V معرف \mathbf{C}^n یا \mathbf{R}^n باشد

۱. برای $x, y, z \in V$

$$(x, y+z) = (x, y) + (x, z) \quad (x+y, z) = (x, z) + (y, z)$$

۲. برای هر $x, y \in V$

$$(\alpha x, y) = \alpha(x, y)$$

و برای $V = C^n$ و $\alpha \in C$

$$(x, \alpha y) = \bar{\alpha}(x, y)$$

۳. در C^n ، $(x, y) = \overline{(y, x)}$ ؛ و در R^n ، $(x, y) = (y, x)$.

۴. برای هر $x \in V$

$$(x, x) \geq 0$$

و $(x, x) = 0$ ، اگر و تنها اگر $x = 0$.

۵. برای هر $x, y \in V$

$$|(x, y)|^2 \leq (x, x)(y, y) \quad (۸.۱.۷)$$

که نابرابری کوشی - شوارتس خواننده می‌شود، و عیناً از راه اثبات (۳.۴.۴) فصل ۴ ثابت می‌شود. با استفاده از نرم اقلیدسی می‌توانیم بنویسیم

$$|(x, y)| \leq \|x\|_2 \|y\|_2 \quad (۹.۱.۷)$$

۶. برای هر $x, y \in V$

$$\|x + y\|_2 \leq \|x\|_2 + \|y\|_2 \quad (۱۰.۱.۷)$$

این رابطه نابرابری مثلثی است. برای یک تعبیر هندسی، توضیحات قبلی در بخش ۱.۴ فصل ۴ را برای نرم $\|f\|_\infty$ روی $C[a, b]$ ببینید. برای اثبات (۱۰.۱.۷)، طرز به دست آوردن (۴.۴.۴) را در فصل چهار ملاحظه کنید.

۷. برای هر ماتریس مربعی از مرتبه n و برای هر $x, y \in C^n$

$$(Ax, y) = (x, A^*y) \quad (۱۱.۱.۷)$$

قبلاً حاصلضرب داخلی برای معرفی طول اقلیدسی به کار رفته بود، ولی حداقل در فضاهایی که مجموعه اسکالر R باشد، برای تعبیری از زاویه نیز به کار می‌رود.

تعریف ۱. برای x و y در R^n ، زاویه بین x و y با رابطه زیر تعریف می‌شود

$$\mathcal{A}(x, y) = \cos^{-1} \left[\frac{(x, y)}{\|x\|_2 \|y\|_2} \right]$$

توجه نمایید که به موجب نابرابری (۹.۱.۷) کوشی-شوارتس، شناسه بین ۱- و ۱ است. این تعریف را می‌توان به‌طور ضمنی به شکل زیر نوشت

$$(x, y) = \|x\|_2 \|y\|_2 \cos(\theta) \quad (۱۲.۱.۷)$$

که یک فرمول آشنا از کاربرد حاصلضرب نقطه‌یی در R^2 و R^3 است.

۲. دو بردار x و y متعامدند اگر و تنها اگر $(x, y) = 0$. انگیزه این اصطلاح از (۱۲.۱.۷) بوده است. اگر $\{x^{(1)}, \dots, x^{(n)}\}$ یک پایه برای C^n یا R^n باشد و اگر $(x^{(i)}, x^{(j)}) = 0$ برای تمام مقادیر $j \neq i$, $n \leq i, j \leq n$ آنگاه $\{x^{(1)}, \dots, x^{(n)}\}$ را پایه متعامد خوانیم. اگر تمام بردارهای پایه دارای طول اقلیدسی ۱ باشند، پایه را یکا متعامد خوانند.

۳. ماتریس مربعی U را یکانی گویند اگر

$$U^*U = UU^* = I$$

اگر ماتریس U حقیقی باشد به جای یکانی آن را متعامد خوانند. سطرها [یا ستونها] یک ماتریس یکانی مرتبه n ، یک پایه یکا متعامد برای C^n تشکیل می‌دهند، و همچنین است برای ماتریسهای متعامد و R^n .

مثال ۱. زاویه بین بردارهای $y = (3, 2, 1)$ ، $x = (1, 2, 3)$ با رابطه زیر داده می‌شود:

$$A = \cos^{-1} \left[\frac{1}{14} \right] \doteq 77.5^\circ \text{ رادیان}$$

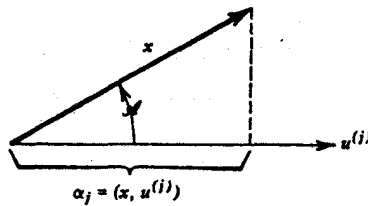
۲. ماتریسهای

$$U_1 = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \quad U_2 = \begin{bmatrix} 1 & i \\ \sqrt{2} & \sqrt{2} \\ i & 1 \\ \sqrt{2} & \sqrt{2} \end{bmatrix}$$

یکانی‌اند، و اولی متعامد است.

برای فضای برداری R^n یا C^n ، داشتن یک پایه یکا متعامد مفید است، زیرا در این صورت به آسانی می‌توان هر بردار دلخواه را به مؤلفه‌های آن در امتداد بردارهای پایه تجزیه کرد. به بیان دقیقتر، گیریم $\{u^{(1)}, \dots, u^{(n)}\}$ پایه یکا متعامد V باشد و $x \in V$. با استفاده از این پایه

$$x = \alpha_1 u^{(1)} + \dots + \alpha_n u^{(n)}$$



شکل ۱.۷ نمایش (۱۵.۱.۷)

ضرایب $\alpha_1, \dots, \alpha_n$ یکتا هستند. برای پیدا کردن α_j ، حاصلضرب داخلی x و $u^{(j)}$ را تشکیل می‌دهیم و سپس با استفاده از ویژگیهای یکا متعامدی پایه خواهیم داشت

$$\begin{aligned} (x, u^{(j)}) &= \alpha_1(u^{(1)}, u^{(j)}) + \dots + \alpha_n(u^{(n)}, u^{(j)}) \\ &= \alpha_j \end{aligned} \quad (۱۳.۱.۷)$$

بنابراین

$$x = \sum_{j=1}^n (x, u^{(j)}) u^{(j)} \quad (۱۴.۱.۷)$$

شکل ۱.۷ یک تعبیر هندسی از این رابطه را نشان می‌دهد. با استفاده از (۱۳.۱.۷)،

$$\begin{aligned} \alpha_j &= (x, u^{(j)}) = \|x\| \|u^{(j)}\| \cos(\angle(x, u^{(j)})) \\ &= \|x\| \cos(\angle(x, u^{(j)})) \end{aligned} \quad (۱۵.۱.۷)$$

بنابراین ضریب α_j طول تصویر قائم x روی محوری است که با $u^{(j)}$ مشخص می‌شود. فرمول (۱۴.۱.۷) تعمیم تجزیه بردار x با استفاده از پایه استاندارد $\{e^{(1)}, \dots, e^{(n)}\}$ است، که قبلاً تعریف شده بود.

مثال گیریم $V = \mathbb{R}^2$ ، و پایه یکا متعامد

$$u^{(1)} = \left(\frac{1}{2}, \frac{\sqrt{3}}{2} \right) \quad u^{(2)} = \left(-\frac{\sqrt{3}}{2}, \frac{1}{2} \right)$$

را در نظر می‌گیریم. در این صورت برای یک بردار داده شده $x = (x_1, x_2)$ ، می‌توان نشان داد که

$$x = \alpha_1 u^{(1)} + \alpha_2 u^{(2)}$$

$$\alpha_1 = (x, u^{(1)}) = \frac{x_1 + x_2 \sqrt{3}}{2} \quad \alpha_2 = (x, u^{(2)}) = \frac{x_2 - x_1 \sqrt{3}}{2}$$

برای مثال

$$(1, 0) = \frac{1}{2}u^{(1)} - \frac{\sqrt{3}}{2}u^{(2)}$$

۲.۷ ویژه بردارها و شکلهای متعارف (کانونی) ماتریسها

عدد λ ، مختلط یا حقیقی، یک ویژه مقدار ماتریس مربعی A است اگر یک بردار $x \in C^n$ ، $x \neq 0$ وجود داشته باشد به طوری که

$$Ax = \lambda x \quad (1.2.7)$$

بردار x را ویژه بردار متناظر با ویژه مقدار λ نامند. با توجه به قضیه ۲.۷ و احکام (۳) و (۵)، λ یک ویژه مقدار A است اگر و تنها اگر

$$\det(A - \lambda I) = 0 \quad (2.2.7)$$

این معادله را معادله مشخصه A خوانند و برای تحلیل آن تابع زیر را معرفی می‌کنیم

$$f_A(\lambda) \equiv \det(A - \lambda I)$$

اگر A از مرتبه n باشد، $f_A(\lambda)$ یک چندجمله‌ی دقیقاً از مرتبه n است که چندجمله‌ی مشخصه A نامیده می‌شود. برای اثبات اینکه این یک چندجمله‌ی است این دترمینان را برحسب دترمینانهای جزئی متوالیاً بسط می‌دهیم، تا عبارت زیر به دست آید

$$f_A(\lambda) = \det(A - \lambda I)$$

$$= \det \begin{bmatrix} a_{11} - \lambda & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} - \lambda & \dots & a_{2n} \\ \vdots & & \ddots & \\ a_{n1} & \dots & \dots & a_{nn} - \lambda \end{bmatrix}$$

$$= (a_{11} - \lambda)(a_{22} - \lambda) \dots (a_{nn} - \lambda)$$

(جملاتی از درجه نایزگتر از $n - 2$)

$$f_A(\lambda) = (-1)^n \lambda^n + (-1)^{n-1} (a_{11} + \dots + a_{nn}) \lambda^{n-1}$$

(جملاتی از درجه نایزگتر از $n - 2$)

$$(3.2.7)$$

همچنین توجه نمایید که جمله ثابت عبارت است از

$$f_A(0) = \det(A) \quad (۴.۲.۷)$$

با توجه به ضریب λ^{n-1} ، تعریف می‌کنیم

$$\text{trace}(A) = a_{11} + a_{22} + \dots + a_{nn} \quad (۵.۲.۷)$$

که اغلب یک کمیت مورد علاقه در مطالعه A است.

چون $f_A(\lambda)$ از درجه n است، دقیقاً n ویژه مقدار برای A موجود است، به شرطی که ریشه‌های چندگانه را برحسب درجه چندگانی آنها به حساب آوریم. هر ماتریس حداقل یک زوج ویژه مقدار-ویژه بردار دارد و ماتریس A از مرتبه $n \times n$ ، حداکثر n ویژه مقدار متمایز خواهد داشت.

مثال ۱. چندجمله‌یی مشخصه

$$A = \begin{bmatrix} 2 & 1 & 0 \\ 1 & 3 & 1 \\ 0 & 1 & 2 \end{bmatrix}$$

عبارت است از

$$f_A(\lambda) = -\lambda^3 + 7\lambda^2 - 14\lambda + 8$$

ویژه‌مقدارها عبارت‌اند از، $\lambda_1 = 1$ ، $\lambda_2 = 2$ ، $\lambda_3 = 4$ و ویژه‌بردارهای متناظر آنها چنین‌اند

$$u^{(1)} = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \quad u^{(2)} = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \quad u^{(3)} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$$

توجه نمایید که این ویژه‌بردارها برهم عمود و بنابراین خطی مستقل‌اند. چون بُعد R^3 (و C^3) برابر ۳ است، این ویژه‌بردارها یک پایه متعامد برای R^3 (و C^3) تشکیل می‌دهند. این نکته توضیحی برای قضیه ۴.۷ است که بعداً بیان خواهد شد.

۲. برای ماتریس

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad f_A(\lambda) = (1 - \lambda)^3$$

و سه ویژه بردار خطی مستقل برای ویژه مقدار $\lambda = 1$ وجود دارد، برای مثال

$$[1, 0, 0]^T \quad [0, 1, 0]^T \quad [0, 0, 1]^T$$

تمام ویژه بردارهای دیگر ترکیبهای خطی از این سه بردارند.

۳. برای ماتریس

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} \quad f_A(\lambda) = (1 - \lambda)^3$$

ماتریس A فقط یک ویژه بردار خطی مستقل برای ویژه مقدار $\lambda = 1$ یعنی

$$x = [1, 0, 0]^T$$

و مضارب آن دارد.

چندگانگی جبری یک ویژه مقدار ماتریس A ، چندگانگی ریشه $f_A(\lambda)$ است و چندگانگی هندسی آن بیشترین عدد ویژه بردارهای خطی مستقل آن ماتریس متناظر با ویژه مقدار آن است. مجموع چندگانگیهای جبری ویژه بردارهای ماتریس A از مرتبه $n \times n$ ، نسبت به اختلالات جزئی در A ، ثابت، یعنی n ، است. ولی مجموع چندگانگیهای هندسی ممکن است با اختلالات جزئی، تغییرات بزرگی داشته باشد و این امر موجب می شود که محاسبه عددی ویژه بردارها اغلب یک مسأله مشکل باشد. همچنین چندگانگیهای جبری و هندسی لازم نیست مساوی باشند، همان گونه که مثالهای قبل نشان می دهند.

تعریف گیریم A و B ماتریسهای مربعی هم مرتبه باشند. در این صورت A متشابه با B است اگر یک ماتریس ناسنگین P وجود داشته باشد که

$$B = P^{-1}AP \quad (۶.۲.۷)$$

توجه نمایید که این رابطه یک رابطه متقارن است زیرا

$$A = Q^{-1}BQ \quad Q = P^{-1}$$

رابطه (۶.۲.۷) را می‌توان چنین تعبیر کرد که بگوییم A و B نمایشهای یک تبدیل خطی T از V به V ، $[V = R^n$ یا $C^n]$ هستند ولی نسبت به پایه‌های مختلف V . ماتریس P ماتریس تبدیل پایه خوانده می‌شود، و این ماتریس دو نمایش مختلف یک بردار $x \in V$ را نسبت به دو پایه‌ای که مورد استفاده قرار گرفته به هم مربوط می‌کند. [آنتون (۱۹۸۴)، بخش ۵.۵] یا نوبل (۱۹۶۹)، بخش ۵.۱۴] را برای جزئیات بیشتر ببینید.

ما اکنون چند ویژگی ساده ماتریسهای مشابه و ویژه مقدرهای آنها را بیان می‌کنیم.

۱. اگر A و B مشابه باشند، آنگاه $f_A(\lambda) = f_B(\lambda)$. برای اثبات این حکم از (۶.۲.۷) استفاده می‌کنیم تا نشان دهیم

$$\begin{aligned} f_B(\lambda) &= \det(B - \lambda I) = \det[P^{-1}(A - \lambda I)P] \\ &= \det(P^{-1}) \det(A - \lambda I) \det(P) = f_A(\lambda) \end{aligned}$$

زیرا

$$\det(P) \det(P^{-1}) = \det(PP^{-1}) = \det(I) = 1$$

۲. ویژه‌مقدرهای ماتریسهای مشابه A و B دقیقاً یکی هستند و یک تناظر یک به یک بین این ویژه‌مقدارها وجود دارد. اگر $Ax = \lambda x$ ، آنگاه با استفاده از

$$\begin{aligned} P^{-1}AP(P^{-1}x) &= \lambda P^{-1}x \\ Bz &= \lambda z \quad z = P^{-1}x \end{aligned} \quad (۷.۲.۷)$$

بدیهی است که $z \neq 0$ ، زیرا در غیر این صورت x صفر می‌شود. همچنین برای هر ویژه‌بردار داده‌شده z از B ، این استدلال را می‌توان معکوس کرد تا ویژه‌بردار متناظر $x = Pz$ برای A تولید شود.

۳. چون $f_A(\lambda)$ تحت تبدیلهای تشابهی ناورداست، ضرایب $f_A(\lambda)$ نیز تحت چنین تبدیلهای تشابهی، ناوردا خواهند بود. به‌ویژه برای ماتریسهای مشابه A و B ،

$$\text{trace}(A) = \text{trace}(B) \quad \det(A) = \det(B) \quad (۸.۲.۷)$$

شکلهای متعارف اکنون به ذکر چندین شکل متعارف مهم ماتریسها می‌پردازیم. این شکلها بین ساختار یک ماتریس و ویژه‌مقدارها و ویژه‌بردارهای آن ارتباط برقرار می‌کنند، و در کاربردهای گوناگون عرضه‌های دیگر ریاضیات و علوم مورد استفاده قرار می‌گیرند.

قضیه ۳.۷ (شکل نرمال شور) گیریم A دارای مرتبه n با عناصری از C باشد. در این صورت یک ماتریس یکانی وجود دارد به طوری که

$$T \equiv U^* A U \quad (۹.۲.۷)$$

یک ماتریس بالا مثلثی است.

چون T مثلثی است و چون $U^* = U^{-1}$

$$f_A(\lambda) = f_T(\lambda) = (\lambda - t_{11}) \dots (\lambda - t_{nn}) \quad (۱۰.۲.۷)$$

و بنابراین ویژه‌مقدارهای A عناصر قطری T هستند.

برهان برهان به استقراء روی n مرتبه A ، انجام می‌گیرد. قضیه برای $n = 1$ با استفاده از $U = [1]$ درست است. گیریم که قضیه برای تمام ماتریسهای مرتبه $n \leq k - 1$ درست باشد، و ثابت می‌کنیم که قضیه برای تمام ماتریسهای از مرتبه $n = k$ درست است.

گیریم λ_1 یک ویژه‌مقدار A باشد، و $u^{(1)}$ ویژه‌بردار متناظر آن باشد با $\|u^{(1)}\|_2 = 1$. یک پایهٔ یکانی متعامد برای C^k می‌گیریم که $u^{(1)}$ اولین بردار پایه باشد و آن را $\{u^{(1)}, \dots, u^{(k)}\}$ می‌نامیم. ماتریس P_1 را با

$$P_1 = [u^{(1)}, u^{(2)}, \dots, u^{(k)}]$$

تعریف می‌کنیم که به شکل افزایشی با ستونهای متعامد $u^{(1)}, \dots, u^{(k)}$ نوشته شده است. پس $P_1^* P_1 = I$ و بنابراین $P_1^{-1} = P_1^*$. تعریف می‌کنیم

$$B_1 = P_1^* A P_1$$

ادعا می‌کنیم

$$B_1 = \begin{bmatrix} \lambda_1 & \alpha_2 & \dots & \alpha_k \\ \circ & & & \\ \vdots & & A_2 & \\ \circ & & & \end{bmatrix}$$

که در آن A_1 از مرتبه $k-1$ است، $\alpha_1, \dots, \alpha_k$ اعدادی هستند. برای اثبات این ادعا از ضرب ماتریسهای افرازشده استفاده می‌نماییم:

$$AP_1 = A[u^{(1)}, \dots, u^{(k)}] = [Au^{(1)}, \dots, Au^{(k)}]$$

$$= [\lambda_1 u^{(1)}, v^{(2)}, \dots, v^{(k)}] \quad v^{(j)} = Au^{(j)}$$

$$B_1 = P_1^* AP_1 = [\lambda_1 P_1^* u^{(1)}, P_1^* v^{(2)}, \dots, P_1^* v^{(k)}]$$

چون $P_1^* P_1 = I$ داریم $P_1^* u^{(1)} = e^{(1)} = [1, 0, \dots, 0]^T$ بنابراین

$$B_1 = [\lambda_1 e^{(1)}, w^{(2)}, \dots, w^{(k)}] \quad w^{(j)} = P_1^* v^{(j)}$$

که دارای شکل مطلوب است.

با فرض استقرای، یک ماتریس یکانی \hat{P}_1 از مرتبه $k-1$ موجود است به طوری که

$$\hat{T} = \hat{P}_1^* A_1 \hat{P}_1$$

یک ماتریس بالا مثلثی از مرتبه $k-1$ است. تعریف می‌کنیم،

$$P_1 = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & & & \\ \vdots & & \hat{P}_1 & \\ 0 & & & \end{bmatrix}$$

پس P_1 یکانی است و

$$P_1^* B_1 P_1 = \begin{bmatrix} \lambda_1 & \gamma_2 & & \gamma_k \\ 0 & & & \\ \vdots & & \hat{P}_1^* A_1 \hat{P}_1 & \\ 0 & & & \end{bmatrix} = \begin{bmatrix} \lambda_1 & \gamma_2 & \dots & \gamma_k \\ 0 & & & \\ \vdots & & \hat{T} & \\ 0 & & & \end{bmatrix} \equiv T$$

یک ماتریس بالامثلثی است. بنابراین

$$T = P_1^* B_1 P_1 = P_1^* P_1^* A P_1 P_1 = (P_1 P_1)^* A (P_1 P_1)$$

$$T = U^* A U \quad U = P_1 P_1$$

و به‌سادگی معلوم می‌شود که U یکانی است. پس، استقرا و برهان کامل است. ■

مثال برای ماتریس

$$A = \begin{bmatrix} ۰۲ & ۰۶ & ۰ \\ ۱۶ & -۰۲ & ۰ \\ -۱۶ & ۱۶ & ۳۰ \end{bmatrix}$$

ماتریسهای قضیه و (۹.۲.۷) چنین‌اند

$$T = \begin{bmatrix} ۱ & ۰ & -۱ \\ ۰ & ۳ & ۲ \\ ۰ & ۰ & -۱ \end{bmatrix} \quad U = \begin{bmatrix} ۰۶ & ۰ & -۰۸ \\ ۰۸ & ۰ & ۰۶ \\ ۰ & ۱۰ & ۰ \end{bmatrix}$$

این راه معمولی محاسبه ویژه‌مقدارها نیست، و فقط برای توضیح قضیه باید در نظر گرفته شود. این قضیه فقط به‌عنوان یک ابزار نظری به‌کار می‌رود، نه یک ابزار محاسباتی.

با استفاده از (۸.۲.۷) و (۹.۲.۷)

$$\text{trace}(A) = \lambda_1 + \lambda_2 + \dots + \lambda_n \quad \det(A) = \lambda_1 \lambda_2 \dots \lambda_n \quad (۱۱.۲.۷)$$

که $\lambda_1, \dots, \lambda_n$ ویژه‌مقدارهای A هستند که باید عناصر قطری T را تشکیل دهند. به‌عنوان یک کاربرد خیلی مهمتر، قضیه معروف زیر را داریم.

قضیه ۴.۷ (قضیه محوره‌های اصلی) گیریم A یک ماتریس اریتمی از مرتبه n باشد، یعنی $A^* = A$. پس A دارای n ویژه‌مقدار $\lambda_1, \dots, \lambda_n$ است که الزاماً متمایز نیستند، و n ویژه‌بردار متناظر u_1, \dots, u_n دارند که یک پایهٔ یکا متعامد C^n را تشکیل می‌دهند. اگر A حقیقی باشد، ویژه‌بردارهای u_1, \dots, u_n را می‌توان حقیقی گرفت که یک پایهٔ یکا متعامد برای R^n تشکیل خواهند داد. بالاخره یک ماتریس یکانی U وجود دارد به‌طوری که

$$U^*AU = D \equiv \text{diag}[\lambda_1, \dots, \lambda_n] \quad (۱۲.۲.۷)$$

یک ماتریس قطری با عناصر قطری $\lambda_1, \dots, \lambda_n$ است. اگر A حقیقی هم باشد، U را می‌توان متعامد گرفت.

برهان با توجه به قضیه ۳.۷، یک ماتریس یکانی U موجود است که

$$U^*AU = T$$

که T بالامتثلی است. ترانزاده مزدوج دو طرف را تشکیل می‌دهیم تا به دست آوریم

$$T^* = (U^*AU)^* = U^*A^*(U^*)^* = U^*AU = T$$

چون T^* پایین مثلثی است، باید داشته باشیم

$$T = \text{diag}[\lambda_1, \dots, \lambda_n]$$

همچنین، $T^* = T$ شامل مزدوجهای مختلط تمام عناصر T است، و بنابراین همه عناصر قطری T باید حقیقی باشند.

U را به صورت زیر می‌نویسیم

$$U = [u^{(1)}, \dots, u^{(n)}]$$

پس $T = U^*AU$ ایجاب می‌کند که $AU = UT$

$$A[u^{(1)}, \dots, u^{(n)}] = [u^{(1)}, \dots, u^{(n)}] \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix}$$

$$[Au^{(1)}, \dots, Au^{(n)}] = [\lambda_1 u^{(1)}, \dots, \lambda_n u^{(n)}]$$

و

$$Au^{(j)} = \lambda_j u^{(j)} \quad j = 1, \dots, n \quad (۱۳.۲.۷)$$

چون ستونهای U یکا متعامدند و چون بعد C^n برابر n است، پس این بردارها باید یک پایه یکا متعامد برای C^n تشکیل دهند. ما از ذکر برهان نتایجی که از حقیقی بودن A ناشی می‌شوند صرف نظر می‌کنیم. پس برهان کامل است. ■

مثال از مثال ۱ در این بخش، ماتریس

$$A = \begin{bmatrix} 2 & 1 & 0 \\ 1 & 3 & 1 \\ 0 & 1 & 2 \end{bmatrix}$$

ویژه بردارها و شکلهای متعارف (کانونی) ماتریسها ۵۴۱

دارای ویژه مقدرهای $\lambda_1 = 1$, $\lambda_2 = 2$ و $\lambda_3 = 4$ و ویژه بردارهای یکا متعامد متناظر زیر است

$$u^{(1)} = \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \quad u^{(2)} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \quad u^{(3)} = \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$$

اینها یک پایهٔ یکا متعامد برای R^3 یا C^3 تشکیل می دهند.

یک شکل متعارف دو می وجود دارد که اخیراً در حل مسائل جبر خطی عددی، به ویژه در حل دستگاه معادلات خطی فرامعین، بسیار مهم شده است. این دستگاهها از برازش داده های تجربی با استفاده از شیوه های خطی کمترین مربعات به دست می آیند. [گلوب و وان لون^۱ (۱۹۸۳)، فصل ۶) و لوسن^۲ و هانسن^۳ (۱۹۷۴) را ببینید].

قضیهٔ ۵.۷ (تجزیهٔ تکین - مقدار) گیریم A از مرتبهٔ $n \times m$ باشد. در این صورت ماتریسهای یکانی U و V از مرتبه های به ترتیب m و n وجود دارند به طوری که ماتریس

$$V^*AU = F \quad (۱۴.۲.۷)$$

یک ماتریس مستطیلی «قطری» از مرتبهٔ $n \times m$ به شکل زیر است

$$F = \begin{bmatrix} \mu_1 & & & & & & & & \\ & \mu_2 & & & & & & & \\ & & \dots & & & & & & \\ & & & \mu_r & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & \dots & \end{bmatrix} \quad (۱۵.۲.۷)$$

اعداد μ_1, \dots, μ_r مقادیر تکین A نامیده می شوند. این مقادیر همگی حقیقی و مثبت اند و آنها را می توان به شکلی مرتب کرد که

$$\mu_1 \geq \mu_2 > \dots \geq \mu_r > 0 \quad (۱۶.۲.۷)$$

که r رتبهٔ ماتریس A است.

برهان ماتریس مربعی A^*A از مرتبه m را در نظر می‌گیریم. این ماتریس یک ماتریس ارمیتی است، و در نتیجه قضیه ۴.۷ را می‌توان برای آن به‌کار برد. ویژه‌مقدارهای A^*A همه حقیقی، به‌علاوه همگی نامنفی‌اند. برای پی‌بردن به آن، فرض می‌کنیم

$$A^*Ax = \lambda x \quad x \neq 0$$

پس

$$(x, A^*Ax) = (x, \lambda x) = \lambda \|x\|_2^2$$

$$(x, A^*Ax) = (Ax, Ax) = \|Ax\|_2^2$$

$$\lambda = \left(\frac{\|Ax\|_2^2}{\|x\|_2^2} \right) \geq 0$$

از این قضیه همچنین ثابت می‌شود که

$$A^*Ax = 0, x \in C^m \quad \text{اگر و تنها اگر} \quad Ax = 0 \quad (17.2.7)$$

با توجه به قضیه ۳.۷، یک ماتریس یکانی U از مرتبه $m \times m$ موجود است به طوری که

$$U^*A^*AU = \text{diag}[\lambda_1, \dots, \lambda_r, 0, \dots, 0] \quad (18.2.7)$$

که همه λ_i ها، $1 \leq i \leq r$ ، مخالف صفر و همگی مثبت‌اند. چون A^*A دارای مرتبه m است، اندیس نابزرگتر از m است. مقادیر تکین زیر را معرفی می‌کنیم

$$\mu_i = \sqrt{\lambda_i} \quad i = 1, \dots, r \quad (19.2.7)$$

U می‌تواند طوری انتخاب شود که ترتیب (۱۶.۲.۷) به‌دست آید. با استفاده از ماتریس قطری مرتبه m

$$D = \text{diag}[\mu_1, \dots, \mu_r, 0, \dots, 0]$$

(۱۸.۲.۷) را می‌توانیم به شکل زیر بنویسیم

$$(AU)^*(AU) = D^2 \quad (20.2.7)$$

گیریم $W = AU$. پس از (۳۰.۲.۷)، $W^*W = D^2$. اگر W را به شکل

$$W = [W^{(1)}, \dots, W^{(m)}] \quad W^{(j)} \in C^n$$

بنویسیم داریم

$$(W^{(j)}, W^{(j)}) = \begin{cases} \mu_j^2 & 1 \leq j \leq r \\ 0 & j > r \end{cases} \quad (21.2.7)$$

و

$$(W^{(i)}, W^{(j)}) = 0 \quad \text{اگر } i \neq j \quad (22.2.7)$$

از (۲۱.۲.۷) نتیجه می‌شود، $W^{(j)} = 0$ اگر $j > r$. و از (۲۲.۲.۷) نتیجه می‌شود که r ستون اول W عناصر متعامد در C^n ، و بنابراین خطی مستقل‌اند. و این امر ایجاب می‌کند که $r \leq n$.
تعریف می‌کنیم

$$V^{(j)} = \frac{1}{\mu_j} W^{(j)} \quad j = 1, \dots, r \quad (23.2.7)$$

این یک مجموعهٔ یکا متعامد در C^n است. اگر $r < n$ ، آنگاه $V^{(1)}, \dots, V^{(r)}$ را طوری انتخاب می‌کنیم که $\{V^{(1)}, \dots, V^{(n)}\}$ یک پایهٔ یکا متعامد در C^n باشد. تعریف می‌کنیم

$$V = [V^{(1)}, \dots, V^{(n)}] \quad (24.2.7)$$

به‌سادگی معلوم می‌شود که V یک ماتریس یکانی از مرتبهٔ $n \times n$ است و مستقیماً می‌توان تحقیق کرد که $VF = W$ ، که به‌صورت (۱۵.۲.۷) است. بنابراین

$$VF = AU$$

که (۱۴.۲.۷) را ثابت می‌کند. برهان r بودن رتبهٔ (A) و به‌دست آوردن سایر ویژگیهای تجزیهٔ مقدار تکین به مسئلهٔ ۱۹ واگذار شده است. تجزیهٔ تکین - مقدار در فصل ۹، در حل کمترین مربعات دستگاههای خطی فرا معین به‌کار برده شده است. ■

برای عرضهٔ اساسیترین شکل متعارف، نماد زیر را وارد می‌کنیم. ماتریس $n \times n$ زیر را تعریف

می‌کنیم

$$J_n(\lambda) = \begin{bmatrix} \lambda & 1 & 0 & \dots & 0 \\ 0 & \lambda & 1 & & \\ \cdot & & \cdot & & \\ \cdot & & & \cdot & 1 \\ \cdot & & & & \cdot \\ 0 & \cdot & \cdot & \cdot & \lambda \end{bmatrix} \quad n \geq 1 \quad (25.2.7)$$

که $J_n(\lambda)$ دارای ویژه مقدار λ با چندگانگی جبری n و چندگانگی هندسی ۱ است. این $J_n(\lambda)$ بلوک ژوردان نامیده می‌شود.

قضیه ۴.۷ (شکل متعارف ژوردان) گیریم A از مرتبه n باشد. آنگاه یک ماتریس ناتکین P وجود دارد که برای آن

$$P^{-1}AP = \begin{bmatrix} J_{n_1}(\lambda_1) & & & 0 \\ & J_{n_2}(\lambda_2) & & \\ & & \ddots & \\ 0 & & & J_{n_r}(\lambda_r) \end{bmatrix} \quad (26.2.7)$$

ویژه‌مقدارهای $\lambda_1, \dots, \lambda_r$ لازم نیست متمایز باشند. اگر A ارمیتی باشد، قضیه ۴.۷ ایجاب می‌کند $n_1 = n_2 = \dots = n_r = 1$ ، زیرا در این حالت مجموع چندگانگیهای هندسی باید برابر n ، مرتبه ماتریس A ، باشد.

اغلب مناسبتر است که (۲۶.۲.۷) را به شکل

$$P^{-1}AP = D + N$$

$$D = \text{diag}[\lambda_1, \dots, \lambda_r] \quad (27.2.7)$$

بنویسیم، که هر λ_i, n_i بار، بر قطر D ظاهر می‌شود. درایه‌های ماتریس N همه صفرند جز احتمالاً ۱‌های واقع بر بالای قطر. این ماتریس یک ماتریس پوچتوان است و به عبارت دقیقتر، در رابطه زیر صدق می‌کند

$$N^n = 0 \quad (28.2.7)$$

قضیه شکل متعارف ژوردان به سادگی قابل اثبات نیست، و برای مطالعه بیشتر این موضوع پر بار، خوانندگان را به کتابهای جبر خطی زیادی که نوشته شده ارجاع می‌دهیم [برای مثال فرانکلین (۱۹۶۸، فصل ۵)؛ هالموس (۱۹۵۸، بخش ۵۸)؛ یا، نوبل (۱۹۶۹، فصل ۱۱) را ببینید].

۳.۷ نرمهای برداری و ماتریسی

نرم اقلیدسی $\|x\|_2$ قبلاً معرفی شده است و اغلب آن را برای اندازه‌گیری بزرگی بردار به کار می‌برند. ولی موارد زیادی هستند که اندازه‌گیری بزرگی بردار از راههای دیگری مناسبتر است. از این رو مفهوم کلی نرم بردار را تعریف می‌کنیم.

تعریف گیریم V یک فضای برداری باشد و $N(x)$ تابعی حقیقی - مقدار که در V تعریف شده است. در این صورت $N(x)$ یک نرم است اگر:

$$(N1) \text{ برای همه مقادیر } x \in V, N(x) \geq 0, \text{ و } N(x) = 0 \text{ اگر و تنها اگر } x = 0.$$

$$(N2) \text{ برای همه مقادیر } x \in V \text{ و تمام اسکالرهاى } \alpha, N(\alpha x) = |\alpha| N(x)$$

$$(N3) \text{ برای همه مقادیر } x, y \in V, N(x + y) \leq N(x) + N(y)$$

$\|x\| = N(x)$ نماد معمولی است. نماد $N(x)$ به کار برده شده تا تأکید شود که نرم، تابعی است با حوزه V و برد اعداد حقیقی نامنفی. فاصله x از y را با $\|x - y\|$ تعریف می‌کنیم. نتایج ساده این تعریف، نابرابری مثلثی در شکل دیگر آن

$$\|x - z\| \leq \|x - y\| + \|y - z\|$$

و معکوس نابرابری مثلثی زیر است.

$$\| \|x\| - \|y\| \| \leq \|x - y\| \quad x, y \in V \quad (1.3.7)$$

مثال ۱. p -نرم را برای $1 \leq p < \infty$ چنین تعریف می‌کنیم

$$\|x\|_p = \left[\sum_{j=1}^n |x_j|^p \right]^{1/p} \quad x \in C^n \quad (2.3.7)$$

۲. نرم ماکسیم چنین است

$$\|x\|_\infty = \max_{1 \leq j \leq n} |x_j| \quad x \in C^n \quad (3.3.7)$$

انگیزه استفاده از زیرنمایه ∞ برای نرم، نتیجه مسأله ۲۳ است.

۳. برای فضای برداری $V = C[a, b]$ ، نرمهای تابعی $\|f\|_2$ و $\|f\|_\infty$ به ترتیب در فصلهای ۴ و ۱ معرفی شده‌اند.

مثال بردار $x = (1, 0, -1, 2)$ را در نظر می‌گیریم. آنگاه

$$\|x\|_1 = 4 \quad \|x\|_2 = \sqrt{6} \quad \|x\|_\infty = 2$$

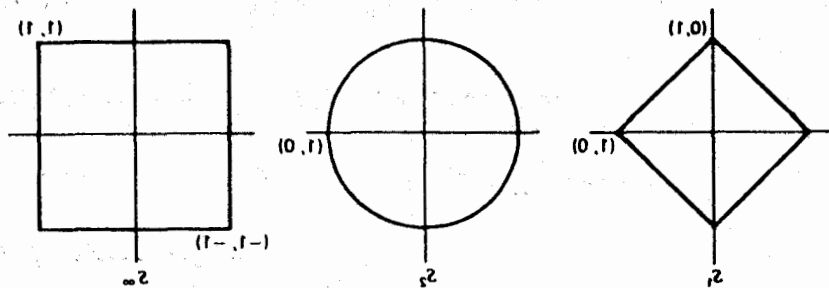
اثبات نرم بودن $\|\cdot\|_p$ برای یک p کلی، بدیهی نیست. اثبات حالت‌های $p = 1$ و $p = \infty$ ساده است و $\|\cdot\|_2$ در بخش ۱.۴ بررسی شده است. ولی برای $1 < p < \infty$ و $p \neq 2$ مشکل است نشان دهیم که $\|\cdot\|_p$ در نابرابری مثلثی صدق می‌کند. این مطلب برای ما مسأله مهمی نیست زیرا حالت‌های عمده مورد نظر ما $p = 1, 2, \infty$ هستند. برای آنکه یک شهود هندسی از این نرمها به دست دهیم، دایره‌های واحد

$$S_p = \{x \in \mathbf{R}^2 \mid \|x\|_p = 1\} \quad p = 1, 2, \infty \quad (۴.۳.۷)$$

را در شکل ۲.۷ رسم کرده‌ایم.

اکنون چند قضیه مربوط به نرمهای مختلف را ثابت می‌کنیم. با قضیه زیر در مورد پیوستگی $N(x) = \|x\|$ ، به صورت تابعی از آغاز می‌کنیم.

لم گیریم $N(x)$ یک نرم در C^n (یا \mathbf{R}^n) باشد. آنگاه $N(x)$ یک تابع پیوسته از x_1, x_2, \dots, x_n ، مؤلفه‌های بردار x است.



شکل ۲.۷ دایره واحد S_p با استفاده از نرم برداری $\|\cdot\|_p$

برهان می‌خواهیم نشان دهیم که

$$x_i = y_i \quad i = 1, 2, \dots, n$$

ایجاب می‌کند که

$$N(x) = N(y)$$

با استفاده از معکوس نابرابری مثلثی (۱.۳.۷)

$$|N(x) - N(y)| \leq N(x - y) \quad x, y \in C^n$$

با توجه به (۱.۱.۷)، تعریف پایه‌استانده $\{e^{(1)}, \dots, e^{(n)}\}$ را برای C^n به خاطر می‌آوریم. پس

$$x - y = \sum_{j=1}^n (x_j - y_j) e^{(j)}$$

$$N(x - y) \leq \sum_{j=1}^n |x_j - y_j| N(e^{(j)}) \leq \|x - y\|_\infty \sum_{j=1}^n N(e^{(j)})$$

$$|N(x) - N(y)| \leq c \|x - y\|_\infty \quad c = \sum_{j=1}^n N(e^{(j)}) \quad (۵.۳.۷)$$

و این رابطه اثبات را کامل می‌کند. ■

باید توجه داشت که این رابطه همچنین ثابت می‌کند که برای هر نرم برداری N در C^n یک $c > 0$ وجود دارد به طوری که

$$N(x) \leq c \|x\|_\infty \quad x \in C^n \text{ هر برای} \quad (۶.۳.۷)$$

تنها کافی است که در (۵.۳.۷)، y را مساوی صفر قرار دهیم. قضیه زیر عکس این قضیه است.

قضیه ۷.۷ (هم‌ارزی نرمها) گیریم N و M نرمهایی در C^n یا R^n باشند. در این صورت ثابتی مانند $c_1 > 0$ و c_2 وجود دارند به طوری که

$$c_1 M(x) \leq N(x) \leq c_2 M(x) \quad x \in V \text{ هر برای} \quad (۷.۳.۷)$$

برهان کافی است حالتی را در نظر بگیریم که N دلخواه باشد و $M(x) = \|x\|_\infty$. ترکیب دو چنین حکمی قضیه کلی را به دست خواهد داد. بنابراین می‌خواهیم نشان دهیم که ثابت‌هایی چون c_1 و c_2 وجود دارند به طوری که

$$c_1 \|x\|_\infty \leq N(x) \leq c_2 \|x\|_\infty \quad (۸.۳.۷)$$

یا هم‌ارز با آن

$$c_1 \leq N(z) \leq c_2 \quad z \in S \text{ هر } \quad (۹.۳.۷)$$

که در آن S مجموعه تمام نقاط z در C^n است به طوری که $\|z\|_\infty = 1$. نامساوی سمت راست (۹.۳.۷) بلافاصله از (۶.۳.۷) نتیجه می‌شود.

توجه کنید که S یک مجموعه بسته و کراندار در C^n است و N یک تابع پیوسته در S . یک قضیه استاندارد در حساب دیفرانسیل و انتگرال پیشرفته داریم که N به ماکسیمم و مینیمم خود در S در نقاطی از S می‌رسد، یعنی ثابت‌های c_1, c_2 و نقاط z_1, z_2 در S وجود دارند به طوری که

$$c_1 = N(z_1) \leq N(z) \leq N(z_2) = c_2 \quad z \in S \text{ هر}$$

واضح است که $c_1, c_2 \geq 0$. و اگر $c_1 = 0$ ، آنگاه $N(z_1) = 0$. ولی در این صورت $z_1 = 0$ که با ساختمان S که در آن $\|z\|_\infty = 1$ ، متناقض است. این مطلب (۹.۳.۷) را ثابت می‌کند و برهان قضیه کامل می‌شود. توجه: این قضیه برای فضاهای نامتناهی - بعد تعمیم نمی‌یابد. ■

بسیاری از روشهای عددی برای مسائلی که متضمن دستگاههای خطی هستند، دنباله‌ای از بردارهای $\{x^{(m)} \mid m \geq 0\}$ تولید می‌کنند و ما می‌خواهیم راجع به همگرایی این دنباله به یک بردار x صحبت کنیم.

تعریف دنباله بردارهای $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}, \dots\}$ در C^n یا R^n ، به بردار x همگرا خوانده می‌شود اگر و تنها اگر

$$\|x - x^{(m)}\| \rightarrow 0 \quad \text{وقتی} \quad m \rightarrow \infty$$

توجه نمایید که انتخاب نرم تصریح نشده است. برای فضاهای متناهی بعد، مهم نیست که چه نرمی به کار برده شده است. گیریم M و N دو نرم بر C^n باشند. در این صورت از (۷.۳.۷)،

$$c_1 M(x - x^{(m)}) \leq N(x - x^{(m)}) \leq c_2 M(x - x^{(m)}) \quad m \geq 0$$

و $M(x - x^{(m)})$ به صفر همگرا می شود اگر و تنها اگر $N(x - x^{(m)})$ به صفر همگرا شود. بنابراین $x^{(m)}$ با نرم M به x همگرا می شود اگر و تنها اگر با نرم N همگرا باشد. این یک نتیجه مهمی است و برای فضاهای نامتناهی بعد درست نیست.

نرمهای ماتریسی مجموعه تمام ماتریسهای $n \times n$ با درایه های مختلط را می توان هم ارز با فضای برداری C^{n^2} در نظر گرفت که یک عمل ضربی خاص به فضای برداری اضافه شده است. بنابراین نرم ماتریسی باید در سه شرط معمول N_1 تا N_3 صدق کند. به علاوه، می خواهیم دو شرط دیگر هم داشته باشد.

تعریف نرم ماتریسی در $N_3 - N_1$ و شرایط زیر صدق می کند

$$\| AB \| \leq \| A \| \| B \| \quad (N_4)$$

(N5) معمولاً فضای برداری که با آن کار می کنیم، \mathbb{R}^n یا $C^n = V$ ، یک نرم برداری دارد، که آن را $\|x\|_v$ می نامیم. می خواهیم که نرمهای ماتریسی و برداری سازگار باشند:

$$\| Ax \|_v \leq \| A \| \| x \|_v \quad x \in V \text{ و هر } A$$

مثال گیریم A ماتریس $n \times n$ باشد و $\| \cdot \|_2$ و $\| \cdot \|_v = \| \cdot \|$. در این صورت برای $x \in C^n$ ، با استفاده از نامساوی (۸.۱.۷) - کوشی - شوارتس:

$$\| Ax \|_2 = \left[\sum_{i=1}^n \left| \sum_{j=1}^n a_{ij} x_j \right|^2 \right]^{1/2} \leq \left[\sum_{i=1}^n \left\{ \sum_{j=1}^n |a_{ij}|^2 \right\} \left\{ \sum_{j=1}^n |x_j|^2 \right\} \right]^{1/2}$$

پس

$$\| Ax \|_2 \leq F(A) \| x \|_2 \quad F(A) = \left[\sum_{i,j=1}^n |a_{ij}|^2 \right]^{1/2} \quad (۱۰.۳.۷)$$

$F(A)$ را نرم فروبنیوس A خوانند. ویژگی N_5 با استفاده از (۱۰.۳.۷) مستقیماً ثابت می شود. ویژگیهای N_1 تا N_3 برقرارند زیرا $F(A)$ نرم اقلیدسی در C^{n^2} است. اثبات N_4 باقی می ماند. با استفاده از نابرابری کوشی - شوارتس

$$\begin{aligned}
 F(AB) &= \left[\sum_{i,j=1}^n \left| \sum_{k=1}^n a_{ik} b_{kj} \right|^2 \right]^{1/2} \\
 &\leq \left[\sum_{i,j=1}^n \left\{ \sum_{k=1}^n |a_{ik}|^2 \right\} \left\{ \sum_{k=1}^n |b_{kj}|^2 \right\} \right]^{1/2} \\
 &= F(A)F(B)
 \end{aligned}$$

پس $F(A)$ یک نرم ماتریسی، سازگار با نرم اقلیدسی است.

معمولاً وقتی یک فضای برداری با یک نرم $\| \cdot \|_v$ داده شده است نرم ماتریسی وابسته به آن با رابطه زیر تعریف می‌شود

$$\| A \| = \sup_{x \neq 0} \frac{\| Ax \|_v}{\| x \|_v} \quad (۱۱.۳.۷)$$

این نرم را اغلب نرم عملگر می‌خوانند و با تعریفی که دارد در N_0 صدق می‌کند:

$$\| Ax \|_v \leq \| A \| \| x \|_v \quad x \in C^n \quad (۱۲.۳.۷)$$

برای یک ماتریس A ، نرم عملگر که بر اثر نرم برداری $\| x \|_p$ القا شده با نماد $\| A \|_p$ نمایش داده می‌شود. مهمترین موارد در جدول ۱.۷ داده شده‌اند، و طرز پیدا کردن آنها بعداً داده شده است. تعریف زیر را برای تعریف $\| A \|_2$ نیاز داریم

تعریف گیریم A یک ماتریس دلخواه باشد. طیف A مجموعه همه ویژه‌مقدارهای A است که با $\sigma(A)$ نمایش داده می‌شود. شعاع طیفی ماکسیمم اندازه این ویژه‌مقدارهاست، و با نماد زیر بیان می‌شود

$$r_\sigma(A) = \max_{\lambda \in \sigma(A)} | \lambda | \quad (۱۳.۳.۷)$$

جدول ۱.۷ نرمهای برداری و نرمهای عملگر ماتریسی متناظر

نرم برداری	نرم ماتریسی
$\ x \ _1 = \sum_{i=1}^n x_i $	$\ A \ _1 = \max_{1 \leq j \leq n} \sum_{i=1}^n a_{ij} $
$\ x \ _2 = \left[\sum_{j=1}^n x_j ^2 \right]^{1/2}$	$\ A \ _2 = \sqrt{r_\sigma(A^*A)}$
$\ x \ _\infty = \max_{1 \leq i \leq n} x_i $	$\ A \ _\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n a_{ij} $

برای آنکه نشان دهیم (۱۱.۳.۷) در حالت کلی یک نرم است، ابتدا نشان می‌دهیم که متناهی است. از قضیه ۷.۷، ثابتهای $c_1, c_2 > 0$ با

$$c_1 \|x\|_2 \leq \|x\|_v \leq c_2 \|x\|_2 \quad x \in \mathbb{C}^n$$

وجود دارند. بنابراین

$$\frac{\|Ax\|_v}{\|x\|_v} \leq \frac{c_2 \|Ax\|_2}{c_1 \|x\|_2} \leq \frac{c_2}{c_1} F(A)$$

که ثابت می‌کند $\|A\|$ متناهی است.

در اینجا جالب است که به معنای هندسی $\|A\|$ توجه نماییم

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|_v}{\|x\|_v} = \sup_{x \neq 0} \left\| A \left(\frac{x}{\|x\|_v} \right) \right\|_v = \sup_{\|z\|_v=1} \|Az\|_v$$

با توجه به اینکه اگر $\|z\|_v \leq 1$ انتخاب شود این سوپرمم تغییر نخواهد کرد،

$$\|A\| = \sup_{\|z\|_v \leq 1} \|Az\|_v \quad (۱۴.۳.۷)$$

گیریم

$$B = \{z \in \mathbb{C}^n \mid \|z\|_v \leq 1\}$$

گوی واحد نسبت به $\|\cdot\|_v$ باشد. پس

$$\|A\| = \sup_{z \in B} \|Az\|_v = \sup_{w \in A(B)} \|w\|_v$$

که $A(B)$ نگاشت B است، وقتی A بر آن اثر کرده باشد. بنابراین $\|A\|$ تأثیر A را بر گوی واحد اندازه‌گیری می‌نماید، و اگر $\|A\| > 1$ ، کشیدگی ماکسیمم گوی B را وقتی A بر آن عمل کرده است نشان می‌دهد.

برهان در زیر ثابت می‌کنیم که نرم عملگر $\|A\|$ یک نرم ماتریسی است.

۱. واضح است که $\|A\| \geq 0$ ، و اگر $A = 0$ ، آنگاه $\|A\| = 0$. برعکس، اگر $\|A\| = 0$ ،

آنگاه $\|Ax\|_v = 0$ برای جميع مقادیر x . بنابراین برای تمام مقادیر x ، $Ax = 0$ ، که ایجاب

می‌کند $A = 0$.

۲. گیریم α یک اسکالر باشد. آنگاه

$$\begin{aligned} \|\alpha A\| &= \sup_{\|x\|_v \leq 1} \|\alpha Ax\|_v = \sup_{\|x\|_v \leq 1} |\alpha| \|Ax\|_v \\ &= |\alpha| \sup_{\|x\|_v \leq 1} \|Ax\|_v = |\alpha| \|A\| \end{aligned}$$

۳. برای هر $x \in \mathbb{C}^n$

$$\|(A+B)x\|_v = \|Ax+Bx\|_v \leq \|Ax\|_v + \|Bx\|_v$$

زیرا $\|\cdot\|_v$ یک نرم است. با استفاده از ویژگی (۱۲.۳.۷)

$$\|(A+B)x\|_v \leq \|A\| \|x\|_v + \|B\| \|x\|_v$$

$$\frac{\|(A+B)x\|_v}{\|x\|_v} \leq \|A\| + \|B\|$$

که ایجاب می‌کند

$$\|A+B\| \leq \|A\| + \|B\|$$

۴. برای هر $x \in \mathbb{C}^n$ ، (۱۲.۳.۷) را به‌کار می‌بریم تا به‌دست آوریم

$$\|(AB)x\|_v = \|A(Bx)\|_v \leq \|A\| \|Bx\|_v \leq \|A\| \|B\| \|x\|_v$$

$$\frac{\|ABx\|_v}{\|x\|_v} \leq \|A\| \|B\|$$

که ایجاب می‌کند

$$\|AB\| \leq \|A\| \|B\|$$

اکنون نتایج جدول ۱.۷ را به‌طور مفصل‌تری شرح می‌دهیم

مثال ۱. از نرم برداری زیر استفاده می‌کنیم

$$\|x\|_1 = \sum_{j=1}^n |x_j| \quad x \in \mathbb{C}^n$$

پس

$$\|Ax\|_1 = \sum_{i=1}^n \left| \sum_{j=1}^n a_{ij} x_j \right| \leq \sum_{i=1}^n \sum_{j=1}^n |a_{ij}| \|x_j\|$$

با تغییر ترتیب مجموعیابی، می‌توانیم جمعووندها را مجزاً کنیم

$$\|Ax\|_1 \leq \sum_{j=1}^n |x_j| \sum_{i=1}^n |a_{ij}|$$

گیریم

$$c = \text{Max}_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}| \quad (15.3.7)$$

پس

$$\|Ax\|_1 \leq c \|x\|_1$$

و بنابراین

$$\|A\|_1 \leq c$$

برای اینکه ثابت کنیم این یک تساوی است x ی پیدا می‌کنیم که:

$$\frac{\|Ax\|_1}{\|x\|_1} = c$$

گیریم k اندیس ستونی باشد که برای آن ماکسیم (15.3.7) حاصل شده است. فرض می‌کنیم $x = e^{(k)}$ k امین بردار یکه باشد. پس $\|x\|_1 = 1$ و

$$\|Ax\|_1 = \sum_{i=1}^n \left| \sum_{j=1}^n a_{ij} x_j \right| = \sum_{i=1}^n |a_{ik}| = c$$

این ثابت می‌کند که برای نرم برداری $\|\cdot\|_1$ ، نرم عملگر چنین است

$$\|A\|_1 = \text{Max}_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}| \quad (16.3.7)$$

این رابطه اغلب نرم ستونی خوانده می‌شود.

۲. برای C^n با نرم $\|x\|_\infty$ ، نرم عملگر چنین است

$$\|A\|_\infty = \text{Max}_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \quad (17.3.7)$$

این نرم سطری A خوانده می‌شود. اثبات این فرمول به مسأله ۲۵ واگذار شده است، گرچه مشابه با اثبات $\|A\|_1$ است.

۳. نرم $\|x\|_2$ در C^n را به‌کار می‌بریم. از (۱۰.۳.۷) نتیجه می‌گیریم که

$$\|A\|_2 \leq F(A) \quad (18.3.7)$$

در حالت کلی تساوی برقرار نیست. برای مثال، با $A = I$ ، ماتریس همانی، (۱۰.۳.۷) و (۱۱.۳.۷) را به‌کار می‌بریم تا به‌دست آوریم

$$F(I) = \sqrt{n} \quad \|I\|_2 = 1$$

ثابت می‌کنیم

$$\|A\|_2 = \sqrt{r_\sigma(A^*A)} \quad (19.3.7)$$

همان‌گونه که قبلاً در جدول ۱.۷ بیان کردیم. ماتریس A^*A ارمیتی است و همه ویژه‌مقدارهای آن همان‌گونه که در اثبات قضیه ۵.۷ نشان داده شد، نامنفی‌اند؛ گیریم که ویژه‌مقدارهای آن

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$$

باشند که برحسب چندگانگی آنها شمرده شده‌اند، و گیریم $u^{(1)}, \dots, u^{(n)}$ ویژه‌بردارهای متناظر با آنها باشند که همچون یک پایهٔ یکا متعامد C^n مرتب شده‌اند.

برای یک $x \in C^n$

$$\|Ax\|_2^2 = (Ax, Ax) = (x, A^*Ax)$$

x را به صورت زیر می‌نویسیم

$$x = \sum_{j=1}^n \alpha_j u^{(j)} \quad \alpha_j \equiv (x, u^{(j)}) \quad (20.3.7)$$

پس

$$A^*Ax = \sum_{j=1}^n \alpha_j A^*A u^{(j)} = \sum_{j=1}^n \alpha_j \lambda_j u^{(j)}$$

و

$$\begin{aligned} \|Ax\|_2^2 &= \left(\sum_{i=1}^n \alpha_i u^{(i)}, \sum_{j=1}^n \alpha_j \lambda_j u^{(j)} \right) = \sum_{j=1}^n \lambda_j |\alpha_j|^2 \\ &\leq \lambda_1 \sum_{j=1}^n |\alpha_j|^2 = \lambda_1 \|x\|_2^2 \end{aligned}$$

که از (۲۰.۳.۷) برای محاسبه $\|x\|_2$ استفاده کرده‌ایم. بنابراین

$$\|A\|_2^2 \leq \lambda_1$$

تساوی با توجه به $x = u^{(1)}$ به دست می‌آید، لذا $\|x\|_2 = 1$ و

$$\|Ax\|_2^2 = (x, A^*Ax) = (u^{(1)}, \lambda_1 u^{(1)}) = \lambda_1$$

رابطه (۱۹.۳.۷) از اینجا ثابت می‌شود، زیرا $\lambda_1 = r_\sigma(A^*A)$ می‌توان نشان داد که AA^* و A^*A ویژه‌مقدارهای ناصفر متساوی دارند (مسئله ۱۹ را ببینید)، بنابراین؛ $r_\sigma(AA^*) = r_\sigma(A^*A)$ که شکل دیگری از فرمول (۱۹.۳.۷) است. این رابطه همچنین ثابت می‌کند که

$$\|A\|_2 = \|A^*\|_2 \quad (21.3.7)$$

این تساوی برای نرمهای ماتریسی قبلی صادق نیست.

نسبتاً به آسانی می‌توان نشان داد که اگر A ارمیتی باشد، آنگاه

$$\|A\|_2 = r_\sigma(A) \quad (22.3.7)$$

برهان درستی این رابطه به مسئله ۲۷ ارجاع شده است.

مثال ماتریس زیر را در نظر می‌گیریم

$$A = \begin{bmatrix} 1 & -2 \\ -3 & 4 \end{bmatrix}$$

پس

$$\|A\|_1 = 6 \quad \|A\|_2 = \sqrt{15 + \sqrt{221}} \doteq 5.46 \quad \|A\|_\infty = 7$$

به‌عنوان مثالی از نابرابری (۲۳.۳.۷) در قضیه زیر، داریم

$$r_\sigma(A) = \frac{5 + \sqrt{33}}{2} \doteq 5.37 < \|A\|_2$$

قضیه ۸.۷ را بگیریم A یک ماتریس دلخواه مربعی باشد. پس برای هر نرم عملگر ماتریسی،

$$r_\sigma(A) \leq \|A\| \quad (23.3.7)$$

بعلاوه، اگر $\varepsilon > 0$ داده شده باشد، یک نرم عملگر ماتریسی که با $\|\cdot\|_\varepsilon$ نمایش داده می‌شود وجود دارد به طوری که

$$\|A\|_\varepsilon \leq r_\sigma(A) + \varepsilon \quad (24.3.7)$$

برهان برای اثبات (23.3.7)، گیریم $\|\cdot\|$ یک نرم ماتریسی سازگار با نرم برداری $\|\cdot\|_v$ باشد. فرض می‌کنیم λ ویژه‌مقداری در $\sigma(A)$ باشد که برای آن

$$|\lambda| = r_\sigma(A)$$

و گیریم x یک ویژه بردار متناظر با آن باشد، $\|x\|_v = 1$. پس

$$r_\sigma(A) = |\lambda| = \|\lambda x\|_v \leq \|Ax\|_v \leq \|A\| \|x\|_v = \|A\|$$

که (23.3.7) را ثابت می‌کند.

برهان (24.3.7) یک ساختمان غیرنمایان است، و در آیزکسن و کلر (1966، ص 12) داده شده است. ■

فرع زیر یک نتیجه ساده ولی مهم قضیه 8.7 است.

فرع برای یک ماتریس مربعی A ، $r_\sigma(A) < 1$ اگر و تنها اگر، به‌ازای یک نرم عملگر ماتریسی داشته باشیم، $\|A\| < 1$. ■

این فرع را می‌توان برای اثبات قضیه 9.7 در بخش بعد به‌کار برد، ولی ما ترجیح می‌دهیم که از شکل متعارف ژوردان که در قضیه 6.7 داده شده استفاده کنیم. نتیجه (22.3.7) و قضیه 8.7 نشان می‌دهند که $r_\sigma(A)$ در واقع یک نرم ماتریسی است، و این فرع برای تحلیل نرخ همگرایی بعضی روشهای بارستی که در فصل 8 برای حل دستگاه معادلات خطی داده شده است به‌کار برده خواهد شد.

4.7 قضایای همگرایی و اختلال

قضایای زیر چارچوبی نظری هستند که ما بعداً با استفاده از آنها، به تحلیلهای خطاها در روشهای عددی حل دستگاه معادلات خطی می‌پردازیم.

قضیه ۹.۷. اگر A ماتریس مربعی از مرتبه n باشد. در این صورت A^m به ماتریس صفر همگرا خواهد بود اگر و تنها اگر $r_\sigma(A) < 1$.

برهان قضیه ۶.۷ را به عنوان ابزار اصلی به کار می‌بریم. گیریم J شکل متعارف زوردان برای A باشد،

$$P^{-1}AP = J$$

پس

$$A^m = (PJP^{-1})^m = PJ^mP^{-1} \quad (۱.۴.۷)$$

و $A^m \rightarrow 0$ اگر و تنها اگر $J^m \rightarrow 0$. با توجه به (۲۷.۲.۷) و (۲۸.۲.۷) یادآور می‌شویم که J را می‌توان به شکل

$$J = D + N$$

نوشت که در آن

$$D = \text{diag}[\lambda_1, \dots, \lambda_n]$$

شامل ویژه‌مقدارهای J (و A) است، و N ماتریسی است که برای آن

$$N^n = 0$$

بنابراین

$$J^m = (D + N)^m = \sum_{j=0}^m \binom{m}{j} D^{m-j} N^j$$

و با استفاده از $N^j = 0$ برای $j \geq n$

$$J^m = \sum_{j=0}^n \binom{m}{j} D^{m-j} N^j \quad (۲.۴.۷)$$

توجه داریم که توانهای D در رابطه زیر صدق می‌کنند

$$m - j \geq m - n \rightarrow \infty \quad \text{وقتی} \quad m \rightarrow \infty \quad (۳.۴.۷)$$

ما به حدود زیر نیاز داریم: برای هر مقدار مثبت $c < 1$ و هر $r \geq 0$

$$\lim_{m \rightarrow \infty} m^r c^m = 0 \quad (۴.۴.۷)$$

این رابطه را می‌توان با قاعده هوییتال از راه حساب دیفرانسیل انتگرال مقدماتی ثابت کرد. در (۲.۴.۷)، بدون توجه به بزرگی m ، تعداد $n + 1$ جمله ثابت، وجود دارند، و ما می‌توانیم همگرایی J^m را با بررسی همگرایی هر یک از جملات آن بررسی کنیم. فرض کنید که $r_\sigma(k) < 1$ می‌دانیم که هر $|\lambda_i| < 1$ ، $i = 1, 2, \dots, n$ و برای هر نرم ماتریسی

$$\left\| \binom{m}{j} D^{m-j} N^j \right\| \leq \frac{m^j}{j!} \|N\|^j \|D^{m-j}\|$$

با استفاده از نرم سطری، عبارت فوق با مقدار زیر کراندار می‌شود

$$\frac{1}{j!} \|N\|_\infty^j m^j [r_\sigma(A)]^{m-j}$$

که با استفاده از (۳.۴.۷) و (۴.۴.۷) به‌ازای $0 \leq j \leq n$ ، وقتی $m \rightarrow \infty$ ، به صفر می‌گراید. تا اینجا نیمی از قضیه ثابت شده است، یعنی اگر $r_\sigma(A) < 1$ ، آنگاه J^m و A^m با توجه به (۱.۴.۷)، وقتی $m \rightarrow \infty$ ، به صفر همگرا هستند.

فرض کنید $r_\sigma(A) \geq 1$ ، گیریم λ یک ویژه مقدار A با $|\lambda| \geq 1$ ، و ویژه بردار متناظر با آن، $x \neq 0$ ، باشند. پس

$$A^m x = \lambda^m x$$

و روشن است که وقتی $m \rightarrow \infty$ این عبارت به صفر همگرا نخواهد بود. بنابراین ممکن نیست که $A^m \rightarrow 0$ ، زیرا که لازم می‌آید $A^m x \rightarrow 0$ ؛ و اثبات کامل می‌شود. ■

قضیه ۱۰.۷ (سری هندسی) گیریم A یک ماتریس مربعی باشد. اگر $r_\sigma(A) < 1$ ، آنگاه $(I - A)^{-1}$ وجود دارد و می‌توان آن را با سری همگرای زیر بیان کرد

$$(I - A)^{-1} = I + A + A^2 + \dots + A^m + \dots \quad (5.4.7)$$

به عکس، اگر سری (۵.۴.۷) همگرا باشد، آنگاه $r_\sigma(A) < 1$.

برهان فرض می‌کنیم $r_\sigma(A) < 1$ وجود $(I - A)^{-1}$ را با اثبات عبارت هم‌ارز (۳) از قضیه ۲.۷ نشان خواهیم داد. فرض می‌کنیم

$$(I - A)x = 0$$

در این صورت $Ax = x$ ، و از اینجا لازم می‌آید که اگر $x \neq 0$ ، یک ویژه‌مقدار A باشد. ولی ما فرض کردیم $r_\sigma(A) < 1$ پس باید داشته باشیم $x = 0$ که به نتیجه اثبات وجود $(I - A)^{-1}$ می‌رسیم.

ما به اتحاد زیر نیاز داریم

$$(I - A)(I + A + A^2 + \dots + A^m) = I - A^{m+1} \quad (6.4.7)$$

که برای هر ماتریس A درست است. اگر طرفین را در $(I - A)^{-1}$ ضرب کنیم خواهیم داشت:

$$I + A + \dots + A^m = (I - A)^{-1}(I - A^{m+1})$$

طرف چپ دارای حد است اگر طرف راست حد داشته باشد. طبق قضیه ۹.۷، $r_\sigma(A) < 1$ ایجاب می‌کند $A^{m+1} \rightarrow 0$ وقتی $m \rightarrow \infty$. بنابراین نتیجه (۵.۴.۷) به دست می‌آید.

به عکس، فرض می‌کنیم که این سری هندسی همگراست و آن را با

$$B = I + A + A^2 + \dots + A^m + \dots$$

نشان می‌دهیم. پس $B - AB = B - BA = I$ ، و بنابراین $I - A$ معکوس دارد، و B معکوس آن است. از طرفین رابطه (۶.۴.۷) حد می‌گیریم؛ طرف چپ دارای حد $(I - A)B = I$ است، و بنابراین حد طرف راست هم باید I باشد. ولی لازمه آن این است که

$$A^{m+1} \rightarrow 0 \quad m \rightarrow \infty \text{ وقتی}$$

به موجب قضیه ۹.۷ باید داشته باشیم $r_\sigma(A) < 1$.

قضیه ۱۱.۷ گیریم A یک ماتریس مربعی باشد. اگر نرم ماتریسی عملگری، $\|A\| < 1$ ، آنگاه $(I - A)^{-1}$ وجود دارد و دارای بسط هندسی (۵.۴.۷) است. به علاوه

$$\|(I - A)^{-1}\| \leq \frac{1}{1 - \|A\|} \quad (7.4.7)$$

برهان چون $\|A\| < 1$ ، از (۲۳.۳.۷) و قضیه ۸.۷ نتیجه می‌شود که $r_\sigma(A) < 1$ غیر از (۷.۴.۷) بقیه حکمها از قضیه ۱۰.۷ نتیجه می‌شوند. برای (۷.۴.۷) فرض می‌کنیم

$$B_m = I + A + \dots + A^m$$

از (۶.۴.۷) نتیجه می‌شود،

$$\begin{aligned} B_m &= (I - A)^{-1}(I - A^{m+1}) \\ (I - A)^{-1} - B_m &= (I - A)^{-1}[I - (I - A^{m+1})] \\ &= (I - A)^{-1}A^{m+1} \end{aligned} \quad (۸.۴.۷)$$

با استفاده از عکس ناپرابری مثلثی،

$$\begin{aligned} \left| \| (I - A)^{-1} \| - \| B_m \| \right| &\leq \| (I - A)^{-1} - B_m \| \\ &< \| (I - A)^{-1} \| \| A \|^{m+1} \end{aligned}$$

چون طرف راست وقتی $m \rightarrow \infty$ به صفر میل می‌کند، داریم

$$\| B_m \| \rightarrow \| (I - A)^{-1} \| \quad \text{وقتی} \quad m \rightarrow \infty \quad (۹.۴.۷)$$

با توجه به تعریف B_m و ویژگیهای نرم ماتریسی،

$$\begin{aligned} \| B_m \| &\leq \| I \| + \| A \| + \dots + \| A \|^m \\ &= \frac{1 - \| A \|^{m+1}}{1 - \| A \|} \leq \frac{1}{1 - \| A \|} \end{aligned}$$

که از ترکیب با (۹.۴.۷) برهان (۷.۴.۷) پایان می‌یابد.

قضیه ۱۲.۷ گیریم A و B ماتریسهای مربعی هم مرتبه باشند. فرض می‌کنیم A نانکین باشد و

$$\| A - B \| < \frac{1}{\| A^{-1} \|} \quad (۱۰.۴.۷)$$

پس B نیز نانکین است،

$$\| B^{-1} \| \leq \frac{\| A^{-1} \|}{1 - \| A^{-1} \| \| A - B \|} \quad (۱۱.۴.۷)$$

و

$$\| A^{-1} - B^{-1} \| \leq \frac{\| A^{-1} \|^2 \| A - B \|}{1 - \| A^{-1} \| \| A - B \|} \quad (۱۲.۴.۷)$$

برهان به اتحاد زیر توجه کنید

$$B = A - (A - B) = A[I - A^{-1}(A - B)] \quad (۱۳.۴.۷)$$

با استفاده از قضیه ۱۱.۷، ماتریس $[I - A^{-1}(A - B)]$ ، بر اساس نامساوی (۱۰.۴.۷)، ناتکین است که ایجاب می‌کند

$$\|A^{-1}(A - B)\| \leq \|A^{-1}\| \|A - B\| < ۱$$

چون B حاصلضرب ماتریسهای ناتکین است، پس خود نیز ناتکین است و

$$B^{-1} = [I - A^{-1}(A - B)]^{-1} A^{-1}$$

با گرفتن نرم از دو طرف و استفاده از قضیه ۱۱.۷، کران (۱۱.۴.۷) به دست می‌آید. برای اثبات از (۱۲.۴.۷)

$$A^{-1} - B^{-1} = A^{-1}(B - A)B^{-1}$$

استفاده می‌کنیم. باز از دو طرف نرم می‌گیریم و (۱۱.۴.۷) را به کار می‌بریم. ■

این قضیه از جهات مختلف مهم است. ولی در حال حاضر، مطلب مورد نظر ما بیان این نکته است که تمام اختلالهای به اندازه کافی کوچک یک ماتریس ناتکین، ماتریسهای ناتکین می‌دهند.

مثال قضیه ۱۱.۷ را با ملاحظهٔ وارونپذیری ماتریس زیر توضیح می‌دهیم

$$A = \begin{bmatrix} ۴ & ۱ & ۰ & ۰ & \dots & ۰ \\ ۱ & ۴ & ۱ & ۰ & & ۰ \\ ۰ & ۱ & ۴ & ۱ & & \vdots \\ \vdots & & & \ddots & & \\ & & & & ۱ & ۴ & ۱ \\ ۰ & \dots & ۰ & ۱ & ۴ \end{bmatrix}$$

A را مجدداً به صورت زیر می‌نویسیم

$$A = 4(I + B)$$

$$B = \begin{bmatrix} 0 & \frac{1}{4} & 0 & 0 & \dots & 0 \\ \frac{1}{4} & 0 & \frac{1}{4} & 0 & & 0 \\ 0 & & & \ddots & & \vdots \\ \vdots & & & & & \frac{1}{4} \\ 0 & \dots & & & & 0 \end{bmatrix}$$

با استفاده از نرم سطری (۱۷.۳.۷)، $\|B\|_{\infty} = \frac{1}{4}$. بنابراین طبق قضیه ۱۱.۷، $(I+B)^{-1}$ وجود دارد و با توجه به (۷.۴.۷)،

$$\|(I+B)^{-1}\|_{\infty} \leq \frac{1}{1 - \frac{1}{4}} = 2$$

لذا A^{-1} موجود است، $A^{-1} = \frac{1}{4}(I+B)^{-1}$ و

$$\|A^{-1}\|_{\infty} \leq \frac{1}{4}$$

با استفاده از نرم سطری و نابرابری (۲۳.۳.۷)،

$$r_{\sigma}(A) \leq 6 \quad r_{\sigma}(A^{-1}) \leq \frac{1}{4}$$

چون ویژه‌مقدارهای A^{-1} عکس ویژه‌مقدارهای A هستند (مسأله ۲۷ را ببینید)، و چون A ارمیتی است ویژه‌مقدارهای آن حقیقی‌اند، کران زیر را خواهیم داشت:

$$2 \leq |\lambda| \leq 6 \quad \lambda \in \sigma(A) \quad \text{جميع مقادير}$$

برای کرانه‌های بهتر در این حالت، قضیه دایره‌گرگورین^۱ فصل ۹ را ببینید.

بحث در آثار خواندنی

موضوع این فصل جبر خطی است، که به‌خصوص برای استفاده در به‌دست آوردن و تحلیل

روشهای جبر خطی عددی برگزیده شده است. کتابهای آنتون (۱۹۸۴) و استرنگ (۱۹۸۰) کتابهای درسی در سطح مقدماتی برای جبرخطی دوره کارشناسی هستند. کتاب فرانکلین (۱۹۶۸) یک کتاب مقدماتی در سطح بالاتری برای نظریه ماتریسهاست و کتاب هالموس (۱۹۵۸) یک متن درسی معروفی در جبر خطی مجرد است. کتاب نوبل (۱۹۶۹) یک کتاب درسی جبر خطی کاربردی در زمینه‌های گسترده است. معرفی مبانی جبرخطی در کتابهای فاده‌یوا (۱۹۵۹)، گلوب و ون‌لون^۱ (۱۹۸۲)، پارلت^۲ (۱۹۸۰)، استوارت (۱۹۷۳) و ویلکینسن (۱۹۶۵) آمده است، همه این کتابها کلاً برای جبر خطی عددی نوشته شده‌اند. برای اطلاعات نظری بیشتر با جزئیات بیشتر و در سطح بالاتر، کارهای کلاسیک گانتماخر^۳ (۱۹۶۰) و هاوس هولدر (۱۹۶۵) را ببینید. مراجع بیشتر در فهرست مراجع فصلهای ۸ و ۹ داده شده‌اند.

مراجع

- Anton, H. (1984). *Elementary Linear Algebra*, 4th ed. Wiley, New York.
- Fadeeva, V. (1959). *Computational Methods of Linear Algebra*. Dover, New York.
- Franklin, J. (1968). *Matrix Theory*. Prentice-Hall, Englewood Cliffs, N.J.
- Gantmacher, F. (1960). *The Theory of Matrices*, vols. I and II. Chelsea, New York.
- Golub, G., and C. Van Loan (1983). *Matrix Computations*. Johns Hopkins Press, Baltimore.
- Halmos, P. (1958). *Finite-Dimensional Vector Spaces*. Van Nostrand, Princeton, N.J.
- Householder, A. (1965). *The Theory of Matrices in Numerical Analysis*. Ginn (Blaisdell), Boston.
- Isaacson, E., and H. Keller (1966). *Analysis of Numerical Methods*. Wiley, New York.
- Lawson, C., and R. Hanson (1974). *Solving Least Squares Problems*. Prentice-Hall, Englewood Cliffs, N.J.
- Noble, B. (1969). *Applied Linear Algebra*. Prentice-Hall, Englewood Cliffs, N.J.
- Parlett, B. (1980). *The Symmetric Eigenvalue Problem*. Prentice-Hall, Englewood Cliffs, N.J.
- Stewart, G. (1973). *Introduction to Matrix Computations*. Academic Press, New York.
- Strang, G. (1980). *Linear Algebra and Its Applications*, 2nd ed. Academic Press, New York.

1. Van Loan

2. Parlett

3. Gantmacher

Wilkinson, J. (1965). *The Algebraic Eigenvalue Problem*. Oxford Univ. Press, Oxford, England.

مسائل

۱. تعیین کنید که آیا مجموعه‌های بردارهای زیر وابسته‌اند یا مستقل.

(الف) $(1, 2, -1, 3), (3, -1, 1, 1), (1, 9, -5, 11)$

(ب) $(1, 1, 0), (0, 1, 1), (1, 0, 1)$

۲. گیریم A و B و C ماتریس‌هایی به ترتیب از مرتبه $m \times n$, $n \times p$ و $p \times q$ باشند.

(الف) قانون شرکتپذیری $(AB)C = A(BC)$ را ثابت کنید

(ب) ثابت کنید $(AB)^T = B^T A^T$

۳. (الف) دو ماتریس مربعی A و B بنویسید که $AB \neq BA$

(ب) ماتریس‌های مربعی A و B با درایه‌های غیرصفر پیدا کنید که $AB = 0$, $BA \neq 0$.

۴. گیریم A ماتریسی از مرتبه $m \times n$ و r و c به ترتیب معرف رتبه‌های سطری و ستونی A باشند.

ثابت کنید که $r = c$. راهنمایی: برای راحتی، فرض کنید که اولین r سطر A مستقل‌اند و بقیه سطرها

وابسته به این r سطرند و به همین ترتیب برای اولین c ستون A . گیریم \hat{A} معرف ماتریس $r \times n$

حاصل از حذف $m - r$ سطر آخر A باشد و \hat{r} و \hat{c} به ترتیب رتبه‌های سطری و ستونی \hat{A} باشند.

روشن است که $\hat{r} = r$. همچنین ستونهای \hat{A} عناصری از C^r هستند که دارای بعد r است، و

بنابراین باید داشته باشیم $\hat{c} \leq r$. نشان دهید که $\hat{c} = c$ ، و بنابراین ثابت می‌شود که $c \leq r$. نابرابری

$r \leq c$ از همین استدلال برای A^T نتیجه می‌شود و این دو نابرابری با هم نتیجه می‌دهند که $r = c$.

۵. هم‌ارزی احکام (۱) - (۴) و (۶) قضیه ۲.۷ را ثابت کنید.

راهنمایی: قضیه ۱.۷، نتیجه مسأله ۴، و تجزیه (۶.۱.۷) را به کار برید.

۶. گیریم

$$f_n(x) = \det \begin{bmatrix} x & 1 & 0 & \dots & 0 \\ 1 & x & 1 & 0 & 0 \\ 0 & 1 & x & 1 & \vdots \\ \vdots & & & \ddots & 0 \\ & & & & 1 \\ 0 & \dots & 0 & 1 & x \end{bmatrix}$$

و مرتبه ماتریس n باشد. همچنین تعریف می‌کنیم $f_0(x) \equiv 1$ (الف) نشان دهید

$$f_{n+1}(x) = xf_n(x) - f_{n-1}(x) \quad n \geq 1$$

(ب) نشان دهید

$$f_n(x) = S_n \left(\frac{x}{\sqrt{2}} \right) \quad n \geq 0$$

که $S_n(x)$ چندجمله‌بی نوع دوم چبیشف از درجه n است (مسئله ۲۴ فصل ۲ را ببینید).
۷. گیریم A یک ماتریس مربعی از مرتبه n با درایه‌های حقیقی باشد. تابع

$$q(x_1, \dots, x_n) = (Ax, x) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j \quad x \in \mathbf{R}^n$$

صورت درجه دوم حاصل از A نامیده می‌شود. این عبارت یک چندجمله‌بی درجه دومی است با n متغیر x_1, \dots, x_n و در ماکسیم‌سازی یا مینیم‌سازی تابع n متغیره ظاهر می‌شود.

(الف) ثابت کنید که اگر A یک ماتریس پادمتقارن باشد، آنگاه $q(x) \equiv 0$.

(ب) برای یک ماتریس مربعی A در حالت کلی، $A_1 = \frac{1}{2}(A + A^T)$ و $A_2 = \frac{1}{2}(A - A^T)$

را تعریف می‌کنیم. پس $A = A_1 + A_2$. نشان دهید که A_1 متقارن است و به‌ازای جمیع مقادیر $x \in \mathbf{R}^n$ ، $(Ax, x) = (A_1x, x)$. این رابطه نشان می‌دهد که، بدون از دست‌دادن کلیت، ماتریس ضرایب A را برای یک صورت درجه دوم همیشه می‌توان متقارن فرض کرد.

۸. بردارهای متعامد زیر داده شده‌اند

$$u^{(1)} = (1, 2, -1) \quad u^{(2)} = (1, 1, 3)$$

یک بردار سوم $u^{(3)}$ پیدا کنید که $\{u^{(1)}, u^{(2)}, u^{(3)}\}$ یک پایه \mathbf{R}^3 باشد. این پایه را به صورت پایه‌ای یکا در آورید.

۹. برای بردارستونی $w \in \mathbf{C}^n$ با $\|w\|_2 = \sqrt{w^*w} = 1$ ، ماتریس $n \times n$ زیر را تعریف می‌کنیم

$$A = I - 2ww^*$$

(الف) برای حالت خاص $w = \left[\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right]^T$ ، ماتریس A را پیدا کنید. تحقیق کنید که این ماتریس متقارن و متعامد است.

(ب) نشان دهید که در حالت کلی، چنین ماتریسهایی ارمیتی و یکانی هستند.

۱۰. گیریم W یک زیرفضای \mathbf{R}^n باشد. برای $x \in \mathbf{R}^n$ تعریف می‌کنیم

$$\rho(x) = \inf_{y \in W} \|x - y\|_2$$

گیریم $\{u_1, \dots, u_m\}$ یک پایهٔ یکا متعامد W با بعد m باشد. این پایه را به یک پایهٔ یکا متعامد $\{u_1, \dots, u_m, \dots, u_n\}$ تمام \mathbf{R}^n گسترش دهید.
(الف) نشان دهید

$$\rho(x) = \left[\sum_{j=m+1}^n |(x, u_j)|^2 \right]^{1/2}$$

که به‌طور یکتا در

$$y = Px \quad P = \sum_{j=1}^m u_j u_j^T$$

به این‌فیمم می‌رسد. نقطهٔ y تصویر قائم x روی W خوانده می‌شود.

(ب) نشان دهید $P^2 = P$. چنین ماتریسی را ماتریس تصویر خوانند.

(ج) نشان دهید $P^T = P$

(د) نشان دهید برای هر $x, z \in \mathbf{R}^n$ ، $(Px, z - Pz) = 0$.

(ه) نشان دهید که به‌ازای جمیع مقادیر $x \in \mathbf{R}^n$ ، $\|x\|_2^2 = \|Px\|_2^2 + \|x - Px\|_2^2$. این شکلی از قضیه فیثاغورس است.

۱۱. ویژه‌مقدارها و ویژه‌بردارهای ماتریسهای زیر را به‌دست آورید.

$$\begin{bmatrix} 1 & 1 \\ -1 & 3 \end{bmatrix} \quad (\text{ب}) \quad \begin{bmatrix} 1 & 6 \\ 1 & 1 \end{bmatrix} \quad (\text{الف})$$

۱۲. گیریم $y \neq 0$ در \mathbf{R}^n باشد و $A = yy^T$ را یک ماتریس $n \times n$ تعریف می‌کنیم. نشان دهید که $\lambda = 0$ ویژه‌مقداری دقیقاً از چندگانگی $n - 1$ است. تنها ویژه‌مقدار ناصفر آن چیست؟

۱۳. گیریم U یک ماتریس یکانی $n \times n$ باشد.

(الف) نشان دهید برای هر $x \in \mathbf{C}^n$ ، $\|Ux\|_2 = \|x\|_2$. با استفاده از این ویژگی نشان دهید که فاصلهٔ x از y برابر فاصلهٔ Ux از Uy است، که نشان می‌دهد تبدیلات یکانی بر \mathbf{C}^n فاصلهٔ بین نقاط را حفظ می‌کنند.

(ب) گیریم U متعامد باشد، نشان دهید

$$(Ux, Uy) = (x, y) \quad x, y \in \mathbf{R}^n$$

این رابطه نشان می‌دهد که تبدیلات متعامد \mathbf{R}^n ، زاویهٔ بین خطوط را هم، آن‌گونه که در (۱۲.۱.۷) تعریف شد، حفظ می‌نمایند.

(ج) نشان دهید که همهٔ ویژه‌مقدارهای یک ماتریس یکانی طول واحد دارند.

۱۴. گیریم A یک ماتریس ارمیتی از مرتبهٔ n باشد. آن را معین مثبت خوانند اگر و تنها اگر، برای هر $x \neq 0$ در \mathbf{C}^n ، $(Ax, x) > 0$. نشان دهید که A معین مثبت است اگر و تنها اگر همهٔ ویژه‌مقدارهای آن حقیقی و مثبت باشند.

راهنمایی: قضیهٔ ۴.۷ را به‌کار برید و برای بیان یک بردار دلخواه $x \in \mathbf{R}^n$ ، (Ax, x) را برحسب یک پایه از ویژه‌بردارها بسط دهید.

۱۵. گیریم A حقیقی و متقارن باشد و ویژه‌مقدارهای $\lambda_1, \dots, \lambda_n$ برحسب چندگانگی آنها تکرار شده باشند. با استفاده از یک پایهٔ یکا متعامد از ویژه‌بردارها نشان دهید که صورت درجهٔ دوم مسألهٔ $q(x) = (Ax, x)$ ، $x \in \mathbf{R}^n$ ، را می‌توان به شکل ساده‌تر زیر بدل کرد.

$$q(x) = \sum_{j=1}^n \alpha_j^2 \lambda_j$$

با $\{\alpha_j\}$ که از روی x تعیین شده است. با استفاده از این رابطه، نمودارهای ممکن

$$(Ax, x) = \text{ثابت}$$

را وقتی A از مرتبهٔ ۳ باشد پیدا کنید.

۱۶. فرض می‌کنیم A حقیقی، متقارن، معین مثبت و از مرتبهٔ n باشد. تعریف می‌کنیم

$$f(x) = \frac{1}{2} x^T A x - b^T x \quad x, b \in \mathbf{R}^n$$

نشان دهید که مینیمم یکتای $f(x)$ از حل $Ax = b$ به‌ازای $\alpha = A^{-1}b$ به‌دست می‌آید.

۱۷. گیریم $f(x)$ تابعی حقیقی مقدار از \mathbf{R}^n باشد، و فرض می‌کنیم $f(x)$ سه بار پیوسته مشتق‌پذیر نسبت به مؤلفه‌های x باشد. تعمیم قضیهٔ تیلر، ۵.۱ فصل ۱ به n متغیر را برای به‌دست آوردن

$$f(x) = f(\alpha) + (x - \alpha)^T \nabla f(\alpha) + \frac{1}{2} (x - \alpha)^T H(\alpha) (x - \alpha) + O(\|x - \alpha\|^2)$$

به‌کار برید. در اینجا

$$\nabla f(x) = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right]^T$$

گرادیان f و

$$H(x) = \left[\frac{\partial^2 f(x)}{\partial x_i \partial x_j} \right] \quad 1 \leq i, j \leq n$$

ماتریس هسی $f(x)$ است. جمله آخر نشان می‌دهد که به‌ازای مقادیر x نزدیک به α جملات باقیمانده از مضربی از $\|x - \alpha\|^3$ کوچک‌ترند.

اگر بخواهیم α ماکسیمم یا مینیمم موضعی باشد، آنگاه یک شرط لازم برای آن، برقراری $\nabla f(x) = 0$ است. فرض می‌کنیم $\nabla f(\alpha) = 0$ ، نشان دهید α یک مینیمم موضعی اکید (یا یکتای) $f(x)$ است اگر و تنها اگر $H(\alpha)$ معین مثبت باشد، [توجه کنید که $H(x)$ همیشه متقارن است].

۱۸. رابطه (۲۸.۲.۷) را ثابت کنید.

۱۹. نمادی را که در قضیه ۵.۷ در تجزیه تکین - مقدار ماتریس A به‌کار برده شد به‌خاطر آورید. (الف) نشان دهید μ_1, \dots, μ_r ویژه‌مقدارهای ناصفر A^*A و AA^* هستند که به‌ترتیب با ویژه‌بردارهای $U^{(1)}, \dots, U^{(r)}$ و $V^{(1)}, \dots, V^{(r)}$ متناظرند. بردار $U^{(j)}$ ستون j ام U را مشخص می‌کند، همین‌طور است برای $V^{(j)}$ و V .

(ب) نشان دهید $AU^{(j)} = \mu_j V^{(j)}$ ، $A^*V^{(j)} = \mu_j U^{(j)}$ ، $1 \leq j \leq r$.

(ج) ثابت کنید $r = \text{rank}(A)$.

۲۰. برای هر چندجمله‌یی $p(x) = b_0 + b_1x + \dots + b_mx^m$ و برای هر ماتریس مربعی A تعریف می‌کنیم

$$p(A) = b_0I + b_1A + \dots + b_mA^m$$

گیریم A ماتریسی باشد که شکل متعارف ژوردان آن ماتریس قطری

$$P^{-1}AP = D = \text{diag}[\lambda_1, \dots, \lambda_n]$$

است. $f_A(\lambda)$ چندجمله‌یی مشخصه A است، ثابت کنید $f_A(A) = 0$. (این قضیه، قضیه کیلی-هاملتن است. و برای همه ماتریسهای مربعی درست است، نه فقط برای ماتریسهایی که دارای شکل قطری متعارف ژوردان هستند).

راهنمایی: برای ساده‌کردن $f_A(A)$ ، $A = PDP^{-1}$ را به‌کار ببرید.

۲۱. روابط زیر را ثابت کنید

$$\|x\|_\infty \leq \|x\|_1 \leq n \|x\|_\infty \quad (\text{الف})$$

$$\|x\|_{\infty} \leq \|x\|_2 \leq \sqrt{n} \|x\|_{\infty} \quad (\text{ب})$$

$$\|x\|_2 \leq \|x\|_1 \leq \sqrt{n} \|x\|_2 \quad (\text{ج})$$

۲۲. گیریم A ماتریس حقیقی نانتکینی از مرتبه m ، و $\|\cdot\|_v$ یک نرم برداری در \mathbf{R}^n باشد. تعریف می‌کنیم

$$\|x\|_* = \|Ax\|_v \quad x \in \mathbf{R}^n$$

نشان دهید $\|\cdot\|_*$ یک نرم برداری روی \mathbf{R}^n است.

۲۳. نشان دهید

$$\lim_{p \rightarrow \infty} \left[\sum_{j=1}^n |x_j|^p \right]^{1/p} = \max_{1 \leq i \leq n} |x_i| \quad x \in \mathbf{C}^n$$

این رابطه استفاده از نماد $\|x\|_{\infty}$ را برای طرف راست موجه می‌سازد.

۲۴. برای هر نرم ماتریسی نشان دهید که (الف) $\|I\| \geq 1$ ، و (ب) $\|A^{-1}\| \geq (1/\|A\|)$.

برای هر نرم عملگر، از (۱۱.۳.۷) فوراً نتیجه می‌شود که $\|I\| = 1$.

۲۵. فرمول (۱۷.۳.۷) را برای نرم عملگر ماتریسی $\|A\|_{\infty}$ به دست آورید.

۲۶. یک نرم برداری در \mathbf{R}^n را با

$$\|x\| = \frac{1}{n} \sum_{j=1}^n |x_j| \quad x \in \mathbf{R}^n$$

تعریف می‌کنیم. نرم عملگر ماتریسی مربوط به این نرم برداری چیست؟

۲۷. گیریم A یک ماتریس مربعی از مرتبه $n \times n$ باشد.

(الف) اگر ویژه‌مقدارها و ویژه بردارهای A داده شده باشند آنها را برای ماتریسهای زیر پیدا کنید.

$$(۱) A^m, m \geq 2, (۲) A^{-1}, \text{ با فرض نانتکینی } A, (۳) A + cI, \text{ ثابت } c.$$

(ب) ثابت کنید اگر A ارمیتی باشد، آنگاه $\|A\|_2 = r_{\sigma}(A)$.

(ج) برای A دلخواه و U یکانی از همان مرتبه، نشان دهید $\|AU\|_2 = \|UA\|_2 = \|A\|_2$.

۲۸. گیریم A ماتریس مربعی از مرتبه $n \times n$ باشد.

(الف) نشان دهید برای هر ماتریس یکانی U ، $F(AU) = F(UA) = F(A)$.

(ب) اگر A ارمیتی باشد، نشان دهید که

$$F(A) = \sqrt{\lambda_1^2 + \dots + \lambda_n^2}$$

که $\lambda_1, \dots, \lambda_n$ ویژه مقدارهای A هستند که برحسب چندگانگی تکرار شده اند. به علاوه

$$\frac{1}{\sqrt{n}} F(A) \leq \|A\|_2 \leq F(A)$$

۲۹. نماد قضیه ۵.۷ را به یاد آورده نشان دهید

$$\|A\|_2 = \mu_1 \quad F(A) = \sqrt{\mu_1^2 + \dots + \mu_r^2}$$

۳۰. گیریم A از مرتبه $n \times n$ باشد. نشان دهید

$$| \text{اثر}(A) | \leq nr_\sigma(A)$$

اگر A متقارن و معین مثبت باشد، نشان دهید

$$\text{اثر}(A) \geq r_\sigma(A)$$

۳۱. نشان دهید که سری نامتناهی

$$I + A + \frac{A^2}{2!} + \dots + \frac{A^n}{n!} + \dots$$

برای هر ماتریس مربعی A همگراست. مجموع سری را با e^A نمایش می دهیم

$$\text{الف) اگر } A = P^{-1}BP \text{، نشان دهید } e^A = P^{-1}e^B P$$

ب) گیریم $\lambda_1, \dots, \lambda_n$ ویژه مقدارهای A برحسب چندگانگی تکرار شده باشند، نشان دهید

ویژه مقدارهای e^A عبارت اند از $e^{\lambda_1}, \dots, e^{\lambda_n}$.

۳۲. ماتریس زیر را در نظر می گیریم

$$A = \begin{bmatrix} 6 & 1 & 1 & 0 & \dots & 0 \\ 1 & 6 & 1 & 1 & 0 & \vdots \\ 1 & 1 & 6 & 1 & 1 & 0 \\ 0 & 1 & 1 & 6 & 1 & 1 & 0 & 0 \\ \vdots & & & & & & & 1 \\ & & & & & 1 & 1 & 6 & 1 \\ 0 & \dots & & & 0 & 1 & 1 & 6 \end{bmatrix}$$

نشان دهید A ناتکین است. یک کران بالا برای $\|A^{-1}\|_\infty$ و $\|A^{-1}\|_2$ پیدا کنید.
 ۳۳. در به دست آوردن درونیاب برازای درجه ۳ در بخش ۷.۳ از فصل ۳، لازم بود که دستگاه خطی $AM = D$ از (۲۱.۷.۳) را حل کنیم که در آن

$$A = \begin{bmatrix} \frac{h_1}{3} & \frac{h_1}{6} & 0 & \dots & 0 \\ \frac{h_1}{6} & \frac{h_1 + h_2}{3} & \frac{h_2}{6} & & \vdots \\ \vdots & & \ddots & & \\ & & \frac{h_{m-1}}{6} & \frac{h_{m-1} + h_m}{3} & \frac{h_m}{6} \\ 0 & \dots & 0 & \frac{h_m}{6} & \frac{h_m}{3} \end{bmatrix}$$

و همه h_i ها، $i = 1, \dots, m$ مثبت اند. با استفاده از یک یا چند قضیه بخش ۴.۷ نشان دهید A ناتکین است. به علاوه، کرانهای

$$\frac{1}{6} \text{Min}(h_i) \leq |\lambda| \leq \text{Max } h_i \quad \lambda \in \sigma(A)$$

را برای ویژه مقدارهای A به دست آورید.

۳۴. گیریم A یک ماتریس مربعی با $A^m = 0$ برای $m \geq 2$ باشد. نشان دهید که $I - A$ نانکین است. چنین ماتریس A را پوچتوان خوانند.



حل عددی دستگاههای معادلات خطی

دستگاههای معادلات خطی در زمینه‌های وسیعی، هم مستقیماً در مدل‌سازی وضعیتهای فیزیکی و هم غیرمستقیم در حل عددی سایر مدل‌های ریاضی ظاهر می‌شوند. این کاربردها تقریباً در تمام زمینه‌های فیزیکی، زیست‌شناسی و علوم اجتماعی صورت می‌گیرند. به علاوه، دستگاههای خطی در مباحث زیر دخالت دارند: نظریهٔ بهینه‌سازی، حل دستگاه معادلات ناخطی، تقریب توابع، حل عددی مسائل مقدار مرزی برای معادلات دیفرانسیل معمولی، معادلات دیفرانسیل جزئی و معادلات انتگرالی، استنباط آماری؛ و مسائل متعدد دیگر. به دلیل اهمیت بسیار گستردهٔ دستگاههای خطی، پژوهشهای فراوانی در حل عددی آنها به عمل آمده است. الگوریتمهای عالی برای عادیترین نوع مسائل دستگاههای خطی ایجاد شده و بعضی از آنها در این فصل تعریف، تحلیل و توضیح داده شده‌اند. عادیترین نوع مسأله، حل دستگاه خطی مربعی

$$Ax = b$$

از مرتبهٔ متوسط، با ضرایبی است که اغلب ناصفرند. این‌گونه دستگاههای خطی، از هر مرتبه، چگال خوانده می‌شوند. برای چنین دستگاههایی، ماتریس ضرایب A معمولاً باید در حافظهٔ اصلی رایانه نگهداری شود تا بتوان دستگاه را به‌طور کارا حل کرد، و بنابراین محدودیتهای حافظه در اغلب

رایانه‌ها، مرتبه دستگاه را محدود می‌نماید. با کاهش سریع هزینه حافظه رایانه‌ها، دستگاه‌های خطی خیلی بزرگی را می‌توان در بعضی ماشینها جا داد، ولی انتظار می‌رود برای اغلب ماشینهای کوچکتر، کرانه‌های بالای عملی از مرتبه‌ای به بزرگی ۱۰۰ تا ۵۰۰ برسد. اغلب الگوریتمها برای حل چنین دستگاههای چگال، براساس حذف گاوسی ساخته شده‌اند که در بخش ۱.۸ تعریف شده‌اند. از جنبه نظری، این یک روش مستقیم است که اگر خطای گردکردن نادیده گرفته شود، پس از یک تعداد متناهی مرحله محاسباتی، جواب درست از حذف گاوسی به دست می‌آید. جرح و تعدیل برای بهبود رفتار خطا در حذف گاوسی، شکلهای متنوع بعضی دسته‌های خاص ماتریسها، و تحلیلهای خطا در بخش ۲.۸ تا ۵.۸ داده شده‌اند.

نوع مهم دیگر مسأله، حل $Ax = b$ است، وقتی که A مربعی، با درایه‌های تنگ، و از مرتبه بزرگ باشد. ماتریس تنگ ماتریسی است که بیشتر درایه‌های آن صفر باشند. چنین دستگاههایی به شکلهای گوناگونی پیدا می‌شوند، ولی ما بحث خود را به آنهایی که الگویی ساده و معلوم برای عناصر غیرصفر دارند، محدود می‌سازیم. این دستگاهها معمولاً در حل عددی معادلات دیفرانسیل جزئی ظاهر می‌شوند، و یک مثال در بخش ۸.۸ داده شده است. به علت مرتبه بزرگ اغلب دستگاههای معادلات خطی تنگ، گاهی به بزرگی 10^5 یا بیشتر، دستگاه خطی را نمی‌توان با روشهایی مستقیم مانند حذف گاوسی حل کرد. روشهای بارستی، روشهای ارجح‌اند و این روشها در بخش ۶.۸ تا ۹.۸ معرفی شده‌اند.

برای حل دستگاههای مربعی چگال از مرتبه متوسط، اغلب مراکز رایانه‌یی مجموعه برنامه‌هایی دارند که می‌توان برای انواع مسائل به‌کاربرد. دانشجویان باید با این برنامه‌ها در مرکز رایانه دانشگاه خود آشنا شده و آنها را برای بیشتر روشن کردن مطالب این فصل به‌کار برند. یک بسته نرم‌افزاری عالی به نام LINPACK ساخته شده است، که در دون‌گارا^۱ و همکاران (۱۹۷۹)، توضیح داده شده است. این نرم‌افزار به‌طور گسترده در دسترس است و ما باز هم در این فصل به آن مراجعه خواهیم کرد.

۱.۸ حذف گاوسی

حذف گاوسی نام رسمی روش حل دستگاههای معادلات خطی با حذف تدریجی مجهولات و تبدیل آنها به دستگاههای مرتبه پایینتر است، این روشی است که بیشتر در دبیرستان یا در درس جبر خطی دوره کارشناسی (که اغلب با تولید شکل سطری-پلکانی ماتریس بستگی دارد) تدریس می‌شود. در اینجا یک تعریف دقیق حذف گاوسی داده شده است، که برای کاربرد در رایانه و تحلیل خطای گردکردن که در محاسبه آن پیش می‌آید، لازم است.

برای حل $Ax = b$ ، آن را به یک دستگاه هم‌ارز $Ux = g$ که U یک ماتریس بالامتثالی است تبدیل می‌کنیم. این دستگاه با فرایند جایگذاری پسر و به‌سادگی حل می‌شود. دستگاه خطی اصلی را با $A^{(1)}x = b^{(1)}$ نشان می‌دهیم

$$A^{(1)} = [a_{ij}^{(1)}] \quad b^{(1)} = [b_1^{(1)}, \dots, b_n^{(1)}]^T \quad 1 \leq i, j \leq n$$

که در آن n مرتبه دستگاه است. این دستگاه را به شکل یک دستگاه مثلثی $Ux = g$ بدل می‌کنیم بدین نحو که با افزودن مضاربی از یک معادله به معادله دیگر مجهولی را از معادله دوم حذف می‌نماییم. از عملیات سطری دیگر در تبدیلیهایی که در بخشهای آینده آمده‌اند، استفاده شده است. برای آنکه ارائه مطلب ساده باشد، بعضی فرضهای فنی را در تعریف الگوریتم به‌کار می‌بریم؛ این فرضها در بخش بعد حذف می‌شوند.

الگوریتم حذف گاوسی

مرحله ۱ فرض می‌کنیم $a_{11}^{(1)} \neq 0$. ضریبهای سطری را چنین تعریف می‌کنیم

$$m_{i1} = \frac{a_{i1}^{(1)}}{a_{11}^{(1)}} \quad i = 2, 3, \dots, n$$

این ضریبها برای حذف جمله x_1 از معادلات ۲ تا n به‌کار می‌روند. تعریف می‌کنیم

$$\begin{aligned} a_{ij}^{(2)} &= a_{ij}^{(1)} - m_{i1}a_{1j}^{(1)} & i, j &= 2, \dots, n \\ b_i^{(2)} &= b_i^{(1)} - m_{i1}b_1^{(1)} & i &= 2, \dots, n \end{aligned}$$

همچنین سطرهای اول A و b دست‌نخورده باقی می‌مانند و در ستون اول $A^{(1)}$ ، همه عناصر زیر قطر، صفر گذاشته می‌شود.

دستگاه $A^{(2)}x = b^{(2)}$ بدین شکل است

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ \circ & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ \vdots & \vdots & & \\ \circ & a_{n2}^{(2)} & \dots & a_{nn}^{(2)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(2)} \\ \vdots \\ b_n^{(2)} \end{bmatrix}$$

حذف مجهولات را با رفتن به ستونهای ۲ و ۳ و غیره ادامه می‌دهیم. و این وضع به‌طور کلی در زیر بیان شده است.

مرحله k : گیریم $1 \leq k \leq n-1$. فرض می‌کنیم $A^{(k)}x = b^{(k)}$ ساخته شده و x_1, x_2, \dots, x_{k-1} در مراحل پی‌درپی حذف شده‌اند و $A^{(k)}$ دارای شکل زیر است

$$\begin{bmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ \circ & a_{22}^{(2)} & & a_{2n}^{(2)} \\ \vdots & & \ddots & \vdots \\ \circ & \dots & \circ & a_{kk}^{(k)} & \dots & a_{kn}^{(k)} \\ \vdots & & & \vdots & & \vdots \\ \circ & \dots & \circ & a_{nk}^{(k)} & \dots & a_{nn}^{(k)} \end{bmatrix}$$

فرض می‌کنیم $a_{kk}^{(k)} \neq 0$. ضریبها را چنین تعریف می‌کنیم

$$m_{ik} = \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} \quad i = k+1, \dots, n \quad (1.1.8)$$

از این ضریبها برای حذف مجهول x_k از معادلات $k+1$ تا n استفاده می‌کنیم. تعریف می‌کنیم

$$\begin{aligned} a_{ij}^{(k+1)} &= a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} \\ b_i^{(k+1)} &= b_i^{(k)} - m_{ik} b_k^{(k)} \quad i, j = k+1, \dots, n \end{aligned} \quad (2.1.8)$$

سطرهای قبلی 1 تا k دست‌نخورده می‌مانند و در ستون k ام زیر عنصر قطر، همه جا صفرگذارده می‌شود. با ادامه این شیوه، پس از $n-1$ مرحله، $A^{(n)}x = b^{(n)}$ را به دست می‌آوریم:

$$\begin{bmatrix} a_{11}^{(1)} & \dots & a_{1n}^{(1)} \\ \circ & & \\ \vdots & \ddots & \vdots \\ \circ & \dots & a_{nn}^{(n)} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1^{(1)} \\ \vdots \\ b_n^{(n)} \end{bmatrix}$$

برای راحتی در نمادگذاری، می‌گیریم $U = A^{(n)}$ و $g = b^{(n)}$. دستگاه $Ux = g$ یک دستگاه بالا مثلثی و به‌سادگی قابل حل است. ابتدا داریم

$$x_n = \frac{g_n}{u_{nn}}$$

$$x_k = \frac{1}{u_{kk}} \left[g_k - \sum_{j=k+1}^n u_{kj} x_j \right] \quad k = n-1, n-2, \dots, 1 \quad (3.1.8)$$

بدین صورت الگوریتم حذف گاوسی کامل می‌شود.

مثال دستگاه خطی زیر را حل کنید

$$\begin{aligned} x_1 + 2x_2 + x_3 &= 0 \\ 2x_1 + 2x_2 + 3x_3 &= 3 \\ -x_1 - 3x_2 &= 2 \end{aligned} \quad (4.1.8)$$

برای ساده کردن نمادگذاری، توجه داریم که مجهولات x_1 و x_2 تا آخرین مرحله، در الگوریتم وارد نمی‌شوند. بنابراین دستگاه خطی بالا را با ماتریس افزوده زیر نشان می‌دهیم:

$$[A|b] = \left[\begin{array}{ccc|c} 1 & 2 & 1 & 0 \\ 2 & 2 & 3 & 3 \\ -1 & -3 & 0 & 2 \end{array} \right]$$

عملیات سطری در این ماتریس افزوده اجرا و مجهولات در آخرین مرحله داده شده‌اند. در نمودار زیر، ضرایب در کنار پیکان، متناظر با تغییراتی که ایجاد می‌نمایند، داده شده‌اند

$$\begin{aligned} \left[\begin{array}{ccc|c} 1 & 2 & 1 & 0 \\ 2 & 2 & 3 & 3 \\ -1 & -3 & 0 & 2 \end{array} \right] & \xrightarrow{\substack{m_{21}=2 \\ m_{31}=-1}} \left[\begin{array}{ccc|c} 1 & 2 & 1 & 0 \\ 0 & -2 & 1 & 3 \\ 0 & -1 & 1 & 2 \end{array} \right] \\ & \downarrow m_{22}=\frac{1}{2} \\ [U|g] & \equiv \left[\begin{array}{ccc|c} 1 & 2 & 1 & 0 \\ 0 & -2 & 1 & 3 \\ 0 & 0 & \frac{1}{2} & \frac{1}{2} \end{array} \right] \end{aligned}$$

از حل $Ux = g$ داریم

$$x_3 = 1 \quad x_2 = -1 \quad x_1 = 1$$

تجزیه مثلثی ماتریس خیلی بجاست که ضرایب m_{ij} را حفظ کنیم، زیرا اغلب $Ax = b$ را با یک A ولی با بردارهای مختلف b حل می‌کنیم. در رایانه، عناصر $a_{ij}^{(k+1)}$ ، $j \geq i$ همیشه در حافظه $a_{ij}^{(k)}$ ذخیره می‌شوند. عناصر زیر قطر صفر شده‌اند، و این یک حافظه مناسب برای عناصر m_{ij} در اختیار می‌گذارد. m_{ij} را در محلی که ابتدا برای ذخیره a_{ij} ، $i > j$ ، به‌کار برده شده بود ذخیره می‌نماییم.

باز یک دلیل دیگر وجود دارد که به ضرایب m_{ij} به‌عنوان عناصر یک ماتریس نگاه می‌کنیم. ابتدا ماتریس پایین مثلثی زیر را معرفی می‌کنیم

$$L = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ m_{21} & 1 & 0 & \dots & 0 \\ \vdots & & & & \vdots \\ m_{n1} & m_{n2} & \dots & 0 & 1 \end{bmatrix}$$

قضیه ۱.۸ اگر L و U ماتریسهای پایین مثلثی و بالا مثلثی باشند که قبلاً در استفاده از حذف گاوسی تعریف شده‌اند، آنگاه

$$A = LU \quad (5.1.8)$$

برهان این اثبات اساساً یک محاسبه جبری است که با استفاده از تعاریف (۱.۸) و (۲.۱.۸) صورت می‌گیرد. برای تجسم درایه‌های $(LU)_{ij}$ ، از فرمول برداری زیر استفاده می‌نماییم

$$(LU)_{ij} = [m_{i1}, \dots, m_{i, i-1}, 1, 0, \dots, 0] \begin{bmatrix} u_{1j} \\ \vdots \\ u_{ij} \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

برای $j \leq i$,

$$\begin{aligned}(LU)_{ij} &= m_{i1}u_{1j} + m_{i2}u_{2j} + \dots + m_{i,i-1}u_{i-1,j} + u_{i,j} \\ &= \sum_{k=1}^{i-1} m_{ik}a_{kj}^{(k)} + a_{ij}^{(i)} \\ &= \sum_{k=1}^{i-1} [a_{ij}^{(k)} - a_{ij}^{(k+1)}] + a_{ij}^{(i)} \\ &= a_{ij}^{(1)} = a_{ij}\end{aligned}$$

برای $j > i$

$$\begin{aligned}(LU)_{ij} &= m_{i1}u_{1j} + \dots + m_{ij}u_{jj} \\ &= \sum_{k=1}^{j-1} m_{ik}a_{kj}^{(k)} + m_{ij}a_{jj}^{(j)} \\ &= \sum_{k=1}^{j-1} [a_{ij}^{(k)} - a_{ij}^{(k+1)}] + a_{ij}^{(j)} \\ &= a_{ij}^{(1)} = a_{ij}\end{aligned}$$

و اثبات کامل می‌شود.

تجزیه (۸.۱.۵) نتیجه مهمی است که از آن در بسط انواع روشهای حذف گاوسی برای رده‌های خاصی از ماتریسها، استفاده وسیعی می‌شود. ولی در حال حاضر ما فقط فرع زیر را مورد نظر قرار می‌دهیم.

فرع با همان ماتریسهای A, L و U در قضیه ۸.۸،

$$\begin{aligned}\det(A) &= u_{11}u_{22} \dots u_{nn} \\ &= a_{11}^{(1)} a_{22}^{(2)} \dots a_{nn}^{(n)}\end{aligned}$$

برهان طبق قاعده حاصلضرب دترمینانها

$$\det(A) = \det(L) \det(U)$$

چون L و U مثلثی هستند، دترمینان آنها برابر حاصلضرب عناصر قطری آنهاست. و نتیجه مطلوب به سادگی حاصل می شود، زیرا $\det(L) = 1$

مثال برای دستگاه (۴.۱.۸) از مثال قبلی،

$$L = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & \frac{1}{2} & 1 \end{bmatrix} \quad U = \begin{bmatrix} 1 & 2 & 1 \\ 0 & -2 & 1 \\ 0 & 0 & \frac{1}{2} \end{bmatrix}$$

به سادگی تحقیق می شود که $A = LU$. همچنین $\det(A) = \det(U) = -1$.

شمارش عملیات. برای تحلیل تعداد عملیات لازم در حل $Ax = b$ با استفاده از حذف گاوسی، پیدایش L و U را از A و تبدیل b به g و بالاخره پیدا کردن جواب x را به طور جداگانه در نظر می گیریم.

۱. محاسبه L و U . در مرحله ۱، $(n-1)$ تقسیم برای محاسبه مضارب m_{i1} ، $2 \leq i \leq n$ به کار برده شده بود. سپس $(n-1)^2$ ضرب و $(n-1)^2$ جمع به کار رفته بود تا عناصر جدید $a_{ij}^{(1)}$ به دست آیند. شمارش را می توانیم به همین شیوه برای سایر مراحل ادامه دهیم.

نتایج در جدول ۱.۸ خلاصه شده است. مقدار حاصل جمع هر ستون با استفاده از اتحادهای زیر به دست آمده است

$$\sum_{j=1}^p j = \frac{p(p+1)}{2} \quad \sum_{j=1}^p j^2 = \frac{p(p+1)(2p+1)}{6} \quad p \geq 1$$

معمولاً، تعداد ضرب و تقسیمهاست که روی هم رفته به عنوان تعداد عملیات در روش حذف گاوسی به کار می رود. در رایانه های قدیمی، جمع بسیار سریعتر از ضرب و تقسیم انجام می شد و

جدول ۱.۸ شمارش عملیات برای تجزیه یک ماتریس به LU

مرحله k	جمع	ضرب	تقسیم
۱	$(n-1)^2$	$(n-1)^2$	$n-1$
۲	$(n-2)^2$	$(n-2)^2$	$n-2$
\vdots	\vdots	\vdots	\vdots
$n-1$	۱	۱	۱
حاصل جمع	$\frac{n(n-1)(2n-1)}{6}$	$\frac{n(n-1)(2n-1)}{6}$	$\frac{n(n-1)}{2}$

بنابراین در محاسبه هزینه بسیاری از الگوریتمها جمعها نادیده گرفته می‌شدند. ولی، در رایانه‌های جدید، زمانهای انجام جمع، ضرب و تقسیم از لحاظ اندازه خیلی به هم نزدیک‌اند. برای سهولت در نمادگذاری، فرض می‌کنیم $MD(\cdot)$ و $AS(\cdot)$ به ترتیب تعداد ضرب و تقسیمها، و تعداد جمع و تفریقها برای محاسبه کمیت داخل پرانتزها باشند.

برای تجزیه A به LU ، داریم

$$\begin{aligned} MD(LU) &= \frac{n(n^2 - 1)}{3} \doteq \frac{n^3}{3} \\ AS(LU) &= \frac{n(n-1)(2n-1)}{6} \doteq \frac{n^3}{3} \end{aligned} \quad (6.1.8)$$

برآوردهای نهایی برای مقادیر بزرگتر n معتبرند.

۲. تبدیل b به $g = b^{(n)}$

$$\begin{aligned} MD(g) &= (n-1) + (n-2) + \dots + 1 = \frac{n(n-1)}{2} \\ AS(g) &= \frac{n(n-1)}{2} \end{aligned} \quad (7.1.8)$$

۳. حل $Ux = g$

$$MD(x) = \frac{n(n+1)}{2} \quad AS(x) = \frac{n(n-1)}{2} \quad (8.1.8)$$

۴. حل $Ax = b$. مراحل ۱ تا ۳ را ترکیب می‌کنیم تا به دست آوریم

$$\begin{aligned} MD(LU, x) &= \frac{n^3}{3} + n^2 - \frac{n}{3} \doteq \frac{1}{3}n^3 \\ AS(LU, x) &= \frac{n(n-1)(2n+5)}{6} \doteq \frac{1}{3}n^3 \end{aligned} \quad (9.1.8)$$

تعداد جمعها همیشه در همان حدود تعداد ضرب و تقسیمهاست، و لذا از این پس ما فقط آخری را در نظر می‌گیریم. به اولین چیزی که باید توجه کنیم این است که حل $Ax = b$ ، در مقایسه با عمل ظاهراً ساده‌ای مثل ضرب دو ماتریس $n \times n$ ، بسیار ناچیز است. ضرب دو ماتریس نیاز به n^3 عمل دارد و حل $Ax = b$ فقط در حدود $\frac{1}{3}n^3$ عمل در برخواهد داشت.

دوم اینکه، هزینه اصلی حل $Ax = b$ در تجزیه $A = LU$ است. وقتی این را پیدا کردیم، فقط n^2 عمل دیگر برای حل $Ax = b$ لازم است. وقتی یک بار $Ax = b$ حل شده باشد حل دستگاههای اضافی دیگر با همان ماتریس ضرایب نسبتاً ناچیز خواهد بود، به شرطی که تجزیه LU ذخیره شده باشد.

بالاخره، روش حذف گاوسی بسیار ناچیزتر از قاعده کرامر است، که از دترمینان استفاده می‌کند و اغلب در درس جبر خطی تدریس می‌شود [برای مثال، بخش ۴.۲ آنتون (۱۹۸۴) را ببینید]. اگر دترمینان‌ها در قاعده کرامر از راه بسط به دترمینان‌های جزئی محاسبه شوند، تعداد عملیات $(n+1)!$ خواهد شد. برای $n = 10$ ، حذف گاوسی 43° عمل و قاعده کرامر 3991680° عمل خواهد داشت. باید تأکید کنیم که قاعده کرامر یک وسیله محاسبه عملی نیست و فقط به عنوان یک وسیله ریاضیات نظری باید به آن بنگریم.

۵. وارون A . وارون A ، معمولاً لازم نیست ولی می‌توان با روش حذف گاوسی آن را پیدا کرد. پیدا کردن A^{-1} هم‌ارز حل معادله $AX = I$ است که X یک ماتریس $n \times n$ مجهول است. اگر X و I را بر حسب ستونهای آنها بنویسیم.

$$X = [x^{(1)}, \dots, x^{(n)}] \quad I = [e^{(1)}, \dots, e^{(n)}]$$

آنگاه حل $AX = I$ هم‌ارز حل n دستگاه

$$Ax^{(1)} = e^{(1)}, \dots, Ax^{(n)} = e^{(n)} \quad (10.1.8)$$

است که همگی دارای یک ماتریس ضرایب A هستند. با استفاده از (۱) - (۳)

$$MD(A^{-1}) = \frac{3}{4}n^3 - \frac{n}{3} = \frac{4}{3}n^3$$

محاسبه A^{-1} چهار برابر حل $Ax = b$ برای یک بردار تنهای b ، هزینه خواهد داشت، نه n بار که ممکن است در ابتدا به نظر آید. با توجه دقیق به جزئیات فرایند وارون‌سازی، و استفاده از شکل خاص بردارهای سمت راست $e^{(1)}, \dots, e^{(n)}$ ، ممکن است هزینه دقیقاً به

$$MD(A^{-1}) = n^3 \quad (11.1.8)$$

عمل کاهش یابد.

ولی، باز هم در بسیاری از موارد، اتلاف وقت است که در حل $Ax = b$ وارون A ، یعنی A^{-1} ، را به دست آوریم. و هیچ مزیتی در ذخیره A^{-1} به جای تجزیه LU برای حل دستگاههای دیگر $Ax = b$ وجود ندارد. در هر دو حالت تعداد ضرب و تقسیمهای لازم برای حل $Ax = b$ دقیقاً n^2 است.

۲.۸ محورگیری و مقیاس دهی در حذف گاوسی

در هر مرحله از فرایند حذف در بخش اخیر، عنصر محور مربوط $a_{kk}^{(k)}$ را ناصفر فرض کردیم. برای حذف این فرض، هر مرحله از فرایند حذف را، با تعویض سطرها با هم و قراردادن یک عنصر ناصفر در جای محور، آغاز می‌کنیم. اگر چنین عنصری موجود نباشد، این ماتریس باید تکین باشد که با فرض تناقض دارد.

ولی، این کافی نیست که فقط بخواهیم عنصر محور ناصفر باشد. اغلب یک عنصر صفر به علت خطاهای گردکردن در محاسبه آن ناصفر می‌شود. استفاده از چنین عنصری به عنوان عنصر محور به خطاهای بزرگی در محاسبات بعدی در ماتریس منجر خواهد شد. برای اجتناب از این امر و به دلایل دیگری مربوط به انتشار خطاهای گردکردن، ما محورگیری جزئی و محورگیری کلی را معرفی می‌نماییم.

تعریف ۱. محورگیری جزئی. برای $1 \leq k \leq n-1$ ، در فرایند حذف گاوسی در مرحله k ام، می‌گیریم

$$c_k = \text{Max}_{k \leq i \leq n} |a_{ik}^{(k)}| \quad (1.2.8)$$

گیریم i کوچکترین اندیس سطری باشد، $i \geq k$ ، که برای آن ماکسیم c_k به دست آمده است. اگر $i > k$ ، آنگاه سطرهای k و i در A و b را با هم تعویض کرده و به مرحله k ام فرایند حذف می‌پردازیم. همه ضریبها اکنون در

$$|m_{ik}| \leq 1 \quad i = k+1, \dots, n \quad (2.2.8)$$

صدق می‌کنند و این امر به جلوگیری از بسیار زیاد شدن تفاوت اندازه‌های عناصر $A^{(k)}$ کمک خواهد کرد، و بنابراین امکان خطاهای کاهش بزرگ در ارقام بامعنی را کمتر خواهد نمود.

۲. محورگیری کلی. تعریف می‌کنیم

$$c_k = \text{Max}_{k \leq i, j \leq n} |a_{ij}^{(k)}|$$

سطرهای A و b و ستونهای A را تعویض می‌کنیم تا عنصری که ماکسیم c_k را دارد در جای محور قرار گیرد. باید توجه کرد که با تعویض یک ستون، ترتیب مجهولات عوض می‌شود. پس از اتمام حذف و فرایند جایگذاری پسر، این تعویض باید عکس شود.

ثابت شده است که محورگیری کلی موجب می‌شود که خطای گردکردن در حذف گاوسی، در مقایسه با موقعی که از هیچ محورگیری استفاده نمی‌شود، با سرعت کم و قابل قبولی انتشار یابد. نتایج

نظری در استفاده از محورگیری جزئی کاملاً به آن خوبی نیست، ولی تقریباً در تمام مسائل عملی، رفتار خطا اساساً مانند رفتار خطا در محورگیری کلی است. با مقایسه زمانهای عملیات، محورگیری کلی یک شیوه گرانتری است و بنابراین، در بیشتر الگوریتم‌های عملی، محورگیری جزئی به کار می‌رود. از این پس وقتی کلمه محورگیری را به کار می‌بریم منظورمان محورگیری جزئی است. کل مسئله انتشار خطای گرد کردن در روش حذف گاوسی توسط ج. ه. ویلکینسن عمیقاً تحلیل شده است [مثلاً (صفحات ۲۰۹-۲۲۰) ویلکینسن (۱۹۵۶) را ببینید]، و بعضی از نتیجه‌گیریهایی او در بخش ۴.۸ ارائه شده است.

مثال تأثیر به کار بردن محورگیری را با حل دستگاه زیر نشان می‌دهیم:

$$0.729x + 0.81y + 0.9z = 0.6867$$

$$x + y + z = 0.8338$$

$$1.331x + 1.21y + 1.1z = 1.000 \quad (3.2.8)$$

جواب دقیق که تا چهار رقم معنی دار گرد شده چنین است:

$$x = 0.2245 \quad y = 0.2814 \quad z = 0.3279 \quad (4.2.8)$$

از حساب اعشاری با ممیز شناور، با چهار رقم اعشاری، برای حل دستگاه خطی استفاده خواهد شد. دلیل استفاده از این حساب نشان دادن اثر کار با یک تعداد متناهی رقم است، در حالی که بتوان اندازه ارقام را معقول نگهداشت. نماد ماتریس افزوده برای نمایش دستگاه (۳.۲.۸) به کار برده می‌شود، درست مانند آنچه در مثال قبل (۴.۱.۸) عمل شد.

۱. حل بدون محورگیری

$$\left[\begin{array}{ccc|c} 0.729 & 0.81 & 0.9 & 0.6867 \\ 1.000 & 1.000 & 1.000 & 0.8338 \\ 1.331 & 1.21 & 1.1 & 1.000 \end{array} \right]$$

$$\downarrow \begin{array}{l} m_{21} = 1.372 \\ m_{31} = 1.826 \end{array}$$

$$\left[\begin{array}{ccc|c} 0.729 & 0.81 & 0.9 & 0.6867 \\ 0 & -0.111 & -0.235 & -0.1084 \\ 0 & -0.269 & -0.543 & -0.2540 \end{array} \right]$$

$$\begin{array}{c} \downarrow m_{32} = 2,422 \\ \left[\begin{array}{ccc|c} 0,7290 & 0,8100 & 0,9000 & 0,6867 \\ 0,0 & -0,1110 & -0,2350 & -0,1084 \\ 0,0 & 0,0 & 0,2640 & 0,008700 \end{array} \right] \end{array}$$

جواب چنین است

$$x = 0,2251 \quad y = 0,2790 \quad z = 0,3295 \quad (5.2.8)$$

۲. حل با محورگیری. برای نشان دادن تعویض سطرها i و j ، نماد $r_i \leftrightarrow r_j$ را به کار خواهیم برد.

$$\left[\begin{array}{ccc|c} 0,7290 & 0,8100 & 0,9000 & 0,6867 \\ 1,0000 & 1,0000 & 1,0000 & 0,8338 \\ 1,331 & 1,210 & 1,100 & 1,000 \end{array} \right]$$

$$\begin{array}{c} \downarrow \\ r_1 \leftrightarrow r_2 \quad \begin{array}{l} m_{21} = 0,7512 \\ m_{31} = 0,5277 \end{array} \end{array}$$

$$\left[\begin{array}{ccc|c} 1,331 & 1,210 & 1,100 & 1,000 \\ 0,0 & 0,09090 & 0,1736 & 0,08250 \\ 0,0 & 0,1473 & 0,2975 & 0,1390 \end{array} \right]$$

$$\begin{array}{c} \downarrow \\ r_2 \leftrightarrow r_3 \quad m_{32} = 0,6171 \end{array}$$

$$\left[\begin{array}{ccc|c} 1,331 & 1,210 & 1,100 & 1,000 \\ 0,0 & 0,1473 & 0,2975 & 0,1390 \\ 0,0 & 0,0 & -0,1000 & -0,003280 \end{array} \right]$$

جواب چنین است:

$$x = 0,2246 \quad y = 0,2812 \quad z = 0,3280 \quad (6.2.8)$$

خطا در (۵.۲.۸) برحسب اینکه کدام مؤلفهٔ جواب در نظر گرفته شود، از هفت تا شانزده برابر بزرگتر از خطا در (۶.۲.۸) است. نتایج (۶.۲.۸) یک رقم معنی دار بیش از نتایج (۵.۲.۸) دارند. این نکته اثر مثبت استفاده از محورگیری را در رفتار خطا، در روش حذف گاوسی، نشان می دهد.

محورگیری نتیجهٔ تجزیه به عوامل (۵.۱.۸) را که در قضیهٔ ۱.۸ داده شده عوض می کند. ولی

در یک شکل اصلاح شده، نتیجه باز هم درست است. اگر تعویض سطرها که بر اثر محورگیری القا شده است، قبل از عمل حذف در A انجام شده باشد، دیگر محورگیری لزومی نخواهد داشت. تعویض سطرها در A را می توان با ضرب قدامی (premultiplication) A در ماتریس جایگشت مناسب P ، به صورت PA به دست آورد. در این صورت حذف گاوسی روی PA به

$$LU = PA \quad (۷.۲.۸)$$

می انجامد که U ماتریس بالا مثلثی حاصل از فرایند حذف با محورگیری است. ماتریس پایین مثلثی L را می توان با استفاده از ضرایب در حذف گاوسی با محورگیری، ساخت. ما از ذکر جزئیات صرف نظر می کنیم، زیرا ساخت عملی مهم نیست.

مثال با توجه به مثال قبل که با محورگیری حل شد ماتریسهای زیر را تشکیل می دهیم

$$L = \begin{bmatrix} ۱۰۰۰ & ۰ & ۰ \\ ۰.۵۴۷۷ & ۱۰۰۰ & ۰ \\ ۰.۷۵۱۳ & ۰.۶۱۷۱ & ۱۰۰۰ \end{bmatrix} \quad U = \begin{bmatrix} ۱.۳۳۱ & ۱.۲۱۰ & ۱.۱۰۰ \\ ۰ & ۰.۱۴۷۳ & ۰.۲۹۷۵ \\ ۰ & ۰ & -۰.۱۰۰۰ \end{bmatrix}$$

از ضرب این دو

$$LU = \begin{bmatrix} ۱.۳۳۱ & ۱.۲۱۰ & ۱.۱۰۰ \\ ۰.۷۲۸۹ & ۰.۸۱۰۰ & ۰.۹۰۰۰ \\ ۱.۰۰۰ & ۱.۰۰۰ & ۱.۰۰۰ \end{bmatrix} = PA \quad P = \begin{bmatrix} ۰ & ۰ & ۱ \\ ۱ & ۰ & ۰ \\ ۰ & ۱ & ۰ \end{bmatrix}$$

نتیجه PA ماتریس A است که در آن ابتدا، سطرها ۱ و ۳، و سپس سطرها ۲ و ۳ تعویض شده اند. این مثالی از (۷.۲.۸) است.

مقیاس دهی به طور تجربی مشاهده شده است که اگر عناصر ماتریس ضرایب A از نظر اندازه زیاد متفاوت باشند، خطاهای کاهش در ارقام بامعنی بزرگی احتمالاً وارد محاسبه می شوند و انتشار خطاهای گردکردن بدتر می شود. برای اجتناب از این مسأله، معمولاً به ماتریس A مقیاس می دهیم تا عناصر کمتر متفاوت باشند. این کار معمولاً از ضرب سطرها و ستونها در ثابتهای مناسب انجام می گیرد. موضوع مقیاس دهی فعلاً به خوبی معلوم نیست، به ویژه اینکه چگونه می توان تضمین کرد که اثر خطاهای گردکردن در حذف گاوسی، با چنین مقیاس دهی، کوچک می شوند. از تجربه محاسباتی به نظر می رسد که اغلب تمام سطرها باید طوری مقیاس دهی شوند که از نظر اندازه

تقریباً مساوی شوند. به علاوه، به تمام ستونها می‌توان مقیاس داد تا از نظر قدرمطلق تقریباً به یک اندازه درآیند. این عمل اخیر، هم‌ارز است با اینکه به مؤلفه‌های مجهول، یعنی x_i ها مقیاس دهیم؛ و آن را می‌توان اغلب این طور تعبیر نمود که بگوییم x_i ها باید با واحدهایی قابل مقایسه با هم، اندازه‌گیری شوند.

هیچ شیوه‌ای از پیش تعیین شده‌ای برای انتخاب عاملهای مقیاس‌دهی در دست نیست که بتواند تنها برای اطلاع از A و b ، اثر انتشار خطای گردکردن را همیشه کاهش دهد. آنچه در استوارت (۱۹۷۷) در استفاده کلی از مقیاس‌دهی آمده است، آن‌گونه که در پاراگراف قبلی توضیح داده شد، تا اندازه‌ای قابل انتقاد است. او پیشنهاد می‌نماید که فاکتورهای مقیاس‌دهی طوری انتخاب شوند که در ماتریس با مقیاسهای جدید، خطاهای دریاها تا حدی برابر باشند. وقتی گردکردن تنها منبع خطا باشد، این عمل به یک شیوه‌ی مقیاس‌دهی می‌انجامد که تمام عناصر را تقریباً به یک اندازه می‌کند. در برنامه‌های LIN-PACK، مقیاس‌دهی دخالت داده نمی‌شود، ولی یک شیوه به همان طریق استوارت توصیه می‌شود. [برای بحث مفصلتر دونگارا و همکاران (۱۹۷۹، صص ۱۷-۱۱۲) را ببینید. برای بحثهای دیگر مقیاس‌دهی فورسایت ومولر (۱۹۶۷، فصل ۱۱) و گلوب و ون لون (۱۹۸۳، صص ۷۲-۷۴) را ببینید].

اگر B معرف نتیجه مقیاس‌دهی سطری و ستونی A باشد، آنگاه

$$B = D_1 A D_2$$

که D_1 و D_2 ماتریسهای قطری هستند با دریا‌های متشکل از ثابتهای مقیاس‌دهی. برای حل $Ax = b$ مشاهده می‌کنیم که

$$D_1 A D_2 (D_2^{-1} x) = D_1 b$$

بنابراین جواب x را از حل

$$Bz = D_1 b \quad x = D_2 z \quad (۸.۲.۸)$$

به دست می‌آوریم. بقیه بحث، به مقیاس‌دهی سطری محدود شده است، زیرا نوعی از آن تا حدی مورد موافقت همگان قرار گرفته است.

معمولاً سعی می‌کنیم ضرایب را طوری انتخاب نماییم که داشته باشیم

$$\max_{1 \leq j \leq n} |b_{ij}| \doteq 1 \quad i = 1, \dots, n \quad (۹.۲.۸)$$

که $B = [b_{ij}]$ نتیجه مقیاس‌دهی A است. ساده‌ترین راه این است که چنین تعریف کنیم

$$s_i = \max_{1 \leq j \leq n} |a_{ij}| \quad i = 1, \dots, n$$

$$b_{ij} = \frac{a_{ij}}{s_i}, \quad j = 1, \dots, n \quad (۱۰.۲.۸)$$

ولی چون این روش خطای گردکردن اضافه‌تری برای هر عنصر ماتریس ضرایب ایجاد می‌نماید، از دو روش دیگر بیشتر استفاده می‌شوند.

۱. تبدیل مقیاس با استفاده از پایهٔ عددی رایانه. گیریم β معرّف پایه‌ای باشد که در حساب رایانه‌یی به‌کار برده می‌شود، برای مثال، $\beta = ۲$ در ماشینهای دودویی. گیریم r_i کوچکترین عدد صحیحی باشد که برای آن $\beta^{r_i} \geq s_i$. ماتریس مقیاس‌دهی شدهٔ B را به شکل زیر تعریف می‌کنیم

$$b_{ij} = \frac{a_{ij}}{\beta^{r_i}} \quad i, j = 1, \dots, n \quad (۱۱.۲.۸)$$

هیچ گردکردنی در تعریف b_{ij} وجود ندارد، فقط تغییر نما در شکل با ممیز شناور a_{ij} رخ می‌دهد. مقادیر B در رابطهٔ زیر صدق می‌کنند

$$\beta^{-1} < \text{Max}_{1 \leq j \leq n} |b_{ij}| \leq 1$$

و بنابراین (۹.۲.۸) به‌خوبی برقرار می‌گردد.

۲. مقیاس‌دهی ضمنی. چنانچه محورگیری در حذف گاوسی به‌کار رفته باشد، استفاده از مقیاس‌دهی معمولاً انتخاب عناصر محورگیری را عوض می‌کند. و فقط به دلیل همین تغییر عناصر محورگیری است که نتایج در حذف گاوسی عوض می‌شوند. این توجه را مدیون ف. باوئر هستیم [در صفحهٔ ۳۸ فورسایت و مولر (۱۹۶۷)] که می‌گوید اگر مقیاس‌دهی (۱۱.۲.۸) به‌کار برده شود، و اگر عناصر محورگیری همانهایی انتخاب شوند که در حل $Ax = b$ به‌کار رفته‌اند، آنگاه حل (۸.۲.۸) دقیقاً همان مقدار محاسبه‌شدهٔ x را خواهد داد. بنابراین تنها انتخاب عناصر محورگیری است که اهمیت دارد. برای مقیاس‌دهی ضمنی، کار با ماتریس A را ادامه می‌دهیم، ولی برای تعیین عنصر محورگیری در مرحلهٔ k -ام الگوریتم حذف گاوسی c_k را با

$$c_k = \text{Max}_{k \leq i \leq n} \left| \frac{a_{ik}^{(k)}}{s_i^{(k)}} \right| \quad (۱۲.۲.۸)$$

تعریف می‌نماییم که به‌جای تعریف (۱۰.۲.۸) که در تعریف محورگیری جزئی استفاده شده، به‌کار می‌رود. کوچکترین اندیس $i \geq k$ را انتخاب می‌کنیم که c_k در (۱۲.۲.۸) را به‌دست دهد، و اگر $i \neq k$ ، سطرهای i و k را تعویض می‌کنیم. سپس الگوریتم حذف بخش ۱.۸ را مانند قبل ادامه می‌دهیم. به‌نظر می‌آید که این شکل مقیاس‌دهی در بین الگوریتمهای منتشر شده که در آنها از مقیاس‌دهی استفاده شده از همه متداولتر است.

الگوریتمی برای حذف گاوسی ابتدا الگوریتمی به نام *Factor* برای تجزیه مثلثی ماتریس A خواهیم داد. در این الگوریتم، حذف گاوسی با محورگیری جزئی توأم با مقیاس دهی ضمنی (۱۰.۲.۸) و (۱۲.۲.۸)، به کار برده می شود. سپس یک الگوریتم به نام *Solve* برای استفاده از نتایج *Factor*، در حل دستگاه خطی $Ax = b$ خواهیم داد. دلیل تجزیه فرایند حذف به این دو مرحله این است که اغلب می خواهیم چندین دستگاه $Ax = b$ را با یک ماتریس A ، ولی با مقادیر مختلف b ، حل کنیم.

الگوریتم $Factor(A, n, Pivot, det, ier)$

۱. ملاحظات: A یک ماتریس $n \times n$ است که با استفاده از تجزیه LU به عوامل صورت می گیرد. حذف گاوسی همراه با محورگیری جزئی و مقیاس دهی ضمنی در سطرها به کار برده می شود. پس از تکمیل الگوریتم، ماتریس بالامثلثی U ، در قسمت بالایی مثلث A ذخیره می شود؛ و ضریبهای (۱.۱.۸) که L را زیر قطرش تشکیل می دهد در جاهای متناظر در A ذخیره می شوند. بردار محور شامل ثابت تمام تعویضهای سطری خواهد بود. اگر $Pivot(k) = k$ ، آنگاه هیچ تعویضی در مرحله k ام فرایند حذف، صورت نگرفته است. ولی اگر $Pivot(k) = i \neq k$ ، آنگاه در مرحله k ام فرایند حذف، سطرها i و k باهم تعویض شده اند. متغیر det در خروج شامل مقدار دترمینان A خواهد بود. متغیر ier یک شاخص خطاست. اگر $ier = 0$ این زیربرنامه به نحو رضایتبخشی انجام شده است. ولی اگر $ier = 1$ ، ماتریس A تکین است، بدین معنی که تمام عناصر محورگیری ممکن در مرحله ای از فرایند حذف، صفر بوده اند. در این حالت تمام محاسبات متوقف شده و از زیربرنامه خارج شده ایم. هیچ اقدامی برای کنترل دقت تجزیه محاسبه شده A به عمل نیامده است و A ممکن است بدون اینکه متوجه شده باشیم تقریباً تکین بوده باشد.

$$2. \quad det := 1.$$

$$3. \quad i = 1, \dots, n, s_i = \max_{1 \leq j \leq n} |a_{i,j}|.$$

۴. برای $k = 1, \dots, n - 1$ اعمال را تا مرحله ۱۶ انجام دهید.

$$5. \quad c_k := \max_{k \leq i \leq n} \left| \frac{a_{ik}}{s_i} \right|.$$

۶. i را کوچکترین زیرنمایه $i \geq k$ بگیرید که ماکسیم آن در مرحله ۵ حاصل شده است.

$$Pivot(k) := i.$$

۷. اگر $c_k = 0$ ، آنگاه $ier := 1$ ، $det := 0$ و از الگوریتم خارج شوید.

۸. اگر $i = k$ ، به مرحله ۱۱ بروید.

$$9. \quad det := -det.$$

۱۰. a_{kj} را با $a_{i,j}$ برای $j = k, \dots, n$ عوض کنید و s_k را با s_i .

۱۱. برای $i = k + 1, \dots, n$ تا مرحله ۱۴ کار را ادامه دهید.

$$a_{ik} := m_i := \frac{a_{ik}}{a_{kk}} \quad ۱۲$$

$$j = k + 1, \dots, n \quad a_{ij} := a_{ij} - m_i a_{kj} \quad ۱۳$$

۱۴. حلقه را در i ختم کنید.

$$\det := a_{kk} \cdot \det \quad ۱۵$$

۱۶. حلقه را در k ختم کنید.

۱۷. $ier := 0$ و از الگوریتم خارج شوید.

الگوریتم $Solve(A, n, b, Pivot)$

۱. ملاحظات: باین الگوریتم دستگاه خطی $Ax = b$ حل می‌شود. فرض بر این است که ماتریس

اصلی A با استفاده از الگوریتم $Factor$ تجزیه شده است که تعویض سطرها در $Pivot$ ثبت شده است. هنگام خروج، جواب در b ذخیره می‌شود. ماتریس A و بردار محور دست‌نخورده می‌مانند.

۲. برای $k = 1, 2, \dots, n - 1$ تا مرحله ۵ جلو بروید.

۳. اگر $pivot(k) \neq k$ آنگاه $i := pivot(k)$ و b_i و b_k را با هم عوض کنید.

$$i = k + 1, \dots, n \quad b_i := b_i - a_{ik} b_k \quad ۴$$

۵. حلقه را در k قطع کنید.

$$b_n := b_n / a_{nn} \quad ۶$$

۷. برای $i = n - 1, \dots, 1$ تا مرحله ۹ عمل کنید.

$$b_i := \frac{1}{a_{ii}} \left\{ b_i - \sum_{j=i+1}^n a_{ij} b_j \right\} \quad ۸$$

۹. حلقه را در i ختم کنید.

۱۰. از الگوریتم خارج شوید.

مثال قبلی (۳.۲.۸) را به عنوان یک توضیح به‌کار می‌بریم. استفاده از تبدیل مقیاس ضمنی در این حالت نیاز به تغییر انتخاب عناصر محور، در محورگیری جزئی، ندارد. الگوریتمهایی شبیه $Solve$ و $Factor$ در منابع چندی داده شده‌اند [فورسایت و مولر (۱۹۶۷)] و فصل اول دون‌گارا و همکاران (۱۹۷۹) را برای نسخه‌های تکمیل‌شده این الگوریتمها ببینید. در برنامه‌های LINPACK اطلاعاتی در مورد شرط یا پایداری مسئله $Ax = b$ و دقت جواب محاسبه شده به دست می‌آید.

یک جنبه مهم برنامه‌های LINPACK استفاده از زیرروالهای اساسی جبر خطی^۱ (BLAS) است. این زیر برنامه‌ها عملیاتی ساده از قبیل حاصلضرب نقطه‌ای دو بردار، یا جمع ضربی از یک بردار با بردار دیگر را انجام می‌دهند. برنامه‌های LINPACK این BLASها را به جای خیلی از حلقه‌های داخلی یک روش به‌کار می‌برند. در صورت تمایل BLAS را می‌توان برای هر رایانه اپتیم نمود؛ بنابراین، با حفظ استقلال برنامه منبع از ماشین می‌توان اجرای برنامه‌های اصلی LINPACK را نیز بهبود بخشید. برای یک بحث کاملتر در مورد BLAS به لاوسن^۲ و همکاران (۱۹۷۹) مراجعه کنید.

۳.۸ صورت‌های دیگر حذف گاوسی

صورت‌های متفاوت زیادی از حذف گاوسی وجود دارند. بعضی صورت‌ها، اصلاحات یا ساده‌سازیهایی هستند که براساس ویژگی‌های خاص رده‌هایی از ماتریسها صورت می‌گیرند، مثل، ماتریسهای متقارن، ماتریسهای معین مثبت. صورت‌های دیگر، راه‌های بازنویسی حذف گاوسی، به شکل فشرده‌تر، گاهی برای استفاده از فن‌های بخصوصی برای کاهش خطا هستند. ما فقط بعضی از این صورت‌ها را بررسی می‌کنیم و سپس برای بقیه، مراجعی ذکر خواهیم کرد.

روش گاوس-یوردان این روند عیناً از نظر حذف معمولی از جمله امکان استفاده از محورگیری و مقیاس‌دهی است. تفاوت این روند در حذف مجهول در معادلات بالای قطر و معادلات پایین قطر است. در مرحله k ام الگوریتم حذف، عنصر محور را مانند قبل انتخاب می‌نماییم. سپس تعریف می‌کنیم

$$a_{kj}^{(k+1)} = \frac{a_{kj}^{(k)}}{a_{kk}^{(k)}} \quad j = k, \dots, n$$

$$b_k^{(k+1)} = \frac{b_k^{(k)}}{a_{kk}^{(k)}}$$

مجهول x_k را در معادلات بالا و پایین معادله k ام حذف می‌نماییم. تعریف می‌کنیم

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - a_{ik}^{(k)} a_{kj}^{(k+1)}$$

$$b_i^{(k+1)} = b_i^{(k)} - a_{ik}^{(k)} b_k^{(k+1)} \quad (1.3.8)$$

برای $j = k, \dots, n$ و $i = 1, \dots, n$ و $i \neq k$. روش گاوس-یوردان هم‌ارز با استفاده از

صورت ساده شده پلکانی - سطری در کتابهای جبرخطی است. [برای مثال، صفحات ۸-۹، آنتون (۱۹۸۴) را ببینید].

در این روند، ماتریس افزوده $[A | b]$ به $[I | b^{(n)}]$ بدل می‌شود، به طوری که در پایان روند حذف پیشین $x = b^{(n)}$. حل $Ax = b$ با این روش به تعداد

$$\frac{n(n-1)^2}{2} \doteq \frac{n^3}{2} \quad (۲.۳.۸)$$

ضرب و تقسیم نیاز دارد. این تعداد ۵۰ درصد بیشتر از روش حذف معمولی است؛ در نتیجه برای حل دستگاههای خطی معمولاً نباید از روش گاوس - یوردان استفاده کرد. ولی، این روش را می‌توان برای تولید یک برنامه وارونیابی ماتریس که از مینیمم حافظه استفاده می‌کند، به کار برد. با بهره‌گیری از مزیت ویژه ساختار خاص طرف راست معادله $AX = I$ ، با روش گاوس - یوردان می‌توان جواب $X = A^{-1}$ را فقط با n محل حافظه اضافی، به جای n^2 محل حافظه اضافی معمولی، محاسبه کرد. محورگیری جزئی و مقیاس‌دهی ضمنی را کماکان می‌توان به کار برد.

روشهای فشرده ممکن است مستقیماً از ماتریس A به تجزیه LU ی آن رسید و این کار را می‌توان با محورگیری جزئی و مقیاس‌دهی توأم کرد. اگر امکان محورگیری را موقتاً ندیده بگیریم، نتیجه

$$A = LU \quad (۳.۳.۸)$$

ما را به مجموعه‌ای از فرمول‌های بازگشتی برای عناصر L و U می‌رساند. چنانچه فقط بخواهیم که L و U به ترتیب پایین مثلثی و بالامثلثی باشند، یک نایکتایی در انتخاب L و U وجود خواهد داشت. اگر A ناتکین باشد، و اگر دو تجزیه داشته باشیم

$$A = L_1 U_1 = L_2 U_2 \quad (۴.۳.۸)$$

آنگاه

$$L_2^{-1} L_1 = U_2 U_1^{-1} \quad (۵.۳.۸)$$

معکوس و حاصلضربهای ماتریسهای پایین مثلثی بازهم پایین مثلثی‌اند و همینطور است برای ماتریسهای بالامثلثی. طرفهای چپ و راست (۵.۳.۸) به ترتیب پایین مثلثی و بالامثلثی هستند. بنابراین باید با یک ماتریس قطری، که آن را D می‌نامیم، برابر باشند، و

$$L_1 = L_2 D \quad U_1 = D^{-1} U_2 \quad (۶.۳.۸)$$

انتخاب D مستقیماً در گرو انتخاب عناصر قطری L و U است و وقتی آنها انتخاب شدند، D به طور یکتا معین می‌شود.

اگر بخواهند تمام عناصر قطری L برابر یک باشند، تجزیه $A = LU$ همان تجزیه‌ای است که در حذف گاوسی، در بخش ۱.۸ به دست آمده است. روش فشردهٔ متناظر، فرمولهای صریحی برای l_{ij} و u_{ij} به دست می‌دهد که به روش دولیتل^۱ معروف شده است. اگر بخواهیم تمام عناصر قطری U یک باشند، روش فشردهٔ متناظر برای محاسبهٔ $A = LU$ روش کروت^۲ نامیده می‌شود. فقط یک ماتریس قطری ضربی آن را از روش دولیتل متمایز می‌سازد. برای الگوریتمی که از الگوریتم کروت برای تجزیهٔ (۳.۳.۸) با محورگیری جزئی و مقیاس‌دهی ضمنی استفاده می‌کند، برنامهٔ *unsymdet* ویلکینسن و راینس (۱۹۷۱، صص ۹۳-۱۱۰) را ببینید. در بعضی موردها، روش کروت نسبت به روش معمولی دولیتل ارجحیت دارد.

امتیاز اصلی فرمول فشرده آن است که عناصر l_{ij} و u_{ij} همه در حاصلضربهای داخلی ظاهر می‌شوند، همان‌گونه که در فرمولهای (۱۴.۳.۸) - (۱۵.۳.۸) در تجزیهٔ ماتریس متقارن معین مثبت، در زیر دیده می‌شود. این حاصلضربهای داخلی و شاید تقسیم پایانی را می‌توان با استفاده از حساب با دقت مضاعف گردآوری کرد و سپس با یک دقت ساده گرد کرد. این طریق محاسبهٔ حاصلضربهای داخلی در فصل یک، قبل از فرمول خطای (۱۹.۵.۱) مورد بحث قرار گرفته است. با این استفادهٔ محدود از حساب با دقت مضاعف، می‌توان دقت عاملهای L و U را افزایش داد و این کار با روش حذف معمولی ممکن نخواهد بود مگر آنکه تمام عملیات و ذخیره‌سازی با حساب دقت مضاعف انجام شوند [برای بحث کامل این روشهای فشرده ویلکینسن (۱۹۶۵، صص ۱۲۸-۲۲۱) و گلوب و ون‌لون (۱۹۸۳، بخش ۱۰۵) را ببینید].

روش چولسکی^۳ گیریم A یک ماتریس متقارن، معین مثبت از مرتبهٔ n باشد. ماتریس A معین مثبت است اگر به ازای تمام مقادیر $x \in \mathbf{R}^n$ و $x \neq 0$

$$(Ax, x) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j > 0 \quad (۷.۳.۸)$$

بعضی از ویژگیهای ماتریسهای معین مثبت در مسألهٔ ۱۴ فصل ۷ و مسائل ۹ و ۱۱ و ۱۲ این فصل داده شده‌اند. ماتریسهای متقارن معین مثبت در کاربردهای متنوع و گسترده‌ای ظاهر می‌شوند.

برای یک چنین ماتریس A ، تجزیه به عوامل بسیار مناسبی وجود دارد که می‌توان بدون احتیاج به محورگیری و مقیاس‌دهی آن را انجام داد. این تجزیه روش چولسکی خوانده می‌شود،

که می‌گوئید می‌توانیم یک ماتریس حقیقی پایین مثلثی L پیدا کنیم به طوری که

$$A = LL^T \quad (۸.۳.۸)$$

در این روش به جای n^2 محل حافظه برای L فقط به $\frac{1}{2}n(n+1)$ محل حافظه نیاز است و تعداد عملیات، به جای $\frac{1}{2}n^3$ در حدود $\frac{1}{6}n^3$ است، که در تجزیه معمولی مورد نیاز بود.

برای اثبات امکان (۸.۳.۸)، یک راه به دست آوردن L بر پایه استقرا را بیان می‌کنیم. فرض می‌کنیم که نتیجه برای ماتریسهای متقارن معین مثبت از مرتبه ناپزرگتر از $n-1$ درست باشد. نشان می‌دهیم که نتیجه برای تمام این‌گونه ماتریسهای مرتبه n نیز درست است. ماتریس مطلوب L از مرتبه n را به صورت

$$L = \begin{bmatrix} \hat{L} & \circ \\ \gamma^T & x \end{bmatrix}$$

می‌نویسیم که \hat{L} یک ماتریس مربعی از مرتبه $n-1$ است، $\gamma \in \mathbb{R}^{n-1}$ و x یک اسکالر است. L باید طوری انتخاب شود که در $A = LL^T$ صدق کند:

$$\begin{bmatrix} \hat{L} & \circ \\ \gamma^T & x \end{bmatrix} \begin{bmatrix} \hat{L}^T & \gamma \\ \circ & x \end{bmatrix} = A = \begin{bmatrix} \hat{A} & c \\ c^T & d \end{bmatrix} \quad (۹.۳.۸)$$

که \hat{A} از مرتبه $n-1$ ، $c \in \mathbb{R}^{n-1}$ و $d = a_{nn}$ حقیقی است. چون (۷.۳.۸) برای A درست است x_n را صفر می‌گیریم، تا گزاره مشابهی برای \hat{A} به دست آوریم که نشان دهد \hat{A} نیز معین مثبت و متقارن است. به علاوه، $d > 0$ ، زیرا کافی است در (۷.۳.۸) بگیریم $x_1 = x_2 = \dots = x_{n-1} = 0$ و $x_n = 1$. از ضرب در (۹.۳.۸)، \hat{L} را با فرض استقرا که باید در

$$\hat{L}\hat{L}^T = \hat{A} \quad (۱۰.۳.۸)$$

صدق کند انتخاب می‌کنیم، سپس γ را از حل

$$\hat{L}\gamma = c \quad (۱۱.۳.۸)$$

به دست می‌آوریم. γ به دست می‌آید زیرا از $[\det[\hat{L}]]^2 = \det(\hat{A})$ نتیجه می‌شود که \hat{L} نانتکین است. بالاخره، x باید در رابطه زیر صدق کند

$$\gamma^T \gamma + x^2 = d \quad (۱۲.۳.۸)$$

برای آنکه به مثبت بودن x^T پی ببریم، دترمینان طرفین (۹.۳.۸) را تشکیل می‌دهیم و نتیجه می‌گیریم

$$[\det(\tilde{L})]^2 x^T = \det(A) \quad (۱۳.۳.۸)$$

چون دترمینان A حاصلضرب ویژه مقادیرهای A است، و چون ویژه مقادیرهای ماتریس متقارن معین مثبت، مثبت‌اند (مسئله ۱۴ فصل ۷ را ببینید)، $\det(A)$ مثبت است. همچنین با فرض استقرای، \tilde{L} حقیقی است. بنابراین x^T در (۱۳.۳.۸) مثبت است، و x را ریشهٔ دوم مثبت آن می‌گیریم. چون رابطه (۸.۳.۸) به‌طور نمایان برای ماتریس از مرتبهٔ $n = ۱$ درست است، پس اثبات تجزیهٔ (۸.۳.۸) کامل شده است. برای یک روش دیگر، بخش ۲.۵، گلوب و ون لون (۱۹۸۳) را ببینید. یک راه ساختن عملی L را می‌توان بر مبنای (۹.۳.۸) - (۱۲.۳.۸) انجام داد، ولی ما روشی را ارائه می‌دهیم که بر پایهٔ پیدا کردن مستقیم عناصر L نهاده شده است. گیریم $L = [l_{ij}]$ ، $l_{ij} = 0$ برای $i > j$. ساختن L را بدین شکل آغاز می‌کنیم که سطر اول L را در ستون اول L^T ضرب می‌نماییم تا به دست آوریم

$$l_{11}^2 = a_{11}$$

چون A معین مثبت است، $a_{11} > 0$ و $l_{11} = \sqrt{a_{11}}$. سطر دوم L را در دو ستون اول L^T ضرب می‌کنیم تا به دست آوریم

$$l_{21}l_{11} = a_{21} \quad l_{21}^2 + l_{22}^2 = a_{22}$$

باز هم می‌توانیم دستگاه فوق را نسبت به مجهولهای l_{21} و l_{22} حل کنیم. به‌طور کلی^۱ برای $i = 1, 2, \dots, n$

$$l_{ij} = \frac{a_{ij} - \sum_{k=1}^{j-1} l_{ik}l_{jk}}{l_{jj}} \quad j = 1, \dots, i-1 \quad (۱۴.۳.۸)$$

$$l_{ii} = \left[a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2 \right]^{1/2} \quad (۱۵.۳.۸)$$

مقدار زیر رادیکال همان جملهٔ x^T در دستور به‌دست آمده قبلی (۱۲.۳.۸) است و l_{ii} حقیقی و مثبت است. برای برنامه‌هایی که روش چولسکی را اجرا می‌کنند،

۱. توجه کنید به‌ازای $i = 1, 2, \dots, n$ ، $l_{i1} = \frac{a_{i1}}{l_{11}}$.

دونگارا و همکاران (۱۹۷۹، فصل ۳) و ویلکینسن و راینش^۱ (۱۹۷۱، صص ۱۰-۳۰) را ببینید.

به حاصلضربهای داخلی در (۱۴.۳.۸) و (۱۵.۳.۸) توجه نمایید. این حاصلضربهای داخلی را می‌توان با دقت مضاعف محاسبه کرد. و تعداد خطاهای گردکردن را مینیمم نمود، در این صورت خطا در عناصر l_{ij} خیلی کمتر از موقعی خواهد بود که با دقت ساده حساب شوند. همچنین توجه نمایید که عناصر L نسبت به A همچنان کراندار باقی می‌مانند، زیرا (۱۵.۳.۸) با استفاده از رابطه زیر یک کران برای عناصر سطر i ام به دست می‌دهد

$$l_{i1}^2 + \dots + l_{ii}^2 = a_{ii} \quad (۱۶.۳.۸)$$

با یک اصلاح جزئی در (۸.۳.۸) می‌توان از به‌کار بردن ریشه‌های دوم در (۱۵.۳.۸) ی روش چولسکی دوری کرد. یک ماتریس قطری D و یک ماتریس پایین‌مثلثی \tilde{L} که عناصر قطری آن همه ۱ باشند، پیدا می‌کنیم به طوری که

$$A = \tilde{L}D\tilde{L}^T \quad (۱۷.۳.۸)$$

این تجزیه به عوامل را می‌توان با همان تعداد عملیات روش چولسکی، در حدود $\frac{1}{2}n^3$ بدون هیچ ریشه‌گیری انجام داد. برای بحث بیشتر و یک برنامه، ویلکینسن و راینش (۱۹۷۱، صص ۱۰-۳۰) را ببینید.

مثال ماتریس مرتبه ۳ ی هیلبرت در زیر را در نظر می‌گیریم

$$A = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \end{bmatrix} \quad (۱۸.۳.۸)$$

در تجزیه چولسکی،

$$L = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{\sqrt{3}} & 0 \\ \frac{1}{3} & \frac{1}{2\sqrt{3}} & \frac{1}{\sqrt{5}} \end{bmatrix}$$

و در (۱۷.۳.۸)،

$$\bar{L} = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ \frac{1}{3} & 1 & 1 \end{bmatrix} \quad D = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \frac{1}{12} & 0 \\ 0 & 0 & \frac{1}{180} \end{bmatrix}$$

برای بسیاری از دستگاههای خطی در کاربردها، ماتریس ضرایب A نواری است یعنی برای یک مقدار کوچک مثبت m

$$a_{ij} = 0 \quad |i - j| > m \quad \text{اگر} \quad (۱۹.۳.۸)$$

الگوریتمهای قبلی در چنین حالتی ساده خواهند شد و صرفه جویی قابل ملاحظه‌ای در زمان محاسبات به وجود خواهد آمد. برای چنین الگوریتم‌هایی وقتی A متقارن و معین مثبت باشد برنامه‌های LINPACK در دونگارا و همکاران (۱۹۷۹، فصل ۴) را ببینید.

دستگاههای سه قطری ماتریس $A = [a_{ij}]$ سه قطری است اگر

$$a_{ij} = 0 \quad |i - j| > 1 \quad \text{برای} \quad (۲۰.۳.۸)$$

که شکل زیر را به دست می‌دهد

$$A = \begin{bmatrix} a_1 & c_1 & 0 & 0 & \dots & 0 \\ b_2 & a_2 & c_2 & 0 & & \vdots \\ 0 & b_3 & a_3 & c_3 & & \\ & & & \ddots & & \\ \vdots & & & & b_{n-1} & a_{n-1} & c_{n-1} \\ 0 & \dots & & 0 & b_n & a_n \end{bmatrix} \quad (۲۱.۳.۸)$$

ماتریسهای سه قطری در کاربردهای متنوعی ظاهر می‌شوند. دستگاه خطی (۲۲.۷.۳) در توابع برازا در بخش ۷.۳ از فصل ۳ را به خاطر آورید. به علاوه، بسیاری از روشهای عددی حل مسائل مقدار مرزی در معادلات دیفرانسیل معمولی و معادلات دیفرانسیل جزئی متضمن حل دستگاههای سه قطری‌اند. تقریباً تمام این کاربردها ماتریسهای سه قطری به دست می‌دهند که تجزیه LU ی

آنها را می‌توان بدون محورگیری انجام داد و در نتیجه برای آنها افزایش زیاد خطا پیش نمی‌آید. فرض‌های دقیق برای A در قضیه ۲.۸ در ذیل داده شده است.

با توجه به تجزیه $A = LU$ بدون محورگیری، می‌بینیم که بیشتر عناصر L و U صفرند و برای تجزیه به فرمول کلی زیر می‌رسیم:

$$A = LU = \begin{bmatrix} \alpha_1 & \circ & \dots & \circ \\ b_2 & \alpha_2 & \circ & \vdots \\ \circ & b_2 & \alpha_2 & \\ \vdots & & \ddots & \\ \circ & \dots & b_n & \alpha_n \end{bmatrix} \begin{bmatrix} 1 & \gamma_1 & \circ & \dots & \circ \\ \circ & 1 & \gamma_2 & \circ & \vdots \\ \vdots & & \ddots & \ddots & \\ \circ & \dots & & 1 & \gamma_{n-1} \\ \circ & & & \circ & 1 \end{bmatrix}$$

می‌توانیم از ضرب دو ماتریس، یک راه بازگشتی، برای محاسبه $\{\alpha_i\}$ و $\{\gamma_i\}$ به‌دست آوریم

$$\alpha_1 = a_1 \quad \alpha_1 \gamma_1 = c_1$$

$$a_i = \alpha_i + b_i \gamma_{i-1} \quad i = 2, \dots, n \quad (22.3.8)$$

$$\alpha_i \gamma_i = c_i \quad i = 2, 3, \dots, n-1$$

این معادلات را می‌توان حل کرد و به‌دست آورد

$$\alpha_1 = a_1 \quad \gamma_1 = \frac{c_1}{\alpha_1}$$

$$a_i = a_i - b_i \gamma_{i-1} \quad \gamma_i = \frac{c_i}{\alpha_i} \quad i = 2, 3, \dots, n-1 \quad (23.3.8)$$

$$\alpha_n = a_n - b_n \gamma_{n-1}$$

برای حل $LUx = f$ ، می‌گیریم $Ux = z$ و $Lz = f$. پس

$$z_1 = \frac{f_1}{\alpha_1} \quad z_i = \frac{f_i - b_i z_{i-1}}{\alpha_i} \quad i = 2, 3, \dots, n \quad (24.3.8)$$

$$x_n = z_n \quad x_i = z_i - \gamma_i x_{i+1} \quad i = n-1, n-2, \dots, 1$$

مقادیر ثابت در (۲۳.۳.۸) را می‌توان برای استفاده بعدی در حل دستگاه خطی $Ax = f$ ، برای هر f طرف راست دستگاه که خواسته شده باشد، ذخیره نمود.

با شمردن فقط ضربها و تقسیمها، تعداد عملیات در محاسبه L و U برابر $2n - 2$ است؛ برای

حل $Ax = f$ ، $2 - 3n$ عمل دیگر لازم است. بنابراین به $5n - 4$ عمل برای حل $Ax = f$

در اولین بار نیاز داریم، برای هر f دیگر، با همان A ، فقط $2 - 3n$ عمل لازم خواهد بود. این روشی بسیار سریع است. برای توضیح، توجه نمایید که A^{-1} معمولاً چگال است و بیشتر درایه‌های آن ناصفرند؛ بنابراین محاسبه $x = A^{-1}f$ به n^2 عمل نیاز دارد. در بسیاری از کاربردها، n ممکن است بزرگتر از 1000 باشد و بنابراین یک صرفه‌جویی قابل توجهی در استفاده از (۲۳.۳.۸) - (۲۴.۳.۸) در مقایسه با سایر روشهای حل دستگاه، وجود دارد.

برای توجیه تجزیه قبلی A ، بخصوص نشان دادن اینکه تمام ضرایب $\alpha_i \neq 0$ ، قضیه زیر را داریم:

قضیه ۲.۸ فرض می‌کنیم که ضرایب $\{a_i, b_i, c_i\}$ در (۲۱.۳.۸) در شرایط زیر صدق می‌کنند:

$$|a_1| > |c_1| > 0 \quad 1$$

$$i = 2, \dots, n-2, b_i, c_i \neq 0, |a_i| \geq |b_i| + |c_i| \quad 2$$

$$|a_n| > |b_n| > 0 \quad 3$$

در این صورت A نانگین است و

$$|\gamma_i| < 1 \quad i = 1, \dots, n-1$$

$$|a_i| - |b_i| < |\alpha_i| < |a_i| + |b_i| \quad i = 2, \dots, n$$

برهان برای اثبات، (صفحه ۵۷) آیزکسون و کالر (۱۹۶۶) را ببینید. ملاحظه می‌کنید که کران آخری نشان می‌دهد که $|c_i| > |\alpha_i|$ ، برای $i = 2, \dots, n-2$. بنابراین ضرایب L و U کراندار باقی می‌مانند، و هیچ مقسوم‌علیه‌ی که تقریباً صفر باشد به کار نمی‌رود مگر به دلیل خطای گردکردن. شرط مخالف صفر بودن c_i و b_i اساسی نیست. برای مثال، اگر یکی از b_i ها مساوی صفر باشد دستگاه خطی را می‌توان به دو دستگاه جدید تجزیه کرد، یکی از مرتبه $i-1$ و دیگری از مرتبه $i+1-n$. برای مثال اگر

$$A = \begin{bmatrix} a_1 & c_1 & 0 & 0 \\ b_2 & a_2 & c_2 & 0 \\ 0 & 0 & a_3 & c_3 \\ 0 & 0 & b_4 & a_4 \end{bmatrix}$$

آنگاه $Ax = f$ را با تبدیل به دو دستگاه خطی زیر حل می‌کنیم

$$\begin{bmatrix} a_3 & c_3 \\ b_4 & a_4 \end{bmatrix} \begin{bmatrix} x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} f_3 \\ f_4 \end{bmatrix} \quad \begin{bmatrix} a_1 & c_1 \\ b_2 & a_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 - c_2 x_3 \end{bmatrix}$$

و اثبات کامل می‌شود.

مثال ماتریس ضرایب درونیایی برآزای (۲۲.۷.۳) از فصل ۳ را در نظر می‌گیریم. در آن ماتریس h_i را ثابت می‌گیریم و سپس از $h/6$ در هر سطر فاکتور می‌گیریم. اگر توجه خود را به ماتریس مرتبه چهار محدود کنیم، ماتریس حاصل چنین خواهد شد:

$$A = \begin{bmatrix} 2 & 1 & 0 & 0 \\ 1 & 4 & 1 & 0 \\ 0 & 1 & 4 & 1 \\ 0 & 0 & 1 & 2 \end{bmatrix}$$

با استفاده از روش (۲۳.۳.۸)، تجزیه LU به شکل زیر خواهد بود

$$L = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 1 & \frac{7}{2} & 0 & 0 \\ 0 & 1 & \frac{26}{7} & 0 \\ 0 & 0 & 1 & \frac{45}{26} \end{bmatrix} \quad U = \begin{bmatrix} 1 & \frac{1}{2} & 0 & 0 \\ 0 & 1 & \frac{2}{7} & 0 \\ 0 & 0 & 1 & \frac{7}{26} \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

که مثال را کامل می‌کند. و باید اشاره کرد که حل مسأله درونیایی برآزای درجه سوم که در بخش ۷.۳ فصل ۳، توضیح داده شد مشکل نیست.

۴.۸ تخلیل خطا

تخلیل خطای روشهای حل $Ax = b$ را با بررسی پایداری جواب x نسبت به اختلالات کوچک b ی سمت راست معادله، آغاز می‌کنیم. طرح کلی بخش ۶.۱ فصل یک را دنبال، و به‌ویژه، ضریب وضعیت در (۶.۶.۱) را مطالعه می‌کنیم.

گیریم $Ax = b$ ، از مرتبه n ، به‌طور یکتا حلپذیر باشد، و جواب مسأله اختلال یافته زیر را در نظر می‌گیریم

$$A\tilde{x} = b + r \quad (1.4.8)$$

فرض می‌کنیم $e = \tilde{x} - x$ ، و از تفاضل $Ax = b$ از رابطه بالا به دست می‌آوریم

$$Ae = r \quad e = A^{-1}r \quad (2.4.8)$$

برای بررسی پایداری $Ax = b$ ، مانند (۶.۶.۱)، می‌خواهیم کمیت

$$\frac{\|e\|}{\|x\|} \div \frac{\|r\|}{\|b\|} \quad (۳.۴.۸)$$

را کراندار کنیم وقتی r تمام عناصری از \mathbf{R}^n را که نسبت به b کوچک‌اند، اختیار می‌کند. از (۲.۴.۸) نرم می‌گیریم تا حاصل شود

$$\|r\| \leq \|A\| \|e\| \quad \|e\| \leq \|A^{-1}\| \|r\|$$

نابرابری اول را بر $\|A\| \|x\|$ و نابرابری دوم را بر $\|x\|$ تقسیم می‌نماییم، پس

$$\frac{\|r\|}{\|A\| \|x\|} \leq \frac{\|e\|}{\|x\|} < \frac{\|A^{-1}\| \|r\|}{\|x\|}$$

نرم ماتریسی در اینجا، همان نرم عملگر ماتریسی است که توسط نرم برداری القا شده است. با استفاده از کرانهای

$$\|b\| \leq \|A\| \|x\| \quad \|x\| \leq \|A^{-1}\| \|b\|$$

به دست می‌آوریم

$$\frac{1}{\|A\| \|A^{-1}\|} \cdot \frac{\|r\|}{\|b\|} \leq \frac{\|e\|}{\|x\|} \leq \|A\| \|A^{-1}\| \cdot \frac{\|r\|}{\|b\|} \quad (۴.۴.۸)$$

با توجه به (۳.۴.۸)، این نتیجه معرفی ضریب وضعیت A :

$$\text{Cond}(A) = \|A\| \|A^{-1}\| \quad (۵.۴.۸)$$

را توجیه می‌کند. برای هر A ی داده شده، انتخابهایی برای b و r وجود دارند به طوری که یکی از نابرابریها در (۴.۴.۸) را به برابری تبدیل کنند. این دلیل دیگری برای معرفی ضریب وضعیت A است وقتی (۳.۴.۸) در نظر گرفته می‌شود. اثبات را به مسئله ۲۰ واگذار می‌نماییم.

کمیت ضریب وضعیت A با نرمی که به کار برده شود تغییر می‌نماید، ولی همیشه دارای کران پایین ۱ است، زیرا

$$1 \leq \|I\| = \|AA^{-1}\| \leq \|A\| \|A^{-1}\| = \text{cond}(A)$$

اگر ضریب وضعیت نزدیک ۱ باشد، با توجه به (۴.۴.۸) معلوم می‌شود که اختلال نسبی کوچک در b به اختلال نسبی کوچک در جواب x می‌انجامد. ولی اگر $\text{cond}(A)$ بزرگ باشد، (۴.۴.۸)

بیان می‌کند که ممکن است اختلالات نسبی کوچکی برای b وجود داشته باشند که اختلالات نسبی بزرگی در x ایجاد نمایند.

چون (۵.۴.۸) با انتخاب نرم تغییر می‌کند، گاهی تعریف دیگری برای ضریب وضعیت به کار می‌بریم که مستقل از نرم است. از قضیه ۸.۷، فصل ۷،

$$\text{cond}(A) \geq r_\sigma(A)r_\sigma(A^{-1})$$

چون ویژه‌مقدارهای A^{-1} ، معکوس ویژه‌مقدارهای A هستند، نتیجه زیر را خواهیم داشت

$$\text{cond}(A) \geq \frac{\text{Max}_{\lambda \in \sigma(A)} |\lambda|}{\text{Min}_{\lambda \in \sigma(A)} |\lambda|} \equiv \text{cond}(A)_* \quad (۶.۴.۸)$$

که $\sigma(A)$ معرف مجموعه همه ویژه‌مقدارهای A است.

مثال دستگاه خطی زیر را در نظر می‌گیریم

$$7x_1 + 10x_2 = b_1$$

$$5x_1 + 7x_2 = b_2 \quad (۷.۴.۸)$$

برای ماتریس ضرایب

$$A = \begin{bmatrix} 7 & 10 \\ 5 & 7 \end{bmatrix}, \quad A^{-1} = \begin{bmatrix} -7 & 10 \\ 5 & -7 \end{bmatrix}$$

گیریم ضریب وضعیت در (۵.۴.۸) وقتی از نرم ماتریس $\|\cdot\|_p$ تولید شده است. با $\text{cond}(A)_p$ نمایش داده شود در این مثال

$$\text{cond}(A)_1 = \text{cond}(A)_\infty = (17)(17) = 289$$

$$\text{cond}(A)_2 \doteq 223 \quad \text{cond}(A)_* \doteq 198$$

این ضرایب وضعیت همگی بیان می‌کنند که (۷.۴.۸) ممکن است نسبت به تغییرات b طرف راست معادله حساس باشد. برای نشان دادن این امکان، حالت خاص زیر را در نظر می‌گیریم

$$7x_1 + 10x_2 = 1$$

$$5x_1 + 7x_2 = 0.7$$

که دارای جواب

$$x_1 = 0, \quad x_2 = 0.1$$

است. برای دستگاه اختلال یافته

$$7\tilde{x}_1 + 10\tilde{x}_2 = 1.01$$

$$5\tilde{x}_1 + 7\tilde{x}_2 = 0.69$$

جواب چنین است

$$\tilde{x}_1 = -0.17 \quad \tilde{x}_2 = 0.22$$

تغییرات نسبی در x ، در مقایسه با تغییرات نسبی در b ی طرف راست معادله، بسیار بزرگ‌اند.

یک دستگاه خطی که جوابش نسبت به تغییرات نسبی جزئی b ی طرف راست معادله ناپایدار باشد بدوضع خوانده می‌شود. دستگاه قبلی (۷.۴.۸)، تا حدی بدوضع است، به‌ویژه وقتی که فقط در حل آن از حساب سه یا چهار رقم با ممیز شناور استفاده شده باشد. ضرایب وضعیت $\text{cond}(A)$ و $\text{cond}(A)_*$ شاخصهای نسبتاً خوبی برای نشان‌دادن بدوضع هستند. از آنجا که این اختلالها با عامل 10^0 بزرگ می‌شوند. احتمال دارد که دقت در جواب با یک رقم کمتر به دست آید.

به‌طور کلی، اگر $\text{cond}(A)_*$ بزرگ باشد، برای بعضی مقادیر b دستگاه $Ax = b$ نسبت به تغییرات r در b ، خیلی حساس خواهد بود. گیریم λ_l و λ_u ویژه مقدارهای A باشند به‌طوری که

$$|\lambda_l| = \text{Min}_{\lambda \in \sigma(A)} |\lambda| \quad |\lambda_u| = \text{Max}_{\lambda \in \sigma(A)} |\lambda|$$

و بنابراین

$$\text{cond}(A)_* = \left| \frac{\lambda_u}{\lambda_l} \right| \quad (۸.۴.۸)$$

گیریم x_l و x_u ویژه بردارهای متناظر باشند با $\|x_u\|_\infty = \|x_l\|_\infty = 1$. پس

$$Ax = \lambda_u x_u$$

دارای جواب $x = x_u$ است. و دستگاه

$$A\tilde{x} = \lambda_u x_u + \lambda_l x_l = \lambda_u \left[x_u + \frac{1}{\text{cond}(A)_*} x_l \right]$$

$$\tilde{x} = x_u + x_l$$

اگر $\text{cond}(A)_*$ بزرگ باشد، طرف راست فقط یک اختلال نسبی کوچک دارد؛

$$\frac{\|r\|_\infty}{\|b\|_\infty} = \frac{1}{\text{cond}(A)_*} \quad (۹.۴.۸)$$

ولی جواب، اختلال نسبی بسیار بزرگتر زیر را دارد

$$\frac{\|\tilde{x} - x\|_\infty}{\|x\|_\infty} = \frac{\|x_l\|_\infty}{\|x_u\|_\infty} = 1 \quad (۱۰.۴.۸)$$

دستگاههایی وجود دارند که در عمل بدوضع نیستند ولی ضرایب وضعیت قبلی برای آنها کاملاً بزرگاند. برای مثال در ماتریس

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 10^{-10} \end{bmatrix}$$

ضرایب وضعیت $\text{cond}(A)_p$ و $\text{cond}(A)_*$ همگی برابر 10^{10} هستند. ولی معمولاً این ماتریس بدوضع تلقی نمی‌شود. مشکل در استفاده از نرم‌ها برای اندازه‌گیری تغییرات بردار است، به جای آنکه به هر مؤلفه جداگانه توجه شود. اگر مقیاس‌دهی در ماتریس ضرایب و بردار مجهول اجرا شود، این مسأله معمولاً پیش نمی‌آید و در آن صورت ضرایب وضعیت معمولاً پیشگوی دقیقی برای بدوضع می‌شوند. در یک توجیه نهایی برای استفاده از $\text{cond}(A)$ به عنوان ضریب وضعیت، قضیه زیر را می‌آوریم.

قضیه ۳.۸ (گاستینل^۱) گیریم A ماتریس نانکینی از مرتبه n باشد و $\| \cdot \|$ معرف یک نرم ماتریس عملگر در این صورت

$$\frac{1}{\text{cond}(A)} = \text{Min} \left\{ \frac{\|A - B\|}{\|A\|} \mid B \text{ یک ماتریس تکین} \right\} \quad (۱۱.۴.۸)$$

که $\text{cond}(A)$ همان است که در (۵.۴.۸) تعریف شده است.

برهان کاهان^۲ (۱۹۶۶، ص ۷۷۵) را ببینید.

این قضیه بیان می‌کند که از لحاظ خطای نسبی، A را می‌توان به خوبی با یک ماتریس تکین تقریب زد اگر و تنها اگر $\text{cond}(A)$ خیلی بزرگ باشد. و از دید ما، ماتریس تکین B در حدّ اعلاّی

جدول ۲.۸ ضرایب وضعیت برای ماتریس هیلبرت

n	$\text{cond}(H_n)_*$	n	$\text{cond}(H_n)_*$
۳	$5,24E + 2$	۷	$4,75E + 8$
۴	$1,55E + 4$	۸	$1,53E + 10$
۵	$4,77E + 5$	۹	$4,93E + 11$
۶	$1,50E + 7$	۱۰	$1,60E + 13$

بدو وضعی است. با ویژه بردار متناظر با ویژه مقدار $\lambda = 0$ اختلالاتی ناصفر برای جواب وجود دارند که متناظر با اختلال صفر در b ، طرف راست معادله هستند. مهمتر از آن، مقادیری از b وجود دارند که برای آنها $Bx = b$ دیگر حلپذیر نیست.

ماتریس هیلبرت ماتریس هیلبرت از مرتبه n چنین تعریف می شود

$$H_n = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \dots & \frac{1}{n} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & & \frac{1}{n+1} \\ \vdots & & & & \vdots \\ \frac{1}{n} & \frac{1}{n+1} & \dots & & \frac{1}{2n-1} \end{bmatrix} \quad (12.4.8)$$

این ماتریس به طور طبیعی در حل مسأله تقریب با کمترین مربعات پیوسته ظاهر می شود. اصل آن در نزدیک به انتهای بخش ۳.۴ از فصل ۴ آمده، همراه با دستگاه خطی حاصل که در (۱۴.۳.۴) داده شده است. همان طور که در بخش ۳.۴ اشاره شد و به دنبال (۹.۶.۱) در بخش ۴.۱ فصل ۱ نشان داده شد، ماتریس هیلبرت بسیار بدوضع است و با افزایش n بدو وضعتر می شود. بدین لحاظ این، یک مثال عددی مورد پسندی، برای بررسی برنامه های حل دستگاههای معادلات خطی، و تعیین حدود کارایی برنامه هنگام برخورد با مسائل بدوضع، بوده است. جدول ۲.۸ ضرایب وضعیت $\text{cond}(H_n)_*$ را برای چند مقدار n نشان می دهد. ماتریس معکوس $[\alpha_{ij}^{(n)}] = H_n^{-1}$ به صورت صریح معلوم است:

$$\alpha_{ij}^{(n)} = \frac{(-1)^{i+j} (n+i-1)!(n+j-1)!}{(i+j-1)[(i-1)!(j-1)!]^2 (n-i)!(n-j)!} \quad 1 \leq i, j \leq n \quad (13.4.8)$$

برای اطلاع بیشتر از H_n ، از جمله یک فرمول مجانبی برای $\text{Cond}(H_n)_*$ ، گرگوری^۱ و کارنی^۲ (۱۹۶۹، صص ۳۳-۳۸ و ۶۶-۷۳) را ببینید.

اگرچه ماتریس هیلبرت به عنوان مثال زیاد به کار برده می شود، ولی باید مواظب باشیم و ببینیم که جواب درست چیست. گیریم \bar{H}_n معرف H_n پس از ورود به حساب متناهی یک رایانه باشد. برای یک برنامه وارون ماتریس، نتایج برنامه باید با \bar{H}_n^{-1} مقایسه شود، نه با H_n^{-1} ؛ این دو ماتریس وارون ممکن است کاملاً متفاوت باشند. برای مثال، اگر از حساب ممیز شناور چهار رقم اعشاری با گردکردن استفاده کنیم، آنگاه

$$\bar{H}_3 = \begin{bmatrix} 1.000 & 0.5000 & 0.3333 \\ 0.5000 & 0.3333 & 0.2500 \\ 0.3333 & 0.2500 & 0.2000 \end{bmatrix} \quad (14.4.8)$$

گردکردن فقط در بسط $\frac{1}{3}$ به شکل کسر اعشاری صورت گرفته است. در این صورت

$$H_3^{-1} = \begin{bmatrix} 9.000 & -36.00 & 30.00 \\ -36.00 & 192.0 & -180.0 \\ 30.00 & -180.0 & 180.0 \end{bmatrix} \quad (15.4.8)$$

$$\bar{H}_3^{-1} = \begin{bmatrix} 9.062 & -36.32 & 30.30 \\ -36.32 & 193.7 & -181.6 \\ 30.30 & -181.6 & 181.5 \end{bmatrix}$$

هر برنامه وارون ماتریس، که برای \bar{H}_3 به کار برده شده باشد جوابش باید با \bar{H}_3^{-1} مقایسه شود، نه با H_3^{-1} . ما به این مسأله بعداً در بخش ۵.۸ باز می گردیم.

کرانه های خطا اکنون آثار خطای گردکردن روی جواب \hat{x} از $Ax = b$ را که با حذف گاوسی به دست آمده است، در نظر می گیریم. مطلب را با ارائه قضیه ای که خطا را کراندار می کند، وقتی b و A به اندازه کوچکی تغییر یافته اند، آغاز می کنیم. این خود قضیه مفیدی است، و برای تحلیل خطای حذف گاوسی که بعداً خواهد آمد، لازم است.

قضیه ۴.۸ دستگاه $Ax = b$ را، وقتی که A ناتکین است، در نظر می گیریم. فرض می کنیم δA و δb اختلالات A و b باشند و

$$\|\delta A\| < \frac{1}{\|A^{-1}\|} \quad (16.4.8)$$

پس $A + \delta A$ ناتکین است. و اگر δx را به طور ضمنی با رابطه

$$(A + \delta A)(x + \delta x) = b + \delta b \quad (17.4.8)$$

تعریف کنیم، آنگاه

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \frac{\|\delta A\|}{\|A\|}} \cdot \left\{ \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right\} \quad (18.4.8)$$

برهان ابتدا باید توجه داشت که δA معرّف ماتریسی است که در رابطه (۱۶.۴.۸) صدق می‌کند، نه یک ثابت δ که در ماتریس A ضرب شده است، و همچنین است برای δb و δx . با استفاده از (۱۶.۴.۸) ناتیکنی $A + \delta A$ بلافاصله از قضیه ۱۲.۷ فصل ۷ نتیجه می‌شود. با توجه به (۱۱.۴.۷)،

$$\|(A + \delta A)^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \|\delta A\|} \quad (19.4.8)$$

از حل (۱۷.۴.۸) نسبت به δx و استفاده از $Ax = b$

$$(A + \delta A)\delta x + Ax + (\delta A)x = b + \delta b$$

$$\delta x = (A + \delta A)^{-1}[\delta b - (\delta A)x]$$

با استفاده از (۱۹.۴.۸) و تعریف (۵.۴.۸) برای $\text{cond}(A)$

$$\|\delta x\| \leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \frac{\|\delta A\|}{\|A\|}} \cdot \left\{ \frac{\|\delta b\|}{\|b\|} + \|x\| \frac{\|\delta A\|}{\|A\|} \right\}$$

طرفین را بر $\|x\|$ تقسیم و از $\|x\| \leq \|A\| \|b\|$ استفاده می‌کنیم تا (۱۸.۴.۸) حاصل شود. ■

تحلیل اثر خطاهای گردکردن در حذف گاوسی منسوب به ج. ه. ویلکینسن است و می‌توان آن را در ویلکینسن (۱۹۶۳، صص ۹۴-۹۹)، (۱۹۶۵، صص ۲۰۹-۲۱۶)، فورسایت و مولر (۱۹۶۷، فصل ۲۱)، و گلوب و ون لون (۱۹۸۳، فصل ۴) یافت. گیریم \hat{x} جواب محاسبه شده $Ax = b$ باشد. محاسبه مستقیم اثرات گردکردن x در هر مرحله، به عنوان وسیله‌ای برای به دست آوردن یک کران $\|x - \hat{x}\|$ ، دشوار است. گرچه نمایان نیست، ولی ساده‌تر آن است که \hat{x} و الگوریتم حذف را در نظر گرفته به عقب برگردیم و نشان دهیم که \hat{x} جواب دقیق دستگاه

$$(A + \delta A)\hat{x} = b$$

است که در آن می‌توان کرانه‌های δA را تعیین کرد. این شیوه به تحلیل پسرو خطا معروف است. در این صورت می‌توانیم قضیه قبلی ۴.۸ را برای کراندارکردن $\|x - \hat{x}\|$ به‌کار ببریم. در قضیه زیر، نرم ماتریس $\|A\|_\infty$ ، نرم سطری (۱۷.۳.۷) است که بر اثر نرم برداری $\|x\|_\infty$ القا شده است.

قضیه ۵.۸ گیریم A از مرتبه n و ناکین باشد، و فرض می‌کنیم از محورگیری جزئی یا کلی در فرایند حذف استفاده شده باشد. تعریف می‌کنیم

$$\rho = \frac{1}{\|A\|_\infty} \max_{1 \leq i, j, k \leq n} |a_{ij}^{(k)}| \quad (20.4.8)$$

گیریم u واحدگرد کردن در رایانه مورد استفاده باشد [(۱۱.۲.۱) و (۱۲.۲.۱)] را برای تعریف u ببینید، ۱. ماتریسهای L و U که با حذف گاوسی محاسبه شده‌اند در روابط زیر صدق می‌کنند:

$$LU = A + E$$

$$\|E\|_\infty \leq n^2 \rho \|A\|_\infty u \quad (21.4.8)$$

۲. جواب تقریبی \hat{x} ، جواب تقریبی $Ax = b$ که با استفاده از حذف گاوسی محاسبه شده در رابطه زیر صدق می‌کند

$$(A + \delta A)\hat{x} = b \quad (22.4.8)$$

با

$$\frac{\|\delta A\|_\infty}{\|A\|_\infty} \leq \{1.0 \cdot 1(n^2 + 3n^2)\rho u\} \quad (23.4.8)$$

۳. با استفاده از قضیه ۴.۸

$$\frac{\|x - \hat{x}\|_\infty}{\|x\|_\infty} \leq \frac{\text{cond}(A)_\infty}{1 - \text{cond}(A)_\infty \frac{\|\delta A\|_\infty}{\|A\|_\infty}} [1.0 \cdot 1(n^2 + 3n^2)\rho u] \quad (24.4.8)$$

برهان اثبات ۱ و ۲ در (فصل ۲۱) فورسایت و مولر (۱۹۶۷) داده شده است. صورتهای مختلف این قضایا در (فصل ۴) گلوب و ونلون (۱۹۸۳) آمده است. ■

به روش تجربی، کران (۲۴.۴.۸)، به علت حذف خطاهای گردکردن با اندازه و علامت متغیر، بسیار بزرگ است. به گفته ویلکینسن (۱۹۶۳، ص ۱۰۸)، یک کران تجربی بهتر در اغلب حالات

$$\frac{\|\delta A\|_\infty}{\|A\|_\infty} \leq nu \quad (25.4.8)$$

است. فرمول (۲۴.۴.۸) اهمیت اندازه $\text{cond}(A)$ را نشان می‌دهد.

کمیت ρ در کرانها را می‌توان در خلال فرایند حذف محاسبه کرد و می‌توان آن را از پیش نیز کراندار کرد. برای محورگیری کامل، یک کران از پیش تعیین شده چنین است

$$\rho \leq 1.8n^{(\ln n)/4} \quad n \geq 1$$

و حدس زده‌اند که برای یک مقدار c ، $\rho \leq cn$. برای محورگیری جزئی، یک کران از پیش تعیین شده 2^{n-1} است و مثالهای غیرعادی دالّ بر امکان این امر وجود دارد. با این همه، در تمام مطالعات تجربی تا به امروز، ρ با یک عدد نسبتاً کوچک مستقل از n کراندار شده‌است. به لحاظ اختلاف در کرانهای نظری برای ρ ، گاهی محورگیری کلی را ترجیح می‌دهند. ولی، در عمل، رفتار خطا با محورگیری جزئی به همان خوبی محورگیری کلی است. به علاوه، محورگیری کلی، نیاز به مقایسه‌های خیلی بیشتری در هر مرحله از فرایند حذف خواهد داشت. نتیجتاً محورگیری جزئی روشی است که در همه برنامه‌های جدید رایانه‌یی حذف گاوسی به‌کار می‌رود.

یکی از مهمترین نتایج تحلیل بالا این است که نشان می‌دهد حذف گاوسی یک فرایند بسیار پایدار است، فقط مشروط بر آنکه ماتریس A خیلی بدوضع نباشد. از نظر تاریخی، پژوهشگران در اوایل دهه ۱۹۵۰، از پایداری حذف گاوسی برای دستگاههای بزرگ، مثلاً $n \geq 10^6$ مطمئن نبودند، ولی این مسأله اکنون حل شده است.

اندازه مانده در جواب محاسبه شده \hat{x} ، یعنی

$$r = b - A\hat{x} \quad (26.4.8)$$

گاهی به اشتباه به اندازه خطای $x - \hat{x}$ ارتباط داده شده است. درواقع، خطای \hat{x} ممکن است بزرگ و r کوچک باشد و این حالت معمولاً در مسائل بدوضع به‌وجود می‌آید. از (۲۶.۴.۸) و $Ax = b$ ،

$$r = A(x - \hat{x})$$

$$x - \hat{x} = A^{-1}r \quad (27.4.8)$$

بنابراین اگر A^{-1} عناصر بزرگی داشته باشد، $x - \hat{x}$ ممکن است بسیار بزرگتر از r شود.

در عمل، مانده r حتی برای مسائل بدوضع خیلی کوچک است. برای اینکه ببینیم چرا این امر اتفاق می‌افتد، از (۲۲.۴.۸) استفاده می‌کنیم تا به‌دست آوریم

$$r = (\delta A)\hat{x}$$

$$\| r \|_{\infty} \leq \| \delta A \|_{\infty} \| \hat{x} \|_{\infty}$$

$$\frac{\| r \|_{\infty}}{\| A \|_{\infty} \| \hat{x} \|_{\infty}} \leq \frac{\| \delta A \|_{\infty}}{\| A \|_{\infty}} \quad (28.4.8)$$

کرانه‌های $\| \delta A \| / \| A \|$ در (۲۳.۴.۸) یا در (۲۵.۴.۸) مستقل از بدوضعی یا خوش‌وضعی مسئله هستند. بنابراین $\| r \|_{\infty}$ نسبت به $\| \hat{x} \| \| A \|$ ، معمولاً کوچک است. جمله اخیر اغلب نزدیک به $\| b \|$ یا به اندازه آن است، زیرا $b = Ax$ ، پس $\| r \|$ نسبت به $\| b \|$ کوچک خواهد بود. به عنوان آخرین نکته در مورد اندازه مانده، مسائلی هستند که در آنها نیاز به اینکه $x - \hat{x}$ کوچک باشد، نیست مهم آن است که فقط r کوچک باشد، در چنین حالتی، بدوضعی معنای خود را نخواهد داشت. کرانه‌های (۱۸.۴.۸) و (۲۴.۴.۸) اهمیت $\text{cond}(A)$ را در تعیین خطا، نشان می‌دهند. به طور کلی اگر $\text{cond}(A) \doteq 10^m$ ، برای یک $m \geq 0$ ، آنگاه در حدود m رقم از دقت در محاسبه \hat{x} از تعداد ارقامی که در حساب رایانه‌یی به کار رفته، از دست خواهد رفت. بنابراین اندازه‌گیری $\| A^{-1} \| \| A \| = \text{cond}(A)$ مطلوب خواهد بود. محاسبه $\| A \|$ آسان و ارزان است و $\| A^{-1} \|$ مسئله اصلی در محاسبه $\text{cond}(A)$ است. محاسبه A^{-1} نیاز به n^3 عمل دارد و این راه بسیار گرانی برای محاسبه $\| A^{-1} \|$ است. یک روش ارزانتر، با $O(n^2)$ عمل در بسته نرم‌افزاری LINPACK وجود دارد.

برای هر دستگاه $Ay = d$

$$y = A^{-1}d$$

$$\| y \| \leq \| A^{-1} \| \| d \|$$

$$\| A^{-1} \| \geq \frac{\| y \|}{\| d \|} \quad (29.4.8)$$

می‌خواهیم d را طوری انتخاب کنیم تا این نسبت، در حد ممکن بزرگ باشد. می‌نویسیم $A = LU$ که LU از حذف گاوسی به دست آمده است. پس حل $Ay = d$ هم‌ارز با حل

$$Lw = d \quad Uy = w$$

است. وقتی $Lw = d$ را حل می‌کنیم، d را طوری پیدا می‌کنیم که w هرچه بیشتر بزرگ و $\| d \|_{\infty} = 1$ نگهداشته شود. آنگاه $Uy = w$ را نسبت به y حل می‌کنیم. این عمل در (۲۹.۴.۸)، کرانی بهتر از d بی که به طور تصادفی انتخاب شده به دست خواهد داد. یک الگوریتم برای انتخاب d در گلوب و ون‌لون (۱۹۸۳، ص ۷۷) داده شده است. الگوریتم LINPACK تعمیمی پیچیده‌تر از الگوریتم قبلی است. برای توضیح گلوب و ون‌لون (۱۹۸۳، ص ۷۸) را ببینید.

کرانه‌های پسین خطا با کرانه‌های خطا برای وارون حساب شده C از ماتریس داده شده A شروع می‌کنیم. ماتریس مانده را به شکل زیر تعریف می‌کنیم

$$R = I - CA$$

قضیه ۶.۸ اگر $\|R\| < 1$ ، آنگاه A و C ناتکین هستند و

$$\frac{\|R\|}{\|A\| \|C\|} \leq \frac{\|A^{-1} - C\|}{\|C\|} \leq \frac{\|R\|}{1 - \|R\|} \quad (30.4.8)$$

برهان چون $\|R\| < 1$ ، با استفاده از قضیه ۱۱.۷ فصل ۷، $I - R$ ناتکین است، و

$$\|(I - R)^{-1}\| \leq \frac{1}{1 - \|R\|}$$

ولی

$$I - R = CA \quad (31.4.8)$$

$$\det(I - R) = \det(CA) = \det(C) \det(A) \neq 0$$

و بنابراین $\det(A)$ و $\det(C)$ هر دو مخالف صفرند. این نشان می‌دهد که A و C هر دو ناتکین‌اند.

برای کران پایین در (30.4.8)،

$$R = I - CA = (A^{-1} - C)A$$

$$\|R\| \leq \|A^{-1} - C\| \|A\|$$

و تقسیم بر $\|A\| \|C\|$ قضیه را ثابت می‌کند. برای کران بالا، (31.4.8) ایجاب می‌کند که

$$(I - R)^{-1} = A^{-1}C^{-1}$$

$$A^{-1} = (I - R)^{-1}C \quad (32.4.8)$$

برای خطا در C ،

$$A^{-1} - C = (I - CA)A^{-1} = RA^{-1} = R(I - R)^{-1}C$$

$$\|A^{-1} - C\| \leq \frac{\|R\| \|C\|}{1 - \|R\|}$$

که برهان کامل می‌شود. ■

جنبه نظری این قضیه معمولاً مهمتر از جنبه عملی آن است. همان‌گونه که قبلاً در بخش ۱.۸ اشاره شد، ماتریسهای وارون لازم نیست برای حل دستگاه خطی حساب شوند. و در نتیجه، به ندرت نیاز واقعی به کران از نوع قبل، خواهد بود. استثنای عمده وقتی است که C تقریبی باشد که به وسیله‌ای غیر از حذف گاوسی، بیشتر از راههای نظری، به دست آمده باشد. در آن صورت، چنین وارونهای تقریبی، برای حل $Ax = b$ ، با روش تصحیح مانده (۳.۵.۸)، که در بخش بعد توضیح داده خواهد شد، به کار برده می‌شوند. در چنین حالتی، کران (۳.۴.۸) می‌تواند اطلاعات سودمندی درباره C به ما بدهد.

فرع گیریم A ، C و R همانهایی باشند که در قضیه ۶.۸ داده شده‌اند. فرض می‌کنیم \hat{x} یک جواب تقریبی برای $Ax = b$ باشد و $r = b - A\hat{x}$ در این صورت

$$\|x - \hat{x}\| \leq \frac{\|Cr\|}{1 - \|R\|} \quad (۳۳.۴.۸)$$

برهان داریم

$$\begin{aligned} r &= b - A\hat{x} = Ax - A\hat{x} = A(x - \hat{x}) \\ x - \hat{x} &= A^{-1}r = (I - R)^{-1}Cr \end{aligned} \quad (۳۴.۴.۸)$$

با استفاده از (۳۲.۴.۸) که در معادله آخر به کار رفته است. اگر از معادله اخیر نرم بگیریم (۳۳.۴.۸) به دست می‌آید.

معلوم شده است که این کران (۳۳.۴.۸) کاملاً دقیق است، به ویژه در مقایسه با بعضی از کرانهایی که اغلب به کار می‌روند. برای یک بحث کامل درباره کرانهای محاسبه پذیر خطا، از جمله تعدادی مثال، ارد ولینچ^۱ (۱۹۷۵) را ببینید.

پیدا کردن کران خطای (۳۳.۴.۸) نسبتاً گران تمام می‌شود. اگر فرض کنیم که \hat{x} از حذف گاوسی به دست آمده است، $n^3/3$ عمل برای محاسبه \hat{x} و تجزیه A به LU به کار رفته است. برای محاسبه $C \equiv A^{-1}$ از راه حذف، حداقل $n^3/2$ عمل اضافی لازم است، به دست آوردن CA نیاز به n^3 عمل ضرب دارد، و به دست آوردن Cr نیاز به n^2 ضرب. بنابراین کران خطا حداقل یک افزایش پنج برابر عمل در بر خواهد داشت. معمولاً ترجیح داده می‌شود که خطا از حل تقریبی معادله خطای

$$A(x - \hat{x}) = r$$

با استفاده از LU که قبلاً ذخیره شده برآورد شود. این امر نیاز به n^2 عمل برای محاسبه r و n^2 عمل دیگر برای حل دستگاه خطی دارد. مگر وقتی که ماتریس $R = I - CA$ دارای نرم نزدیک به یک باشد، این روش یک خطای بسیار قابل قبولی به دست می‌دهد. این مطلب در بخش بعد دنبال و توضیح داده شده است.

۵.۸ روش تصحیح مانده

فرض می‌کنیم که $Ax = b$ حل شده، جواب تقریبی $\hat{x} \equiv x^{(0)}$ به دست آمده است. همچنین گیریم تجزیه LU و تعداد تمام تعویضهای سطری یا ستونی ذخیره شده است. کمیت زیر را حساب می‌کنیم

$$r^{(0)} = b - Ax^{(0)} \quad (۱.۵.۸)$$

تعریف می‌کنیم $e^{(0)} = x - x^{(0)}$. در این صورت مانند قبل در (۳۴.۴.۸)،

$$Ae^{(0)} = r^{(0)}$$

این دستگاه را با تجزیه LU ی ذخیره شده حل کرده جواب تقریبی نتیجه را $\hat{e}^{(0)}$ می‌نامیم. یک جواب تقریبی جدید

$$x^{(1)} = x^{(0)} + \hat{e}^{(0)} \quad (۲.۵.۸)$$

را برای $Ax = b$ تعریف می‌کنیم این فرایند را می‌توانیم تکرار و $x^{(2)}, \dots$ را، با خطایی که مرتب کاهش می‌یابد، حساب کنیم. محاسبه $r^{(0)}$ ، نیاز به n^2 عمل، و محاسبه $\hat{e}^{(0)}$ ، به n^2 عمل دیگر نیاز دارد. بنابراین محاسبه مقادیر بهبودیافته $x^{(1)}, x^{(2)}, \dots$ در مقایسه با محاسبه مقدار اصلی $x^{(0)}$ ، گران نخواهد بود. این روش را روش بهبود بارستی یا روش تصحیح مانده نیز می‌نامند.

به دست آوردن مقادیر دقیق برای $r^{(0)}$ بسیار اهمیت دارد. چون جواب تقریبی $Ax = b$ است، معمولاً در محاسبه اش خطاهای کاهش دقت خواهد داشت، وقتی $Ax^{(0)}$ و b تقریباً تا آخرین درجه دقت حساب ماشین حساب شده‌اند، بنابراین در به دست آوردن مقادیر دقیق برای $r^{(0)}$ ، معمولاً باید به سمت حساب با دقت بالاتر برویم. اگر از همان درجه دقتی که برای حساب LU و $x^{(0)}$ به کار برده شده برای محاسبه $r^{(0)}$ استفاده کنیم، این نتیجه غیردقیق $r^{(0)}$ معمولاً به $\hat{e}^{(0)}$ ی منجر خواهد شد که یک تقریب ضعیف برای $e^{(0)}$ خواهد بود. در حساب با دقت ساده،

$r^{(0)}$ را با دقت مضاعف محاسبه می‌کنیم. ولی اگر محاسبات قبلاً با دقت مضاعف انجام شده باشد، معمولاً مشکل است که به دقت بالاتر برسیم.

مثال دستگاه $Ax = b$ را با $A = \bar{H}_3$ که در (۱۴.۴.۸) آمده با حساب ممیز شناور با چهار رقم اعشاری و گرد کردن حل کنید. برای طرف راست از

$$b = [1, 0, 0]^T$$

استفاده می‌کنیم. جواب درست اولین ستون \bar{H}_3^{-1} است، که با توجه به (۱۵.۴.۸)، تا چهار رقم معنی‌دار چنین است

$$x = [9, 062, -36, 32, 30, 30]^T$$

با استفاده از حذف و محورگیری جزئی،

$$x^{(0)} = [8, 968, -35, 77, 29, 77]^T$$

مانده $r^{(0)}$ با حساب با دقت مضاعف به دست آمده و سپس تا چهار رقم با معنی گرد شده است. مقادیر به دست آمده چنین‌اند

$$r^{(0)} = [-0, 0005341, -0, 0004359, -0, 0005344]^T$$

از حل دستگاه $Ae^{(0)} = r^{(0)}$ با تجزیه ذخیره‌شده LU ،

$$\hat{e}^{(0)} = [0, 09216, -0, 5442, 0, 5239]^T$$

$$x^{(1)} = [9, 060, -36, 31, 30, 29]^T$$

از تکرار این عملیات،

$$r^{(1)} = [-0, 00006570, -0, 00003770, -0, 00001980]^T$$

$$\hat{e}^{(1)} = [0, 0001707, -0, 01300, 0, 01241]^T$$

$$x^{(2)} = [9, 062, -36, 32, 30, 30]^T$$

بردار $x^{(2)}$ تا چهار رقم با معنی دقیق است. ضمناً توجه کنید که $x^{(1)} - x^{(0)} = \hat{e}^{(0)}$ یک پیشگوی دقیق برای خطای $e^{(0)}$ در $x^{(0)}$ است.

فرمولها را می‌توان بسط داده برآورد نمود که چند بارست برای پیدا کردن جواب x با دقت کامل، لازم است. برای بحث در مورد مطالب مربوطه و الگوریتم‌هایی که از این روش استفاده می‌کنند دون‌گارا و همکاران (۱۹۷۹، ص ۱۹)، فورسایت و مولر (۱۹۶۷، فصلهای ۱۳، ۱۶، ۱۷)، گلوب و ون‌لون (۱۹۸۳) و ویلکینسن و راینش (۱۹۷۱، صص ۹۳-۱۱۰) را ببینید.

روش دیگر تصحیح مانده مواردی وجود دارند که می‌توانیم یک وارون تقریبی C را برای ماتریس داده‌شده A محاسبه نماییم. این عمل معمولاً با بررسی دقیق ساختار A ، و سپس استفاده از تکنیکهای متنوع تقریب در برآورد A^{-1} انجام می‌گیرد. بدون توجه به منشأ C ، نشان خواهیم داد که چگونه آن را به گونه‌ای بارستی برای حل $Ax = b$ به کار می‌بریم.

گیریم $x^{(0)}$ یک حدس اولیه باشد و تعریف می‌کنیم $r^{(0)} = b - Ax^{(0)}$. مانند قبل، $A(x - x^{(0)}) = r^{(0)}$ را به‌طور ضمنی با رابطه زیر تعریف می‌کنیم

$$x^{(1)} - x^{(0)} = Cr^{(0)}$$

به‌طور کلی، تعریف می‌کنیم

$$r^{(m)} = b - Ax^{(m)} \quad x^{(m+1)} = x^{(m)} + Cr^{(m)} \quad m = 0, 1, 2, \dots \quad (3.5.8)$$

اگر C یک تقریب خوبی برای A^{-1} باشد، همان‌گونه که در تحلیل زیر نشان داده خواهد شد این بارست به‌سرعت همگرا می‌شود.

ابتدا یک فرمول بازگشتی برای خطا به‌دست می‌آوریم:

$$\begin{aligned} x - x^{(m+1)} &= x - x^{(m)} - Cr^{(m)} = x - x^{(m)} - C[b - Ax^{(m)}] \\ &= x - x^{(m)} - C[Ax - Ax^{(m)}] \\ x - x^{(m+1)} &= (I - CA)(x - x^{(m)}) \end{aligned} \quad (4.5.8)$$

با استقرا

$$x - x^{(m)} = (I - CA)^m(x - x^{(0)}) \quad m \geq 0 \quad (5.5.8)$$

اگر برای یک نرم ماتریسی،

$$\|I - CA\| < 1 \quad (6.5.8)$$

آنگاه با استفاده از نرم برداری وابسته

$$\|x - x^{(m)}\| \leq \|I - CA\|^m \|x - x^{(0)}\| \quad (۷.۵.۸)$$

و این رابطه با هر انتخابی برای $x^{(0)}$ هرگاه $m \rightarrow \infty$ به صفر همگرا خواهد شد. به طور کلیتر، برای هر انتخاب $x^{(0)}$ ، $x^{(m)}$ به x همگرا خواهد بود اگر و تنها اگر

$$(I - CA)^m \rightarrow 0 \quad m \rightarrow \infty \quad \text{وقتی}$$

و با توجه به قضیه ۹.۷ از فصل ۷، این عبارت هم‌ارز است با

$$r_\sigma(I - CA) < 1 \quad (۸.۵.۸)$$

برای شعاع خاص $I - CA$. این رابطه را ممکن است، حتی وقتی که (۶.۵.۸) برای یک نرم ماتریسی معمولی برقرار باشد، ثابت کرد. همچنین توجه نمایید که

$$I - AC = A(I - CA)A^{-1}$$

و بنابراین $I - CA$ و $I - AC$ ماتریسهای مشابه‌اند و دارای ویژه‌مقدارهای برابرند. اگر

$$\|I - AC\| < 1 \quad (۹.۵.۸)$$

آنگاه (۸.۵.۸) درست است حتی اگر (۶.۵.۸) درست نباشد، و بازم هم همگرایی برقرار است. عبارت (۴.۵.۸) نشان می‌دهد که نرخ همگرایی $x^{(m)}$ به x خطی است:

$$\|x - x^{(m+1)}\| \leq c \|x - x^{(m)}\| \quad m \geq 0 \quad (۱۰.۵.۸)$$

با مقدار مجهول $c < 1$. ثابت c اغلب از راه محاسبه با رابطه زیر برآورد می‌شود

$$c = \text{Max} \frac{\|x^{(m+2)} - x^{(m+1)}\|}{\|x^{(m+1)} - x^{(m)}\|} \quad (۱۱.۵.۸)$$

که ماکسیمم روی بعضی یا تمام بارستهایی که محاسبه شده، گرفته می‌شود. این رابطه دقیق نیست و از فرمول زیر حاصل شده است

$$x^{(m+2)} - x^{(m+1)} = (I - CA)(x^{(m+1)} - x^{(m)}) \quad (۱۲.۵.۸)$$

برای اثبات این تساوی، تنها (۴.۵.۸) را به‌کار می‌بریم و فرمولها را برای مقادیر پیاپی m از هم کم می‌کنیم.

اگر فرض کنیم $(۱۰.۵.۸)$ برای بارستهایی که محاسبه می‌کنیم معتبر باشد، و اگر یک برآوردی برای c داشته باشیم، آنگاه می‌توانیم یک کران خطا بسازیم.

$$\begin{aligned} \|x^{(m+1)} - x^{(m)}\| &= \| [x - x^{(m)}] - [x - x^{(m+1)}] \| \\ &\geq \|x - x^{(m)}\| - \|x - x^{(m+1)}\| \\ &\geq \|x - x^{(m)}\| - c \|x - x^{(m)}\| \\ \|x - x^{(m)}\| &\leq \frac{1}{1-c} \|x^{(m+1)} - x^{(m)}\| \\ \|x - x^{(m+1)}\| &\leq \frac{c}{1-c} \|x^{(m+1)} - x^{(m)}\| \quad (۱۳.۵.۸) \end{aligned}$$

برای بارستهایی که همگرایی آهسته دارند [با $c \approx 1$]، این کران مهم است، زیرا در این صورت $\|x^{(m+1)} - x^{(m)}\|$ ممکن است بسیار کوچکتر از $\|x - x^{(n)}\|$ باشد. همچنین به دست آوردن کران در بخش ۵.۲ از فصل دوم را به یاد آورید. یک کران مشابه، $(۵.۵.۲)$ ، برای خطا در یک روش همگرایی خطی به دست آمده بود.

مثال تعریف می‌کنیم $A(\varepsilon) = A_0 + \varepsilon B$ ، با

$$A_0 = \begin{bmatrix} 2 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 2 \end{bmatrix} \quad B = \begin{bmatrix} 0 & 1 & 1 \\ -1 & 0 & 1 \\ -1 & -1 & 0 \end{bmatrix}$$

به‌عنوان یک وارون تقریبی برای $A(\varepsilon)$ ، از

$$A(\varepsilon)^{-1} \doteq C = A_0^{-1} = \begin{bmatrix} \frac{2}{3} & -\frac{1}{3} & \frac{1}{3} \\ -\frac{1}{3} & 1 & -\frac{1}{3} \\ \frac{1}{3} & -\frac{1}{3} & \frac{2}{3} \end{bmatrix}$$

استفاده می‌کنیم. دستگاه $A(\varepsilon)x = b$ را می‌توان با روش تصحیح مانده $(۳.۵.۸)$ حل کرد. برای تحلیل همگرایی

$$\begin{aligned} I - CA(\varepsilon) &= I - A_0^{-1}[A_0 + \varepsilon B] = -\varepsilon A_0^{-1}B \\ &= -\varepsilon \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ -\frac{1}{3} & 0 & \frac{1}{3} \\ -\frac{1}{3} & -\frac{1}{3} & -\frac{1}{3} \end{bmatrix} \end{aligned}$$

همگرایی وقتی تضمین می شود که

$$\|I - CA(\varepsilon)\|_{\infty} = |\varepsilon| < 1$$

و از (۴.۵.۸)

$$\|x - x^{(m+1)}\|_{\infty} \leq |\varepsilon| \|x - x^{(m)}\|_{\infty} \quad m \geq 0$$

موارد زیادی از نوع این مثال وجود دارند. ممکن است لازم باشد دستگاههای خطی در شکل کلی $A(\varepsilon)x = b$ را برای هر ε نزدیک به صفر حل کنیم. برای صرفه جویی در وقت، یا $A(0)^{-1}$ یا LU ، تجزیه $A(0)$ ، را به دست می آوریم. در این صورت این را به عنوان یک وارون تقریبی برای $A(\varepsilon)$ به کار می بریم، و $A(\varepsilon)x = b$ را با استفاده از روش تصحیح مانده حل می کنیم.

۶.۸ روشهای بارستی

همان گونه که در مقدمه این فصل اشاره شد، بسیاری از دستگاهها خیلی بزرگتر از آن هستند که با روشهای مستقیم، بر پایه حذف گاوسی، حل شوند. برای این دستگاهها، روشهای بارستی، اغلب تنها راه حل ممکن و همچنین در بسیاری از حالات، سریعتر از حذف هستند. وسیعترین زمینه برای کاربرد روشهای بارستی دستگاههای خطی است که از حل عددی معادلات دیفرانسیل جزئی حاصل می شوند. دستگاههای از مراتب 10^3 تا 10^5 غیر معمول نیستند، ولی تقریباً همه ضرایب دستگاه صفرند. به عنوان یک مثال از چنین مسائلی، حل عددی معادله بواسن است که در بخش ۸.۸ مطالعه شده است. خواننده، در صورتی که بخواهد می تواند آن بخش را با این بخش با هم بخواند.

گذشته از بزرگ بودن، دستگاههای خطی که باید حل شوند، اغلب چندین ویژگی مهم دیگر دارند. معمولاً تنگ هستند، یعنی درصد کوچکی از ضرایب ناصفرند. ضرایب ناصفر به شیوه خاصی در A ظاهر می شوند، و معمولاً فرمول ساده ای برای تولید ضرایب a_{ij} ، در صورت لزوم، وجود دارد تا مجبور به ذخیره کردن آنها نباشیم. به عنوان یک نتیجه این ویژگیها، توجه به فضای حافظه برای بردارهای x و b ممکن است بسیار مهمتر از رعایت فضای حافظه برای A باشد. ماتریسهای A اغلب دارای ویژگیهای خاصی هستند، که در این بخش و دو بخش دیگر مورد بحث قرار خواهند گرفت.

موضوع را با تعریف و تحلیل دو روش بارستی کلاسیک آغاز می کنیم؛ به دنبال آن، یک چارچوب مجرد کلی برای مطالعه روشهای بارستی ارائه می کنیم. ویژگیهای خاص دستگاه خطی $Ax = b$ ، هنگام تدوین روش بارستی برای حل آن، بسیار مهم اند. نتایج این بخش فقط آغاز طرح یک روش برای هر زمینه خاص کاربردی خواهند بود.

روش گاوس-ژاکوبی (جایگزینی همزمان) دستگاه $Ax = b$ را به شکل زیر بازنویسی می‌نماییم

$$x_i = \frac{1}{a_{ii}} \left\{ b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j \right\} \quad i = 1, 2, \dots, n \quad (۱.۶.۸)$$

با فرض $a_{ii} \neq 0$. بارست را این طور تعریف می‌نماییم

$$x_i^{(m+1)} = \frac{1}{a_{ii}} \left\{ b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(m)} \right\} \quad i = 1, \dots, n \quad m \geq 0 \quad (۲.۶.۸)$$

و فرض می‌کنیم که حدسهای اولیه $x_i^{(0)}$, $i = 1, \dots, n$ داده شده‌اند. صورتهای دیگری از این روش وجود دارند. برای مثال، بسیاری از مسائل طبیعتاً به صورت

$$(I - B)x = b$$

داده شده‌اند، بنابراین معمولاً ابتدا بارست

$$x^{(m+1)} = b + Bx^{(m)} \quad m \geq 0 \quad (۳.۶.۸)$$

را در نظر می‌گیریم. اولین تحلیل خطا به (۲.۶.۸) محدود می‌شود، ولی همین تحلیل را می‌توان برای (۳.۶.۸) به کار برد.

برای تحلیل همگرایی گیریم $e^{(m)} = x - x^{(m)}$ با کم کردن (۲.۶.۸) از (۱.۶.۸)

داریم:

$$e_i^{(m+1)} = - \sum_{\substack{j=1 \\ j \neq i}}^n \frac{a_{ij}}{a_{ii}} e_j^{(m)} \quad i = 1, \dots, n \quad m \geq 0 \quad (۴.۶.۸)$$

$$\| e_i^{(m+1)} \| \leq \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{ij}}{a_{ii}} \right| \| e^{(m)} \|_{\infty}$$

تعریف می‌کنیم

$$\mu = \max_{1 \leq i \leq n} \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{ij}}{a_{ii}} \right| \quad (۵.۶.۸)$$

پس

$$|e_i^{(m+1)}| \leq \mu \|e^{(m)}\|_\infty$$

و چون سمت راست مستقل از i است،

$$\|e^{(m+1)}\|_\infty \leq \mu \|e^{(m)}\|_\infty \quad (۶.۶.۸)$$

اگر $\mu < ۱$ ، وقتی m با یک نرخ خطی که با μ کراندار می شود به بینهایت میل کند، $e^{(m)} \rightarrow 0$ و

$$\|e^{(m)}\|_\infty \leq \mu^m \|e^{(0)}\|_\infty \quad (۷.۶.۸)$$

برای آنکه $\mu < ۱$ درست باشد، ماتریس A باید غالب قطری باشد، یعنی، باید در شرط زیر صدق کند.

$$\sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| < |a_{ii}| \quad i = 1, 2, \dots, n \quad (۸.۶.۸)$$

چنین ماتریسهایی در کاربردهای متعددی پیدا می شوند و اغلب ماتریس وابسته آنها تنگ است. برای یک نتیجه کلیتر، (۴.۶.۸) را به شکل

$$e^{(m+1)} = M e^{(m)} \quad m \geq 0 \quad (۹.۶.۸)$$

می نویسیم.

$$M = - \begin{bmatrix} 0 & \frac{a_{۱۲}}{a_{۱۱}} & \dots & \frac{a_{۱n}}{a_{۱۱}} \\ \frac{a_{۲۱}}{a_{۲۲}} & 0 & \frac{a_{۲۳}}{a_{۲۲}} & \dots & \frac{a_{۲n}}{a_{۲۲}} \\ \vdots & & \ddots & & \vdots \\ \frac{a_{n۱}}{a_{nn}} & & \dots & & 0 \end{bmatrix}$$

به استقرا

$$e^{(m)} = M^m e^{(0)} \quad (۱۰.۶.۸)$$

اگر بخواهیم وقتی $m \rightarrow \infty$ ، $e^{(m)}$ مستقل از انتخاب $x^{(0)}$ (و بنابراین مستقل از $e^{(0)}$) به 0 میل کند، لازم و کافی است که

$$M^m \rightarrow 0 \quad m \rightarrow \infty \quad \text{وقتی}$$

۶۲۰ حل عددی دستگاههای معادلات خطی

یا هم‌ارز با آن، بنابر قضیه ۹.۷ از فصل ۷،

$$r_\infty(M) < 1 \quad (11.6.8)$$

شرط $\mu < 1$ تنها شرط کوچکتر از یک بودن نرم سطری M ، $\|M\|_\infty < 1$ است و این (۱۱.۶.۸) را ایجاب می‌کند. ولی اکنون می‌بینیم که اگر برای هر عملگر نرم ماتریسی، $\|M\| < 1$ ، آنگاه $e^{(m)} \rightarrow 0$.

مثال حل $Ax = b$ را به روش گاوس-ژاکوبی با ماتریسهای زیر در نظر می‌گیریم

$$A = \begin{bmatrix} 10 & 3 & 1 \\ 2 & -10 & 3 \\ 1 & 3 & 10 \end{bmatrix} \quad b = \begin{bmatrix} 14 \\ -5 \\ 14 \end{bmatrix} \quad x^{(0)} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad (12.6.8)$$

اگر معادله i ام را نسبت به x_i حل کنیم، داریم $x = g + Mx$

$$M = \begin{bmatrix} 0 & -0.3 & -0.1 \\ 0.2 & 0 & 0.3 \\ -0.1 & -0.3 & 0 \end{bmatrix} \quad g = \begin{bmatrix} 1.4 \\ 0.5 \\ 1.4 \end{bmatrix}$$

جواب درست $x = [1, 1, 1]^T$ است. برای واری همگرایی، توجه داریم که $\|M\|_\infty = 0.5$

و $\|M\|_1 = 0.6$ بنابراین

$$\|e^{(m+1)}\|_\infty \leq 0.5 \|e^{(m)}\|_\infty \quad m \geq 0 \quad (13.6.8)$$

حکم مشابهی برای $\|e^{(m+1)}\|_1$ برقرار است. بنابراین همگرایی تضمین می‌شود و خطاها در هر بارست اقلاباً با ضریب $\frac{1}{2}$ کاهش می‌یابند. نتایج عددی واقعی در جدول ۳.۸ داده شده‌اند و مؤید نتیجه (۱۳.۶.۸) هستند. ستون آخر عبارت است از:

$$\text{Ratio} \equiv \frac{\|e^{(m)}\|_\infty}{\|e^{(m-1)}\|_\infty} \quad (14.6.8)$$

این رابطه نشان می‌دهد که ممکن است همگرایی از یک مرحله به مرحله دیگر تغییر کند در حالی که در (۱۳.۶.۸) یا به‌طور کلیتر در (۶.۶.۸) صدق می‌کند.

جدول ۳.۸ نتایج عددی برای روش گاوس-ژاکوبی

m	$x_1^{(m)}$	$x_2^{(m)}$	$x_3^{(m)}$	$\ e^{(m)}\ _\infty$	نسبت
۰	۰	۰	۰	10^0	
۱	۱٫۴	۰٫۵	۱٫۴	0.5	0.5
۲	۱٫۱۱	۱٫۲۰	۱٫۱۱	0.2	0.4
۳	0.929	1.055	0.929	0.071	0.36
۴	0.9906	0.9645	0.9906	0.0355	0.50
۵	1.01159	0.9953	1.01159	0.01159	0.33
۶	1.000251	1.005795	1.000251	0.005795	0.50

روش گاوس-زایدل^۱ (جایگزینی پی در پی) با استفاده از (۱.۶.۸)، تعریف می‌کنیم

$$x_i^{(m+1)} = \frac{1}{a_{ii}} \left\{ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(m+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(m)} \right\} \quad i = 1, 2, \dots, n \quad (15.6.8)$$

هر مؤلفه جدید $x_i^{(m+1)}$ بلافاصله در محاسبه مؤلفه بعدی به‌کار گرفته می‌شود. این امر برای محاسبه رایانه‌یی بسیار مناسب است، زیرا مقدار جدید بلافاصله در جایی که مقدار قدیم قرار داشته ذخیره می‌شود، و این کار جای لازم برای ذخیره‌سازی را مینیمم می‌نماید. روش گاوس-زایدل در مقایسه با روش گاوس-ژاکوبی، به نصف جا برای ذخیره x نیاز دارد.

برای تحلیل خطا، (۱۵.۶.۸) را از (۱.۶.۸) کم می‌کنیم.

$$e_i^{(m+1)} = - \sum_{j=1}^{i-1} \frac{a_{ij}}{a_{ii}} e_j^{(m+1)} - \sum_{j=i+1}^n \frac{a_{ij}}{a_{ii}} e_j^{(m)} \quad i = 1, 2, \dots, n \quad (16.6.8)$$

تعریف می‌کنیم

$$\alpha_i = \sum_{j=1}^{i-1} \left| \frac{a_{ij}}{a_{ii}} \right|, \quad \beta_i = \sum_{j=i+1}^n \left| \frac{a_{ij}}{a_{ii}} \right| \quad i = 1, \dots, n$$

با $\alpha_1 = \beta_n = 0$. از همان تعریف (۵.۶.۸) برای μ که برای روش ژاکوبی به‌کار رفت استفاده می‌کنیم،

$$\mu = \text{Max}_{1 \leq i \leq n} (\alpha_i + \beta_i)$$

فرض می‌کنیم $\mu < 1$. سپس تعریف می‌کنیم

$$\eta = \text{Max}_{1 \leq i \leq n} \frac{\beta_i}{1 - \alpha_i} \quad (17.6.8)$$

از (۱۶.۶.۸)،

$$|e_i^{(m+1)}| \leq \alpha_i \|e^{(m+1)}\|_\infty + \beta_i \|e^{(m)}\|_\infty \quad i = 1, \dots, n \quad (۱۸.۶.۸)$$

گیریم k اندیسی باشد که برای آن

$$\|e^{(m+1)}\|_\infty = |e_k^{(m+1)}|$$

در این صورت با $i = k$ در (۱۸.۶.۸)،

$$\|e^{(m+1)}\|_\infty \leq \alpha_k \|e^{(m+1)}\|_\infty + \beta_k \|e^{(m)}\|_\infty$$

$$\|e^{(m+1)}\|_\infty \leq \frac{\beta_k}{1 - \alpha_k} \|e^{(m)}\|_\infty$$

و بنابراین

$$\|e^{(m+1)}\|_\infty \leq \eta \|e^{(m)}\|_\infty \quad (۱۹.۶.۸)$$

چون برای هر i ،

$$(\alpha_i + \beta_i) - \frac{\beta_i}{1 - \alpha_i} = \frac{\alpha_i[1 - (\alpha_i + \beta_i)]}{1 - \alpha_i} \geq \frac{\alpha_i}{1 - \alpha_i} [1 - \mu] \geq 0$$

داریم

$$\eta \leq \mu < 1 \quad (۲۰.۶.۸)$$

که ترکیب آن با (۱۹.۶.۸)، همگرایی $e^{(m)} \rightarrow 0$ را وقتی $m \rightarrow \infty$ نشان می‌دهد. در این روش هم نرخ همگرایی خطی خواهد بود، ولی نسبت به روش ژاکوبی، همگرایی سریعتر است.

مثال دستگاه (۱۲.۶.۸) از مثال قبل را بگیرید و آن را با روش گاوس-زایدل حل کنید. با یک محاسبه ساده از (۱۷.۶.۸) و (۱۲.۶.۸)،

$$\eta = 0.4$$

نتایج عددی در جدول ۴.۸ داده شده‌اند. سرعت همگرایی به مراتب بهتر از سرعت همگرایی روش گاوس-ژاکوبی مثال قبل است که در جدول ۳.۸ داده شده است. به نظر می‌آید که مقادیر نسبت همگرایی حدود ۱۸٪ است.

چارچوب کلی برای روشهای بارستی برای حل $Ax = b$ یک شکل دوباره‌شده A را در نظر می‌گیریم:

$$A = N - P \quad (۲۱.۶.۸)$$

جدول ۴.۸ نتایج عددی برای روش گاوس - زایدل

m	$x_1^{(m)}$	$x_2^{(m)}$	$x_3^{(m)}$	$\ e^{(m)}\ _\infty$	نسبت
۰	۰	۰	۰	۱	
۱	۱٫۴	۰٫۷۸	۱٫۰۲۶	۰٫۴	۰٫۴
۲	۱٫۰۶۳۴۰۰	۱٫۰۲۰۴۸۰	۰٫۹۸۷۵۱۶	$۶٫۳۴E-۲$	۰٫۱۶
۳	۰٫۹۹۵۱۰۴	۰٫۹۹۵۲۷۶	۱٫۰۰۱۹۰۷	$۴٫۹۰E-۳$	۰٫۰۷۷
۴	۱٫۰۰۱۲۲۷	۱٫۰۰۰۸۱۷	۰٫۹۹۹۶۳۲	$۱٫۲۳E-۳$	۰٫۲۵
۵	۰٫۹۹۹۷۹۲	۰٫۹۹۹۸۴۸	۱٫۰۰۰۰۶۶	$۲٫۰۸E-۴$	۰٫۱۷
۶	۱٫۰۰۰۰۳۹	۱٫۰۰۰۰۲۸	۰٫۹۹۹۹۸۸	$۳٫۹۰E-۵$	۰٫۱۹

و $Ax = b$ را به صورت زیر می نویسیم

$$Nx = b + Px \quad (۲۲.۶.۸)$$

ماتریس N به گونه‌ای انتخاب شده است که دستگاه خطی $Nz = f$ ، برای هر f ، «به سادگی حلپذیر» باشد. مثلاً N ممکن است قطری، سه قطری یا مثلثی باشد. روش بارستی را با $x^{(۰)}$ داده شده به صورت زیر تعریف می کنیم

$$Nx^{(m+1)} = b + Px^{(m)} \quad m \geq 0 \quad (۲۳.۶.۸)$$

مثال ۱. روش ژاکوبی

$$N = \text{diag}[a_{11}, a_{22}, \dots, a_{nn}] \quad P = N - A$$

۲. روش گاوس - زایدل

$$N = \begin{bmatrix} a_{11} & 0 & \dots & 0 \\ a_{21} & a_{22} & 0 & \dots \\ \vdots & & \ddots & \vdots \\ a_{n1} & \dots & & a_{nn} \end{bmatrix} \quad (۲۴.۶.۸)$$

برای تحلیل خطا، (۲۳.۶.۸) را از (۲۲.۶.۸) کم می کنیم تا به دست آوریم

$$Ne^{(m+1)} = Pe^{(m)}$$

$$e^{(m+1)} = Me^{(m)} \quad M = N^{-1}P \quad (۲۵.۶.۸)$$

با استقرا

$$e^{(m)} = M^m e^{(0)} \quad m \geq 0 \quad (26.6.8)$$

برای آنکه برای حدس اولیه دلخواه $x^{(0)}$ ، (و بنابراین هر مقدار دلخواه $e^{(0)}$) وقتی n به بینهایت میل می‌کند، $e^{(m)} \rightarrow 0$ ، لازم و کافی است که

$$M^m \rightarrow 0 \quad m \rightarrow \infty \quad \text{وقتی}$$

یا هم‌ارز با آن از قضیه ۹.۷

$$r_\sigma(M) < 1 \quad (27.6.8)$$

این چارچوب کلی برای روشهای بارستی، از آیزکسون و کلا (۱۹۶۶، صص ۶۱-۸۱) اقتباس شده است. شرط (۲۷.۶.۸) قبلاً در (۱۱.۶.۸) برای روش گاوس-ژاکوبی به دست آمده بود. در روش گاوس-زایدل کار با ماتریس $N^{-1}P$ که N با (۲۴.۶.۸) داده شده، دشوارتر است. باید مقادیری از λ را به دست آوریم که برای آن

$$\det(\lambda I - N^{-1}P) = 0$$

یا هم‌ارز آن

$$\det(\lambda N - P) = 0 \quad (28.6.8)$$

برای استفاده در حل عددی معادلات دیفرانسیل جزئی در بخش ۸.۸، تحلیل همگرایی بالا با روش گاوس-زایدل مناسب نیست. ثابتهای μ و η در (۵.۶.۸) و (۱۷.۶.۸) هر دو مساوی ۱ می‌شوند، با اینکه به‌طور تجربی، این روش باز هم همگراست. برای پرداختن به بسیاری از این دستگاهها، اغلب قضیه مهم زیر را به‌کار می‌برند.

قضیه ۷.۸ گیریم A ارمیتی با عناصر قطری مثبت باشد. در این صورت روش گاوس-زایدل برای حل $Ax = b$ ، برای هر انتخاب $x^{(0)}$ همگرا می‌شود، اگر و تنها اگر A معین مثبت باشد.

برهان اثبات در آیزکسون و کلا (۱۹۶۶، صص ۷۰-۷۱) داده شده است. برای تعریف ماتریس معین مثبت، مسأله ۱۴ فصل ۷ را به یاد آورید. این قضیه در بخش ۸.۸ توضیح داده شده است. ■

روشهای دیگر بارستی بهترین روشهای بارستی روشهایی هستند که بر پایه یک اطلاع کامل

از مسأله‌ای که باید حل شود، جنبه‌های خاص آن در طرح برنامه بارستی به حساب آیند. این کار معمولاً با نگاه کردن به شکل ماتریس و منشاء دستگاه خطی صورت می‌گیرد.

ماتریس A ممکن است شکل خاصی داشته باشد که به روش بارستی ساده‌ای منجر شود. به‌عنوان مثال، فرض کنید A یک شکل سه قطری بلوکی دارد:

$$A = \begin{bmatrix} B_1 & C_1 & & \dots & \circ \\ A_2 & B_2 & C_2 & & \\ \vdots & & & \ddots & \vdots \\ \circ & & \dots & A_r & B_r \end{bmatrix} \quad (29.6.8)$$

ماتریسهای A_i, B_i و C_i مربعی از مرتبه m و A از مرتبه $n = rm$ است. به ازای $x, b \in \mathbf{R}^n$ بردارهای x و b را به شکل افزاز شده مثل:

$$x = \begin{bmatrix} x_{(1)} \\ \vdots \\ x_{(r)} \end{bmatrix} \quad b = \begin{bmatrix} b_{(1)} \\ \vdots \\ b_{(r)} \end{bmatrix} \quad x_{(i)}, b_{(i)} \in \mathbf{R}^m$$

می‌نویسیم. در این صورت $Ax = b$ به شکل زیر نوشته می‌شود

$$\begin{aligned} B_1 x_{(1)} + C_1 x_{(2)} &= b_{(1)} \\ A_i x_{(i-1)} + B_i x_{(i)} + C_i x_{(i+1)} &= b_{(i)} \quad 2 \leq i \leq r-1 \\ A_r x_{(r-1)} + B_r x_{(r)} &= b_{(r)} \end{aligned} \quad (30.6.8)$$

فرض می‌کنیم که دستگاههای خطی

$$B_j x_{(j)} = d_{(j)} \quad 1 \leq j \leq r \quad (31.6.8)$$

به‌سادگی، شاید مستقیماً برای تمام $d_{(j)}$ های طرفهای راست حلپذیر باشند. مثلاً اغلب همه B_i ها یک ماتریس ثابت سه قطری T هستند که برای آنها شیوه به‌کاررفته در (۲۰.۳.۸) - (۲۴.۳.۸) را می‌توان به‌کار برد.

می‌توان روشی از نوع روش ژاکوبی برای (۳۰.۶.۸) به‌کار برد

$$\begin{aligned} B_1 x_{(1)}^{(v+1)} &= b_{(1)} - C_1 x_{(2)}^{(v)} \\ B_i x_{(i)}^{(v+1)} &= b_{(i)} - A_i x_{(i-1)}^{(v)} - C_i x_{(i+1)}^{(v)} \quad 2 \leq i \leq r-1 \\ B_r x_{(r)}^{(v+1)} &= b_{(r)} - A_r x_{(r-1)}^{(v)} \end{aligned} \quad (32.6.8)$$

تحلیل همگرایی پیچیده‌تر از تحلیل همگرایی در روشهای گاوس-ژاکوبی و گاوس-زایدل است؛ به بعضی از نتایج در مسأله ۲۹ اشاره شده است. روشهای مشابهی برای دستگاههای خطی که از حل بعضی معادلات دیفرانسیل جزئی به دست می‌آیند، به کار برده شده‌اند، که در بخش ۸.۸ به آنها اشاره شده است.

یک جنبه مهم دیگر حل دستگاههای خطی $Ax = b$ ، نگرستن به منشاء پیدایش آنهاست. در بسیاری از موارد یک معادله دیفرانسیل یا انتگرالی مثل

$$\mathcal{A}x = y \quad (۳۳.۶.۸)$$

داریم که x و y تابع‌اند. اگر این معادله را به صورت گسسته بنویسیم، یک خانواده مسائل

$$A_n x_n = y_n \quad x_n, y_n \in \mathbf{R}^n \quad (۳۴.۶.۸)$$

با A_n از مرتبه n حاصل می‌شود. وقتی $n \rightarrow \infty$ ، جوابهای x_n دستگاه (۳۴.۶.۸)، (به تعبیری) به جواب x (۳۳.۶.۸) میل می‌کند. بنابراین دستگاههای خطی در (۳۴.۶.۸) بستگی نزدیکی با هم دارند. به عنوان مثال، به تعبیری، برای m و n به اندازه کافی بزرگ $A_m^{-1} \cong A_n^{-1}$ ، اگرچه مرتبه ماتریسها مختلف‌اند. به این مطلب می‌توان معنای دقیقتری داد که ما را به روشهای بارستی دیگری هدایت می‌کند که در آنجا، حل دستگاههای بزرگ خطی با حل دستگاههای کوچکتر امکان‌پذیر می‌شود. اخیراً، چنین روشهایی به نام روشهای چندشبکه‌یی به وجود آمده‌اند که کاربردهایی به‌ویژه در حل معادلات دیفرانسیل جزئی دارند [هاکبوش^۱ و تروتنبرگ^۲ (۱۹۸۲)] را ببینید. برای روشهای بارستی در حل معادلات انتگرالی، پیشرفتهای مربوطه ولی متفاوت را در اتکینسن (۱۹۷۶، بخش ۲ فصل ۴) را ببینید. روشهای چندشبکه‌یی، روشهای بارستی بسیار مؤثر و کارا برای معادلات دیفرانسیل و انتگرالی هستند.

۷.۸ پیشگویی خطا و شتاب

با توجه به (۲۵.۶.۸)، رابطه خطای زیر را داریم

$$x - x^{(m+1)} = M(x - x^{(m)}) \quad m \geq 0 \quad (۱.۷.۸)$$

نحوه همگرایی $x^{(m)}$ به x ، بسته به ویژه‌مقدارها و ویژه‌بردارهای M ، ممکن است کاملاً پیچیده باشد. ولی در بیشتر حالات عملی، رفتار خطا خیلی ساده است: اندازه $\|x - x^{(m)}\|_\infty$ در هر مرحله با یک عامل تقریباً ثابت کاهش می‌یابد، و

$$\|x - x^{(m+1)}\|_\infty \leq c \|x - x^{(m)}\|_\infty \quad (۲.۷.۸)$$

که در آن $c < 1$ و خیلی به $r_\sigma(M)$ وابسته است. برای اندازه‌گیری این ثابت c ، با توجه به (۱.۷.۸) داریم

$$x^{(m+1)} - x^{(m)} = e^{(m)} - e^{(m+1)} = M e^{(m-1)} - M e^{(m)}$$

$$x^{(m+1)} - x^{(m)} = M(x^{(m)} - x^{(m-1)}) \quad m \geq 0 \quad (3.7.8)$$

که موجب می‌شود از

$$c \doteq \frac{\|x^{(m+1)} - x^{(m)}\|_\infty}{\|x^{(m)} - x^{(m-1)}\|_\infty} \quad (4.7.8)$$

یا برای اطمینان بیشتر، از ماکسیمم چندین نسبت پی‌درپی از نوع فوق استفاده کنیم. در بسیاری از کاربردها، این نسبت برای مقادیر بزرگ m ، تقریباً ثابت است.

وقتی این ثابت c به دست آمد، و با فرض (۲.۷.۸)، می‌توانیم خطا در $x^{(m+1)}$ را با استفاده از (۱۳.۵.۸) کراندار نماییم.

$$\|x - x^{(m+1)}\|_\infty \leq \frac{c}{1-c} \|x^{(m+1)} - x^{(m)}\|_\infty \quad (5.7.8)$$

این کران وقتی $c \doteq 1$ ، مسأله‌ساز می‌شود و همگرایی کند خواهد بود. در این حالت، تفاضل $\|x^{(m+1)} - x^{(m)}\|_\infty$ ممکن است از خطای واقعی $\|x - x^{(m+1)}\|_\infty$ خیلی کوچکتر باشد.

مثال دستگاه خطی (۵.۸.۸) از بخش ۸.۸ با استفاده از روش گاوس-زایدل حل شده بود. برای هم‌آهنگی با (۵.۸.۸)، بردار مجهول را u گرفته‌ایم. در (۴.۸.۸)، تابع $f = x^2 y^2$ ، و در (۵.۸.۸) تابع $g = 2(x^2 + y^2)$ است. ناحیهٔ جواب $0 \leq x, y \leq 1$ و اندازهٔ شبکه در هر جهت $h = \frac{1}{16}$ است. این مفروضات، یک مرتبهٔ ۲۲۵ برای دستگاه خطی (۵.۸.۸) به دست می‌داد. حدس اولیهٔ $u^{(0)}$ در این بارست، بر پایهٔ درونیابی دوخطی $f = x^2 y^2$ بر ناحیهٔ $0 \leq x, y \leq 1$ نهاده شده است. [۱۷.۸.۸) را ببینید]. نمونه‌ای از نتایج عددی در جدول ۵.۸ داده شده است. ستون نسبت از روی (۴.۷.۸) محاسبه شده و ستون برآورد خطا از روی (۵.۷.۸)؛ و ستون خطا، خطای درست $\|u - u^{(m)}\|_\infty$ است. همان‌گونه که از جدول می‌توان دید، همگرایی کاملاً آهسته است که نیاز به (۵.۷.۸) را به جای مقدار بسیار بزرگتر $\|u^{(m)} - u^{(m-1)}\|_\infty$ توجیه می‌نماید. وقتی $m \rightarrow \infty$ ، مقدار نسبت به ۹۶۲^۰ همگرا می‌شود، و برآورد خطای (۵.۷.۸) یک برآوردگر دقیق برای خطای درست بارستی است.

سرعت همگرایی اکنون ببینیم که چند بارست باید حساب کنیم تا به خطای مطلوب برسیم. و چه موقع روش بارستی در حل $Ax = b$ بر روش حذف گاوسی برتری دارد؟ مقداری از m

جدول ۵.۸ مثالی از بارست گاوس - زایدل

m	$\ u^{(m)} - u^{(m-1)}\ _{\infty}$	نسبت	برآورد خطا	خطا
۲۰	$۱٫۲۰E-۳$	۰٫۹۶۶	$۳٫۴۲E-۲$	$۳٫۰۹E-۲$
۲۱	$۱٫۱۶E-۳$	۰٫۹۶۶	$۳٫۲۴E-۲$	$۲٫۹۸E-۲$
۲۲	$۱٫۱۲E-۳$	۰٫۹۶۵	$۳٫۰۸E-۲$	$۲٫۸۶E-۲$
۲۳	$۱٫۰۸E-۳$	۰٫۹۶۵	$۲٫۹۳E-۲$	$۲٫۷۶E-۲$
۲۴	$۱٫۰۴E-۳$	۰٫۹۶۴	$۲٫۸۰E-۲$	$۲٫۶۵E-۲$
۶۰	$۲٫۶۰E-۴$	۰٫۹۶۲	$۶٫۵۸E-۳$	$۶٫۵۸E-۳$
۶۱	$۲٫۵۰E-۴$	۰٫۹۶۲	$۶٫۳۳E-۳$	$۶٫۳۳E-۳$
۶۲	$۲٫۴۱E-۴$	۰٫۹۶۲	$۶٫۰۹E-۳$	$۶٫۰۹E-۳$

را پیدا می‌کنیم که برای آن

$$\|x - x^{(m)}\|_{\infty} \leq \varepsilon \|x - x^{(0)}\|_{\infty} \quad (۶.۷.۸)$$

که ε ضریب داده‌شده‌ای است که خطای اولیه باید با آن کاهش یابد. تحلیل را بر پایه فرض (۲.۷.۸) می‌گذاریم. معمولاً ثابت c تقریباً با $r_{\sigma}(M)$ برابر است که M همان است که در (۱.۷.۸) داده شده است.

از رابطه (۲.۷.۸) چنین نتیجه می‌شود

$$\|x - x^{(m)}\|_{\infty} \leq c^m \|x - x^{(0)}\|_{\infty} \quad m \geq 0 \quad (۷.۷.۸)$$

بنابراین کوچکترین مقدار m را که برای آن

$$c^m \leq \varepsilon$$

می‌یابیم. از حل این رابطه به دست می‌آید

$$m \geq \frac{-\ln \varepsilon}{R(c)} = m^* \quad R(c) = -\ln c \quad (۸.۷.۸)$$

دو برابرکردن $R(c)$ منجر به نصف‌کردن تعداد بارست‌هایی خواهد شد که باید محاسبه شود. برای آنکه به این نتیجه معنی بیشتری بدهیم، آن را در حل یک دستگاه خطی چگال به روش بارستی به‌کار می‌بریم. فرض می‌کنیم که روش گاوس-ژاکوبی یا گاوس-زایدل برای حل

جدول ۶.۸ مثال از شمارش بارستی

c	$R(c)$	m^*
۰٫۹	۰٫۱۰۵	۱۳۱
۰٫۸	۰٫۲۲۳	۶۲
۰٫۶	۰٫۵۱۱	۲۷
۰٫۴	۰٫۹۱۶	۱۵
۰٫۲	۱٫۶۱	۹

$Ax = b$ ، با دقت ساده در رایانهٔ بزرگ IBM، یعنی، در حدود ۶ رقم با معنی، به‌کار برده شده باشد. فرض می‌کنیم $x^{(0)} = 0$ ، و می‌خواهیم m را طوری پیدا کنیم که

$$\frac{\|x - x^{(m)}\|_{\infty}}{\|x\|_{\infty}} \leq 10^{-6} = \varepsilon \quad (۹.۷.۸)$$

فرض می‌کنیم که A از مرتبهٔ n باشد، تعداد عملها (ضرب و تقسیم) در هر بارست n^2 است. برای به‌دست آوردن نتیجهٔ (۹.۷.۸)، تعداد بارستهای لازم عبارت است از

$$m^* = \frac{6 \ln_e 10}{R(c)}$$

و تعداد عملها چنین است

$$m^* n^2 = (6 \ln_e 10) \frac{n^2}{R(c)}$$

اگر حذف گاوسی برای حل $Ax = b$ با همان دقت به‌کار برده شده باشد، تعداد عملها در حدود $\frac{n^3}{3}$ است. روش بارستی از روش حذف گاوسی کاراتر خواهد بود اگر

$$\begin{aligned} m^* n^2 &< \frac{n^3}{3} \\ m^* &< \frac{n}{3} \end{aligned} \quad (۱۰.۷.۸)$$

مثال ماتریس A از مرتبهٔ $n = 51$ را در نظر می‌گیریم. پس روش بارستی کاراتر است اگر $m^* < 7$. جدول ۶.۸ مقادیر m^* را به ازای مقادیر مختلف c به‌دست می‌دهد. برای $c \leq 0.44$ ، روش بارستی از روش حذف گاوسی کاراتر خواهد بود. و اگر دقتی کمتر از درجهٔ دقت کامل مورد نظر باشد، روش بارستی حتی برای c های بزرگتر نیز کاراتر است. در عمل، معمولاً از یک حدس اولیهٔ $x^{(0)}$ نیز آگاهی داریم که از $x^{(0)} = 0$ بهتر است، که تعداد بارستهای لازم را بازهم کاهش خواهد داد.

استفاده عمده از روشهای بارستی در حل دستگاههای بزرگ تنک است، که در آنها حالت حذف گاوسی اغلب ناممکن است. و حتی وقتی که حذف گاوسی هم ممکن باشد، روش بارستی باز ممکن است ارجح باشد. چند مثال از این دستگاهها در بخش ۸.۸ داده شده است.

روشهای شتابی اکثر روشهای بارستی الگوی منظمی دارند که در آنها خطا کاهش می یابد. این الگو را اغلب می توان برای تسریع همگرایی به کار برد، درست همان گونه که در فصلهای پیشین برای روشهای عددی دیگر انجام گرفته است. به جای آنکه یک نظریه کلی برای شتاب روشهای بارستی در حل $Ax = b$ بیاوریم، فقط شتاب برای روش گاوس-زایدل را شرح خواهیم داد، که یکی از حالتهای عمده مطاب در کاربردهاست.

تعریف (۱۵.۶.۸) در روش گاوس-زایدل را به یاد می آوریم. یک پارامتر شتاب ω معرفی می کنیم، و (۱۵.۶.۸) را با یک تغییر جزئی در نظر می گیریم:

$$z_i^{(m+1)} = \frac{1}{a_{ii}} \left\{ b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(m+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(m)} \right\}$$

$$x_i^{(m+1)} = \omega z_i^{(m+1)} + (1 - \omega) x_i^{(m)} \quad i = 1, \dots, n \quad m \geq 0 \quad (11.7.8)$$

حالت $m = 1$ همان روش گاوس-زایدل معمولی است. منظور از شتاب دادن، انتخاب بهینه یک ترکیب خطی از بارست قبلی و بارست گاوس-زایدل معمولی است. روش (۱۱.۷.۸)، با انتخاب بهینه ω روش *SOR* خوانده می شود که مخفف اصطلاح تاریخی *successive over relaxation* است.

برای اینکه بفهمیم چگونه ω باید انتخاب شود، (۱۱.۷.۸) را به شکل ماتریسی می نویسیم. A را به شکل زیر تجزیه می کنیم

$$A = D + L + U$$

که $D = \text{diag}[a_{11}, \dots, a_{nn}]$ ، L ماتریس پایین مثلثی، U ماتریس بالامثلثی است و قطره های اصلی L و U هر دو صفرند. در این صورت (۱۱.۷.۸) چنین خواهد شد:

$$z^{(m+1)} = D^{-1}[b - Lx^{(m+1)} - Ux^{(m)}]$$

$$x^{(m+1)} = \omega z^{(m+1)} + (1 - \omega)x^{(m)} \quad m \geq 0$$

از حذف $z^{(m+1)}$ و حل نسبت به $x^{(m+1)}$ داریم

$$[I + \omega D^{-1}L]x^{(m+1)} = \omega D^{-1}b + [(1 - \omega)I - \omega D^{-1}U]x^{(m)}$$

برای خطا،

$$e^{(m+1)} = M(\omega)e^{(m)} \quad m \geq 0 \quad (12.7.8)$$

$$M(\omega) = [I + \omega D^{-1}L]^{-1}[(1 - \omega)I - \omega D^{-1}U] \quad (13.7.8)$$

پارامتر ω باید طوری انتخاب شود که $r_\sigma(M(\omega))$ را مینیمم کند، تا اینکه $x^{(m)}$ هرچه سریعتر به x همگرا شود. این مقدار بهینه را ω^* می‌نامیم.

محاسبه ω^* جز در ساده‌ترین حالتها دشوار است، و معمولاً فقط به‌طور تقریبی، با انتخاب چند مقدار برای ω و توجه به اثر آن روی سرعت همگرایی به دست می‌آید. علی‌رغم مسأله محاسبه ω^* ، افزایش به دست آمده در سرعت همگرایی $x^{(m)}$ به x ، بسیار چشمگیر است و سعی در محاسبه ω^* کاملاً ارزش دارد. این موضوع در بخش بعد نشان داده شده است.

مثال ما شتاب (۱۱.۷.۸) در مثال قبلی را برای روش گاوس-زایدل که به دنبال (۵.۷.۸) آمده به کار می‌بریم. مقدار پارامتر بهینه شتاب $\omega^* = 1.6735$ است. یک بحث کامل از روش SOR برای حل دستگاه خطی که در حل معادلات دیفرانسیل جزئی پیش می‌آید در بخش بعدی ارائه شده است. حدس اولیه مانند حدس قبل بوده است. نتایج در جدول ۷.۸ داده شده‌اند.

نتایج نشان می‌دهند که نرخ همگرایی بسیار سریعتر از نرخ همگرایی در روش گاوس-زایدل است. به‌عنوان مثال، با روش گاوس-زایدل، داریم $\|u - u^{(28)}\|_\infty = 9.7 \times 10^{-6}$ ، ولی، همان‌طور که می‌توان از مقادیر نسبت SOR داریم، $\|u - u^{(22)}\|_\infty = 8.71 \times 10^{-6}$. ولی، همان‌طور که می‌توان از مقادیر نسبت ملاحظه کرد، رفتار منظم همگرایی بارستنها را از دست داده‌ایم. مقدار c که در آزمون خطای (۵.۷.۸) به کار برده شده باید با دقتی بیشتر از نسبت c که در جدول به کار برده شده، انتخاب شود. می‌توانید از یک میانگین یا ماکسیمم چندین مقدار پی‌درپی نسبت استفاده نمایید.

۸.۸ حل عددی معادله پواسن

مهمترین کاربرد روشهای بارستی خطی، در دستگاههای خطی بزرگی است که از حل عددی معادلات دیفرانسیل جزئی به روش تفاضلات متناهی، به‌وجود می‌آیند. برای نشان دادن این موضوع، مسأله دیریکله را برای معادله پواسن بر یک مربع واحد در صفحه xy حل می‌کنیم:

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = g(x, y) \quad 0 < x, y < 1 \quad (1.8.8)$$

$$u(x, y) = f(x, y) \quad \text{یک نقطه مرزی}$$

جدول ۷.۸ مثال روش SOR (۱۱.۷.۸)

m	$\ u^{(m)} - u^{(m-1)}\ _\infty$	نسبت	برآورد خطا	خطا
۲۱	$۲.۰۶E-۴$	۰.۶۹۳	$۴.۶۵E-۴$	$۳.۶۴E-۴$
۲۲	$۱.۳۵E-۴$	۰.۶۵۷	$۲.۵۹E-۴$	$۲.۶۵E-۴$
۲۳	$۸.۷۶E-۵$	۰.۶۴۸	$۱.۶۱E-۴$	$۱.۸۷E-۴$
۲۴	$۵.۱۱E-۵$	۰.۵۸۴	$۷.۱۷E-۵$	$۱.۳۹E-۴$
۲۵	$۳.۴۸E-۵$	۰.۶۸۰	$۷.۴۰E-۵$	$۱.۰۶E-۴$
۲۶	$۲.۷۸E-۵$	۰.۸۰۰	$۱.۱۱E-۴$	$۸.۰۴E-۵$
۲۷	$۲.۴۶E-۵$	۰.۸۸۴	$۱.۸۷E-۴$	$۶.۱۵E-۵$
۲۸	$۲.۰۷E-۵$	۰.۸۴۲	$۱.۱۱E-۴$	$۴.۱۶E-۵$

توابع $g(x, y)$ و $f(x, y)$ داده شده‌اند و باید $u(x, y)$ را پیدا کنیم. برای $N > ۱$ تعریف می‌کنیم $h = ۱/N$ و

$$(x_j, y_k) = (jh, kh) \quad 0 \leq j, k \leq N$$

این نقاط را نقاط گرهی یا نقاط شبکه‌یی خوانند (شکل ۱.۸ را ببینید). برای تقریب‌زدن (۱.۸.۸) تقریبهایی را برای مشتقات دوم به‌کار می‌بریم. برای یک تابع چهار بار پیوسته مشتقپذیر $G(x)$ در $[x-h, x+h]$ ، نتایج (۱۷.۷.۵) و (۱۸.۷.۵) از بخش (۷.۵) چنین به‌دست می‌دهند،

$$G''(x) = \frac{G(x+h) - 2G(x) + G(x-h)}{h^2} - \frac{h^2}{12} G^{(4)}(\xi) \quad x-h \leq \xi \leq x+h \quad (۲.۸.۸)$$

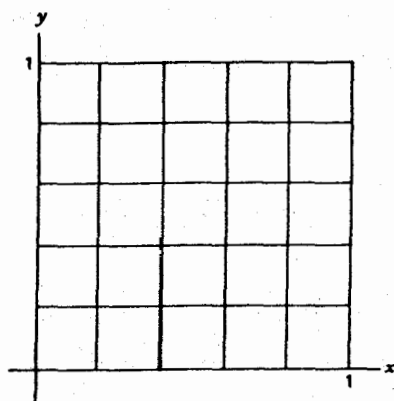
وقتی تقریب فوق را در (۱.۸.۸) در هر نقطه شبکه‌یی داخلی به‌کار ببریم، به‌دست می‌آوریم

$$\frac{u(x_{j+1}, y_k) - 2u(x_j, y_k) + u(x_{j-1}, y_k))}{h^2} + \frac{u(x_j, y_{k+1}) - 2u(x_j, y_k) + u(x_j, y_{k-1}))}{h^2} = g(x_j, y_k) + \frac{h^2}{12} \left\{ \frac{\partial^4 u(\xi, y_k)}{\partial x^4} + \frac{\partial^4 u(x_j, \eta)}{\partial y^4} \right\} \quad (۳.۸.۸)$$

برای مقداری از $1 \leq j, k \leq N-1$ و $y_{k-1} \leq \eta \leq y_{k+1}$ ، $x_{j-1} \leq \xi \leq x_{j+1}$ ، برای تقریب عددی $u_h(x, y)$ مربوط به (۱.۸.۸)، گیریم

$$u_h(x_j, y_k) = f(x_j, y_k) \quad (x_j, y_k) \text{ گرهی مرزی } \quad (۴.۸.۸)$$

در همه نقاط شبکه‌یی داخلی، خطاهای برشی طرف راست در (۳.۸.۸) را حذف و مسأله را



شکل ۱.۸ شبکهٔ تفصیل منتهی

نسبت به جواب تقریبی $u_h(x_j, y_k)$ حل می‌کنیم:

$$u_h(x_j, y_k) = \frac{1}{4} \{u_h(x_{j+1}, y_k) + u_h(x_j, y_{k+1}) + u_h(x_{j-1}, y_k) + u_h(x_j, y_{k-1})\} - \frac{h^2}{4} g(x_j, y_k) \quad 1 \leq j, k \leq N-1 \quad (5.1.8)$$

تعداد معادلات در (۴.۸.۸) - (۵.۸.۸) برابر است با تعداد مجهولات، $(N+1)^2$.

قضیهٔ ۸.۸ برای هر $N \geq 2$ ، دستگاه خطی (۴.۸.۸) - (۵.۸.۸) دارای یک جواب یکتای $\{u_h(x_j, y_k) \mid 0 \leq j, k \leq N\}$ است. اگر $u(x, y)$ ، جواب (۱.۸.۸)، چهاربار پیوسته مشتقپذیر باشد، آنگاه

$$\begin{aligned} \text{Max}_{0 \leq j, k \leq n} |u(x_j, y_k) - u_h(x_j, y_k)| &\leq ch^2 \quad (6.1.8) \\ c &= \frac{1}{24} \left\{ \text{Max}_{0 \leq x, y \leq 1} \left| \frac{\partial^2 u(x, y)}{\partial x^2} \right| + \text{Max}_{0 \leq x, y \leq 1} \left| \frac{\partial^2 u(x, y)}{\partial y^2} \right| \right\} \end{aligned}$$

برهان ۱. حلپذیری یکتای (۴.۸.۸) - (۵.۸.۸) را با استفاده از قضیهٔ ۲.۷، فصل ۷، ثابت می‌کنیم. دستگاه همگن زیر را در نظر می‌گیریم

$$v_h(x_j, y_k) = \frac{1}{4} [v_h(x_{j+1}, y_k) + v_h(x_j, y_{k+1}) + v_h(x_{j-1}, y_k) + v_h(x_j, y_{k-1})] \quad 1 \leq j, k \leq N-1 \quad (7.1.8)$$

$$v_h(x_j, y_k) = 0 \quad \text{یک نقطه مرزی} \quad (8.1.8)$$

با نشان دادن اینکه این دستگاه فقط جواب نمایان $v_h(x_j, y_k) \equiv 0$ را دارد، از قضیه ۲.۷ نتیجه می‌شود که دستگاه ناهمگن (۴.۸.۸) - (۵.۸.۸) فقط یک جواب یکتا دارد.

گیریم

$$\alpha = \max_{0 \leq j, k \leq N} v_h(x_j, y_k)$$

با توجه به (۸.۸.۸)، $\alpha \geq 0$. فرض می‌کنیم $\alpha > 0$. پس یک نقطه گرهی داخلی (\bar{x}_j, \bar{y}_k) باید وجود داشته باشد که در آن این ماکسیمم رخ دهد. ولی با استفاده از (۷.۸.۸)، میانگین مقادیر v_h در چهار نقطه مجاور (\bar{x}_j, \bar{y}_k) است. تنها راهی که میانگین بودن با ماکسیمم بودن (\bar{x}_j, \bar{y}_k) سازگار باشد آن است که v_h در چهار نقطه شبکه‌یی مجاور نیز برابر α باشد. همین استدلال را برای نقاط مجاور هم ادامه می‌دهیم. چون فقط تعدادی متناهی نقطه شبکه‌یی وجود دارند، سرانجام برای یک نقطه مرزی $(x_j, y_h) = \alpha$ داریم $v_h(x_j, y_h) = \alpha$ ولی در این صورت $\alpha > 0$ با شرط (۸.۸.۸) تناقض پیدا می‌کند. بنابراین ماکسیمم $v_h(x_j, y_k)$ برابر صفر است. یک استدلال مشابه نشان می‌دهد که مینیمم $v_h(x_j, y_k)$ نیز برابر صفر است. این نتایج وقتی با هم در نظر گرفته شوند معلوم می‌شود که تنها جواب (۷.۸.۸) - (۸.۸.۸)، $v_h(x_j, y_k) \equiv 0$ است.

۲. برای ملاحظه همگرایی $u_h(x_j, y_k)$ به $u(x_j, y_k)$ ، تعریف می‌کنیم

$$e_h(x_j, y_k) = u(x_j, y_k) - u_h(x_j, y_k)$$

با کم کردن (۵.۸.۸) از (۳.۸.۸) به دست می‌آوریم

$$e_h(x_j, y_k) = \frac{1}{4} [e_h(x_{j+1}, y_k) + e_h(x_j, y_{k+1}) + e_h(x_{j-1}, y_k) + e_h(x_j, y_{k-1})] - \frac{h^2}{12} \left[\frac{\partial^2 u(x_j, y_k)}{\partial x^2} + \frac{\partial^2 u(x_j, y_k)}{\partial y^2} \right] \quad (۹.۸.۸)$$

و از (۴.۸.۸)،

$$e_h(x_j, y_k) = 0 \quad \text{یک نقطه مرزی} \quad (۱۰.۸.۸)$$

با این دستگاه می‌توان مشابه با قسمت (۱) عمل کرد، نتیجه (۶.۸.۸) به دست خواهد آمد. ولی چون این امر برای بحث دستگاههای خطی اساسی نیست، از استدلال آن صرف نظر می‌شود [ایزکسون و کلر (۱۹۶۶)، صص ۴۴۷-۴۵۰] را ببینید. ■

مثال مسأله زیر را حل کنید

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 0 \quad 0 \leq x, y \leq 1 \quad (11.8.8)$$

$$u(0, y) = \cos(\pi y) \quad u(1, y) = e^\pi \cos(\pi y)$$

$$u(x, 0) = e^{\pi x} \quad u(x, 1) = -e^{\pi x}$$

جواب درست $u(x, y) = e^{\pi x} \cos(\pi y)$ است.

نتایج عددی برای چندین مقدار N در جدول ۸.۸ داده شده است. خطا، در همه نقاط شبکه ماکسیمم است، و ستون نسبت عامل کاهش خطای ماکسیمم وقتی h نصف شود را می دهد. از لحاظ نظری، به موجب (۶.۸.۸)، این عامل باید برابر 4^r باشد. نتایج عددی این را تایید می کند.

راه حل بارستی چون روش گاوس - زایدل معمولاً سریعتر از روش گاوس - ژاکوبی است، ما فقط همان روش اولی را در نظر می گیریم. برای $k = 1, 2, \dots, N - 1$ تعریف می کنیم

$$u_h^{(m+1)}(x_j, y_k) = \frac{1}{4} [u_h^{(m)}(x_{j+1}, y_k) + u_h^{(m)}(x_j, y_{k+1}) + u_h^{(m+1)}(x_{j-1}, y_k) + u_h^{(m+1)}(x_j, y_{k-1})] - \frac{h^2}{4} g(x_j, y_k) \quad j = 1, 2, \dots, N - 1 \quad (12.8.8)$$

برای نقاط مرزی، از رابطه زیر استفاده می کنیم

$$u_h^{(m)}(x_j, y_k) = f(x_j, y_k) \quad m \geq 0 \text{ برای هر}$$

مقادیر $u_h^{(m+1)}(x_j, y_k)$ سطر به سطر حساب می شوند، ابتدا، از سطر آخر نقاط شبکه تا سطر اول نقاط شبکه. و در هر سطر، از چپ به راست.

جدول ۸.۸ حل عددی (۱۱.۸.۸)

N	$\ u - u_h\ _\infty$	نسبت
۴	۰٫۱۴۴	
۸	۰٫۰۳۹۰	۳٫۷
۱۶	۰٫۰۱۰۲	۳٫۸
۳۲	۰٫۰۰۲۶۰	۳٫۹
۶۴	۰٫۰۰۰۶۵۴	۴٫۰

برای بارست (۱۲.۸.۸)، تحلیل همگرایی باید بر پایه قضیه ۷.۸ صورت گیرد. به آسانی می توان نشان داد که این ماتریس متقارن است و بنابراین ویژه مقادارها حقیقی اند. به علاوه می توان نشان داد که ویژه مقادارها در بازه $2 < \lambda < 0$ قرار دارند. با توجه به این نکته و مسأله ۱۴، فصل ۷، این ماتریس معین مثبت است. چون همه ضرایب قطری ماتریس مثبت اند، از قضیه ۷.۸ نتیجه می شود که روش گاوس-زایدل همگرا خواهد بود. برای نشان دادن اینکه ویژه مقادارها در بازه $2 < \lambda < 0$ قرار دارند آیزکسون و کالر (۱۹۶۶، صص ۴۴۸-۴۵۹) یا مسأله ۳ فصل ۹ را ببینید.

محاسبه سرعت همگرایی $r_\sigma(M)$ از روی (۲۸.۶.۸) آسان نیست. استدلال کاملاً پیچیده است، و فقط خواننده را به بحث کامل آن در آیزکسون و کالر (۱۹۶۶، صص ۴۶۳-۴۷۰)، شامل مطالب مربوط به روش شتاب گاوس-زایدل ارجاع می دهیم. می توان نشان داد که

$$r_\sigma(M) = 1 - \pi^2 h^2 + O(h^4) \quad (13.8.8)$$

روش گاوس-زایدل همگراست ولی سرعت همگرایی، حتی برای مقدار نسبتاً کوچک h خیلی آهسته است. این مطلب در جدول ۵.۸ بخش ۷.۸ نشان داده شده است. برای شتاب دادن به روش گاوس-زایدل از رابطه زیر استفاده می کنیم:

$$\begin{aligned} v_h^{(m+1)}(x_j, y_k) &= \frac{1}{4} [u_h^{(m)}(x_{j+1}, y_k) + u_h^{(m)}(x_j, y_{k+1}) \\ &\quad + u_h^{(m+1)}(x_{j-1}, y_k) + u_h^{(m+1)}(x_j, y_{k-1})] - \frac{h^2}{4} g(x_j, y_k) \\ u_h^{(m+1)}(x_j, y_k) &= \omega v_h^{(m+1)}(x_j, y_k) \\ &\quad + (1 - \omega) u_h^{(m)}(x_j, y_k) \quad j = 1, \dots, N-1 \end{aligned} \quad (14.8.8)$$

پارامتر بهینه شتاب عبارت است از

$$\begin{aligned} \omega^* &= \frac{2}{1 + \sqrt{1 - \xi^2}} \\ \xi &= 1 - 2 \sin^2 \left(\frac{\pi}{2N} \right) \end{aligned} \quad (15.8.8)$$

نرخ همگرایی متناظر، چنین است.

$$r_\sigma(M(\omega^*)) = \omega^* - 1 = 1 - 2h\pi + O(h^2) \quad (16.8.8)$$

این نرخ همگرایی خیلی بهتر از نرخی است که با (۱۳.۸.۸)، داده شده است. روش شتاب داده شده (۱۴.۸.۸) گاوس-زایدل با مقدار بهینه ω^* در (۱۵.۸.۸)، روش *SOR* خوانده می شود. نام

SOR مخفف *successive over relaxation* است، نامی است بر پایهٔ تعبیر فیزیکی این روش، که ابتدا در پیدا کردن آن به‌کار برده شده است.

مثال مثال قبلی (۱۱.۸.۸) را به‌یاد آورید. این مثال هم با روش گاوس-زایدل حل شده بود و هم با روش SOR. حدس اولیه برای بارست، از فرمول درونیایی دو خطی برای داده‌های مرزی f ، یعنی

$$u_h^{(\circ)}(x, y) = (1-x)f(\circ, y) + xf(1, y) + (1-y)f(x, \circ) + yf(x, 1) \\ - [(1-y)(1-x)f(\circ, \circ) + (1-y)xf(1, \circ) \\ + y(1-x)f(\circ, 1) + xyf(1, 1)] \quad (17.8.8)$$

در تمام نقاط شبکه‌یی داخلی گرفته شده است. آزمون خطا برای توقف بارست چنین بوده است،

$$\max_{1 \leq j, k \leq N-1} |u_h(x_j, y_k) - u_h^{(m)}(x_j, y_k)| \leq \varepsilon$$

که $\varepsilon > 0$ داده شده و طرف راست (۵.۷.۸) برای پیشگویی خطا در بارست به‌کار برده شده است. نتایج عددی برای تعداد لازم بارستها در جدول ۹.۸ داده شده‌اند. روش SOR به بارستهای بسیار کمتری نسبت به روش گاوس-زایدل برای مقادیر کوچکتر h نیاز دارد.

از بخش قبل یادآوری می‌کنیم که تعداد بارستهای لازم، که m^* نامیده شده بود، برای کاهش خطای بارستی با عامل ε ، متناسب است با $1/\ln(c)$ ، که c نسبتی است که خطای بارستی در هر مرحله کاهش می‌یابد. برای روشهای این بخش، c را، $c = r_\sigma(M)$ می‌گیریم. اگر بنویسیم $1 - \delta = r_\sigma(M)$ ، آنگاه

$$\frac{1}{-\ln r_\sigma(M)} = \frac{1}{-\ln(1 - \delta)} \doteq \frac{1}{\delta}$$

جدول ۹.۸ تعداد بارستهای لازم برای حل (۵.۸.۸)

N	ε	گاوس-زایدل	SOR
۸	0.01	۲۵	۱۲
۸	0.001	۴۰	۱۶
۱۶	0.001	۱۴۲	۳۲
۳۲	0.001	۴۹۵	۶۵
۸	0.0001	۵۴	۱۸
۱۶	0.0001	۲۰۱	۳۵
۳۲	0.0001	۷۳۳	۷۱

وقتی δ نصف شود، تعداد بارستهای محاسبه دو برابر می‌شود. در روش گاوس - زایدل

$$\frac{1}{\ln r_{\sigma}(M)} \doteq \frac{1}{\pi^2 h^2}$$

وقتی h نصف شود تعداد بارستهای محاسبات با ضریب ۴ افزایش می‌یابد. در روش SOR

$$\frac{1}{\ln r_{\sigma}(M(\omega^*))} \doteq \frac{1}{2\pi h}$$

و وقتی که h نصف شود، تعداد بارستها، دو برابر می‌شود. این دو نتیجه برای داده‌های $\varepsilon = 10^{-2}$ و $\varepsilon = 10^{-4}$ در جدول ۹.۸ داده شده‌اند.

ملاحظه می‌کنید که دو برابر کردن N در هر روش تعداد معادلاتی را که باید حل شوند ۴ برابر افزایش می‌دهد و لذا کار در هر بارست نیز به همین اندازه افزایش می‌یابد. استفاده از SOR کار را بسیار تقلیل می‌دهد اگرچه باز هم وقتی N بزرگ باشد زیاد است.

۹.۸ روش گرادیان مزدوج

روش بارستی که در این بخش ارائه شد، در دهه ۱۹۵۰ به وجود آمده، ولی شهرت عمده خود را در سالهای اخیر به دست آورده است، به ویژه در حل دستگاههای خطی که از حل عددی معادلات دیفرانسیل جزئی خطی به دست آمده‌اند. نوشته‌ها درباره روش گرادیان مزدوج (روش CG) بسیار گسترده و پیچیده شده‌اند و با سایر مباحث جبر خطی مربوط هستند. بنابراین به لحاظ کمبود جا، می‌توانیم فقط به یک مقدمه کوتاه، شامل تعریف روش گرادیان مزدوج و ذکر بعضی نتایج نظری اساسی مربوط به آن اکتفا کنیم.

روش گرادیان مزدوج با روشهای دیگر این فصل از این نظر تفاوت دارد که این روش برای حل مسائل غیرخطی پایه‌گذاری شده است؛ در واقع، روش گرادیان مزدوج روشی است که برای مینیمم‌سازی توابع غیرخطی مورد استفاده قرار گرفته است. در دستگاه خطی که باید حل شود،

$$Ax = b \quad (۱.۹.۸)$$

فرض می‌شود ماتریس ضرایب A حقیقی، متقارن، و معین مثبت است. حل این دستگاه هم‌ارز مینیمم‌سازی تابع زیر است:

$$f(x) = \frac{1}{2} x^T A x - b^T x \quad x \in \mathbf{R}^n \quad (۲.۹.۸)$$

جواب یکتای x^* معادله $Ax = b$ نیز مینیمم‌کننده یکتای $f(x)$ است، وقتی x در \mathbf{R}^n تغییر می‌کند. برای مشاهده این موضوع، ابتدا نشان می‌دهیم که

$$f(x) = E(x) - \frac{1}{4} b^T x^* \quad (3.9.8)$$

$$E(x) \equiv \frac{1}{4} (x^* - x)^T A (x^* - x)$$

با استفاده از $Ax^* = b$ ، اثبات سراسر است. تفاوت $E(x)$ و $f(x)$ یک مقدار ثابت است، پس هر دو یک مینیمم‌ساز دارند. به موجب معین مثبت بودن A ، $E(x)$ فقط در $x = x^*$ مینیمم می‌شود، و بنابراین $f(x)$ نیز فقط در x^* مینیمم است.

یک روش بارستی معروف برای پیدا کردن مینیمم یک تابع غیرخطی، تندترین کاهش است، که به اختصار در بخش ۱۲.۲ توضیح داده شده است. برای مینیمم کردن $f(x)$ با این روش، فرض می‌کنیم که یک حدس اولیه x_0 داده شده است. مسیری را انتخاب می‌کنیم که در x_0 ، $f(x)$ تندترین کاهش را داشته باشد. این مسیر با $g_0 = -\nabla f(x_0)$ ، گرادیان $f(x)$ در x_0 با علامت منفی، داده می‌شود:

$$g(x_0) \equiv g_0 = b - Ax_0. \quad (4.9.8)$$

سپس، مسأله مینیمم‌سازی یک‌بعدی

$$\text{Min}_{0 \leq \alpha \leq \infty} f(x_0 + \alpha g_0)$$

را حل می‌کنیم و جواب آن را α_1 می‌نامیم. با استفاده از این جواب بارست جدید زیر را تعریف می‌کنیم

$$x_1 = x_0 + \alpha_1 g_0. \quad (5.9.8)$$

این فرایند را به‌طور استقرایی ادامه می‌دهیم. روش تندترین کاهش همگرا می‌شود ولی همگرایی معمولاً خیلی آهسته است. تدبیر بهینه موضعی، با استفاده از امتداد تندترین کاهش، یک تدبیر خوب برای پیدا کردن یک امتداد بهینه در پیدا کردن مینیمم کلی نخواهد بود. در مقایسه، روش گرادیان مزدوج سریعتر خواهد بود، و اگر فرض کنیم خطای گرد کردن وجود ندارد، به بیشتر از n بارست نیاز نخواهد داشت. برای بقیه این بخش، فرض می‌کنیم که حدس اولیه $x_0 = 0$ است. اگر صفر نباشد، همیشه می‌توانیم مسأله کمی تغییر یافته زیر را حل کنیم

$$Az = b - Ax_0.$$

با نشان دادن جواب آن با z^* ، داریم $z^* = x_0 + z^*$. یک حدس اولیه $z_0 = z^*$ متناظر با $x_0 = 0$ در مسأله اصلی خواهد بود. از این پس فرض می‌کنیم $x_0 = 0$.

روشهای امتداد مزدوج. فرض می‌کنیم A ، $n \times n$ باشد، گوییم در \mathbf{R}^n یک مجموعه از بردارهای
 ناصفر p_1, \dots, p_n - مزدوج است اگر

$$p_i^T A p_j = 0 \quad 1 \leq i, j \leq n \quad i \neq j \quad (۶.۹.۸)$$

بردارهای p_j اغلب امتدادهای مزدوج نامیده می‌شوند. با معرفی یک حاصلضرب داخلی و یک
 نرم جدید برای \mathbf{R}^n ، می‌توان تعریف هندسی هم‌ارزی ارائه داد:

$$(x, y)_A = x^T A y \quad (۷.۹.۸)$$

$$\|x\|_A = \sqrt{(x, x)_A} = \sqrt{x^T A x} \quad x \in \mathbf{R}^n$$

شرط (۶.۹.۸) هم‌ارز آن است که خواسته شود p_1, \dots, p_n یک پایه متعامد برای \mathbf{R}^n نسبت به
 ضرب داخلی $(\cdot, \cdot)_A$ باشد. از این رو، همچنین می‌گوییم که $\{p_1, \dots, p_n\}$ - متعامدند اگر
 در (۶.۹.۸) صدق کنند. مشاهده می‌شود که تابع $E(x)$ در (۳.۹.۸) با نرم $\|\cdot\|_A$ چنین است:

$$E(x) = \frac{1}{2} \|x^* - x\|_A^2 \quad (۸.۹.۸)$$

که به نحو روشنتری یک اندازه برای خطای $x - x^*$ است. رابطه بین $\|x\|_A$ و $\|x\|_2$ در
 مسأله ۳۶ بررسی شده است.

اگر یک مجموعه از امتدادهای مزدوج $\{p_1, \dots, p_n\}$ داده شده باشد، حل $Ax = b$
 سراسر است. گیریم

$$x^* = \alpha_1 p_1 + \dots + \alpha_n p_n$$

با استفاده از (۶.۹.۸)،

$$\alpha_k = \frac{p_k^T A x^*}{p_k^T A p_k} = \frac{p_k^T b}{p_k^T A p_k} \quad k = 1, \dots, n. \quad (۹.۹.۸)$$

ما این فرمول را برای x^* به‌کار می‌بریم تا روش بارستی امتداد مزدوج را معرفی کنیم.
 گیریم $x_0 = 0$

$$x_k = \alpha_1 p_1 + \dots + \alpha_k p_k \quad 1 \leq k \leq n \quad (۱۰.۹.۸)$$

مانده x_k در $Ax = b$ را به شکل زیر معرفی می‌کنیم:

$$r_k = b - A x_k = -\nabla f(x_k)$$

به سادگی دیده می شود که $r_0 = b$ و

$$x_k = x_{k-1} + \alpha_k p_k \quad r_k = r_{k-1} - \alpha_k A p_k \quad k = 1, \dots, n \quad (11.9.8)$$

به ازای $k = n$ داریم $x_n = x^*$ و $r_n = 0$ و x_k و r_n ممکن است به ازای $k < n$ برابر x^* باشد.

لم ۱ بردار r_k بر p_1, \dots, p_k متعامد است، یعنی، $1 \leq i \leq k$ ، $r_k^T p_i = 0$.
برهان مسأله ۳۷ را ببینید.

لم ۲ (الف) مسأله مینیم سازی

$$\text{Min}_{-\infty < \alpha < \infty} f(x_{k-1} + \alpha p_k)$$

با $\alpha = \alpha_k$ به طور یکتا حل می شود و مقدار مینیم $f(x_k)$ به دست می آید.

(ب) گیریم $\mathcal{S}_k = \text{Span}\{p_1, \dots, p_k\}$ ، زیر فضای k -بعدی باشد که از $\{p_1, \dots, p_k\}$ به وجود می آید. در این صورت مسأله

$$\text{Min}_{x \in \mathcal{S}_k} f(x)$$

به طور یکتا با $x = x_k$ حل می شود و مقدار مینیم، $f(x_k)$ را به دست می دهد.

برهان (الف): $\varphi(\alpha) \equiv f(x_{k-1} + \alpha p_k)$ را بسط می دهیم:

$$\varphi(\alpha) = f(x_{k-1}) + \alpha p_k^T A x_{k-1} + \frac{1}{2} \alpha^2 p_k^T A p_k - \alpha b^T p_k$$

جمله $p_k^T A x_{k-1}$ برابر صفر است، زیرا $x_{k-1} \in \mathcal{S}_{k-1}$ و p_k ، A -متعامد بر \mathcal{S}_{k-1} است. از

حل $\varphi'(\alpha) = 0$ ، $\alpha = \alpha_k$ نتیجه می شود که برهان را کامل می کند.

(ب) برای هر $h \in \mathcal{S}_k$ ، تابع $f(x_k + h)$ را بسط می دهیم:

$$\begin{aligned} f(x_k + h) &= f(x_k) + h^T A x_k + \frac{1}{2} h^T A h - h^T b \\ &= f(x_k) + \frac{1}{2} h^T A h - h^T r_k \end{aligned}$$

با توجه به لم ۱ و فرض $h \in \mathcal{S}_k$ نتیجه می شود که $h^T r_k = 0$. بنابراین چون A معین مثبت است،

$$f(x_k + h) = f(x_k) + \frac{1}{2} h^T A h \geq f(x_k)$$

اگر $h = 0$ گرفته شود، مینیم به طور یکتا در \mathcal{S}_k به دست می آید که (ب) را ثابت می کند. ■

لم ۲ یک ویژگی بهینگی برای روشهای امتدادهای مزدوج را به دست می دهد که با (۱۱.۹.۸) و (۹.۹.۸) تعریف شده اند. مسأله، دانستن چگونگی انتخاب امتدادهای مزدوج $\{p_j\}$ است. انتخابهای زیادی وجود دارند که بعضی به روشهای معلوم برای حل مستقیم $Ax = b$ می انجامند [استوارت^۱ (۱۹۷۳b) را ببینید].

روش گرادیان مزدوج. ما راهی ارائه می دهیم که همزمان امتدادهای $\{p_k\}$ و بارستهای $\{x_k\}$ را تولید کند. برای اولین امتداد p_1 ، امتداد تندترین کاهش را به کار می بریم:

$$p_1 = -\nabla f(x_0) = r_0 = b \quad (12.9.8)$$

زیرا $x_0 = 0$. یک ساختمان استقرایی برای بقیه امتدادها می دهیم. فرض می کنیم x_0, \dots, x_k همراه با امتدادهای p_1, \dots, p_k تولید شده باشند. یک جهت جدید p_{k+1} باید انتخاب شود که با p_1, \dots, p_k, A - مزدوج باشد. همچنین فرض می کنیم $x_k \neq x^*$ و بنابراین $r_k \neq 0$ در غیر این صورت جواب x^* به دست آمده و کار تمام شده است.

به موجب لم ۱، r_k بر \mathcal{S}_k متعامد است، و بنابراین r_k به \mathcal{S}_k تعلق ندارد. r_k را برای تولید p_{k+1} ، با انتخاب یک مولفه r_k ، به کار می بریم. به دلایلی که ذکر آن در اینجا دشوار است، کافی است که

$$p_{k+1} = r_k + \beta_{k+1} p_k \quad (13.9.8)$$

را در نظر بگیریم. پس شرط $p_k^T A p_{k+1} = 0$ ایجاب می کند که

$$\beta_{k+1} = -\frac{p_k^T A r_k}{p_k^T A p_k} \quad (14.9.8)$$

مخرج مخالف صفر است زیرا A معین مثبت است و $p_k \neq 0$. می توان نشان داد [لوتنبرگر (۱۹۸۴)، ص ۲۴۵] که این تعریف در رابطه زیر نیز صدق می کند

$$p_j^T A p_{k+1} = 0 \quad j = 1, 2, \dots, k-1 \quad (15.9.8)$$

که نشان می دهد $\{p_1, \dots, p_{k+1}\}, A$ - مزدوج است.

روش گرادیان مزدوج عبارت است از انتخاب $\{p_k\}$ از (۱۲.۹.۸) - (۱۴.۹.۸) و $\{x_k, r_k\}$ از (۱۱.۹.۸) و (۹.۹.۸). اگر از خطاهای گردکردن صرف نظر شود، این روش با n بارست یا کمتر

همگرا می‌شود. سرعت واقعی همگرایی تا اندازه زیادی با ویژه‌مقدارهای A تغییر می‌کند. تحلیل خطای روش گرادیان مزدوج بر پایه قضیهٔ بهینگی زیر استوار است:

قضیهٔ ۹.۸. بارستهای $\{x_k\}$ در روش گرادیان مزدوج در رابطهٔ زیر صدق می‌کنند

$$\|x^* - x_k\|_A = \min_{\deg(q) < k} \|x^* - q(A)b\|_A \quad (۱۶.۹.۸)$$

برهان برای یک چندجمله‌یی $q(\lambda)$ ، نماد $q(A)$ یک عبارت ماتریسی است که در آن به جای هر توان λ^i ، A^i گذاشته شده است. برای مثال

$$q(\lambda) = a_0 + a_1\lambda + a_2\lambda^2 \Rightarrow q(A) = a_0I + a_1A + a_2A^2$$

برهان (۱۶.۹.۸) در لوثنبرگر (۱۹۸۴، ص ۲۴۶) داده شده است. ■

با استفاده از این قضیه، تعدادی از نتیجه‌های خطا داده می‌شوند، که با ویژگیهای ماتریس A تغییر می‌یابند. برای مثال گیریم که ویژه‌مقدارهای A چنین باشند

$$0 < \lambda_1 \leq \dots \leq \lambda_n \quad (۱۷.۹.۸)$$

که برحسب چندگانگی تکرار شده‌اند، و گیریم v_1, \dots, v_n یک پایهٔ یکا متعامد مربوط به ویژه‌بردارها باشند، می‌نویسیم

$$x^* = \sum_{j=1}^n c_j v_j \quad b = Ax^* = \sum_{j=1}^n c_j \lambda_j v_j \quad (۱۸.۹.۸)$$

پس

$$q(A)b = \sum_{j=1}^n c_j \lambda_j q(\lambda_j) v_j$$

$$\|x^* - q(A)b\|_A = \left[\sum_{j=1}^n c_j^2 \lambda_j [1 - \lambda_j q(\lambda_j)]^2 \right]^{1/2} \quad (۱۹.۹.۸)$$

هر انتخاب چندجمله‌یی $q(\lambda)$ از درجهٔ کمتر از k یک کران برای $\|x^* - x_k\|_A$ به دست می‌دهد. یکی از بهترین کرانهایی که می‌شناسیم

$$\|x^* - x_k\|_A \leq 2 \left[\frac{1 - \sqrt{c}}{1 + \sqrt{c}} \right]^k \|x^*\|_A \quad (۲۰.۹.۸)$$

است که در آن $c = \lambda_1 / \lambda_n$ ، معکوس ضریب وضعیت $\text{cond}(A)$ است. این یک کران عادی است، که یک همگرایی ضعیف برای مسائل بدوضع به دست می دهد. طرح کلی اثبات آن در لوتنبرگ (۱۹۸۴، ص ۲۵۸، مسأله ۱۰) آمده است. کرانه‌های دیگر را می توان بر پایه رفتار ویژه مقادیرهای $\{\lambda_i\}$ و ضرایب c_i در (۱۸.۹.۸) به دست آورد. در بسیاری از کاربردها که λ_i زیاد تغییر می نماید، اغلب اتفاق می افتد که c_j برای λ_j ی کوچکتر کاملاً نزدیک صفر باشد. در این صورت فرمول (۱۹.۹.۸) را می توان دستکاری کرد تا یک کران بهبود یافته نسبت به (۲۰.۹.۸) به دست آید. در موارد دیگر، ممکن است تعدادی از ویژه مقادیرها حول یک تعداد از مقادیر یکی شوند و در این صورت (۱۹.۹.۸) را می توان برای نشان دادن همگرایی با یک k ی کوچکتر به کار برد. برای قضایای دیگر، لوتنبرگ (۱۹۸۴، ص ۳۵۰)، جنینگز^۱ (۱۹۷۷)، و وان در اسلویس^۲ و وان در هورست^۳ (۱۹۸۶) را ببینید.

فرمولهای $\{\alpha_j, \beta_j\}$ در تعریف روش CG را بازمه ممکن است بهتر کرد تا فرمولهای ساده تر و کاراتری به دست دهند. ما آنها را به صورت زیر می آوریم.

الگوریتم $\text{CG}(A, b, x, n)$

۱. ملاحظه: در این الگوریتم جواب $Ax = b$ با استفاده از روش گرادیان مزدوج حساب می شود.

$$2. \quad p_0 := 0, \quad r_0 := b, \quad x_0 := 0$$

۳. برای $k = 0, \dots, n-1$ ، تا مرحله ۷ جلو بروید.

۴. اگر $r_k = 0$ ، بگیری $x = x_k$ و خارج شوید.

۵. برای $k = 0$ ، $\beta_1 := 0$ و

$$\beta_{k+1} := r_k^T r_k / r_{k-1}^T r_{k-1}, \quad k > 0$$

$$p_{k+1} := r_k + \beta_{k+1} p_k$$

$$6. \quad \alpha_{k+1} := r_k^T r_k / p_{k+1}^T A p_{k+1}$$

$$x_{k+1} := x_k + \alpha_{k+1} p_{k+1}$$

$$r_{k+1} := b - A x_{k+1}$$

۷. حلقه را در k تمام کنید.

۸. $x := x_n$ خارج شوید.

در این الگوریتم مسائل مربوط به حساب دقت متناهی در نظر گرفته نمی شود. برای یک الگوریتم کاملتر به ویلکینسن و راینش (۱۹۷۱، ۵۷-۶۹) مراجعه کنید.

بحث ما درباره روش گرادیان مزدوج خیلی نزدیک روشی است که در لوتنبرگر (۱۹۸۴، فصل ۸) آمده است. برای روش دیگری با انگیزه هندسی بیشتر گلوب و ون لون (۱۹۸۳، بخش ۲.۱۰) را ببینید. اینان همچنین به منابع وسیعی از نوشته‌ها ارجاع داده‌اند.

مثال به عنوان یک آزمون ساده، ماتریس مرتبه ۵ زیر را به کار می‌بریم

$$A = \begin{bmatrix} 5 & 4 & 3 & 2 & 1 \\ 4 & 5 & 4 & 3 & 2 \\ 3 & 4 & 5 & 4 & 3 \\ 2 & 3 & 4 & 5 & 4 \\ 1 & 2 & 3 & 4 & 5 \end{bmatrix} \quad (21.9.8)$$

کوچکترین و بزرگترین ویژه‌مقادیرها به ترتیب، $\lambda_1 = 0.5484$ و $\lambda_5 = 17.1778$ هستند. برای کران خطای (۲۰.۹.۸)، $c = 0.31925$ ، و

$$\frac{1 - \sqrt{c}}{1 + \sqrt{c}} = 0.697$$

برای دستگاه خطی می‌گیریم

$$b = [7.9380, 12.9763, 17.3057, 19.4332, 18.4196]^T$$

که به جواب درست

$$x^* = [-0.3227, 0.3544, 1.1010, 1.5705, 1.6897]^T$$

می‌رسیم. نتایج با استفاده از روش گرادیان مزدوج، همراه با کران خطای (۲۰.۹.۸) در جدول ۱۰.۸ داده شده است. همان‌گونه که قبلاً بیان شد، کران (۲۰.۹.۸) خیلی عادی است.

همان‌گونه که انتظار می‌رفت، مانده‌ها کاهش یافته‌اند. ولی توجه به طریقی که امتدادهای $\{p_k\}$ ساخته شده‌اند، ایجاب می‌کند که به دست آوردن امتدادهای دقیق p_k برای k های بزرگتر، به علت تعداد کمتر ارقام دقیق مانده‌های r_k ، احتمالاً مشکل شود. برای بحث در این مورد، گلوب و ون لون (۱۹۸۳، ص ۳۷۳) را ببینید، که شامل بیشتری از نوشته‌ها در باب این مسأله نیز هست.

جدول ۱۰.۸ مثال از روش گرادیان مزدوج

k	$\ r_k\ _\infty$	$\ x - x_k\ _\infty$	$\ x - x_k\ _A$	کران (۲۰.۹.۸)
۱	۴٫۲۷	$۸٫۰۵E - ۱$	۲٫۶۲	۱۲٫۷
۲	$۸٫۹۸E - ۲$	$۷٫۰۹E - ۲$	$۱٫۳۱E - ۱$	۸٫۸۳
۳	$۲٫۷۵E - ۳$	$۳٫۶۹E - ۳$	$۴٫۷۸E - ۳$	۶٫۱۵
۴	$۷٫۵۹E - ۵$	$۱٫۳۸E - ۴$	$۱٫۶۶E - ۴$	۴٫۲۹
۵	$\doteq 0$	$\doteq 0$	$\doteq 0$	۲٫۹۹

روش گرادیان مزدوج مشروط کران (۲۰.۹.۸) نشان می‌دهد یا به نظر می‌آید که ایجاب می‌کند بارستهای CG حتی برای روشهایی با ضریب وضعیت معمولی چون $\text{cond}(A)_2 = 1/c = 10^0$ خیلی آهسته همگرا شوند. برای افزایش نرخ همگرایی، یا حداقل تضمین یک نرخ سریع همگرایی، مسئله $Ax = b$ را به یک مسئله هم‌ارز با ضریب وضعیت کوچکتر تبدیل می‌کنند. کران (۲۰.۹.۸) کوچکتر خواهد شد، و انتظار می‌رود که بارستها با سرعت بیشتری همگرا شوند.

برای یک ماتریس ناکین Q ، $Ax = b$ را به وسیله

$$(Q^{-1}AQ^{-T})(Q^T x) = Q^{-1}b \quad (۲۲.۹.۸)$$

با $Q^{-T} \equiv (Q^T)^{-1}$ تبدیل می‌کنیم. می‌نویسیم

$$\tilde{A} = Q^{-1}AQ^{-T}, \quad \tilde{x} = Q^T x, \quad \tilde{b} = Q^{-1}b \quad (۲۳.۹.۸)$$

در این صورت (۲۲.۹.۸) به صورت $\tilde{A}\tilde{x} = \tilde{b}$ درآید. ماتریس Q باید طوری انتخاب شود که $\text{cond}(\tilde{A})_2$ به‌طور چشمگیری از $\text{cond}(A)_2$ کوچکتر باشد. روش واقعی CG دقیقاً برای حل $\tilde{A}\tilde{x} = \tilde{b}$ به‌کار گرفته نمی‌شود، بلکه الگوریتم CG کمی تعدیل می‌شود. برای الگوریتم حاصل، وقتی ماتریس Q متقارن باشد، گلوب و وِن‌لون (۱۹۸۳، ص ۳۷۴) را ببینید.

پیدا کردن Q نیاز به تحلیل دقیق مسئله اصلی $Ax = b$ و درک ساختار A برای انتخاب Q دارد. از (۲۳.۹.۸) داریم:

$$A = Q\tilde{A}Q^T$$

که \tilde{A} باید طوری انتخاب شود که اندازه ویژه‌مقدارها نزدیک به ۱ باشد. مثلاً اگر \tilde{A} نزدیک به ماتریس همانی I باشد، آنگاه $A \doteq QQ^T$. این تجزیه را می‌توان با تجزیه مثلثی چولسکی انجام داد. عملهای تقریبی چولسکی گاهی در تعریف شرط گزارها به‌کار می‌روند. برای آشنایی با مسئله انتخاب شرط گزارها، گلوب و وِن‌لون (۱۹۸۳، بخش ۳.۱۰) و اکسلسون^۱ (۱۹۸۵) را ببینید.

بحث در آثار خواندنی

منابعی که تأثیر زیادی در ارائه حذف گاوسی و مباحث دیگر این فصل دارند کتابهای فورسایت و مولر (۱۹۶۷) و گلوب و ون لون (۱۹۸۳)، آیزکسن و کلر (۱۹۶۶) و ویلکینسن (۱۹۶۳، ۱۹۶۵) همراه با مقاله کاهان (۱۹۶۶) هستند. برداشتهای کلی بسیار خوب دیگر در کنت و دبور^۱ (۱۹۸۰)، نوبل (۱۹۶۹)، رایس^۲ (۱۹۸۱) و استوارت (۱۹۷۳a) دیده می‌شوند. درآمد مقدماتی تری در آنتون (۱۹۸۴) و استرنگ (۱۹۸۰) داده شده‌اند.

بهترین برنامه‌ها برای حل مستقیم صورتهای کلی و خاص دستگاههای خطی با اندازه‌های کوچک تا متوسط، مبتنی بر صورتهایی هستند که در بسته‌های LINPACK آمده و در دونگارا و همکاران (۱۹۷۹) تشریح شده است. این برنامه‌ها کاملاً قابل اجرا هستند و در هر دو حساب با دقت ساده و دقت مضاعف هم برای حساب با اعداد حقیقی و هم حساب با اعداد مختلط موجودند. در کنار حل این دستگاهها، با این برنامه‌ها می‌توان ضریب وضعیت ماتریس مورد نظر را نیز برآورد کرد. برنامه‌های معادلات خطی در IMSL و NAG انواع مختلف و پیشرفته همان برنامه‌های LINPACK هستند. جنبه دیگر LINPACK استفاده از زیر روالهای اساسی جبرخطی (BLAS) است. اینها زیر برنامه‌هایی در سطح پایین‌اند که عملیات اساسی برداری مانند حاصلضرب عددی دو بردار و مجموع دو بردار را انجام می‌دهند. اینها به زبان فورتن، به صورت قسمتی از LINPACK موجودند؛ ولی با ارائه کد زبان اسمبلی آنها، اغلب می‌توان کارایی برنامه‌های اصلی LINPACK را به طور قابل ملاحظه‌ای بهبود بخشید. برای یک بحث کلیتر در مورد BLAS، لاوسن و همکاران (۱۹۷۹) را ببینید. برنامه‌های LINPACK به فراوانی در دسترس‌اند و در پیشرفت برنامه‌های معادلات جبری در سایر بسته‌های نرم‌افزاری مؤثر بوده‌اند.

نوشته‌های زیادی در مورد حل دستگاههای خطی که از حل عددی معادلات دیفرانسیل جزئی به دست می‌آیند وجود دارد. برای متون کلی در حل عددی معادلات دیفرانسیل جزئی، بیرکهاف و لینچ (۱۹۸۴)، فورسایت و وازوو^۳ (۱۹۶۰)، گلا دول^۴ و ویت^۵ (۱۹۷۹)، لاپیدوس و پیندر^۶ (۱۹۸۱) و ریشتمایر و مورتون^۷ (۱۹۶۷) را ببینید. برای منتهایی که به روشهای بارستی کلاسیک برای حل دستگاههای خطی حاصل از حل عددی معادلات دیفرانسیل جزئی اختصاص یافته، هاگمن^۸ و یانگ (۱۹۸۱) و وارگا^۹ (۱۹۶۲) را ببینید. برای شیوه‌های دیگر که اخیراً مورد توجه قرار گرفته اسوارتس تراوبر^{۱۰} (۱۹۸۴)، اسوارتس تراوبر و سویت (۱۹۷۹)، جورج و لیو^{۱۱} (۱۹۸۱) و هاکیوش و تروتنبرگ^{۱۲} (۱۹۸۲) را ببینید.

- | | | | |
|----------------------|-------------------------------|-------------------------|-------------|
| 1. Conte and de Boor | 2. Rice | 3. Wasow | 4. Gladwell |
| 5. Wait | 6. Lapidus and Pinder | 7. Richtmyer and Morton | |
| 8. Hageman | 9. Varga | 10. Swarztrauber | |
| 11. George and Liu | 12. Hackbusch and Trottenberg | | |

حل عددی معادلات دیفرانسیل جزئی منبع درصد زیادی از دستگاههای خطی تنک است که در عمل حل شده‌اند. ولی، دستگاههای خطی تنک از مرتبه بزرگ در کاربردهای دیگر نیز وجود دارند [مثلاً داف^۱ (۱۹۸۱) را ببینید]. روشهای گوناگون زیادی برای حل دستگاههای بزرگ تنک وجود دارند، که در باره بعضی از آنها در بخشهای ۶.۸ تا ۸.۸ بحث کردیم. روشهای مستقیم و بارستی دیگر، که به ساختار ماتریس بستگی دارند، در دسترس‌اند. برای یک نمونه از تحقیقات جاری در این زمینه بسیار فعال، بررسی داف (۱۹۷۷)، گزارش بیورک^۲ و همکاران (۱۹۸۱)، داف (۱۹۸۱)، داف و استوارت (۱۹۷۹)، اوانز^۳ (۱۹۸۵) و کتابهای درسی جورج و لیو (۱۹۸۱) و پیسانتسکی^۴ (۱۹۸۴) را ببینید. بسته‌های نرم‌افزاری چندی برای حل انواع مختلف دستگاههای تنک وجود دارند، که بعضاً در کتابهای قبلی آمده‌اند. برای یک فهرست کلی از بسیاری از بسته‌های موجود، گردآوری هیث^۵ (۱۹۸۲) را ببینید. برای روشهای بارستی در حل دستگاههای مربوط به حل معادلات دیفرانسیل جزئی، در کتابهای وارگا (۱۹۶۲) و هاگمن و یانگ (۱۹۸۱) بسیاری از روشهای کلاسیک مورد بحث قرار گرفته‌اند. معادلات انتگرالی به دستگاههای خطی چگال منجر می‌شوند، و انواع دیگر روشهای بارستی برای حل آنها به‌کار برده شده‌اند. برای بعضی روشهای کاملاً موفق، آنکینسون (۱۹۷۶، فصل ۴) را ببینید. روش گرادیان مزدوج با هستسن^۶ و اشتیفل^۷ (۱۹۵۲) آغاز شده است، و استفاده از آن در حل معادلات دیفرانسیل جزئی و انتگرالی هنوز در حال گسترش است. برای بحثهای مفصلتر مربوط به روش جهت‌های مزدوج برای سایر روشهای عددی هستسن (۱۹۸۰) و استوارت (۱۹۷۳b) را ببینید. برای مراجعه به نوشته‌های تازه، از جمله بحثهایی در روش گرادیان مزدوج مشروط اکسلسون (۱۹۸۵)، اکسلسون و لیندزکوگ^۸ (۱۹۸۶) و (بخشهای ۲.۱۰ و ۳.۱۰) گلوب و ون لون (۱۹۸۳) را ببینید. یک تعمیم برای دستگاههای نامتقارن در آیزن‌اشات^۹ و همکاران (۱۹۸۳) پیشنهاد شده است.

یکی از مهمترین عواملی که جهت پژوهشهای آینده در جبر خطی عددی را مشخص خواهد نمود استفاده روزافزون از رایانه‌های برداری و پردازنده موازی است. ماشینهای برداری، مانند CRAY-2 هنگام عملیات اساسی روی کتیبهای برداری، مانند آنچه که در BLAS از LINPACK به‌کار برده شده، به نحو احسن کار می‌کنند. در سالهای اخیر، زمان دسترسی به استفاده از این ماشینها در شبکه‌های رایانه‌یی جهانی که اخیراً توسعه یافته بسیار افزوده شده است. این موضوع، مقیاس بسیاری از مسائل فیزیکی را که می‌توان به آنها پرداخت، تغییر داده است، و همین امر نیاز بیشتری به برنامه‌های هرچه کارا تر رایانه‌یی را، در حل دسته وسیعی از دستگاههای خطی، موجب شده است. استفاده از رایانه‌های موازی حتی جدیدتر است، و تنها در اواسط تا اواخر دهه ۱۹۸۰ متداول

- | | | | |
|--------------------------|--------------|------------|----------------|
| 1. Duff | 2. Björck | 3. Evans | 4. Pissanetsky |
| 5. Heath | 6. Hestenson | 7. Stiefel | |
| 8. Axelsson and Lindskog | 9. Eisenstat | | |

شده است. معماریهای خیلی متنوعی برای اینگونه ماشینها وجود دارد. بعضی دارای پردازنده‌های چندتایی هستند که از یک حافظه، با طرحهای گوناگون مشترکاً استفاده می‌کنند؛ در بعضی دیگر، هر پردازنده دارای حافظه اختصاصی است که به راههای گوناگون با پردازنده‌های دیگر مرتبط است. رایانه‌های موازی، در مقایسه با رایانه‌هایی که ما مطالعه کرده‌ایم، اغلب به الگوریتمهای عددی کاملاً متفاوتی منجر می‌شوند، الگوریتمهایی که در آن می‌توان از پردازنده‌های همزمان، که در یک مسأله به‌کار می‌آیند استفاده کرد. آثار مکتوب در این موارد کم است، ولی به سرعت در حال افزایش است. به‌عنوان یک بررسی برای حل دستگاههای خطی مرتبط با معادلات دیفرانسیل جزئی، هم در رایانه‌های برداری، هم در رایانه‌های موازی، اورتگا و فویکت (۱۹۸۵) را ببینید. برای یک کتاب درسی برنامه‌ریزی شده برای حل دستگاههای خطی، اورتگا (۱۹۸۷) را ببینید.

مراجع

- Aird, T., and R. Lynch (1975). Computable accurate upper and lower error bounds for approximate solutions of linear algebraic systems, *ACM Trans. Math. Softw.* **1**, 217-231.
- Anton, H. (1984). *Elementary Linear Algebra*, 4th ed. Wiley, New York.
- Atkinson, K. (1976). *A Survey of Numerical Methods for the Solution of Fredholm Integral Equations of the Second Kind*. SIAM Pub., Philadelphia.
- Axelsson, O. (1985). A survey of preconditioned iterative methods for linear systems of algebraic equations, *BIT* **25**, 166-187.
- Axelsson, O., and G. Lindskog (1986). On the rate of convergence of the preconditioned conjugate gradient method, *Numer. Math.* **48**, 499-523.
- Birkhoff, G., and R. Lynch (1984). *Numerical Solution of Elliptic Problems*. SIAM Pub., Philadelphia.
- Björck, A., R. Plemmons, and H. Schneider, Eds. (1981). *Large Scale Matrix Problems*. North-Holland, Amsterdam.
- Concus, P., G. Golub, and D. O'Leary (1984). A generalized conjugate gradient method for the numerical solution of elliptic partial differential equations. In *Studies in Numerical Analysis*, G. Golub, Ed. Mathematical Association of America, pp. 178-198.
- Conte, S., and C. de Boor (1980). *Elementary Numerical Analysis*, 3rd ed. McGraw-Hill, New York.
- Dongarra, J., J. Bunch, C. Moler, and G. Stewart (1979). *Linpack User's Guide*. SIAM Pub., Philadelphia.
- Dorr, F. (1970). The direct solution of the discrete Poisson equations on a

- rectangle, *SIAM Rev.* **12**, 248–263.
- Duff, I. (1977). A survey of sparse matrix research, *Proc. IEEE* **65**, 500–535.
- Duff, I., Ed. (1981). *Sparse Matrices and Their Uses*. Academic Press, New York.
- Duff, I., and G. Stewart, Eds. (1979). *Sparse Matrix Proceedings 1978*. SIAM Pub., Philadelphia.
- Eisenstat, S., H. Elman, and M. Schultz (1983). Variational iterative methods for nonsymmetric systems of linear equations, *SIAM J. Numer. Anal.* **20**, 345–357.
- Evans, D., Ed. (1985). *Sparsity and Its Applications*. Cambridge Univ. Press, Cambridge, England.
- Forsythe, G., and C. Moler (1967). *Computer Solution of Linear Algebraic Systems*. Prentice-Hall, Englewood Cliffs, N.J.
- Forsythe, G., and W. Wasow (1960). *Finite Difference Methods for Partial Differential Equations*. Wiley, New York.
- George, A., and J. Liu (1981). *Computer Solution of Large Sparse Positive Definite Systems*. Prentice-Hall, Englewood Cliffs, N.J.
- Gladwell, I., and R. Wait, Eds. (1979). *A Survey of Numerical Methods for Partial Differential Equations*. Oxford Univ. Press, Oxford, England.
- Golub, G., and C. Van Loan (1983). *Matrix Computations*. Johns Hopkins Press, Baltimore.
- Gregory, R., and D. Karney (1969). *A Collection of Matrices for Testing Computational Algorithms*. Wiley, New York.
- Hackbusch, W., and U. Trottenberg, Eds. (1982). *Multigrid Methods, Lecture Notes Math. 960*. Springer-Verlag, New York.
- Hageman, L., and D. Young (1981). *Applied Iterative Methods*. Academic Press, New York.
- Heath, M., Ed. (1982). *Sparse Matrix Software Catalog*. Oak Ridge National Laboratory, Oak Ridge, Tenn. (Published in connection with the Sparse Matrix Symposium 1982.)
- Hestenes, M. (1980). *Conjugate Direction Methods in Optimization*. Springer-Verlag, New York.
- Hestenes, M., and E. Stiefel (1952). Methods of conjugate gradients for solving linear systems, *J. Res. Nat. Bur. Stand.* **49**, 409–439.
- Isaacson, E., and H. Keller (1966). *Analysis of Numerical Methods*. Wiley, New York.
- Jennings, A. (1977). Influence of the eigenvalue spectrum on the convergence rate of the conjugate gradient method, *J. Inst. Math. Its Appl.* **20**, 61–72.
- Kahan, W. (1966). Numerical linear algebra, *Can. Math. Bull.* **9**, 756–801.
- Lapidus, L., and G. Pinder (1982). *Numerical Solution of Partial Differential*

Equations in Science and Engineering. Wiley, New York.

- Lawson, C., and R. Hanson (1974). *Solving Least Squares Problems*. Prentice-Hall, Englewood Cliffs, N.J.
- Lawson, C., R. Hanson, D. Kincaid, and F. Krogh (1979). Basic linear algebra subprograms for Fortran usage, *ACM Trans. Math. Softw.* 5, 308-323.
- Luenberger, D. (1984). *Linear and Nonlinear Programming*, 2nd ed. Addison-Wesley, Reading, Mass.
- Noble, B. (1969). *Applied Linear Algebra*. Prentice-Hall, Englewood Cliffs, N.J.
- Ortega, J. (1987). *Parallel and Vector Solution of Linear Systems*. Preprint, Univ. of Virginia, Charlottesville.
- Ortega, J., and R. Voigt (1985). Solution of partial differential equations on vector and parallel computers, *SIAM Rev.* 27, 149-240.
- Pissanetzky, S. (1984). *Sparse Matrix Technology*. Academic Press, New York.
- Rice, J. (1981). *Matrix Computations and Mathematical Software*. McGraw-Hill, New York.
- Richtmyer, R., and K. Morton (1967). *Difference Methods for Initial Value Problems*, 2nd ed. Wiley, New York.
- Stewart, G. (1973a). *Introduction to Matrix Computations*. Academic Press, New York.
- Stewart, G. (1973b). Conjugate direction methods for solving systems of linear equations, *Numer. Math.* 21, 284-297.
- Stewart, G. (1977). Research, development, and LINPACK. In *Mathematical Software III*, John Rice (Ed.). Academic Press, New York.
- Stone, H. (1968). Iterative solution of implicit approximations of multidimensional partial differential equations, *SIAM J. Numer. Anal.* 5, 530-558.
- Strang, G. (1980). *Linear Algebra and Its Applications*, 2nd ed. Academic Press, New York.
- Swarztrauber, P. (1984). Fast Poisson solvers. In *Studies in Numerical Analysis*, G. Golub, Ed. Mathematical Association of America, pp. 319-370.
- Swarztrauber, P., and R. Sweet (1979). Algorithm 541: Efficient Fortran subprograms for the solution of separable elliptic partial differential equations, *ACM Trans. Math. Softw.* 5, 352-364.
- Van der Sluis, A., and H. van der Vorst (1986). The rate of convergence of conjugate gradients, *Numer. Math.* 48, 543-560.
- Varga, R. (1962). *Matrix Analysis*. Prentice-Hall, Englewood Cliffs, N.J.
- Wilkinson, J. (1963). *Rounding Errors in Algebraic Processes*. Prentice-Hall, Englewood Cliffs, N.J.
- Wilkinson, J. (1965). *The Algebraic Eigenvalue Problem*. Oxford Univ. Press, Oxford, England.
- Wilkinson, J., and C. Reinsch (1971). *Linear Algebra, Handbook for Automatic Computation*, Vol. 2. Springer-Verlag, New York.

مسائل

۱. دستگاههای $Ax = b$ در زیر را با حذف گاوسی بدون محورگیری حل کنید. مانند (۵.۱.۸).
 $A = LU$ را پیدا کنید.

$$A = \begin{bmatrix} 1 & 1 & -1 \\ 1 & 2 & -2 \\ -2 & 1 & 1 \end{bmatrix} \quad b = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \quad (\text{الف})$$

$$A = \begin{bmatrix} 4 & 3 & 2 & 1 \\ 3 & 4 & 3 & 2 \\ 2 & 3 & 4 & 3 \\ 1 & 2 & 3 & 4 \end{bmatrix} \quad b = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix} \quad (\text{ب})$$

$$A = \begin{bmatrix} 1 & -1 & 1 & -1 \\ -1 & 3 & -3 & 3 \\ 2 & -4 & 7 & -7 \\ -3 & 7 & -10 & 14 \end{bmatrix} \quad b = \begin{bmatrix} 0 \\ 2 \\ -2 \\ 8 \end{bmatrix} \quad (\text{ج})$$

۲. دستگاه خطی زیر را در نظر بگیرید و تحقیق کنید که جواب آن، $x_1 = ۲٫۶$ ، $x_2 = -۳٫۸$ و $x_3 = -۵٫۰$ است.

$$6x_1 + 2x_2 + 2x_3 = -2$$

$$2x_1 + \frac{2}{3}x_2 + \frac{1}{3}x_3 = 1$$

$$x_1 + 2x_2 - x_3 = 0$$

(الف) با استفاده از حساب اعشاری ممیز شناور چهاررقمی، با گرد کردن، دستگاه فوق را با حذف گاوسی بدون محورگیری حل کنید.

(ب) بخش (الف) را با استفاده از محورگیری جزئی تکرار کنید. برای انجام اعمال حساب، به یاد داشته باشید که پس از هر محاسبه اعداد را تا چهاررقم بامعنی گرد کنید، درست همان طور که در رایانه عمل می شود.

۳. (الف) الگوریتمهای *Solve* و *Factor* در بخش ۲.۸، یا برنامه های مشابه داده شده در فورسایت و مولر (۱۹۶۷، فصلهای ۱۶ و ۱۷) (۱۹۶۷) را انجام دهید.

(ب) برای آزمون برنامه، دستگاه $Ax = b$ از مرتبه n را با $A = [a_{ij}]$ که به صورت زیر تعریف شده است حل کنید

$$a_{ij} = \text{Max}(i, j)$$

همچنین $b = [1, 1, \dots, 1]^T$ را تعریف می‌کنیم. جواب درست $x = [0, \dots, 0, (1/n)]^T$ است. این ماتریس از گرگوری و کارنی (۱۹۶۹، ص ۴۲) اقتباس شده است.

۴. حل معادله انتگرالی زیر را با گسسته‌ساختن انتگرال با قاعده انتگرالگیری میانگامی (۱۸.۲.۵) در نظر می‌گیریم

$$\lambda x(s) - \int_0^1 \cos(\pi st)x(t)dt = 1 \quad 0 \leq s \leq 1$$

به عبارت دقیقتر، گیریم $n > 0$ ، $h = 1/n$ ، $t_i = (i - 1/2)h$ به ازای $i = 1, \dots, n$. برای مقادیر تقریبی $x(t_1), \dots, x(t_n)$ دستگاه خطی زیر را حل می‌کنیم

$$\lambda z_i - \sum_{j=1}^n h \cos(\pi t_i t_j) z_j = 1 \quad i = 1, \dots, n$$

این دستگاه خطی را با $(\lambda I - K_n)z = b$ نمایش می‌دهیم که K_n از مرتبه $n \times n$ و به شکل زیر است:

$$(K_n)_{ij} = h \cdot \cos(\pi t_i t_j) \quad b_i = 1 \quad 1 \leq i, j \leq n$$

برای مقدار به اندازه کافی بزرگ n ، $z_i \doteq x(t_i)$ ، $1 \leq i \leq n$. مقدار λ ناصفر است و در اینجا فرض شده است که λ ویژه مقدار K_n نیست.

$(\lambda I - K_n)z = b$ را برای چندین مقدار n ، مثلاً ۲، ۴، ۸، ۳۲، ۶۴، حل کرده، جوابهای بردار z را ثبت کنید. چنانچه ممکن باشد نمودار این جوابها را به دست آورید تا اطلاعاتی راجع به تابع جواب $x(s)$ از معادله انتگرالی اصلی به دست آورید. از $0, 1, 2, 4, \lambda = 4$ استفاده کنید.

۵. (الف) حل $Ax = b$ را با A و b مختلط و A از مرتبه n در نظر می‌گیریم. این مسأله را به حل یک دستگاه حقیقی مربعی از مرتبه $2n$ تبدیل کنید.

راهنمایی: بنویسید $A = A_1 + iA_2$ ، $b = b_1 + ib_2$ ، $x = x_1 + ix_2$ که $A_1, A_2, b_1, b_2, x_1, x_2$ همگی حقیقی‌اند. معادلاتی را که باید x_1, x_2 در آنها صدق کنند تعیین کنید.

(ب) حافظه مورد نیاز و تعداد عملیات روش (الف) را برای حل یک دستگاه مختلط $Ax = b$ تعیین نمایید. این نتایج را با نتایجی که از حل مستقیم $Ax = b$ به روش حذف گاوسی و حساب مختلط به دست می‌آید، مقایسه کنید. به هزینه بیشتر اعمال حساب مختلط توجه نمایید.

۶. گیریم A و B و C به ترتیب ماتریسهای از مراتب $m \times n$ ، $n \times p$ و $p \times q$ باشند. تعداد عملیات محاسبه $A(BC)$ و $(AB)C$ را شمارش کنید. مثالهایی بیاورید که در آنها یک ترتیب محاسبه بر ترتیب محاسبه دیگر رجحان داشته باشد.

۷. (الف) نشان دهید که تعداد ضرب و تقسیمهای روش گاوس - ژوردان در بخش ۳.۸ در حدود $\frac{1}{3}n^3$ است.

(ب) نشان دهید که چگونه روش گاوس - ژوردان، با محورگیری جزئی را می توان برای وارون کردن ماتریس $n \times n$ فقط با $n(n+1)$ جا در حافظه، به کار برد. آیا محورگیری کلی را می توان به کار برد؟
۸. هریک از برنامه های مسأله ۳ (الف) یا روش گاوس - ژوردان را برای وارون کردن ماتریسهای مسائل ۱ و ۳ (ب) به کار برید.

۹. ثابت کنید که اگر $A = LL^T$ ، و L حقیقی و نانکین باشد، A متقارن و معین مثبت است.
۱۰. با استفاده از روش چولسکی، تجزیه $A = LL^T$ را برای ماتریسهای زیر به دست آورید

$$\begin{bmatrix} 15 & -18 & 15 & -3 \\ -18 & 24 & -18 & 4 \\ 15 & -18 & 18 & -3 \\ -3 & 4 & -3 & 1 \end{bmatrix} \quad (\text{ب}) \quad \begin{bmatrix} 225 & -30 & 45 \\ -30 & 50 & -100 \\ 45 & -100 & 340 \end{bmatrix} \quad (\text{الف})$$

۱۱. گیریم $A = LU = LDM$ ، و همه l_{ii} و m_{ii} ها مساوی ۱ باشند و D قطری باشد. به علاوه فرض می کنیم A متقارن است. نشان دهید که $M = L^T$ و بنابراین $A = LDL^T$. نشان دهید که A معین مثبت است اگر و تنها اگر همه d_{ii} ها مثبت باشند.

۱۲. گیریم A حقیقی، متقارن، معین مثبت و از مرتبه n باشد. حل $Ax = b$ با استفاده از حذف گاوسی بدون محورگیری را در نظر می گیریم. منظور ما از این مسأله این است که ثابت کنیم محورها ناصفرند.

(الف) نشان دهید که همه عناصر قطری در $a_{ii} > 0$ صدق می کنند. این امر نشان می دهد که a_{11} را می توان یک عنصر محور گرفت.

(ب) پس از حذف x_1 از معادله های ۲ تا n ، فرض می کنیم ماتریس نتیجه $A^{(2)}$ بدین شکل باشد.

$$A^{(2)} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \circ & & & \\ \vdots & \hat{A}^{(2)} & & \\ \circ & & & \end{bmatrix}$$

نشان دهید که $\hat{A}^{(2)}$ متقارن و معین مثبت است.

این شیوه عمل را می‌توان به‌طور استقرایی در هر مرحله از فرایند حذف ادامه داد و بنابراین وجود محورهای ناصفر در هر مرحله توجیه می‌شود.

راهنمایی: برای اثبات معین مثبت بودن $\hat{A}^{(2)}$ ، ابتدا اتحاد زیر را به‌ازای هر انتخاب x_1, \dots, x_n ثابت کنید

$$\sum_{i,j=2}^n a_{ij}^{(2)} x_i x_j = \sum_{i,j=1}^n a_{ij} x_i x_j - a_{11} \left[x_1 + \sum_{j=2}^n \frac{a_{1j}}{a_{11}} x_j \right]^2$$

سپس x_1 مناسب را اختیار نمایید.

۱۳. به‌عنوان یک راه دیگر برای ارائه یک روش فشرده در به‌دست آوردن تجزیه A به LU ، روش ماتریس جهتدار زیر را در نظر بگیرید. بنویسید

$$A = \begin{bmatrix} \hat{A} & d \\ c^T & \alpha \end{bmatrix} \quad c, d \in \mathbf{R}^{n-1} \quad \alpha \in \mathbf{R}$$

و \hat{A} ماتریس مربع از مرتبه $n - 1$ است. فرض کنید A ناکمین است. به‌عنوان یک مرحله در فرایند استقرا، فرض کنید که $\hat{A} = \hat{L}\hat{U}$ معلوم است. $A = LU$ را به شکل زیر پیدا کنید

$$A = \begin{bmatrix} \hat{L} & \circ \\ m^T & \gamma \end{bmatrix} \begin{bmatrix} \hat{U} & q \\ \circ & \gamma \end{bmatrix} \quad m, q \in \mathbf{R}^{n-1} \quad \gamma \in \mathbf{R}$$

نشان دهید که می‌توان q, m و γ را پیدا کرد و توضیح دهید چگونه. (این روش برای ماتریس اصلی A به‌کار می‌رود، با تجزیه هر زیرماتریس اصلی در گوشه چپ بالا به‌ترتیب صعودی).

۱۴. با استفاده از الگوریتم (۲۳.۳.۸) - (۲۴.۳.۸) برای حل دستگاه‌های سه‌قطری، معادله $Ax = b$ را حل کنید که در آن،

$$A = \begin{bmatrix} 2 & -1 & \circ & \circ & \circ \\ 1 & 2 & -1 & \circ & \circ \\ \circ & 1 & 2 & -1 & \circ \\ \circ & \circ & 1 & 2 & -1 \\ \circ & \circ & \circ & 1 & 2 \end{bmatrix} \quad b = \begin{bmatrix} 3 \\ -2 \\ 2 \\ -2 \\ 1 \end{bmatrix}$$

تحقیق کنید که فرضها و حکمهای قضیه ۲.۸ در این مثال صادق‌اند.

۱۵. ماتریسی سه قطری از مرتبه n به شکل زیر را تعریف می‌کنیم

$$A_n = \begin{bmatrix} 2 & -1 & 0 & \dots & 0 \\ -1 & 2 & -1 & 0 & \\ 0 & -1 & 2 & -1 & \vdots \\ \vdots & & & \ddots & -1 \\ 0 & \dots & -1 & 2 & \end{bmatrix}$$

یک فرمول کلی برای $A_n = LU$ پیدا کنید.

راهنمایی: حالت‌های $n = 2, 3, 5$ را در نظر بگیرید و سپس الگوی کلی را حدس بزنید و درستی آن را بررسی کنید.

۱۶. زیرروالی برای حل دستگاههای سه قطری با استفاده از (۲۳.۳.۸) - (۲۴.۳.۸) بنویسید. درستی آن را با استفاده از مثالهای فصل ۱۴ و ۱۵ امتحان کنید. درگرگوری و کارنی (۱۹۶۹، فصل ۲) نیز تعدادی دستگاههای سه قطری وجود دارند که وارونه‌های درست آنها معلوم‌اند.

۱۷. خانواده‌هایی از دستگاههای خطی $A_k x = b$ وجود دارند که در آنها A_k به طریقی ساده به ماتریس A_{k+1} تبدیل می‌شود و ممکن است ساده‌تر باشد که برای تجزیه A_{k+1} به LU ، از تبدیل تجزیه A_k استفاده شود. به‌عنوان یک مثال که در روش سادگی برای برنامه‌ریزی خطی پیش می‌آید، فرض می‌کنیم $A_1 = [a_1, \dots, a_n]$ و $A_2 = [a_2, \dots, a_{n+1}]$ ، که هر $a_j \in \mathbf{R}^n$ فرض می‌کنیم که $A_1 = L_1 U_1$ معلوم و L_1 پایین‌مثلثی و U_1 بالامثلثی باشد. یک راه ساده پیدا کنید که تجزیه $A_2 = L_2 U_2$ را از تجزیه A_1 به‌دست دهد، با فرض اینکه محورگیری لازم نباشد. راهنمایی: با استفاده از $L_1 u_i = a_i$ ، $1 \leq i \leq n$ ، بنویسید

$$A_2 = L_1 [u_2, u_3, \dots, u_n, L_1^{-1} a_{n+1}] \equiv L_1 \tilde{U}$$

نشان دهید که \tilde{U} را می‌توان به‌سادگی به یک شکل بالامثلثی U_2 تبدیل کرد، و این تبدیل متناظر با برگرداندن L_1 به شکل مطلوب L_2 است. به عبارت دقیقتر، $L_2 = L_1 M^{-1}$ ، $U_2 = M \tilde{U}$ ، نشان دهید که هزینه عمل پیدا کردن $A_2 = L_2 U_2$ برابر $O(n^2)$ است.

۱۸. (الف) ضریبهای وضعیت $\text{cond}(A)_p$ را به‌ازای $p = 1, 2, \infty$ برای ماتریس زیر پیدا کنید

$$A = \begin{bmatrix} 100 & 99 \\ 99 & 98 \end{bmatrix}$$

(ب) ویژه مقدارها و ویژه بردارهای A را پیدا کنید، و آنها را برای روشن کردن تبصره‌هایی که به دنبال (۸.۴.۸) در بخش ۴.۸ آمده‌اند به‌کار برید.

۱۹. ثابت کنید که اگر A یکانی باشد، آنگاه $\text{cond}(A)_* = 1$.

۲۰. نشان دهید که برای هر A ، کران بالای (۴.۸.۸) را با انتخابهای مناسب b و r می‌توان به دست آورد. راهنمایی: به موجب تعریفهای $\|A\|$ و $\|A^{-1}\|$ در بخش ۳.۷، بردارهای \hat{x} و \hat{r} وجود دارند که برای آنها $\|\hat{x}\| \|A\| = \|A\hat{x}\|$ و $\|\hat{r}\| \|A^{-1}\| = \|A^{-1}\hat{r}\|$. با استفاده از این روابط ساختن تساوی در کران بالای (۴.۴.۸) را کامل کنید.

۲۱. ضریب وضعیت $\text{cond}(A)_*$ مربوط به (۶.۴.۸) ممکن است برای ماتریسهای بدوضع A کاملاً کوچک باشد. برای مشاهده این مطلب، ماتریس $n \times n$ در زیر را تعریف می‌کنیم:

$$A_R = \begin{bmatrix} 1 & -1 & -1 & \dots & -1 \\ 0 & 1 & -1 & \dots & -1 \\ \vdots & & \ddots & & \vdots \\ & & & 1 & -1 \\ 0 & \dots & 0 & 0 & 1 \end{bmatrix}$$

به سادگی دیده می‌شود که $\text{cond}(A)_* = 1$. تحقیق کنید که A_n^{-1} با ماتریس بالامثلثی $B = [b_{ij}]$ داده می‌شود که در آن $b_{ii} = 1$,

$$b_{ij} = 2^{j-i-1} \quad i < j \leq n$$

$\text{cond}(A)_\infty$ را محاسبه کنید.

۲۲. مانند بخش ۲.۸، فرض می‌کنیم $H_n = [1/(i+j-1)]$ معرّف ماتریس مرتبه n هیلبرت باشد و گیریم \bar{H}_n ماتریسی باشد که از وارد کردن H_n به رایانه شما در حساب با دقت ساده، به دست آمده است. برای مقایسه H_n^{-1} و \bar{H}_n^{-1} ، ماتریس \bar{H}_n را، با افزودن صفرهایی به جزء اعشاری هر درایه، به یک ماتریس با دقت مضاعف بدل کنید. سپس از یک برنامه رایانه‌یی وارون‌سازی ماتریس با دقت مضاعف استفاده کنید تا \bar{H}_n^{-1} به‌طور عددی محاسبه شود. این کار یک مقدار دقیق برای \bar{H}_n^{-1} با دقت ساده، برای مقادیر کوچک n ، به دست خواهد داد. پس از به دست آوردن \bar{H}_n^{-1} ، آن را با H_n^{-1} که در (۱۳.۴.۸) یا گرگوری و کارنی (۱۹۶۹، صص ۳۴-۳۷) داده شده مقایسه نمایید.

۲۳. با استفاده از برنامه‌های مسأله ۳ یا برنامه‌های SGECS و SGECS از LINPACK، $\bar{H}_n x = b$ را برای چندین مقدار n حل کنید. از $b = [1, -1, 1, -1, \dots]^T$ استفاده کنید و جواب درست را با به‌کار بردن \bar{H}_n^{-1} از مسأله ۲۲ محاسبه کنید. درباره نتایج خود اظهار نظر کنید.

۲۴. با استفاده از روش تصحیح مانده، که در ابتدای بخش ۵.۸ توضیح داده شد، جوابهای دستگاههای خطی مسأله ۲۳ را با استفاده از حساب با دقت ساده به دست آورید. مانده‌ها و تصحیحات را حساب کنید، نرخ کاهش در جملات تصحیح را وقتی n افزایش می‌یابد بررسی کنید. کوشش کنید نتایج خود را توضیح دهید.

۲۵. دستگاه خطی مسأله ۴ را برای حل تقریبی یک معادله انتگرالی در نظر می‌گیریم. گاهی اوقات می‌خواهیم چنین دستگاهی را به‌ازای چندین مقدار λ که بهم نزدیک هستند حل کنیم. برنامه‌ای بنویسید که ابتدا دستگاه را به‌ازای $\lambda_0 = 4$ حل و سپس تجزیه $\lambda_0 I - K$ به LU را ذخیره کند. برای حل $(\lambda I - K)x = b$ با مقادیر دیگر λ نزدیک به λ_0 روش تصحیح مانده (۳.۵.۸) را با $C = [LU]^{-1}$ به‌کار برید. برای مثال، این دستگاه را به‌ازای $\lambda = 4, 1, 4, 5, 10$ حل کنید. در هر حالت بارستها را چاپ کرده و نسبت در (۱۱.۵.۸) را محاسبه کنید. در رفتار بارستها وقتی λ افزایش می‌یابد اظهار نظر کنید.

۲۶. گیریم دستگاه $Ax = b$ با

$$A = \begin{bmatrix} 4 & -1 & 0 & -1 & 0 & 0 \\ -1 & 4 & -1 & 0 & -1 & 0 \\ 0 & -1 & 4 & 0 & 0 & -1 \\ -1 & 0 & 0 & 4 & -1 & 0 \\ 0 & -1 & 0 & -1 & 4 & -1 \\ 0 & 0 & -1 & 0 & -1 & 4 \end{bmatrix} \quad b = \begin{bmatrix} 2 \\ 1 \\ 2 \\ 2 \\ 1 \\ 2 \end{bmatrix}$$

دارای جواب $x = [1, 1, 1, 1, 1, 1]^T$ باشد. دستگاه را با استفاده از روش بارستی گاوس-ژاکوبی و سپس با روش گاوس-زایدل حل کنید. حدس اولیه $x^{(0)} = 0$ را به‌کار برید. به نرخ کاهش خطای بارستی توجه کنید. جواب را با دقت $\varepsilon = 0.0001$ پیدا کنید.

۲۷. گیریم A و B از مرتبه n باشند و A ناتکین باشد. حل دستگاه زیر را در نظر بگیرید

$$Az_1 + Bz_2 = b_1 \quad Bz_1 + Az_2 = b_2$$

با $z_1, z_2, b_1, b_2 \in \mathbb{R}^n$.

(الف) شرایط لازم و کافی برای همگرایی روش بارستی زیر را پیدا کنید

$$Az_1^{(m+1)} = b_1 - Bz_2^{(m)} \quad Az_2^{(m+1)} = b_2 - Bz_1^{(m)} \quad m \geq 0$$

(ب) قسمت (الف) را برای روش بارستی زیر تکرار کنید

$$Az_1^{(m+1)} = b_1 - Bz_2^{(m)} \quad Az_2^{(m+1)} = b_2 - Bz_1^{(m+1)} \quad m \geq 0$$

نرخهای همگرایی دو روش را با هم مقایسه نمایید.

۲۸. برای معادله خطای (۲۵.۶.۸) نشان دهید که $\rho_\sigma(M) < 1$ اگر برای نرم ماتریسی مناسبی، رابطه

$$\|P\| < \frac{1}{2\|A^{-1}\|}$$

برقرار باشد.

۲۹. برای بارست دستگاه سه قطری بلوکی، که در (۳۰.۶.۸) داده شده است، همگرایی را تحت مفروضات زیر ثابت کنید.

$$\|C_1\| < \frac{1}{\|B_1^{-1}\|}; \quad \|A_i\| + \|C_i\| < \frac{1}{\|B_i^{-1}\|},$$

$$2 \leq i \leq r-1; \quad \|A_r\| < \frac{1}{\|B_r^{-1}\|}$$

کران نرخ همگرایی را پیدا کنید.

۳۰. ماتریس A_n از مسأله ۱۵ را به خاطر آورید و دستگاه خطی $A_n x = b$ را در نظر بگیرید. این دستگاه از این نظر مهم است که در تقریب استاندارد تفاضلات متناهی (۳۰.۱۱.۶) در مسأله مقدار مرزی دو نقطه‌یی ظاهر می‌شود

$$y''(x) = f(x, y(x)) \quad \alpha < x < \beta \quad y(\alpha) = a, \quad y(\beta) = a_1$$

این دستگاه همچنین به لحاظ ظاهر شدن در تحلیل روشهای بارستی برای حل گسسته‌سازی معادله پواسن به صورت (۵.۸.۸) نیز اهمیت دارد. در این راستا، روش ژاکوبی را برای حل $A_n x = b$ به‌طور بارستی در نظر بگیرید. همگرایی روش ژاکوبی را با نشان دادن $\rho_\sigma(M) < 1$ برای ماتریس مناسب M ، نشان دهید.

راهنمایی: نتایج مسأله ۶ از فصل ۷ را به‌کار ببرید.

۳۱. روشهای بارستی همگرا ممکن است رفتار غیرعادی داشته باشند، به‌عنوان مثال دستگاه زیر را در نظر بگیرید

$$x^{(k+1)} = b + Ax^{(k)} \quad k \geq 0$$

با

$$A = \begin{bmatrix} \lambda & c \\ 0 & -\lambda \end{bmatrix}, \quad b \in \mathbf{R}^2$$

با فرض $|\lambda| < 1$ ، می‌دانیم که $(I - A)^{-1}$ موجود، و برای تمام مقادیر اولیه x^0 ، $x^{(k+1)} - x^{(k)} \rightarrow x^* - x^{(k)}$ ، فرمولهای صریحی برای A^k ، $x^* - x^{(k)}$ و $x^{(k+1)} - x^{(k)}$ پیدا کنید. با وفق دادن c با λ نشان دهید که امکان دارد $\|x^* - x^{(k)}\|_\infty$ هنگام همگرایی به صفر، به طور تناوبی افزایش و کاهش یابد. به مقادیر متناظر $\|x^{(k+1)} - x^{(k)}\|_\infty$ توجه کنید. برای سادگی در تمام محاسبات $x^{(0)} = 0$ را اختیار کنید.

۳۲. (الف) گیریم C یک وارون تقریبی A باشد. تعریف می‌کنیم، $R_0 = I - AC_0$ و فرض می‌کنیم $\|R_0\| < 1$. روش بارستی زیر را تعریف می‌کنیم

$$C_{m+1} = C_m(I + R_m) \quad R_{m+1} = I - AC_{m+1} \quad m \geq 0$$

که روش بارستی معروفی برای محاسبه وارون A^{-1} است. همگرایی C_m به A^{-1} را ابتدا با مربوط ساختن خطای $A^{-1} - C_m$ به مانده R_m ، نشان دهید. و سپس رفتار مانده R_m را با نشان دادن $R_{m+1} = R_m^2$ ، $m \geq 0$ بررسی کنید.
(ب) C_m را به بسط زیر مربوط سازید

$$A^{-1} = C_0(I - R_0)^{-1} = C_0 \sum_{j=0}^{\infty} R_0^j$$

به رابطه این روش وارون کردن A با روش بارستی (۶.۰.۲) از فصل ۲ در محاسبه $1/a$ به‌ازای مقادیر ناصفر a توجه کنید. همچنین مسأله ۱، فصل ۲، را ملاحظه نمایید.

۳۳. برنامه‌های حل بارستی گسسته‌سازی (۵.۸.۸) معادله یواسون بر مربع واحد را اجرا کنید. برای آنکه وضعیتی داشته باشید که در آن یک جواب درست دستگاه خطی در دست باشد، معادله‌های یواسون را انتخاب کنید که در آنها هیچ خطای گسسته‌سازی در رفتن به (۵.۸.۸) وجود نداشته باشد. این مطلب زمانی درست خواهد بود که خطاهای برشی در (۳.۸.۸) متحداً برابر صفر باشند، مثلاً $u(x, y) = x^2 y^2$.

(الف) (۵.۸.۸) را با روش ژاکوبی حل کنید. به خطای واقعی $\|x - x^{(\nu)}\|_\infty$ در هر بارست و همچنین به $\|x^{(\nu+1)} - x^{(\nu)}\|_\infty$ توجه نمایید. ثابت c در (۲.۷.۸) را برآورد نمایید و کران خطای برآوردشده (۵.۷.۸) را محاسبه کنید. نتیجه را با خطای واقعی بارستی مقایسه نمایید.

(ب) مسأله را با روش گاوس-زایدل مجدداً حل کنید. همچنین نرخ بارستی c را با مقدار برآوردشده (۱۳.۸.۸) مقایسه کنید.

(ج) روش SOR را به‌کار برید، از پارامتر شتاب اپتیمال ω^* در (۱۵.۸.۸) استفاده کنید.

۳۴. الف) گسسته‌سازی معادله پواسن در (۵.۹.۸) را برای معادله زیر تعمیم دهید

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} - c(x, y)u = g(x, y) \quad 0 < x, y < 1$$

با $u = f(x, y)$ که مانند قبل روی مرز، اختیار شده است.

ب) فرض می‌کنیم $0 \leq x, y \leq 1, c(x, y) \geq 0$. قسمت الف) از برهان قضیه ۸.۸ را تعمیم داده نشان دهید که دستگاه خطی قسمت الف) دارای جواب یکتاست.

۳۵. الگوریتم گرادیان مزدوج بخش ۹.۸ را پیاده کنید. آن را برای دستگاههای مسائل ۱ و ۳ و ۴ امتحان کنید. هرگاه ممکن باشد، به منظور امتحان کردن، دستگاههایی با جواب درست معلوم را به کار برید. با استفاده از آن، خطاهای واقعی را در هر بارست محاسبه و بررسی کنید که با چه سرعتی کاهش می‌یابند. برای دستگاه خطی در مسأله ۴، که بر پایه حل یک معادله انتگرالی حاصل شده است، دستگاه را برای چندین مقدار n حل کنید. درباره نتایج خود اظهار نظر کنید.

۳۶. نرم برداری $\|x\|_A$ از (۷.۹.۸) را که A متقارن و معین مثبت است به یاد آورید. گیریم ویژه‌مقدارهای A چنین نشان داده شده باشند:

$$0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

نشان دهید که

$$\sqrt{\lambda_1} \|x\|_2 \leq \|x\|_A \leq \sqrt{\lambda_n} \|x\|_2$$

که هر دو تساوی با انتخابهای مناسبی از x به دست می‌آیند.

راهنمایی: از یک پایه یکا متعامد از ویژه بردارهای A استفاده کنید.

۳۷. لم ۱ را که پس از (۱۱.۹.۸) آمده است ثابت کنید.

راهنمایی: از استقرای ریاضی روی k استفاده کنید. آن را برای $k = 1$ ثابت و سپس فرض کنید که برای $k \leq l$ درست است و آن را برای $k = l + 1$ ثابت کنید. اثبات را به دو بخش بشکنید. (۱): $p_i^T r_{l+1} = 0$ برای $i \leq l$ ، و (۲): $p_{l+1}^T r_{l+1} = 0$.

۳۸. گیریم A متقارن، معین مثبت و از مرتبه $n \times n$ باشد. فرض می‌کنیم $U = \{u_1, \dots, u_n\}$ یک مجموعه از بردارهای ناصفر در R^n باشد. در این صورت اگر U هم یک مجموعه متعامد و هم یک مجموعه A -متعامد باشد، آنگاه برای مقادیر مناسب $\lambda_i > 0, i = 1, \dots, n$ ، $u_i = \lambda_i^{-1} A u_i$ برعکس، همیشه می‌توان یک مجموعه از ویژه بردارهای $\{u_1, \dots, u_n\}$ از A را انتخاب نمود که هم متعامد و هم A -متعامد باشند.

۳۹. گیریم A متقارن، معین مثبت و از مرتبه n باشد. فرض می‌کنیم $\{v_1, \dots, v_n\}$ یک

مجموعه A - متعامد در \mathbf{R}^n باشد که همه v_i ها مخالف صفرند. تعریف می‌کنیم

$$Q_j = \frac{v_j v_j^T A}{v_j^T A v_j} \quad j = 1, \dots, n$$

که ویژگیهای زیر را برای Q_j نشان می‌دهد.

$$Q_j v_j = v_j \text{ و } Q_j v_i = 0 \text{ اگر } i \neq j \quad ۱.$$

$$Q_j^2 = Q_j \quad ۲.$$

$$x, y \in \mathbf{R}^n \text{ مقادیر تمام برای } (x, Q_j y)_A = (Q_j x, y)_A \quad ۳.$$

$$x, y \in \mathbf{R}^n \text{ مقادیر تمام برای } (Q_j x, (I - Q_j) y)_A = 0 \quad ۴.$$

$$x \in \mathbf{R}^n \text{ مقادیر تمام برای } \|x\|_A^2 = \|Q_j x\|_A^2 + \|(I - Q_j)x\|_A^2 \quad ۵.$$

ویژگیهای ۵.۲ می‌گویند که Q_j یک تصویر قائم در فضای برداری \mathbf{R}^n با ضرب داخلی $(\cdot, \cdot)_A$ است. تعریف می‌کنیم

$$S_k = \text{Span}\{v_1, \dots, v_k\}$$

نشان دهید که جواب مسأله مینیم سازی

$$\text{Min}_{y \in S_k} \|x - y\|_A$$

با رابطه زیر داده می‌شود

$$y = \left[\sum_{j=1}^k Q_j \right] x \equiv P_k x$$

ماتریس P_k در ویژگیهای ۵.۲ نیز صدق می‌کند.

مسأله ویژه مقدار ماتریس

ما در این فصل مسأله محاسبه ویژه مقدارها و ویژه بردارهای یک ماتریس مربعی را بررسی می‌کنیم. این مسأله در شماری از زمینه‌ها پیش می‌آید و ماتریسهای حاصل ممکن است شکلهای مختلف داشته باشند. این ماتریسها ممکن است تَنک یا فشرده باشند، ممکن است مرتبه و ساختارهای بسیار متغیر داشته و اغلب متقارن باشند. به علاوه آنچه باید محاسبه شود ممکن است آنقدر تغییر کند که بر انتخاب روشی که باید به کار گرفته شود اثر بگذارد. اگر تنها محاسبه چند ویژه مقدار مورد نظر باشد، روش عددی با روشی که محاسبه تمام ویژه مقدارها خواسته شده باشد متفاوت خواهد بود. مسأله کلی پیدا کردن تمام ویژه مقدارها و ویژه بردارهای یک ماتریس نامتقارن A ممکن است نسبت به اختلالات درایه‌های A ، کاملاً ناپایدار باشد، و این امر طرح روشهای کلی و برنامه‌های رایانه‌یی را مشکلتر می‌سازد. ویژه مقدارهای یک ماتریس متقارن A نسبت به اختلالات A کاملاً پایدارند. این موضوع، همراه با ناپایداری ممکن در ماتریسهای نامتقارن، در بخش ۱.۹ بررسی شده است. به علت پایداری بیشتر مسأله ویژه مقدار ماتریسهای متقارن و به لحاظ پیدایش فراوان آنها، روشهای زیادی به خصوص برای آنها ایجاد شده‌اند. این موضوع تأکید اصلی ما در این فصل خواهد بود، اگرچه روشهایی برای ویژه مقدارهای ماتریسهای نامتقارن نیز مورد بحث قرار گرفته‌اند. معمولاً ویژه مقدارهای ماتریس ابتدا محاسبه می‌شوند و سپس، در صورت لزوم، از آنها در

پیدا کردن ویژه بردارها استفاده می‌شود. استثنای عمده بر این قاعده، روش توانی است که در بخش ۲.۹ توضیح داده شده است و در محاسبه ویژه مقدار غالب منفرد یک ماتریس، مفید است. روال معمولی برای محاسبه ویژه مقدارهای یک ماتریس A شامل دو مرحله است. ابتدا، با استفاده از تبدیلات تشابهی، A به شکل ساده‌تری، که معمولاً برای ماتریسهای متقارن، سه قطری است، بدل می‌شود. و سپس، این ماتریسهای ساده برای محاسبه ویژه مقدارها و در صورت نیاز برای ویژه بردارها، به کار برده می‌شوند. شکل اصلی تبدیلات تشابهی که مورد استفاده قرار می‌گیرند ماتریسهای یکانی یا متعامد خاصی هستند، که در بخش ۳.۹ درباره آنها بحث شده است. برای محاسبه ویژه مقدارهای ماتریس سه قطری متقارن، نظریه دنباله‌های استورم^۱ در بخش ۴.۹ معرفی شده و الگوریتم QR در بخش ۵.۹ مورد بحث قرار گرفته است. هنگامی که ویژه مقدارها محاسبه شده باشند، قویترین تکنیک برای محاسبه ویژه بردارها روش بارست معکوس است. این روش در بخش ۶.۹ بحث و توضیح داده شده است. لازم به تذکر است که ما کلمات متقارن و نامتقارن را عموماً به کار می‌بریم، و حال آنکه معمولاً باید فقط در ارتباط با ماتریسهای حقیقی به کار برده شوند. برای ماتریسهای مختلط همیشه باید به ترتیب واژه‌های ارمیتی و نارمیتی را جایگزین کنیم. بیشتر روشهای عددی که در حال حاضر به کار برده می‌شوند از ۱۹۵۰ به بعد گسترش یافته‌اند. پیاده کردن آنها به صورت برنامه‌های رایانه‌یی چندان آسان نیست، به ویژه برای آن روشهایی که در ماتریسهای نامتقارن به کار برده می‌شوند. از نیمه دهه ۱۹۶۰ الگوریتمهایی برای تعدادی از مسائل ویژه مقداری ماتریس به زبان ALGOL در مجله ریاضیات کاربردی^۲ انتشار یافته‌اند. این برنامه‌ها به نحوی گسترده مورد آزمایش قرار گرفته و براساس آزمونها و نظریه‌های جدید اصلاح شده‌اند. این الگوریتمها در کتاب ویلکینسن و راینش (۱۹۷۱)، قسمت دوم گردآوری شده‌اند. در یک پروژه بخش ریاضیات کاربردی آزمایشگاه ملی آرگون^۳، این برنامه‌ها به زبان فورترن ترجمه شده و آزمایشها و اصلاحات بیشتری در آنها صورت گرفته است. این بسته برنامه‌ها EISPACK نامیده می‌شود و می‌توان آنها را از آزمایشگاه ملی آرگون و منابع دیگر تهیه نمود (به پیوست مراجعه شود). شرح کامل این بسته، شامل تمام برنامه‌ها در اسمیت و همکاران (۱۹۷۶) و گاربو^۴ و همکاران (۱۹۷۷) داده شده است.

۱.۹ جای ویژه مقدار، خطا، و قضایای پایداری

موضوع را با ارائه چند قضیه برای جابجایی و کراندار کردن ویژه مقدارهای یک ماتریس A آغاز می‌نماییم. برای کران بالا، قضیه ۸.۷ از فصل ۷ را یادآوری می‌نماییم که برای هر نرم ماتریسی،

$$\max_{\lambda \in \sigma(A)} |\lambda| \leq \|A\| \quad (1.1.9)$$

نماد $\sigma(A)$ معرف مجموعه تمام ویژه مقدارهای A است. قضیه بعدی یک روش محاسبه‌ای ساده در برآوردهای بهتر جای ویژه مقدارهای A به دست می‌دهد.
برای $A = [a_{ij}]$ از مرتبه n ، تعریف می‌کنیم

$$r_i = \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \quad i = 1, 2, \dots, n \quad (2.1.9)$$

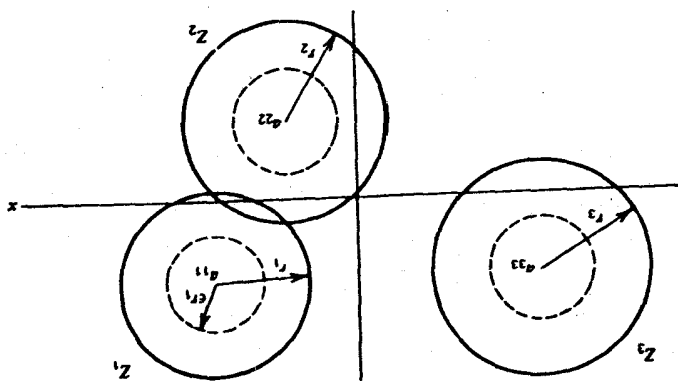
و گیریم Z_i معرف دایره‌ای به مرکز a_{ii} و شعاع r_i در صفحه مختلط باشد:

$$Z_i = \{z \in \mathbf{C} \mid |z - a_{ii}| \leq r_i\} \quad (3.1.9)$$

قضیه ۱.۹ (گرشگورین^۱) گیریم A دارای مرتبه n باشد و λ یک ویژه مقدار A . آنگاه λ به یکی از دایره‌های Z_i تعلق دارد. به علاوه اگر m تا از این دایره‌ها یک مجموعه همبند S تشکیل دهند که از $n - m$ دایره باقیمانده مجزا باشند، آنگاه S دقیقاً m ویژه مقدار A را در بر دارد که بر حسب چندگانی بودن ریشه‌های چندجمله‌یی مشخصه A شمارش شده‌اند.

چون A و A^T ویژه مقدارها و چندجمله‌بیهای مشخصه واحدی دارند، این قضایا، اگر برای تعریف شعاعها در (۲.۱.۹)، از مجموعیابی روی ستونها به جای سطرها، استفاده شود، نیز معتبرند.

برهان شکل ۱.۹ تصویری است از صفحه مختلط که چگونگی دایره را برای یک ماتریس مختلط از مرتبه ۳ نشان می‌دهد. دایره با خط پرنهایی هستند که با (۳.۱.۹) داده شده‌اند و دایره خط چین



شکل ۱.۹ مثال قضیه دایره گرشگورین

بعداً در برهان پیش می‌آیند. به موجب این قضیه، باید یک ویژه مقدار در Z_2 باشد و دو ویژه مقدار در اجتماع Z_1 و Z_2 .

گیریم λ یک ویژه مقدار A باشد و x ویژه بردار متناظر با آن. فرض می‌کنیم k ، زیرنمایه یک مؤلفه x باشد که برای آن

$$\|x_k\| = \max_{1 \leq i \leq n} |x_i| = \|x\|_\infty$$

در این صورت از مؤلفه k ام رابطه $Ax = \lambda x$ نتیجه می‌شود

$$\sum_{j=1}^n a_{kj} x_j = \lambda x_k$$

$$(\lambda - a_{kk}) x_k = \sum_{\substack{j=1 \\ j \neq k}}^n a_{kj} x_j$$

$$\|\lambda - a_{kk}\| \|x_k\| \leq \sum_{\substack{j=1 \\ j \neq k}}^n |a_{kj}| \|x_j\| \leq r_k \|x\|_\infty$$

با حذف $\|x\|_\infty$ از طرفین، اولین قسمت قضیه اثبات می‌شود.

تعریف می‌کنیم

$$D = \text{diag}[a_{11}, a_{22}, \dots, a_{nn}] \quad E = A - D$$

برای $0 \leq \epsilon \leq 1$ تعریف می‌کنیم

$$A(\epsilon) = D + \epsilon E \tag{۴.۱.۹}$$

و ویژه مقدارهای آن را با $\lambda_1(\epsilon), \dots, \lambda_n(\epsilon)$ نشان می‌دهیم. توجه نمایید که $A(1) = A$ که همان ماتریس اصلی است. ویژه مقدارها ریشه‌های چندجمله‌یی مشخصه

$$f_\epsilon(\lambda) \equiv \det[A(\epsilon) - \lambda I]$$

هستند. چون ضرایب $f_\epsilon(\lambda)$ توابع پیوسته از ϵ هستند و چون ریشه‌های هر چندجمله‌یی، توابعی پیوسته از ضرایب آن هستند [هنریچی (۱۹۷۴)، ص ۲۸۱ را ببینید]، پس $\lambda_n(\epsilon), \dots, \lambda_1(\epsilon)$ توابعی پیوسته از ϵ هستند. اگر پارامتر ϵ تغییر کند هر یک از ویژه مقدارها مانند $\lambda_i(\epsilon)$ در صفحه مختلط تغییر می‌نماید و راهی از $\lambda_i(0)$ به $\lambda_i(1)$ رسم می‌کند.

با توجه به قسمت اول قضیه، می‌دانیم که ویژه‌مقادیرهای $\lambda_i(\epsilon)$ در دایره‌های

$$Z_i(\epsilon) = \{z \in \mathbf{C} \mid |z - a_{ii}| \leq \epsilon r_i\} \quad i = 1, \dots, n \quad (5.1.9)$$

جا دارند، r_i مانند قبل با (۲.۱.۹) تعریف می‌شود. مثالهای این دایره‌ها در شکل ۱.۹ با دایره‌های خط‌چین داده شده‌اند. این دایره‌ها وقتی ϵ از ۱ به ۰ می‌رود کاهش می‌یابند و ویژه‌مقادیرهای $\lambda_i(\epsilon)$ باید در داخل این دایره‌ها باقی بمانند. وقتی $\epsilon = 0$ ، ویژه‌مقادیرها

$$\lambda_i(0) = a_{ii}$$

خواهند شد. با توجه به مسیر $\lambda_i(\epsilon)$ ، $0 \leq \epsilon \leq 1$ ، این راه که از $\epsilon = 0$ شروع می‌شود باید در دایره $Z_i(1)$ باقی بماند. بنابراین اگر m تا از این دایره‌های $Z_i(1)$ یک مجموعه همبند S ، مجزا از $n - m$ دایره باقیمانده، تشکیل دهند، آنگاه S باید دقیقاً شامل m ویژه‌مقدار $\lambda_i(1)$ باشد، زیرا شامل m ویژه‌مقدار $\lambda_i(0)$ است. و این قسمت دوم قضیه را اثبات می‌کند.
از آنجا که

$$\det[A - \lambda I] = \det[A - \lambda I]^T = \det[A^T - \lambda I]$$

خواهیم داشت $\sigma(A) = \sigma(A^T)$. لذا این قضیه را برای سطرهای A^T به‌کار می‌بریم تا آن را برای ستونهای A اثبات نماییم؛ و اثبات قضیه کامل می‌شود. ■

از این قضیه می‌توان به راههای مختلف استفاده کرد، ما فقط دو مثال عددی ساده ذکر می‌کنیم.

مثال ماتریس

$$A = \begin{bmatrix} 4 & 1 & 0 \\ 1 & 0 & -1 \\ 1 & 1 & -4 \end{bmatrix}$$

را در نظر می‌گیریم. مطابق با قضیه قبل، ویژه‌مقادیرها باید در دایره‌های

$$|\lambda - 4| \leq 1 \quad |\lambda| \leq 2 \quad |\lambda + 4| \leq 2 \quad (6.1.9)$$

قرار داشته باشند. چون دایره اول از بقیه مجزاست، باید یک ریشه تنها در این دایره واقع باشد. چون ضرایب

$$f(\lambda) = \det | A - \lambda I |$$

حقیقی اند، ویژه مقدارهای مختلط، اگر وجود داشته باشند، باید جفت‌های مزدوج باشند. و این مطلب، با استفاده از (۶.۱.۹)، به سادگی ایجاب می‌کند که یک ویژه مقدار حقیقی در بازه $[3, 5]$ وجود داشته باشد. دو دایره باقیمانده در تنها نقطه $(-2, 0)$ برهم مماس می‌شوند. با استفاده از برهان قبلی، ویژه مقدارهای واقع در این دو دایره باید حقیقی باشند. و با توجه به تعریف $A(\epsilon)$ ، $\epsilon < 1$ در (۴.۱.۹) یک ویژه مقدار در $[-6, -2]$ و یکی در $[-2, 2]$ وجود دارد. چون به آسانی دیده می‌شود که $\lambda = -2$ یک ویژه مقدار نیست، می‌توانیم نتیجه بگیریم که A در هر یک از بازه‌های،

$$[-6, -2), (-2, 2], [3, 5]$$

یک ویژه مقدار دارد. ویژه مقدارهای واقعی عبارت‌اند از:

$$-3.76010, -0.442931, 4.20303$$

مثال ماتریس مرتبه n زیر را در نظر می‌گیریم

$$A = \begin{bmatrix} 4 & 1 & 0 & & \dots & 0 \\ 1 & 4 & 1 & 0 & & \\ 0 & 1 & 4 & 1 & & \vdots \\ \vdots & & & \ddots & & \\ & & & & 1 & 4 & 1 \\ 0 & \dots & & & 0 & 1 & 4 \end{bmatrix} \quad (7.1.9)$$

چون A متقارن است، همه ویژه مقدارهای آن حقیقی‌اند. شعاعهای n_i ی (۲.۱.۹) همگی یا ۱ هستند یا ۲، و مرکز تمام دایره‌ها $a_{ii} = 4$ است. بنابراین به موجب قضیه قبل، ویژه مقدارها باید همگی در بازه $[2, 6]$ واقع باشند. چون ویژه مقدارهای A^{-1} عکس ویژه مقدارهای A هستند، به ازای تمام این ویژه مقدارهای A^{-1} که با μ نمایش می‌دهیم، باید داشته باشیم

$$\frac{1}{6} \leq \mu \leq \frac{1}{2}$$

با استفاده از نرم ماتریسی (۲۲.۳.۷) که از نرم برداری اقلیدسی به دست آمده است داریم:

$$\|A^{-1}\|_2 = r_\sigma(A^{-1}) \leq \frac{1}{\rho}$$

که مستقل از اندازه n است.

کرانها برای ویژه مقدارهای اختلال یافته در یک ماتریس داده شده A ، می خواهیم اختلالی ایجاد و سپس تأثیر آن را بر ویژه مقدارهای A مشاهده کنیم. کرانهای تحلیلی برای اختلالات روی ویژه مقدارها را، براساس اختلالات روی ماتریس A به دست می آوریم. این کرانها همچنین تعریفی برای ضریب وضعیت به ذهن القا می کنند که می توان آن را برای بیان درجه پایداری یا ناپایداری در ویژه مقدارها به کار برد. برای آنکه استدلال را خیلی ساده کنیم، فرض می کنیم که شکل متعارف ژوردان ماتریس A قطری باشد (قضیه ۶.۷ را ببینید). یعنی برای ماتریس نانکینی چون P :

$$P^{-1}AP = \text{diag}[\lambda_1, \dots, \lambda_n] \equiv D \quad (۸.۱.۹)$$

ستونهای P ویژه بردارهای A متناظر با ویژه مقدارهای $\lambda_1, \lambda_2, \dots, \lambda_n$ هستند. ماتریسهایی که در (۸.۱.۹) صدق می کنند مهمترین حالتها در عمل هستند. برای یک بحث مختصر درباره حالتی که در آن شکل متعارف ژوردان قطری نباشد، مبحث آخر این بخش را ملاحظه کنید. همچنین نیازمندیم که ویژگی خاصی را برای نرمهای ماتریسی که به کار برده می شوند بپذیریم. برای هر ماتریس قطری

$$G = \text{diag}[g_1, \dots, g_n]$$

باید داشته باشیم

$$\|G\| = \max_{1 \leq i \leq n} |g_i| \quad (۹.۱.۹)$$

کلیه نرمهای ماتریس عملگر القاشده از نرمهای برداری $\|x\|_p$ ، $1 \leq p \leq \infty$ ، این ویژگی را دارا هستند. اکنون می توانیم قضیه زیر را بیان کنیم.

قضیه ۲.۹ (باوئر-فیکه^۱) گیریم A یک ماتریس با شکل متعارف قطری ژوردان مانند (۸.۱.۹) باشد. و فرض می کنیم که نرم ماتریس در (۹.۱.۹) صدق کند. گیریم $A + E$ اختلال یافته A باشد و λ یک ویژه مقدار $A + E$ در این صورت

$$\min_{1 \leq i \leq n} |\lambda - \lambda_i| \leq \|P\| \|P^{-1}\| \|E\| \quad (۱۰.۱.۹)$$

برهان اگر λ ویژه مقدار A نیز باشد، آنگاه (۱.۱.۹) بالبداهه درست است. پس فرض می‌کنیم $\lambda \neq \lambda_1, \lambda_2, \dots, \lambda_n$ و x یک ویژه بردار $A + E$ متناظر با λ باشد. در این صورت

$$(A + E)x = \lambda x$$

$$(\lambda I - A)x = Ex$$

A را از (۸.۱.۹) به دست آورده در رابطه اخیر می‌گذاریم

$$(\lambda I - PDP^{-1})x = Ex$$

$$(\lambda I - D)(P^{-1}x) = (P^{-1}EP)(P^{-1}x)$$

چون $\lambda I - D$ ، $\lambda \neq \lambda_1, \lambda_2, \dots, \lambda_n$ ناکین است،

$$(\lambda I - D)^{-1} = \text{diag}[(\lambda - \lambda_1)^{-1}, \dots, (\lambda - \lambda_n)^{-1}]$$

پس

$$P^{-1}x = (\lambda I - D)^{-1}(P^{-1}EP)(P^{-1}x)$$

$$\|P^{-1}x\| \leq \|(\lambda I - D)^{-1}\| \|P^{-1}EP\| \|P^{-1}x\|$$

با حذف $\|P^{-1}x\|$ و استفاده از (۹.۱.۹)

$$1 \leq [\text{Max}_i |\lambda - \lambda_i|^{-1}] \|P^{-1}\| \|P\| \|E\|$$

این رابطه هم‌ارز با (۱۰.۱.۹) بوده و برهان کامل می‌شود. ■

فرض اگر A ارمیتی باشد، و اگر $A + E$ یک شکل اختلال یافته A باشد، آنگاه برای هر ویژه مقدار λ ی $A + E$

$$\text{Min}_{1 \leq i \leq n} |\lambda - \lambda_i| \leq \|E\|_2 \quad (۱۱.۱.۹)$$

برهان چون A ارمیتی است ماتریس P را می‌توان به گونه‌ای انتخاب کرد که یکانی باشد. و با استفاده از نرم عملگر (۱۹.۳.۷) که از نرم برداری اقلیدسی القا شده، $\|P\|_2 = \|P^{-1}\|_2 = 1$ (به مثال ۱۳ فصل ۷ نگاه کنید). و در اینجا برهان کامل می‌شود. ■

حکم (۱۱.۱.۹) ثابت می‌کند که اختلالهای کوچک یک ماتریس ارمیتی، چنانچه در مقدمه این فصل گفته شده بود، به اختلالهای کوچک ویژه مقادارها می‌انجامد. باید توجه داشت که خطای نسبی در بعضی یا همه ویژه مقادارها ممکن است بزرگ باشد، و این حالت معمولاً زمانی رخ می‌دهد که اندازه‌های ویژه مقادارهای ماتریس خیلی زیاد تغییر کنند.

مثال ماتریس هیلبرت از مرتبه ۳ را در نظر می‌گیریم،

$$H_3 = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \end{bmatrix} \quad (12.1.9)$$

ویژه مقادارهای آن با هفت رقم بامعنی عبارت‌اند از:

$$\lambda_1 = 1.408319 \quad \lambda_2 = 0.1223271 \quad \lambda_3 = 0.002687340 \quad (13.1.9)$$

اکنون ماتریس اختلال یافته \hat{H}_3 را که معرف ماتریس H_3 تا چهار رقم بامعنی است در نظر می‌گیریم:

$$\hat{H}_3 = \begin{bmatrix} 1.000 & 0.5000 & 0.3333 \\ 0.5000 & 0.3333 & 0.2500 \\ 0.3333 & 0.2500 & 0.2000 \end{bmatrix} \quad (14.1.9)$$

ویژه مقادارهای آن تا ۷ رقم بامعنی عبارت‌اند از:

$$\hat{\lambda}_1 = 1.408294 \quad \hat{\lambda}_2 = 0.1223415 \quad \hat{\lambda}_3 = 0.002664489 \quad (15.1.9)$$

برای تحقیق در صحت (۱۱.۱.۹) برای این حالت، به سادگی می‌توان مقدار

$$\|E\|_2 = r_\sigma(E) = \frac{1}{3} \times 10^{-4} = 0.000033$$

را محاسبه نمود. برای خطاها و خطاهای نسبی در (۱۵.۱.۹)،

$$\lambda_1 - \hat{\lambda}_1 = 0.0000249 \quad \text{Rel}(\hat{\lambda}_1) = 0.0000177$$

$$\lambda_2 - \hat{\lambda}_2 = -0.0000144 \quad \text{Rel}(\hat{\lambda}_2) = -0.000118$$

$$\lambda_3 - \hat{\lambda}_3 = 0.0000229 \quad \text{Rel}(\hat{\lambda}_3) = 0.0085$$

همه خطاها در (۱۱.۱.۹) صدق می‌کنند. ولی خطای نسبی λ_3 در مقایسه با اختلالات نسبی \hat{H}_3 کاملاً چشمگیر است.

برای یک ماتریس نامتقارن A با P یی که در (۸.۱.۹) داده شده، عدد

$$K(A) = \| P \| \| P^{-1} \|$$

ضریب وضعیت مسأله ویژه مقدار برای A خوانده می‌شود. این تعریف براساس کران (۱۰.۱.۹) برای اختلالات ویژه مقدارهای ماتریس اختلال یافته A داده شده است. انتخاب دیگر، که حتی از نظر محاسبه مشکلتر است، استفاده از

$$K(A) = \text{Inf} \| P \| \| P^{-1} \| \quad (۱۶.۱.۹)$$

است که در آن اینفیموم روی تمام ماتریسهای P ، که در رابطه (۸.۱.۹) صدق می‌کنند، و تمام نرمهای ماتریسی که در رابطه (۹.۱.۹) صدق می‌کنند، محاسبه شده است. دلیل در نظر گرفتن ضرایب وضعیت این است که در ماتریسهای نامتقارن A ، اختلالات کوچک E ممکن است به اختلالات نسبی بزرگی در ویژه مقدارهای A منجر شوند.

مثال برای روشن کردن مسائل آسیب‌شناختی که با ماتریسهای نامتقارن ظاهر می‌شوند، در نظر می‌گیریم:

$$A = \begin{bmatrix} ۱۰۱ & -۹۰ \\ ۱۱۰ & -۹۸ \end{bmatrix} \quad A + E = \begin{bmatrix} ۱۰۱ - \epsilon & -۹۰ - \epsilon \\ ۱۱۰ & -۹۸ \end{bmatrix} \quad (۱۷.۱.۹)$$

ویژه مقدارهای A برابر $\lambda = ۱, ۲$ و ویژه مقدارهای $A + E$:

$$\lambda = \frac{۳ - \epsilon \pm \sqrt{۱ - ۸۲۸\epsilon + \epsilon^2}}{۲}$$

هستند. برای درک بهتر به عنوان مثال ϵ را ۰.۰۰۱ انتخاب می‌کنیم. پس

$$A + E = \begin{bmatrix} ۱۰۰.۹۹۹ & -۹۰.۰۰۱ \\ ۱۱۰ & -۹۸ \end{bmatrix}$$

و ویژه مقدارها عبارت‌اند از:

$$\lambda \doteq ۱.۲۹۸, ۱.۷۰۱ \quad (۱۸.۱.۹)$$

از این مسأله نباید نتیجه بگیریم که ماتریسهای نامتقارن بدوضع‌اند. بسیاری حالتها در عمل کاملاً خوش‌وضع هستند. ولی برای نوشتن یک الگوریتم کلی، ما همیشه در پی آن هستیم که هر چه ممکن است حالت‌های بیشتری را پوشش دهیم، و این مثال نشان می‌دهد که این امر احتمالاً برای ردهٔ همهٔ ماتریسهای نامتقارن دشوار است.

برای ماتریسهای متقارن، قضیهٔ (۱۱.۱.۹) را به چندین راه می‌توان بهبود بخشید. یکی توصیف مینیماکس برای ویژه‌مقدارهای ماتریسهای متقارن است. برای بحث در این مورد و کرانه‌های خطای حاصله، به پارلت^۱ (۱۹۸۰، بخش ۲.۱۰) یا ویلکینسن (۱۹۶۵، ص ۱۰۱) مراجعه کنید. به جای آن قضیهٔ زیر را می‌دهیم که برای تحلیل خطای روشهایی که بعداً می‌آیند مفیدتر خواهد بود.

قضیهٔ ۳.۹ (ویلان-هوفمن^۲) گیریم A و E ماتریسهای حقیقی، متقارن و از مرتبهٔ n باشند و $\hat{A} = A + E$ را تعریف می‌کنیم. گیریم λ_i و $\hat{\lambda}_i$ ، $i = 1, \dots, n$ ، به ترتیب ویژه‌مقدارهای A و \hat{A} باشند که به‌طور صعودی مرتب شده‌اند، در این صورت

$$\left[\sum_{j=1}^n (\lambda_j - \hat{\lambda}_j)^2 \right]^{1/2} \leq F(E) \quad (۱۹.۱.۹)$$

که در آن $F(E)$ نرم فروبنیوس E است که در (۱۰.۳.۷) تعریف شده بود.

برهان ویلکینسن (۱۹۶۵، صص ۱۰۴-۱۰۸) را ببینید.

این قضیه برای کراندارکردن اثر خطاهای گردکردن ناشی از تبدیل یک ماتریس متقارن به ماتریس سه قطری بعداً مورد استفاده قرار خواهد گرفت.

یک کران خطای قابل‌محاسبه برای ماتریسهای متقارن گیریم A ماتریس متقارنی باشد که برای آن ویژه‌مقدار تقریبی λ و ویژه‌بردار تقریبی x محاسبه شده است. ماندهٔ

$$\eta = Ax - \lambda x \quad (۲۰.۱.۹)$$

را تعریف می‌کنیم. چون A متقارن است، یک ماتریس یکانی U موجود است که برای آن

$$U^*AU = \text{diag}[\lambda_1, \dots, \lambda_n] \equiv D \quad (۲۱.۱.۹)$$

سپس نشان خواهیم داد که

$$\text{Min}_{1 \leq i \leq n} |\lambda - \lambda_i| \leq \frac{\|\eta\|_2}{\|x\|_2} \quad (22.1.9)$$

با استفاده از (21.1.9)

$$\eta = UDU^*x - \lambda x$$

$$U^*\eta = DU^*x - \lambda U^*x = (D - \lambda I)U^*x$$

اگر λ یک ویژه مقدار A باشد، درستی (22.1.9) روشن است. بنابراین اگر فرض کنیم $\lambda \neq \lambda_1, \dots, \lambda_n$ از کلیت مسأله کاسته نخواهد شد. از این رو $D - \lambda I$ ناتکین بوده، و

$$U^*x = (D - \lambda I)^{-1}U^*\eta$$

$$\|U^*x\|_2 \leq \|(D - \lambda I)^{-1}\|_2 \|U^*\eta\|_2$$

مسأله ۱۳ فصل ۷ را به یاد می آوریم، که ایجاب می کند

$$\|U^*x\|_2 = \|x\|_2 \quad \|U^*\eta\|_2 = \|\eta\|_2$$

پس با استفاده از تعریف نرم ماتریسی، داریم

$$\|x\|_2 \leq [\text{Max}_{1 \leq i \leq n} |(\lambda - \lambda_i)^{-1}|] \|\eta\|_2$$

که با (22.1.9) هم ارز است.

فایده (22.1.9) بعداً در (15.2.9) بخش ۲.۹، با استفاده از یک زوج ویژه مقدار-ویژه بردار تقریبی حاصل از روش توانی نشان داده شده است.

پایداری ویژه مقدارها در ماتریسهای نامستقرن برای اینکه بتوانیم به امکان ناپایداری در مسأله ویژه مقدار ماتریس نامستقرن به طور مؤثر پردازیم، لازم است درک بهتری از ماهیت آن ناپایداری داشته باشیم. برای مثال، یکی از نتایج تحلیل ناپایداری این است که تبدیلات تشابهی یکانی، شرایط مسأله را بدتر نخواهند کرد.

مانند قبل، فرض می کنیم A دارای شکل قطری متعارف ژوردان باشد:

$$P^{-1}AP = \text{diag}[\lambda_1, \dots, \lambda_n] \equiv D \quad (23.1.9)$$

در این صورت $\lambda_1, \dots, \lambda_n$ ویژه مقدهای A بوده و ستونهای P ویژه بردارهای متناظر آنها هستند، که آنها را u_1, \dots, u_n می نامیم. ماتریس P یکتا نیست. برای مثال، اگر F یک

ماتریس قطری نانکین باشد، آنگاه

$$(PF)^{-1}A(PF) = F^{-1}DF = D$$

با انتخاب مناسب F ، ستونهای PF طول واحد خواهند داشت. بنابراین، بی آنکه از کلیت مسأله کاسته شود، فرض می‌کنیم ستونهای P دارای طول واحد باشند:

$$Au_i = \lambda_i u_i \quad u_i^* u_i = 1 \quad i = 1, \dots, n \quad (24.1.9)$$

که

$$P = [u_1, \dots, u_n]$$

ترانزاده مزدوج (۲۳.۱.۹) را می‌گیریم

$$P^* A^* (P^*)^{-1} = D^* = \text{diag}[\bar{\lambda}_1, \dots, \bar{\lambda}_n]$$

که نشان می‌دهد ویژه‌مقدارهای A^* مزدوجهای مختلط ویژه‌مقدارهای A هستند. با نوشتن

$$(P^*)^{-1} = [w_1, \dots, w_n] \quad (25.1.9)$$

داریم

$$A^* w_i = \bar{\lambda}_i w_i \quad i = 1, \dots, n \quad (26.1.9)$$

یا هم‌ارز آن، با تشکیل ترانزاده مزدوج داریم:

$$w_i^* A = \lambda_i w_i^* \quad (27.1.9)$$

این رابطه بیان می‌کند که w_i^* ویژه‌بردار چپ A برای ویژه‌مقدار λ_i است. چون $P^{-1}P = I$ ، و از آنجا که

$$P^{-1} = \begin{bmatrix} w_1^* \\ \vdots \\ w_n^* \end{bmatrix}$$

داریم

$$w_i^* u_j = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases} \quad (28.1.9)$$

این رابطه بیان می‌کند که $\{u_i\}$ یعنی ویژه‌بردارهای A و $\{w_i\}$ یعنی ویژه‌بردارهای A^* یک مجموعهٔ دوپه‌دو متعامد ایجاد می‌کنند.

با نرمالسازی ویژه بردارهای w_i به وسیله

$$v_i = \frac{w_i}{\|w_i\|_2} \quad i = 1, \dots, n$$

عدد حقیقی و مثبت

$$s_i = v_i^* u_i = \frac{\lambda_i}{\|w_i\|_2} \quad (29.1.9)$$

را تعریف می‌کنیم. اکنون ماتریس $(P^*)^{-1}$ را می‌توانیم به شکل

$$(P^*)^{-1} = \left[\frac{v_1}{s_1}, \dots, \frac{v_n}{s_n} \right]$$

بنویسیم. و

$$A^* v_i = \bar{\lambda}_i v_i \quad \|v_i\|_2 = 1 \quad i = 1, \dots, n \quad (30.1.9)$$

اکنون پایداری یک ویژه مقدار ساده λ_k ی A را امتحان می‌کنیم. منظور از ساده بودن این است که λ_k یک ریشه چندگانه چندجمله‌یی مشخصه A نیست. نتایج را می‌توان برای ویژه مقدارهای چندگانه نیز بسط داد ولی ما کاری با آن نداریم. ماتریس اختلال یافته

$$A(\epsilon) = A + \epsilon B \quad \epsilon > 0$$

را برای ماتریس دلخواه B ی مستقل از ϵ ، در نظر می‌گیریم. ویژه مقدارهای $A(\epsilon)$ را با $\lambda_1(\epsilon), \dots, \lambda_n(\epsilon)$ نشان می‌دهیم. پس:

$$P^{-1} A(\epsilon) P = D + \epsilon C \quad C = P^{-1} B P$$

$$c_{ij} = \frac{\lambda_j}{s_i} v_i^* B u_j \quad 1 \leq i, j \leq n \quad (31.1.9)$$

ثابت می‌کنیم که

$$\lambda_k(\epsilon) = \lambda_k + \frac{\epsilon}{s_k} v_k^* B u_k + O(\epsilon^2) \quad (32.1.9)$$

در به دست آوردن این نتیجه از قضیه ۱.۹ گرشگورین، استفاده می‌کنیم. همچنین لازم است توجه کنیم که برای هر ماتریس قطری ناتکین F ،

$$F P^{-1} A(\epsilon) P F^{-1} = D + \epsilon F C F^{-1} \quad (33.1.9)$$

و این تبدیل ویژه مقدارهای $A(\epsilon)$ را عوض نمی‌کند. F را چنین انتخاب می‌کنیم:

$$f_{ii} = \begin{cases} \epsilon\alpha & i = k \\ 1 & i \neq k \end{cases}$$

که در آن α ثابت مثبتی است که باید بعداً معین شود. بیشتر درایه‌های ماتریس (۳۳.۱.۹) عوض نشده‌اند و فقط آنهایی که روی سطر k ام و ستون k ام هستند باید بررسی شوند. این درایه‌ها چنین‌اند.

$$[D + \epsilon FCF^{-1}]_{kj} = \begin{cases} \epsilon^2 \alpha c_{kj} & j \neq k \\ \lambda_k + \epsilon c_{kk} & j = k \end{cases}$$

$$[D + \epsilon FCF^{-1}]_{ik} = \frac{1}{\alpha} c_{ik} \quad i \neq k$$

قضیه ۱.۹ را برای ماتریس (۳۳.۱.۹) به‌کار می‌بریم. مراکز دایره و شعاعهای آنها عبارت‌اند از:

$$\text{مرکز} = \lambda_k + \epsilon c_{kk} \quad r_k = \epsilon^2 \alpha \sum_{j \neq k} |c_{kj}| \quad (34.1.9)$$

$$\text{مرکز} = \lambda_i + \epsilon c_{ii} \quad r_i = \epsilon \sum_{j \neq i, k} |c_{ij}| + \frac{1}{\alpha} |c_{ik}| \quad i \neq k$$

می‌خواهیم α را آنقدر بزرگ و ϵ را آن اندازه کوچک بگیریم که دایره به مرکز $\lambda_k + \epsilon c_{kk}$ را از بقیه دایره جدا سازیم، و از این راه بدانیم که دقیقاً یک ویژه‌مقدار (۳۳.۱.۹) در داخل دایره k قرار دارد. فاصله بین مراکز دایره k و $i \neq k$ دارای کران پایین

$$|\lambda_i - \lambda_k| - \epsilon |c_{ii} - c_{kk}|$$

است، که برای مقادیر کوچک ϵ ، در حدود $|\lambda_i - \lambda_k|$ است. α را طوری می‌گیریم که

$$\frac{1}{\alpha} |c_{ik}| \leq \frac{1}{4} |\lambda_i - \lambda_k| \quad i \neq k \quad (35.1.9)$$

سیس ϵ را طوری می‌گیریم که به‌ازای جميع مقادیر $0 < \epsilon < \epsilon_0$ دایره k هیچ‌یک از دایره دیگر را قطع نکند. این کار شدنی است زیرا λ_k از ویژه‌مقدارهای دیگر λ_i متمایز است، و نابرابری (۳۵.۱.۹)، برقرار است.

با توجه به این ساختمان، قضیه ۱.۹ ایجاب می‌کند که دایره k ام دقیقاً شامل یک ویژه مقدار $A(\epsilon)$ باشد، آن را $\lambda_k(\epsilon)$ می‌نامیم. از (۳۴.۱.۹)،

$$|\lambda_k(\epsilon) - \lambda_k - \epsilon c_{kk}| \leq r_k = O(\epsilon^2)$$

و با به کار گرفتن فرمول مربوط به c_{kk} در (۳۱.۱.۹)، نتیجه مطلوب (۳۲.۱.۹) اثبات می‌شود. با کراندار کردن (۳۲.۱.۹)، به دست می‌آوریم

$$|\lambda_k(\epsilon) - \lambda_k| \leq \frac{\epsilon}{s_k} \|v_k\|_2 \|B\|_2 \|u\|_2 + O(\epsilon^2)$$

و با استفاده از (۲۴.۱.۹) و (۳۰.۱.۹)،

$$|\lambda_k(\epsilon) - \lambda_k| \leq \frac{\epsilon}{s_k} \|B\|_2 + O(\epsilon^2) \quad (36.1.9)$$

عدد s_k ، وقتی ماتریس A به اندازه مقادیر کوچک $E = \epsilon B$ اختلال یابد، به پایداری ویژه مقدار λ_k بستگی نزدیکی پیدا می‌کند. اگر A متقارن بود، $u_k = v_k$ و بنابراین $s_k = 1$ ، که همان نتیجه کیفی را که قبلاً برای ماتریسهای متقارن به دست آوردیم، به دست می‌داد. برای ماتریسهای نامتقارن، اگر s_k کاملاً کوچک باشد، آنگاه اختلال کوچک $E = \epsilon B$ می‌تواند به اختلال بزرگ در ویژه مقدار λ_k منتهی شود. چنین مسائلی بدوضع خوانده می‌شوند.

مثال مثال (۱۷.۱.۹) را به یاد آورید. در این مثال $\epsilon = 10^{-6}$ ،

$$P^{-1}AP = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \quad B = \begin{bmatrix} -1 & -1 \\ 0 & 0 \end{bmatrix}$$

$$P = \begin{bmatrix} 9 & -10 \\ \sqrt{181} & \sqrt{221} \\ 10 & -11 \\ \sqrt{181} & \sqrt{221} \end{bmatrix} \quad P^{-1} = \begin{bmatrix} -11\sqrt{181} & 10\sqrt{181} \\ -10\sqrt{221} & 9\sqrt{221} \end{bmatrix} \quad (37.1.9)$$

اگر نرم سطری را برای برآورد ضریب وضعیت (۱۶.۱.۹) به کار ببریم، آنگاه

$$K(A) \leq \|P\|_\infty \|P^{-1}\|_\infty = 419$$

u_1 و u_2 از ستونهای P ، و بردارهای w_1 و w_2 از ستونهای $(P^{-1})^T$ به دست می‌آیند، (۲۵.۱.۹) را ببینید.

برای محاسبه (۳۶.۱.۹)، وقتی $\lambda_1 = 1$ ،

$$s_1 = v_1^T u_1 = \frac{1}{\|w_1\|_2} = \frac{1}{\sqrt{(221)(181)}} \doteq 0.005$$

و $\|B\|_2 = \sqrt{2}$ از فرمول (۳۶.۱.۹) نتیجه می‌شود

$$|\lambda_1(\epsilon) - \lambda_1| \leq \frac{\sqrt{2}\epsilon}{0.005} + O(\epsilon^2) \doteq 283\epsilon + O(\epsilon^2) \quad (38.1.9)$$

مقدار واقعی خطا $0.298 - 1 = 1.298$ ، $\lambda_1(0.001) - \lambda_1 \doteq 1.298$ است و برآورد قبلی خطا برابر $0.283\epsilon \doteq 0.283$ بود. بنابراین (۳۶.۱.۹) یک برآورد کران معقول برای خطاست.

در بخشهای بعدی، برخی روشهای عددی، ابتدا توسط بعضی تبدیلات تشابهی، ماتریس A را به شکل ساده‌تری بدل می‌کنند. ما می‌خواهیم تبدیلاتی را به‌کار ببریم که اعداد s_k را حتی کوچکتر نسازند، و موجب نشوند که بدو وضعی مسأله بدتر هم بشود. از این دیدگاه، بهترین تبدیلات درکاربرد، تبدیلات متعامد یا یکانی هستند.

گیریم U یکانی باشد و $\hat{A} = U^*AU$. برای یک ویژه‌مقدار ساده λ_k ، گیریم s_k و \hat{s}_k معرف اعداد (۲۹.۱.۹) برای دو ماتریس A و \hat{A} ، باشند. اگر $\{u_i\}$ و $\{v_i\}$ ویژه‌بردارهای A و A^* باشند، آنگاه $\{U^*u_i\}$ و $\{U^*v_i\}$ ویژه‌بردارهای متناظر برای \hat{A} و \hat{A}^* می‌شوند. برای \hat{s}_k ،

$$\begin{aligned} \hat{s}_k &= (U^*v_k)^*(U^*u_k) = v_k^*UU^*u_k = v_k^*u_k \\ &= s_k \end{aligned} \quad (39.1.9)$$

بنابراین پایداری ویژه‌مقدار λ_k نه بهتر می‌شود و نه بدتر. تبدیلات یکانی، طول و زاویه بین بردارها را نیز حفظ می‌کنند (مسأله ۱۳ فصل ۷ را ملاحظه کنید). به‌طور کلی، اعمال ماتریس یکانی روی یک ماتریس داده‌شده A هیچ خللی در وضعیت مسأله ویژه‌مقدار ایجاد نمی‌کند و همین یکی از دلایل اصلی این است که شکل تبدیلات تشابهی را در حل مسأله ویژه‌مقدار ماتریسی ترجیح می‌دهند. روشهای مشابه با روش قبلی در (۲۳.۱.۹) - (۳۶.۱.۹) را، می‌توان برای به‌دست آوردن قضیه پایداری برای ویژه‌بردارهای ویژه‌مقدارهای مجزا به‌کار برد. با استفاده از همین مفروضات درباره $\{\lambda_i\}$ ، $\{u_i \equiv u_i(0)\}$ ، و $\{v_i\}$ ، مسأله ویژه‌بردار زیر را در نظر می‌گیریم

$$(A + \epsilon B)u_k(\epsilon) = \lambda_k(\epsilon)u_k(\epsilon) \quad \|u_k(\epsilon)\| = 1 \quad (40.1.9)$$

که λ_k یک ویژه مقدار ساده A است. لذا

$$u_k(\epsilon) = u_k + \epsilon a_k u_k + \epsilon \sum_{\substack{j=1 \\ j \neq k}}^n \left[\frac{v_j^* B u_k}{(\lambda_k - \lambda_j) s_j} \right] u_j + O(\epsilon^2) \quad (41.1.9)$$

برهان این رابطه به مسأله ۶ احاله شده است.

رابطه (۴۱.۱.۹) نشان می‌دهد که پایداری $u_k(\epsilon)$ به ضرایب وضعیت s_j و نزدیکی λ_k به λ_j بستگی دارد. این امر هشداری بر احتمال ناپایداری است، و مثالهای دیگری از این موضوع در مسأله ۷ داده شده‌اند. یک بررسی عمیقتر از رفتار ویژه بردارها، وقتی A دستخوش اختلالانی باشد، نیاز به بررسی زیرفضاهای ویژه بردار و ارتباط بین آنها را می‌طلبد، خاصه وقتی که ویژه مقدار λ_k ساده نباشد. برای اطلاعات بیشتر در این زمینه، گلوب و ون لون (۱۹۸۳، صص ۲۰۳ تا ۲۰۷ و ۲۷۱ تا ۲۷۵) را ببینید.

ماتریسهای دارای شکل متعارف ناقطری ژوردان از بحث در مسأله ویژه مقدار ماتریسهای که شکل کانونی ژوردان آنها ناقطری بود، احتراز کرده‌ایم. در مسأله ویژه مقدار، مسائلی در مورد ناپایداری وجود دارند که از آنچه در قضیه ۲.۹ داده شده بدترند. و مشکلات عمده‌ای در تعیین یک مبنای صحیح برای ویژه بردارها وجود دارند.

به جای ارائه یک بحث کلی، مشکلات این رده از ماتریسها را با بررسی یک حالت ساده به تفصیل مطالعه می‌کنیم. گیریم

$$A = \begin{bmatrix} 1 & 1 & 0 & & \dots & 0 \\ 0 & 1 & 1 & & 0 & \vdots \\ \vdots & & & \ddots & & \\ & & & & 1 & 1 \\ 0 & & \dots & & 0 & 1 \end{bmatrix} \quad (42.1.9)$$

ماتریسی از مرتبه n باشد. چندجمله‌ی مشخصه آن چنین است:

$$f(\lambda) = (1 - \lambda)^n$$

است و $\lambda = 1$ یک ریشه چندگانه از مرتبه n است. تنها یک مجموعه یک-بعدی از ویژه بردارها وجود دارد که توسط

$$x = [1, 0, \dots, 0]^T \quad (43.1.9)$$

تولید می‌شود.

برای $\epsilon > 0$ ، یک اختلال به شکل زیر برای A ایجاد می‌کنیم

$$A(\epsilon) = \begin{bmatrix} 1 & 1 & 0 & & \dots & 0 \\ 0 & 1 & 1 & & & 0 \\ \vdots & & & \ddots & & \vdots \\ 0 & & & & 1 & 1 \\ \epsilon & 0 & \dots & & 0 & 1 \end{bmatrix}$$

چندجمله‌یی مشخصه آن

$$f_\epsilon(\lambda) = (1 - \lambda)^n - (-1)^n \epsilon$$

خواهد بود. در اینجا n ریشه متمایز

$$\lambda_k(\epsilon) = 1 + \omega_k \epsilon^{1/n} \quad k = 1, \dots, n \quad (44.1.9)$$

وجود دارد که مجموعه $\{\omega_k\}$ ریشه‌های n ام واحد هستند:

$$\omega_k = e^{2\pi ki/n} \quad k = 1, \dots, n$$

برای اختلالهای حاصل در ویژه مقدار A ، داریم

$$|\lambda_k(\epsilon) - \lambda_k| = \epsilon^{1/n} \quad (45.1.9)$$

مثلاً، اگر $n = 10$ و $\epsilon = 10^{-10}$ ، آنگاه

$$|\lambda_k(\epsilon) - \lambda_k| = 0.1 \quad (46.1.9)$$

از قضیه قبلی (۱۰.۱.۹) کرانی به دست آمد که برحسب ϵ خطی بود، و همان‌گونه که از (۴۶.۱.۹) پیداست، (۴۵.۱.۹) از کران قبلی بسیار بدتر است.

چون $A(\epsilon)$ دارای n ویژه مقدار متمایز است، یک مجموعه کامل از n ویژه بردار مستقل خطی نیز دارد، که آنها را $x_1(\epsilon), \dots, x_n(\epsilon)$ می‌نامیم. اولین کاری که باید انجام گیرد، پیدا کردن ارتباط بین این ویژه بردارها با یک ویژه بردار x است که در (۴۳.۱.۹) داده شده است. این مسأله یک مسأله‌ای است مشکل، و همیشه در ماتریسهایی ظاهر می‌شود که شکل زوردان آنها ناقطری است. شکل ماتریسهای A و $A(\epsilon)$ به غایت ساده است، و فقط به مشکلاتی اشاره می‌کنند که ممکن است هنگامی پدید آیند که ماتریس با یک ماتریس قطری مشابه نباشد. در کار عملی، خطاهای گردکردن همیشه تضمین می‌کنند که چنین ماتریسی ویژه مقدارهای متمایز دارد. این مثال از نقطه نظر کیفی، مثال مناسبی برای نشان دادن مشکلاتی است که پیش می‌آیند.

۲.۹ روش توانی

این روش، یک روش کلاسیک، برای پیدا کردن ویژه مقدار غالب و ویژه بردار متناظر با آن، در یک ماتریس است. این، یک روش کلی نیست ولی در بعضی موارد مفید است. مثلاً برای ماتریسهای بزرگ و تنک، که روشهای بخشهای اخیر، به علت محدودیت در حافظه رایانه، نمی توانند به کار برده شوند، اغلب روشی موفق است. به علاوه، روش بارست معکوس، که در بخش ۶.۹ توضیح داده شده، روش توانی است که برای یک ماتریس معکوس مناسب به کار گرفته شده است. و مطالب این بخش، مقدمه ای بر مطالب اخیر است. به کار بردن روش توانی به عنوان یک برنامه همه منظوره رایانه‌یی در مورد رده‌های گوناگون ماتریسهای بزرگ، بسیار مشکل است. ولی می توان آن را برای بعضی رده‌های خاص به آسانی به کار برد. فرض می کنیم A یک ماتریس حقیقی $n \times n$ باشد که شکل متعارف ژوردان آن قطری است. گیریم $\lambda_1, \dots, \lambda_n$ معرّف ویژه مقدارهای A و x_1, \dots, x_n ویژه بردارهای متناظر با آنها باشند که یک پایه برای \mathbb{C}^n تشکیل می دهند. به علاوه فرض می کنیم

$$|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n| \geq 0 \quad (۱.۲.۹)$$

گرچه این حالت، حالتی است کاملاً خاص، ولی توجه اصلی ما به کاربرد روش توانی در این حالت است. بحث مربوط به آن به سادگی قابل تعمیم به حالتی است که یک ویژه مقدار غالب چندگانگی هندسی $r > 1$ داشته باشد. (مسأله ۱۰ را ملاحظه کنید). توجه داشته باشید که لازمه این مفروضات حقیقی بودن λ_1 و x_1 است.

گیریم $z^{(0)}$ یک حدس اولیه حقیقی برای مضر بی از x_1 باشد. اگر هیچ روش منطقی برای انتخاب $z^{(0)}$ در دست نباشد، آنگاه یک مولد عددی تصادفی برای هر مؤلفه انتخاب می کنیم. تعریف می کنیم

$$w^{(m)} = Az^{(m-1)} \quad (۲.۲.۹)$$

گیریم β_m یک مؤلفه $w^{(m)}$ باشد که اندازه آن ماکسیمم است. تعریف می کنیم

$$z^{(m)} = \frac{w^{(m)}}{\beta_m} \quad m \geq 1 \quad (۳.۲.۹)$$

نشان می دهیم که بردارهای $\{z^{(m)}\}$ وقتی $m \rightarrow \infty$ ، کمیت $\|x^{(1)}\| / \|\hat{\sigma}_m x^{(1)}\|$ را با هر یک از مقادیر $\hat{\sigma}_m = \pm 1$ ، تقریب می زنند.

ابتدا نشان می دهیم که:

$$z^{(m)} = \sigma_m \frac{A^m z^{(0)}}{\|A^m z^{(0)}\|_\infty} \quad \sigma_m = \pm 1 \quad m \geq 1 \quad (۴.۲.۹)$$

اولاً، $w^{(1)} = Az^{(0)}$

$$\beta_1 = \sigma_1 \|w^{(1)}\|_\infty = \sigma_1 \|Az^{(0)}\|_\infty \quad \sigma_1 = \pm 1$$

پس

$$z^{(1)} = \frac{w^{(1)}}{\beta_1} = \sigma_1 \frac{Az^{(0)}}{\|Az^{(0)}\|_\infty}$$

در اثبات (۴.۲.۹) برای حالت کلی $m > 1$ از استقرای ریاضی استفاده می‌کنیم. به‌عنوان مثال

برای $m = 2$

$$w^{(2)} = Az^{(1)} = \sigma_1 \frac{A^2 z^{(0)}}{\|Az^{(0)}\|_\infty}$$

$$\beta_2 = \mu \frac{\|A^2 z^{(0)}\|_\infty}{\|Az^{(0)}\|_\infty} \quad \mu = \pm 1$$

$$z^{(2)} = \frac{w^{(2)}}{\beta_2} = \sigma_1 \frac{A^2 z^{(0)}}{\|Az^{(0)}\|_\infty} \div \frac{\mu \|A^2 z^{(0)}\|_\infty}{\|Az^{(0)}\|_\infty} = \sigma_2 \frac{A^2 z^{(0)}}{\|A^2 z^{(0)}\|_\infty}$$

که در آن $\sigma_2 = \sigma_1 \mu$. حالت کلی m اساساً به همین شکل به‌دست می‌آید.

برای بررسی همگرایی $\{z^{(m)}\}$ ، ابتدا $z^{(0)}$ را برحسب ویژه‌بردارهای پایه $\{x_j\}$ بسط می‌دهیم:

$$z^{(0)} = \sum_{j=1}^n \alpha_j x_j$$

فرض می‌کنیم $\alpha_1 \neq 0$ ، که انتخاب تصادفی $z^{(0)}$ معمولاً صحت آن را تضمین می‌کند. همچنین می‌توان نشان داد که α_1 حقیقی است. پس

$$\begin{aligned} A^m z^{(0)} &= \sum_{j=1}^n \alpha_j A_j^m x = \sum_{j=1}^m \alpha_j \lambda_j^m x_j \\ &= \lambda_1^m \left[\alpha_1 x_1 + \sum_{j=2}^n \alpha_j \left(\frac{\lambda_j}{\lambda_1}\right)^m x_j \right] \end{aligned} \quad (5.2.9)$$

با توجه به (۱.۲.۹) داریم

$$\left(\frac{\lambda_j}{\lambda_1}\right)^m \rightarrow 0 \quad m \rightarrow \infty \quad \text{وقتی } 2 \leq j \leq m \quad (6.2.9)$$

با استفاده از این رابطه در (۵.۲.۹) و (۴.۲.۹)، داریم

$$z^{(m)} \doteq \left(\frac{\lambda_1}{|\lambda_1|}\right)^m \frac{\sigma_m \alpha_1 x_1}{|\alpha_1| \|x_1\|_\infty} \equiv \hat{\sigma}_m \cdot \frac{x_1}{\|x_1\|_\infty} \quad (7.2.9)$$

پس $\sigma_m = \pm 1$ ، و معمولاً مستقل از m است. در حالتی که x_1 یک مؤلفهٔ ماکسیمال یکتا داشته

باشد چنین است. در حالتی که x_1 بیش از یک مؤلفهٔ ماکسیمال داشته باشد ممکن است که $\hat{\sigma}_m$

با m تغییر کند (مسأله ۹ را ببینید). نرخ همگرایی در (۷.۲.۹) به $|\lambda_2/\lambda_1|$ بستگی خواهد داشت:

$$\left\| z^{(m)} - \hat{\sigma}_m \frac{x_1}{\|x_1\|_\infty} \right\|_\infty \leq C \left| \frac{\lambda_2}{\lambda_1} \right|^m \quad (۸.۲.۹)$$

زیرا همه نسبت‌های باقیمانده $|\lambda_j/\lambda_1|$ به کران $|\lambda_2/\lambda_1|$ محدود شده است.

برای به دست آوردن یک دنباله از ویژه‌مقدارهای تقریبی، گیریم k اندیس یک مؤلفه ناصفر x_1 باشد. در حالت کلی، k را یک زیرنمایه مؤلفه α_m می‌گیریم و ممکن است با m تغییر کند. تعریف می‌کنیم

$$\lambda_1^{(m)} = \frac{w_k^{(m)}}{z_k^{(m-1)}} \quad m \geq 1 \quad (۹.۲.۹)$$

برای بررسی نرخ همگرایی، از (۴.۲.۹) و (۵.۲.۹) استفاده می‌کنیم:

$$\begin{aligned} \lambda_1^{(m)} &= \left[\sigma_m \cdot \frac{A^m z^{(0)}}{\|A^{m-1} z^{(0)}\|} \right]_k \div \left[\sigma_m \cdot \frac{A^{m-1} z^{(0)}}{\|A^{m-1} z^{(0)}\|} \right]_k \\ &= \frac{\lambda_1^m \left[\alpha_1 x_1 + \sum_{j=2}^n \alpha_j \left(\frac{\lambda_j}{\lambda_1} \right)^m x_j \right]_k}{\lambda_1^{m-1} \left[\alpha_1 x_1 + \sum_{j=2}^n \alpha_j \left(\frac{\lambda_j}{\lambda_1} \right)^{m-1} x_j \right]_k} \end{aligned} \quad (۱۰.۲.۹)$$

$$\lambda_1^{(m)} = \lambda_1 \left[1 + O \left(\left| \frac{\lambda_2}{\lambda_1} \right|^m \right) \right] \quad (۱۱.۲.۹)$$

نرخ همگرایی خطی است و به $|\lambda_2/\lambda_1|$ بستگی دارد. اندیس k معمولاً ثابت انتخاب می‌شود. اگر k را α_m ، اندیس مؤلفه ماکسیمال $w^{(m)}$ بگیریم، و اگر x_1 یک مؤلفه ماکسیمال تنها داشته باشد، آنگاه وقتی $m \rightarrow \infty$ ، مقدار k ثابت می‌شود. با بیش از یک مؤلفه ماکسیمال، k ممکن است تغییر کند همان‌گونه که در مسأله ۹ نشان داده شده است. یک روش دیگر برای تعریف $\lambda_1^{(m)}$ در (کونت و دیور^۱ (۱۹۸۰، ص ۱۹۲) داده شده است، که نیاز به انتخاب یک اندیس مؤلفه خاص ندارد:

$$\lambda_1^{(m)} = \frac{u^T w^{(m)}}{u^T z^{(m-1)}}$$

بردار u باید در $u^T x_1 \neq 0$ صدق کند، و یک انتخاب تصادفی برای u معمولاً کافی است. خطای $\lambda_1^{(m)}$ باز هم در رابطه (۱۱.۲.۹) صدق می‌کند.

جدول ۱.۹ مثالی از روش توانی

m	$z_1^{(m)}$	$z_2^{(m)}$	$z_3^{(m)}$	$\lambda_1^{(m)}$	R_m
۱	۰٫۵۰۰۷۷	۰٫۷۵۰۳۸	۱٫۰۰۰۰	۱۱٫۷۶۲۸۱۳۳	
۲	۰٫۵۲۶۲۶	۰٫۷۶۳۱۳	۱٫۰۰۰۰	۹٫۵۰۳۸۴۹۶	
۳	۰٫۵۲۴۵۹	۰٫۷۶۲۳۰	۱٫۰۰۰۰	۹٫۶۳۱۳۲۳۱	-۰٫۰۵۶۴۳
۴	۰٫۵۲۴۷۰	۰٫۷۶۲۳۵	۱٫۰۰۰۰	۹٫۶۲۲۹۶۷۴	-۰٫۰۶۵۵۵
۵	۰٫۵۲۴۶۹	۰٫۷۶۲۳۵	۱٫۰۰۰۰	۹٫۶۲۳۵۰۸۳	۰٫۰۶۴۷۴

مثال گیریم

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 3 & 4 \\ 3 & 4 & 5 \end{bmatrix} \quad (۱۲.۲.۹)$$

ویژه‌مقدارهای درست آن عبارت‌اند از:

$$\lambda_1 = ۹٫۶۲۳۴۷۵۳۸۳ \quad \lambda_2 = -۰٫۶۲۳۴۷۵۳۸۳۰ \quad \lambda_3 = ۰ \quad (۱۳.۲.۹)$$

یک حدس اولیه $z^{(۰)}$ با یک مولد عددی تصادفی تولید شده بود. پنج بارست اول همراه با نسبت‌های

$$R_m = \frac{\lambda_1^{(m)} - \lambda_1^{(m-1)}}{\lambda_1^{(m-1)} - \lambda_1^{(m-2)}} \quad (۱۴.۲.۹)$$

در جدول ۱.۹ داده شده‌اند. بارست‌های $\lambda_1^{(m)}$ با استفاده از (۹.۲.۹) و $k = ۳$ ، تعریف شده بودند.

براساس یک بحث بعدی، این نسبتها باید λ_2/λ_1 را، وقتی $m \rightarrow \infty$ ، تقریب بزنند.

ما کران خطای قابل محاسبه (۲۲.۱.۹) را که در بخش ۱.۹ به دست آمده بود به کار می‌بریم.

$$\eta = Ax - \lambda x$$

را با

$$x = x^{(۵)} \quad \lambda = \lambda_1^{(۵)}$$

محاسبه می‌کنیم. سپس از (۲۲.۱.۹) به دست می‌آوریم.

$$\text{Min}_i |\lambda_i - \lambda_1^{(۵)}| < \frac{\|\eta\|_2}{\|x\|_2} = ۳٫۳۰ \times ۱۰^{-۵} \quad (۱۵.۲.۹)$$

از مقایسه مستقیم با جواب واقعی λ_1 داریم

$$\lambda_1 - \lambda_1^{(5)} = 0.00000329$$

که نشان می‌دهد (۱۵.۲.۹) در این مورد برآورد خیلی دقیق است.

روشهای شتاب چون برای کاهش خطاهای $\lambda_1^{(m)}$ و $z^{(m)}$ یک الگوی منظم معلومی وجود دارد از این الگو می‌توان برای به دست آوردن روشهای سریعتر همگرایی استفاده کرد.

حالت ۱. انتقال ویژه مقدارها. یک مقدار ثابت b انتخاب می‌کنیم، و به جای محاسبه ویژه مقدارهای A محاسبه ویژه مقدارهای

$$B = A - bI \quad (16.2.9)$$

را قرار می‌دهیم. ویژه مقدارهای B ، $\lambda_i - b$ ، $i = 1, \dots, n$ هستند. b را طوری می‌گیریم که $\lambda_1 - b$ ویژه مقدار غالب B باشد، و b را طوری انتخاب می‌کنیم که نسبت همگرایی را مینیمم کرده‌ایم. به عنوان یک حالت خاص، برای وضوح بیشتر، فرض می‌کنیم که همه ویژه مقدارهای A حقیقی و به گونه‌ای مرتب شده باشند که

$$\lambda_1 > \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_n \quad \lambda_1 > |\lambda_n|$$

در این صورت ویژه مقدار غالب B ، بسته به اندازه b ، ممکن است $\lambda_1 - b$ یا $\lambda_n - b$ باشد. ما می‌خواهیم که b در

$$|\lambda_1 - b| > |\lambda_n - b|$$

صدق کند. نرخ همگرایی عبارت است از

$$\text{Max} \left\{ \left| \frac{\lambda_2 - b}{\lambda_1 - b} \right|, \left| \frac{\lambda_n - b}{\lambda_1 - b} \right| \right\} \quad (17.2.9)$$

اگر دقیقاً به رفتار این دو کسر، وقتی b تغییر می‌کند، بنگریم، خواهیم دید که مینیمم (۷.۲.۹) وقتی حاصل می‌شود که

$$(\lambda_1 - b) - (\lambda_2 - b) = (\lambda_n - b) - [-(\lambda_1 - b)]$$

و

$$b^* = \frac{1}{2}(\lambda_2 + \lambda_n) \quad (18.2.9)$$

انتخاب بهینه b خواهد بود. نسبت همگرایی حاصله چنین است:

$$\frac{\lambda_2 - b^*}{\lambda_1 - b^*} = -\frac{\lambda_n - b^*}{\lambda_1 - b^*} = \frac{\lambda_2 - \lambda_n}{2\lambda_1 - \lambda_2 - \lambda_n} \quad (19.2.9)$$

روشهای تجربی براساس این فرمول و فرمول (۱۱.۲.۹) می‌توانند برای تعیین مقادیر تقریبی b^* به‌کار روند.

تبدیلاتی غیر از (۱۶.۲.۹) را می‌توان برای تبدیل مجموعه ویژه‌مقدارها به گونه‌ای به‌کار برد که همگرایی بازهم سریعتر به‌دست آید. برای بحث بیشتر در مورد این مطالب، ویلکینسن (۱۹۶۵)، صص (۵۷° - ۵۸°) را ملاحظه کنید.

مثال در مثال قبلی (۱۲.۲.۹)، نسبت همگرایی نظری چنین بود

$$\frac{\lambda_2}{\lambda_1} \doteq -0.0648$$

با استفاده از مقدار بهینه b که با (۱۸.۲.۹) داده شده است، و استفاده از (۱۳.۲.۹) با تغییر ترتیب،

$$b^* = \frac{1}{2}(\lambda_2 + \lambda_n) \doteq -0.31174 \quad (20.2.9)$$

ویژه‌مقدارهای $A - bI$ عبارت‌اند از

$$9.93522 \quad 0.31174 \quad -0.31174 \quad (21.2.9)$$

نسبت همگرایی روش توانی برای $A - bI$ چنین است

$$\pm \frac{0.31174}{9.93522} \doteq \pm 0.0314$$

که از نصف قدرمطلق نسبت اصلی کمتر است.

حالت ۲. برون‌یابی این‌کن. شکل همگرایی (۱۱.۲.۹) کاملاً شبیه روشهای ریشه‌یابی به‌طور خطی همگرایی بخش ۵.۲ از فصل ۲ است. با دنبال‌کردن همان بخشی که در بخش ۶.۲ داده شد، برون‌یابی این‌کن را برای تسریع همگرایی $\{\lambda_1^{(m)}\}$ و $\{z^{(m)}\}$ در نظر می‌گیریم. برای استفاده از بحث زیر، باید در (۱.۲.۹) فرض کنیم که

$$|\lambda_2| \neq |\lambda_1| \quad (22.2.9)$$

این فرض ممکن است سست‌تر شده به $|\lambda_2| = |\lambda_1|$ بینجامد، که در نتیجه $\lambda_2 = \lambda_1$ ولی دونسبت همگرایی با اندازه مساوی و مختلف‌العلامه را مجاز نمی‌دانیم (مثال قبلی (۲۱.۲.۹) را برای این مورد، ملاحظه کنید). شیوه این‌کن را نیز می‌توان چنان اصلاح کرد که محدودیت (۲۲.۲.۹) از میان برود.

با (۲۲.۲.۹) و استفاده از (۱۰.۲.۹)،

$$\lambda_1 - \lambda_1^{(m)} \doteq cr^m \quad (23.2.9)$$

که c مقداری است ثابت، r نرخ مجهول همگرایی، و از لحاظ نظری $r = \lambda_2/\lambda_1$. اگر دقیقاً بخش ۶.۲ را دنبال کنیم نتیجه می شود که:

$$R_m = \frac{\lambda_1^{(m)} - \lambda_1^{(m-1)}}{\lambda_1^{(m-1)} - \lambda_1^{(m-2)}} \rightarrow r \quad m \rightarrow \infty \quad \text{وقتی} \quad (24.2.9)$$

و برونمایی ای تکن مقدار بهبودیافته $\hat{\lambda}_1$ را چنین می دهد:

$$\hat{\lambda}_1 = \lambda_1^{(m)} - \frac{[\lambda_1^{(m)} - \lambda_1^{(m-1)}]^2}{[\lambda_1^{(m)} - \lambda_1^{(m-1)}] - [\lambda_1^{(m-1)} - \lambda_1^{(m-2)}]} \quad m \geq 3 \quad (25.2.9)$$

محاسبه مشابهی را برای ویژه بردارهای تقریبی می توان به کار برد تا هر مؤلفه دنباله $\{z^{(m)}\}$ را شتاب دهد، اگرچه ملاحظاتی باید اعمال شوند.

مثال باز هم مثال (۱۲.۲.۹) را در نظر می گیریم. در جدول ۱.۹،

$$R_5 = -0.06474 \doteq \frac{\lambda_2}{\lambda_1} = -0.06479$$

به عنوان یک مثال از (۲۵.۲.۹)، با مقادیر $\lambda_1^{(2)}$ ، $\lambda_1^{(3)}$ و $\lambda_1^{(4)}$ از آن جدول برونمایی می کنیم. پس

$$\hat{\lambda}_1 = 9.6234814 \quad \lambda_1 - \hat{\lambda}_1 = -6.03 \times 10^{-6}$$

در مقایسه با استفاده از مقادیر دقیقتر $\lambda_1^{(5)}$ از همان جدول، $\lambda_1 - \lambda_1^{(5)} = -3.29 \times 10^{-5}$. این امر (هرگاه استفاده از برونمایی از نظر تئوری توجیه پذیر باشد)، ارزش این استفاده را مجدداً نشان می دهد.

حالت ۳. خارج قسمت ریلی-ریتس^۱. هرگاه A متقارن باشد، بهتر است که تقریبهای ویژه مقدار زیر را به کار ببریم

$$\lambda_1^{(m+1)} = \frac{(Az^{(m)}, z^{(m)})}{(z^{(m)}, z^{(m)})} = \frac{(w^{(m+1)}, z^{(m)})}{(z^{(m)}, z^{(m)})} \quad m \geq 0 \quad (26.2.9)$$

ما معمولاً قرارداد استاندارد ضرب داخلی را به کار می بریم:

$$(w, z) = \sum_1^n w_i z_i \quad w, z \in \mathbf{R}^n$$

برای تحلیل این دنباله (۲۶.۲.۹)، توجه کنید که تمام ویژه مقدرهای A حقیقی اند و ویژه بردارهای x_1, \dots, x_n می‌توانند یکا متعامد انتخاب شوند. در این صورت (۲.۲.۹)، (۴.۲.۹)، (۵.۲.۹) همراه با (۲۶.۲.۹) به ما می‌دهند

$$\lambda_1^{(m+1)} = \frac{\sum_{j=1}^n |\alpha_j|^2 \lambda_j^{2m+1}}{\sum_{j=1}^n |\alpha_j|^2 \lambda_j^{2m}}$$

$$\lambda_1^{(m+1)} = \lambda_1 \left[1 + O \left(\left(\frac{\lambda_2}{\lambda_1} \right)^{2m} \right) \right] \quad (27.2.9)$$

نسبت همگرایی $\lambda_1^{(m)}$ به λ_1 برابر $(\lambda_2/\lambda_1)^2$ است، که بهبودی برای نسبت اصلی λ_2/λ_1 در (۲۷.۲.۹) است.

این رویه یک رویه کلاسیک معروف است و جنبه‌های اضافی دیگری در بردارد که در بعضی مسائل دیگر کاربرد دارند. برای بحث بیشتر ویلکینسن (۱۹۶۵، صص ۱۷۲ تا ۱۷۸) را ببینید.

مثال در مثال (۱۲.۲.۹)، ویژه بردار تقریبی $z^{(r)}$ را در (۲۶.۲.۹) به کار می‌بریم. پس

$$\lambda_1^{(r)} = \frac{(Az^{(r)}, z^{(r)})}{(z^{(r)}, z^{(r)})} \doteq 9.623464$$

که به همان دقت $\lambda_1^{(5)}$ است که قبلاً به دست آمده بود.

روش توانی را می‌توان وقتی که یک ویژه مقدار غالب تنها وجود نداشته باشد نیز به کار برد، ولی در آن صورت الگوریتم پیچیده تر خواهد شد. روش توانی را می‌توان برای تعیین ویژه مقدارهایی غیر از ویژه مقدار غالب نیز به کار برد. این عمل متضمن فرایندی است به نام کاهش A برای حذف λ به عنوان یک ویژه مقدار. برای یک بحث کامل در تمام زمینه‌های روش توانی، گلوب و ون لون (۱۹۸۳، صص ۲۰۸-۲۱۸) و ویلکینسن (۱۹۶۵، فصل ۹) را ببینید. گرچه روش توانی یک روش مفید در بعضی شرایط است، ولی باید تأکید کرد که روشهایی که در بخشهای بعدی می‌آیند معمولاً کاراترند. برای یک تغییر همگرایی سریع در روش توانی و خارج قسمت ریلی-ریتس، بارست خارج قسمت ریلی برای ماتریسهای متقارن در پارلت (۱۹۸۰، ص ۷۰) را ببینید.

۳.۹ تبدیلات متعامد با استفاده از ماتریسهای هاوسهولدر

به عنوان یک مرحله در پیدا کردن ویژه مقدارهای یک ماتریس، اغلب استفاده از تبدیلات تشابهی، مسأله را به یک صورت ساده تری بدل می کند. ماتریسهای متعامد رده ماتریسهای خواهند بود که ما برای این تبدیلات به کار می بریم. در (۳۹.۱.۹) نشان داده شده بود که تبدیلات متعامد وضعیت یا پایداری ویژه مقدارهای یک ماتریس نامتقارن را بدتر نمی کنند. همچنین، ماتریسهای متعامد ویژگیهای مطلوب دیگر انتشار خطا را نیز در بر دارند، که یک نمونه از آن بعداً در این بخش داده شده است؛ و به همین دلایل، ماتریسهای خود را به تبدیلاتی که در آنها از ماتریسهای متعامد استفاده می شود محدود می کنیم. این بخش را با نگاهی به رده خاص ماتریسهای متعامدی که به ماتریسهای هاوسهولدر معروف اند آغاز می کنیم. سپس نشان می دهیم که چگونه یک ماتریس هاوسهولدر بسازیم که یک بردار داده شده را به شکل ساده تری بدل کند. با این ساختمان به عنوان یک ابزار، به دو تبدیل یک ماتریس داده شده A توجه می کنیم: (۱) به دست آوردن تجزیه QR آن، و (۲) ساختن یک ماتریس مشابه سه قطری برای حالتی که A متقارن است. این شکلها در بخشهای بعد، در محاسبه ویژه مقدارهای A ، به کار می روند. از نظر نمادگذاری، توجه داشته باشید که اصطلاح متعامد را به ماتریسهای حقیقی، محدود کنیم. ولی خیلی متداول شده است که متعامد را در این زمینه به جای یکانی، در حالت کلی مختلط، به کار می برند، و ما همین نمادگذاری را می پذیریم. خوانندگان باید واژه تعامد را، وقتی برای ماتریسهای مختلط به کار برده می شود، به مفهوم یکانی منظور کنند.

گیریم $w \in \mathbb{C}^*$ با $\|w\|_2 = \sqrt{w^*w} = 1$. تعریف می کنیم

$$U = I - 2ww^* \quad (1.3.9)$$

این حالت کلی ماتریس هاوسهولدر است.

مثال برای $n = 3$ ، می خواهیم

$$w = [w_1, w_2, w_3]^T \quad |w_1|^2 + |w_2|^2 + |w_3|^2 = 1$$

ماتریس U چنین داده می شود:

$$U = \begin{bmatrix} 1 - 2|w_1|^2 & -2w_1\bar{w}_2 & -2w_1\bar{w}_3 \\ -2\bar{w}_1w_2 & 1 - 2|w_2|^2 & -2w_2\bar{w}_3 \\ -2\bar{w}_1w_3 & -2\bar{w}_2w_3 & 1 - 2|w_3|^2 \end{bmatrix}$$

برای حالت خاص

$$w = \left[\frac{1}{3}, \frac{2}{3}, \frac{2}{3} \right]^T$$

داریم

$$U = \frac{1}{9} \begin{bmatrix} 7 & -4 & -4 \\ -4 & 1 & -8 \\ -4 & -8 & 1 \end{bmatrix}$$

ابتدا ثابت می‌کنیم که U ارمیتی و متعامد است. برای نشان دادن ارمیتی بودن

$$\begin{aligned} U^* &= (I - 2ww^*)^* = I^* - 2(ww^*)^* \\ &= I - 2(w^*)^*w^* = I - 2ww^* = U \end{aligned}$$

برای نشان دادن تعامد،

$$\begin{aligned} U^*U &= U^2 = (I - 2ww^*)^2 \\ &= I - 4ww^* + 4(ww^*)(ww^*) \\ &= I \end{aligned}$$

زیرا استفاده از قانون شرکتپذیری و $w^*w = 1$ ایجاب می‌کند که داشته باشیم

$$(ww^*)(ww^*) = w(w^*w)w^* = ww^*$$

ماتریس U در مثال قبل این ویژگیها را نشان می‌دهد. در مسأله ۱۲، یک تعبیر هندسی از تابعخطی $T(x) = Ux$ برای یک ماتریس هاوسهولدر U ارائه می‌دهیم.معمولاً بردارهای w را با چند مؤلفهٔ اول صفر به‌کار می‌بریم،

$$w = [0, \dots, 0, w_r, \dots, w_n]^T = [O_{r-1}, \hat{w}^T]^T \quad (2.3.9)$$

که در آن $\hat{w} \in \mathbb{C}^{n-r+1}$ پس

$$U = \begin{bmatrix} I_{r-1} & O \\ O & I_{n-r+1} - 2\hat{w}\hat{w}^* \end{bmatrix} \quad (3.3.9)$$

ضرب قدامی یک ماتریس A در این U ، $r-1$ سطر اول A را تغییر نمی‌دهد، و ضرب خلفیدر آن، $r-1$ ستون اول A را تغییر نخواهد داد. برای بقیهٔ این بخش ما فرض می‌کنیم تمامماتریسها و بردارها حقیقی‌اند، تا از مقادیر مختلط احتمالی برای w احتراز کنیم.

ماتریسهای هائوسهولدر برای تبدیل بردار ناصفر به بردار جدیدی که قسمت عمده مؤلفه‌های آن صفر است به کار می‌برند. گیریم $b \neq 0$ داده شده است، $b \in \mathbf{R}^n$ ، و بنا به فرض می‌خواهیم ماتریس U به شکل (۱.۳.۹) را به گونه‌ای تولید کنیم که Ub در موضع $r+1$ ام تا m ، برای مقادیر داده شده $r \geq 1$ ، شامل صفر باشد. w را مطابق (۲.۳.۹) انتخاب می‌کنیم. در این صورت $r-1$ عنصر اول b و Ub یکی هستند.

برای ساده کردن کار اخیر می‌گیریم، $m = n - r + 1$

$$w = \begin{bmatrix} O_{r-1} \\ v \end{bmatrix} \quad b = \begin{bmatrix} c \\ d \end{bmatrix}$$

که $v, d \in \mathbf{R}^m$ ، $c \in \mathbf{R}^{r-1}$. در این صورت محدودیت ما در شکل Ub ایجاب می‌کند که $r-1$ مؤلفه اول Ub ، c باشد، و برای مقداری از α ،

$$(I - 2vv^T)d = [\alpha, 0, \dots, 0]^T \quad \|v\|_2 = 1 \quad (۴.۳.۹)$$

چون $I - 2vv^T$ متعامد است، طول d حفظ می‌شود (مسأله ۱۳ فصل ۷)؛ و بنابراین

$$|\alpha| = \|d\|_2 \equiv S$$

$$\alpha = \pm S = \pm \sqrt{d_1^2 + \dots + d_m^2} \quad (۵.۳.۹)$$

تعریف می‌کنیم

$$p = v^T d$$

با توجه به (۴.۳.۹)،

$$d - 2pv = [\alpha, 0, \dots, 0]^T \quad (۶.۳.۹)$$

از ضرب در v^T و استفاده از $\|v\|_2 = 1$ نتیجه می‌شود که

$$p = -\alpha v_1 \quad (۷.۳.۹)$$

از گذاشتن این مقدار در اولین مؤلفه (۶.۳.۹) نتیجه می‌شود

$$d_1 + 2\alpha v_1^2 = \alpha$$

$$v_1^2 = \frac{1}{2} \left[1 - \frac{d_1}{\alpha} \right] \quad (۸.۳.۹)$$

علامت α در (۵.۳.۹) را به شکل زیر انتخاب می‌کنیم

$$\text{sign}(\alpha) = -\text{sign}(d_1) \quad (۹.۳.۹)$$

این انتخاب v_1 را ماکسیمم می‌کند، و از هرگونه خطای ممکن از دست دادن ارقام بامعنی در محاسبه v_1 جلوگیری می‌کند. علامت v_1 مهم نیست. با در دست داشتن v_1 از (۷.۳.۹)، p به دست می‌آید. به (۶.۳.۹) برمی‌گردیم و مؤلفه‌ها را از ۲ تا m می‌گیریم:

$$v_j = \frac{d_j}{\sqrt{p}} \quad j = 2, 3, \dots, m \quad (۱۰.۳.۹)$$

احکام (۵.۳.۹)، (۷.۳.۹) تا (۹.۳.۹)، v را کاملاً مشخص می‌کنند، و بنابراین w و U نیز معین می‌شوند. تعداد محاسبات عبارت است از $2m + 2$ ضرب و تقسیم و دو بار ریشه دوم. در عمل می‌توان از ریشه دوم که v_1 را معین می‌کند احتراز کرد، زیرا وقتی ماتریس ww^T تشکیل شد، ریشه دوم از بین می‌رود. دنباله‌ای از چنین تبدیلات بردار b ، مورد استفاده قرار می‌گیرند تا به طور منظم ماتریسها را به شکلهای ساده‌تری تبدیل کنند.

مثال بردار داده‌شده

$$b = [2, 2, 1]^T$$

را در نظر می‌گیریم. ماتریس U را به گونه‌ای حساب می‌کنیم که Ub در دو موضع آخری خود صفر داشته باشد. برای کمک به ساختمان زیر، بعضی محاسبات میانی نیز آورده شده‌اند. توجه داشته باشید که در این حالت $w = v$ و $b = d$. بنابراین

$$\alpha = -3 \quad v_1 = \sqrt{\frac{5}{6}} \quad p = \sqrt{\frac{15}{2}}$$

$$v_2 = \frac{2}{\sqrt{30}} \quad v_3 = \frac{1}{\sqrt{30}}$$

ماتریس U به شکل زیر داده می‌شود

$$U = \begin{bmatrix} 2 & 2 & 1 \\ -\frac{2}{3} & -\frac{2}{3} & -\frac{1}{3} \\ 2 & 11 & 2 \\ -\frac{2}{3} & \frac{15}{2} & \frac{15}{2} \\ 1 & 2 & 14 \\ -\frac{1}{3} & -\frac{15}{2} & \frac{15}{2} \end{bmatrix}$$

$$Ub = [-3, 0, 0]^T$$

تجزیه یک ماتریس به QR برای ماتریس مفروض A ، نشان می‌دهیم که یک ماتریس متعامد Q و یک ماتریس بالامثلثی R وجود دارند به طوری که

$$A = QR \quad (۱۱.۳.۹)$$

گیریم

$$P_r = I - 2w^{(r)}w^{(r)T} \quad r = 1, \dots, n-1 \quad (۱۲.۳.۹)$$

که در آن $w^{(r)}$ به شکل (۲.۳.۹)، با $r-1$ صفر، در ابتدای آن است. اگر A را برحسب ستونهای آن، $A_{*1}, A_{*2}, \dots, A_{*n}$ بنویسیم، خواهیم داشت

$$P_1 A = [P_1 A_{*1}, \dots, P_1 A_{*n}]$$

P_1 و $w^{(1)}$ را با ساختمان (۵.۳.۹) - (۱۰.۳.۹) با A_{*1} انتخاب می‌کنیم. در این صورت زیرقطر $P_1 A$ در اولین ستون همه صفر خواهند بود.

P_2 را به همین شکل انتخاب می‌کنیم به گونه‌ای که زیرقطر $P_2 P_1 A$ در ستون دوم همه صفر باشند، ابتدا باید توجه کرد که چون $w^{(2)}$ در موضع اول دارای صفر است، و چون $P_1 A$ در ستون اول در زیرقطر صفر است، حاصلضربهای $P_2 P_1 A$ در سطر و ستون اول، یک مقدار دارند. اکنون P_2 و $w^{(2)}$ را مانند قبل در (۵.۳.۹) تا (۱۰.۳.۹) انتخاب می‌کنیم، با b برابر ستون دوم $P_1 A$. با ادامه این کار برای هر ستون A ، یک ماتریس بالامثلثی به شکل زیر خواهیم داشت

$$R \equiv P_{n-1} \dots P_1 A \quad (۱۳.۳.۹)$$

اگر در مرحله r ام ساختمان، تمام عناصر در زیر قطر ستون r ام صفر باشند، آنگاه تنها $P_r = I$ را انتخاب کرده و به مرحله بعدی می‌رویم. برای تکمیل ساختمان، تعریف می‌کنیم

$$Q^T = P_{n-1} \dots P_1 \quad (۱۴.۳.۹)$$

که متعامد است. پس، همان‌گونه که مورد نظر بود، $A = QR$.

مثال ماتریس زیر را در نظر می‌گیریم

$$A = \begin{bmatrix} 4 & 1 & 1 \\ 1 & 4 & 1 \\ 1 & 1 & 4 \end{bmatrix}$$

$$w^{(1)} = [0, 985599, 0, 119573, 0, 119573]^T$$

$$A_2 = P_1 A_1 \begin{bmatrix} -4,24264 & -2,12132 & -2,12132 \\ 0 & 3,62132 & 0,621321 \\ 0 & 0,621321 & 3,62132 \end{bmatrix}$$

$$w^{(2)} = [0, 0, 996393, 0, 0, 848572]^T$$

$$R = P_2 A_2 = \begin{bmatrix} -4,24264 & -2,12132 & -2,12132 \\ 0 & -3,67423 & -1,22475 \\ 0 & 0 & 3,46410 \end{bmatrix}$$

برای تجزیه $A = QR$ ، $Q = P_1 P_2$ را حساب می‌کنیم. ولی در بیشتر حالات این محاسبه برای تولید صریح Q کارا نخواهد بود. ما در این مورد توضیح مختصری خواهیم داد. چون متعامد بودن Q مستلزم $\det(Q) = \pm 1$ داریم.

$$|\det(A)| = |\det(Q) \det(R)| = |\det(R)| = 53,99999$$

این عدد با این واقعیت که ویژه‌مقدارهای A برابر $\lambda = 3, 3, 6$ و حاصلضرب آنها $54 = \det(A)$ است، سازگاری دارد.

بحث در تجزیه QR بجاست که بدانیم تا چه اندازه تجزیه $A = QR$ یکتاست. برای A ی ناکتین، فرض می‌کنیم

$$A = Q_1 R_1 = Q_2 R_2 \quad (15.3.9)$$

پس R_1 و R_2 نیز باید ناکتین باشند، و

$$Q_2^T Q_1 = R_2 R_1^{-1}$$

معکوس یک ماتریس بالامتلی، ماتریسی است بالامتلی و حاصلضرب دو ماتریس بالامتلی یک ماتریس بالامتلی است. بنابراین $R_2 R_1^{-1}$ یک ماتریس بالامتلی است. همچنین، حاصلضرب دو ماتریس متعامد، ماتریسی است متعامد، بنابراین $Q_2^T Q_1$ متعامد است. این حقیقت به ما می‌گوید

که $R_2 R_1^{-1}$ متعامد است. ولی مشکل نیست که نشان دهیم تنها ماتریسهای بالامثلثی متعامد، ماتریسهای قطری اند. برای یک ماتریس قطری D

$$R_2 R_1^{-1} = D$$

چون $R_2 R_1^{-1}$ متعامد است،

$$D^T = I$$

چون فقط با ماتریسهای حقیقی سروکار داریم، عناصر قطر D ، یا 1 هستند یا -1 . از ترکیب این نتایج،

$$Q_2 = Q_1 D \quad R_2 = D R_1 \quad (۱۶.۳.۹)$$

نتیجه می‌گیریم که علامتهای عناصر قطری R در $A = QR$ ، را می‌توان به دلخواه انتخاب کرد، ولی بقیه تجزیه، به‌طور یکتا تعیین می‌شود.

یک موضوع عملی دیگر این است که معین کنیم چگونه ماتریس R در (۱۳.۳.۹) را حساب کنیم. گیریم

$$A_r = P_r A_{r-1} = [I - 2w^{(r)} w^{(r)T}] A_{r-1} \quad r = 1, 2, \dots, n-1 \quad (۱۷.۳.۹)$$

که در آن $A_0 = A$ ، $A_{n-1} = R$. اگر P_r را محاسبه و آن را در A_{r-1} ضرب کنیم تا A_r به‌دست آید، تعداد ضربها برابر خواهد بود با

$$(n-r+1)^2 + \frac{1}{2}(n-r+2)(n-r+1)$$

روش بسیار کاراتری برای محاسبه A_r وجود دارد. (۱۷.۳.۹) را دوباره به شکل زیر می‌نویسیم

$$A_r = A_{r-1} - 2w^{(r)} [w^{(r)T} A_{r-1}] \quad (۱۸.۳.۹)$$

ابتدا $w^{(r)T} A_{r-1}$ را محاسبه می‌کنیم، سپس $w^{(r)} [w^{(r)T} A_{r-1}]$ و A_r را. این محاسبات تقریباً به

$$2(n-r)(n-r+1) + (n-r+1) \quad (۱۹.۳.۹)$$

عمل ضرب نیاز دارد که نشان می‌دهد (۱۹.۳.۹) راه مرتجیحی برای محاسبه A_r ، و بالاخره $R = A_{n-1}$ است. این عملیات، هزینه به‌دست آوردن $w^{(r)}$ را در بر ندارد، که به‌دنبال (۱۰.۳.۹) قبلاً مورد بحث قرار گرفته است.

اگر لازم باشد که ماتریسهای P_1, \dots, P_{n-1} را برای استفاده‌های بعدی ذخیره نماییم، فقط هر ستون $w^{(r)}$ ، $r = 1, \dots, n-1$ را نگه می‌داریم. اولین عنصر ناصفر $w^{(r)}$ را که در موضع r ام است، در یک جای ویژه حافظه ذخیره می‌کنیم، و بقیه عناصر ناصفر $w^{(r)}$ واقع در موضع $r+1$ ام تا n ام، از ستون r ام ماتریس A_r و R را که در زیر قطر واقع‌اند ذخیره می‌نماییم. ماتریس Q در (۱۴.۳.۹) می‌تواند صریحاً محاسبه شود. ولی همان‌گونه که ساختمان (۱۸.۳.۹) نشان می‌دهد، ما برای ضرب در ماتریسهای دیگر، Q را به‌طور صریح نیاز نداریم.

استفاده اصلی از تجزیه A به QR در تعریف روش QR برای محاسبه ویژه‌مقدارهای A خواهد بود که در بخش ۵.۹ عرضه شده است. این تجزیه را برای حل یک دستگاه خطی $Ax = b$ نیز می‌توان به‌کار برد. در اینجا تجزیه مستقیماً به دستگاه هم‌ارز $Rx = Q^T b$ منجر خواهد شد، و چون Q متعامد است، خطای بسیار کمی دخالت خواهد کرد. دستگاه $Rx = Q^T b$ یک بالامثلثی است و می‌توان آن را با استفاده از جایگذاری پسرو از یک طریق پایداری حل کرد. در حالتی که A بدووضع است، این یک راه برتر در حل دستگاه خطی $Ax = b$ می‌تواند باشد. برای بحث از خطاهای موجود در استفاده از تجزیه QR و مقایسه آن با روش حذف گاوسی در حل $Ax = b$ ، به ویلکینسن (۱۹۶۵، صص ۲۴۴ تا ۲۴۹) مراجعه کنید. ما این مبحث را باز هم در بخش ۷.۹، هنگام بحث از حل کمترین مربعات دستگاههای خطی فرامعین، دنبال می‌کنیم.

تبدیل یک ماتریس متقارن به شکل سه‌قطری گیریم A یک ماتریس حقیقی متقارن باشد. برای پیدا کردن ویژه‌مقدارهای A ، معمولاً ابتدا آن را با تبدیلات تشابهی متعامد، به یک شکل سه‌قطری برمی‌گردانند. سپس ویژه‌مقدارهای ماتریس سه‌قطری توسط نظریه دنباله‌های استورم، که در بخش ۴.۹ ارائه شده، یا توسط QR که در بخش ۵.۹ آمده است، محاسبه می‌شوند. برای ماتریسهای متعامد، ماتریسهای (۳.۳.۹) هاوسهولدر را به‌کار می‌بریم.

گیریم

$$P_r = I - 2w^{(r+1)}w^{(r+1)T} \quad r = 1, \dots, n-2 \quad (20.3.9)$$

که در آن $w^{(r+1)}$ مانند (۲.۳.۹) تعریف شده است:

$$w^{(r+1)} = [0, \dots, 0, w_{r+1}, \dots, w_n]^T$$

[به تغییری که در نمادگذاری P_r در (۱۲.۳.۹) هنگام تعریف تجزیه QR ، صورت گرفته توجه

کنید.] ماتریس

$$A_r = P_r^T A P_r = P_r A P_r$$

با A متشابه است، عنصر a_{11} تغییر نکرده است، و A_2 متقارن خواهد بود. $w^{(2)}$ و P_1 را می‌سازیم تا شکل زیر به‌ازای مقداری چون \hat{a}_{21} به‌دست آید

$$P_1 A_{*1} = [a_{11}, \hat{a}_{21}, 0, \dots, 0]^T$$

بردار A_{*1} اولین ستون A است. از (۵.۳.۹) تا (۱۰.۳.۹) با $m = n - 1$ و

$$d = [a_{21}, a_{31}, \dots, a_{n1}]^T$$

استفاده می‌کنیم. برای مثال، از (۸.۳.۹)،

$$w_1^2 = v_1^2 = \frac{1}{2} \left[1 - \frac{a_{21}}{\alpha} \right]$$

$$\alpha = -\text{sign}(a_{21}) \sqrt{a_{21}^2 + \dots + a_{n1}^2}$$

با داشتن P_1 و $P_1 A$ ، از ضرب خلفی در P_1 اولین ستون $P_1 A$ تغییر نخواهد یافت. (خواننده باید این موضوع را تحقیق کند). تقارن A_2 از رابطه زیر نتیجه می‌شود

$$A_2^T = (P_1 A P_1)^T = p_1^T A^T P_1^T = P_1 A P_1 = A_2$$

چون A_2 متقارن است، ساختمان در اولین ستون A ایجاب می‌کند که در موضع سوم تا m ام اولین سطر و ستون A_2 همه صفر باشند.

این روند را ادامه داده، فرض می‌کنیم

$$A_{r+1} = P_r^T A_r P_r \quad r = 1, 2, \dots, n-2 \quad (21.3.9)$$

که در آن $A_1 = A$. سپس P_r را برای ایجاد صفر در موضع $r+2$ تا n از ستون r ام انتخاب می‌کنیم. در محاسبه $P_r A_{r-1}$ به علت شکل خاص P_r ، ستونهای ۱ تا $r-1$ بدون تغییر باقی می‌مانند. بردار $w^{(r+1)}$ را همانند توضیح قبلی برای $w^{(2)}$ برمی‌گزینیم. ماتریس نهایی $T \equiv A_{r-1}$ سه قطری و متقارن است.

$$T = \begin{bmatrix} \alpha_1 & \beta_1 & 0 & \dots & 0 \\ \beta_1 & \alpha_2 & \beta_2 & & \vdots \\ 0 & \beta_2 & \alpha_3 & & \\ \vdots & & & \ddots & \\ 0 & & & & \alpha_{n-1} & \beta_{n-1} \\ 0 & \dots & & & \beta_{n-1} & \alpha_n \end{bmatrix} \quad (22.3.9)$$

این شکل، شکل بسیار مناسبتری برای محاسبه ویژه مقدارهای A خواهد بود و ویژه بردارهای A را می توان به سادگی از ویژه بردارهای T به دست آورد.

T با رابطه زیر به A وابسته است

$$T = Q^T A Q \quad Q = P_1 \dots P_{n-1} \quad (۲۳.۳.۹)$$

مانند قبل در تجزیه QR ، ما به ندرت Q را صریحاً تولید می کنیم، و ترجیح می دهیم که همانند (۱۸.۳.۹)، با هریک از ماتریسهای P_r جداگانه کار کنیم. برای x ، یک ویژه بردار A ، یعنی

$$Ax = \lambda x$$

$$Tz = \lambda z \quad x = Qz \quad (۲۴.۳.۹)$$

اگر یک مجموعه از ویژه بردارهای یکا متعامد $\{z_i\}$ را برای T را تولید کنیم، آنگاه $\{Qz_i\}$ یک مجموعه یکا متعامد از ویژه بردارهای A خواهد شد، زیرا Q طول و زوایا را حفظ می کند (مسأله ۱۳ فصل ۷ را ببینید).

مثال گیریم

$$A = \begin{bmatrix} 1 & 3 & 4 \\ 3 & 1 & 2 \\ 4 & 2 & 1 \end{bmatrix}$$

پس

$$w^{(1)} = \left[0, \frac{2}{\sqrt{5}}, \frac{1}{\sqrt{5}} \right]^T$$

$$P_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & -\frac{3}{5} & -\frac{4}{5} \\ 0 & -\frac{4}{5} & \frac{3}{5} \end{bmatrix}$$

$$T = P_1^T A P_1 = \begin{bmatrix} 1 & -5 & 0 \\ -5 & \frac{73}{25} & \frac{14}{25} \\ 0 & \frac{14}{25} & -\frac{23}{25} \end{bmatrix}$$

برای تجزیه خطای این تبدیل به شکل سه قطری، قضایایی از ویلکینسن (۱۹۶۵) را در ذیل می‌آوریم. گیریم حساب رایانه، ممیز شناور دودویی و گردکردن با t رقم دودویی در جزء کسری باشد. به علاوه فرض می‌کنیم تمام حاصلضربهای داخلی

$$\sum_{i=1}^m a_i b_i$$

که در محاسبات پیش می‌آیند با دقت مضاعف انباشته شده و با دقت ساده هنگام مجموعیایی گرد شده‌اند. این حاصلضربهای داخلی در جاهای گوناگونی در محاسبه T از روی A وجود دارند. گیریم \hat{T} معرّف ماتریس سه قطری متقارن است که از A با حساب رایانه قبلی حساب شده است. فرض می‌کنیم \hat{P}_r ماتریسی باشد که از تبدیل A_{r-1} به A_r حاصل شده است. P_r تبدیل دقیق نظری این ماتریس باشد بدون هیچ خطای گردکردن و $Q = P_1 \dots P_{r-1}$ حاصلضرب دقیق این P_r ها باشد که یک ماتریس متعامد است.

قضیه ۴.۹ فرض می‌کنیم A یک ماتریس حقیقی متقارن از مرتبه n باشد. گیریم \hat{T} یک ماتریس حقیقی متقارن سه قطری باشد که از اعمال تبدیلات تشابهی (۲۰.۳.۹) هاوسهولدر بر A ، مانند (۲۱.۳.۹)، به دست آمده است. فرض می‌کنیم حساب ممیز شناوری که استفاده شده دارای همان ویژگیهای مذکور در پاراگراف قبل است. اگر $\{\lambda_i\}$ و $\{\tau_i\}$ به ترتیب ویژه مقدرهای A و \hat{T} باشند که به طور صعودی مرتب شده‌اند، آنگاه

$$\left[\frac{\sum_{i=1}^n (\tau_i - \lambda_i)^2}{\sum_{i=1}^n \lambda_i^2} \right]^{1/2} \leq c_n 2^{-t} \quad (25.3.9)$$

که در آن

$$c_n = 25(n-1)[1 + (12,36)2^{-t}]^{2n-2}$$

برای مقادیر کوچک و متوسط n ، $c_n \approx 25(n-1)$.

برهان از ویلکینسن (۱۹۶۵، ص ۱۶۱)، با استفاده از ماتریس فروبنیوس با نرم F ،

$$F(\hat{T} - Q^T A Q) \leq 2x(n-1)(1+x)^{2n-2} F(A) \quad (26.3.9)$$

که در آن $x = (۱۲,۳۶)۲^{-t}$. از نتیجه (۱۹.۱.۹) قضیه ۳.۹ ویلانت - هوفمن داریم

$$\left[\sum_{i=1}^n (\tau_i - \lambda_i)^2 \right]^{1/2} \leq F(\hat{T} - Q^T A Q) \quad (۲۷.۳.۹)$$

زیرا ویژه‌مقدارهای A و $Q^T A Q$ یکی هستند. و از مسأله ۲۸ (ب) فصل ۷،

$$F(A) = \left[\sum_{i=1}^n \lambda_i^2 \right]^{1/2}$$

از ترکیب این نتایج، (۲۵.۳.۹) به دست می‌آید.

برای بحث بیشتر از خطا، از جمله حالتی که حاصلضربهای داخلی با دقت مضاعف ذخیره نشده باشند، ویلکینسن (۱۹۶۵، صص ۲۹۷-۲۹۹) را ملاحظه نمایید. قضیه (۲۵.۳.۹) نشان می‌دهد که تبدیل به یک شکل سه‌قطری، عملی است بسیار پایدار که خطای بسیار کم جدیدی در ویژه‌مقدارها به دنبال خواهد داشت.

ماتریسهای متعامد دوران در صفحه رده‌های ماتریسهای متعامد دیگری وجود دارند که می‌توان به جای ماتریسهای هاوسهولدر از آنها استفاده کرد. رده اصلی، مجموعه دورانها در صفحه‌اند، که دوران دستگاه مختصات در صفحه به زاویه داده شده θ در صفحه مختصات یک تعبیر هندسی آن است. برای اعداد صحیح $k, l, 1 \leq k < l \leq n$ ، ماتریس متعامد $n \times n$ ، $R^{(k,l)}$ را با تعویض چهار عنصر ماتریس یکه I_n تعریف می‌کنیم. به‌ازای هر عدد حقیقی θ ، عناصر $R^{(k,l)}$ را چنین تعریف می‌کنیم

$$R_{i,j}^{(k,l)} = \begin{cases} \cos \theta & (i,j) = (k,k) \quad \text{یا} \quad (l,l) \\ \sin \theta & (i,j) = (k,l) \\ -\sin \theta & (i,j) = (l,k) \\ (I_n)_{ij} & \text{برای همه مقادیر دیگر } (i,j) \end{cases} \quad (۲۸.۳.۹)$$

به‌ازای $1 \leq i, j \leq n$.

مثال برای $n = 3$

$$R^{(1,3)} = \begin{bmatrix} \cos \theta & 0 & \sin \theta \\ 0 & 1 & 0 \\ -\sin \theta & 0 & \cos \theta \end{bmatrix}$$

به عنوان یک حالت خاص، θ را مساوی $\frac{\pi}{4}$ می‌گیریم. پس

$$R^{(1,3)} = \begin{bmatrix} \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \\ 0 & 1 & 0 \\ -\frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \end{bmatrix}$$

دورانهای $R^{(k,l)}$ صفحه‌یی را می‌توان برای همان تبدیلیهایی که با ماتریسهای هاوسهولدر انجام می‌شوند، مورد استفاده قرار داد. در بسیاری مواقع، ماتریسهای هاوسهولدر کاراترند، ولی دورانهای صفحه‌یی برای قسمتی از روش QR که در بخش ۵.۹ توضیح داده شده است، کارایی بیشتری دارند. فکر حل مسأله ویژه مقدار ماتریسهای متقارن با تبدیل آن در ابتدا به یک شکل سه قطری را گیونز^۱ در سال ۱۹۵۴ عرضه کرده است. همچنین روشهای محاسبه ویژه مقدارهای یک ماتریس سه قطری را که در بخش بعد توضیح داده می‌شود، او پیشنهاد نموده است. گیونز دورانهای $R^{(k,l)}$ صفحه را به کار برده است و ماتریسهای هاوسهولدر در ۱۹۵۸ توسط هاوسهولدر معرفی شده‌اند. برای بحث بیشتر در مورد ماتریسهای دوران و ویژگیهای آنها، گلوب و ون لون (۱۹۸۳، بخش ۴.۳)، پارلت (۱۹۸۰، بخش ۴.۶)، و همچنین مسائل ۱۵ و ۱۷ (ب) را ملاحظه نمایید.

۴.۹ ویژه مقدارهای یک ماتریس متقارن سه قطری

گیریم T یک ماتریس متقارن سه قطری از مرتبه n ، مانند (۲۲.۳.۹) باشد. چند جمله‌یی مشخصه T را محاسبه و از آن برای پیدا کردن ویژه مقدارهای T استفاده می‌کنیم. برای محاسبه

$$f_n(\lambda) \equiv \det(T - \lambda I) \quad (1.4.9)$$

دنباله

$$f_k(\lambda) = \det \begin{bmatrix} \alpha_1 - \lambda & \beta_1 & \circ & \dots & \circ \\ \beta_1 & \alpha_2 - \lambda & \beta_2 & & \vdots \\ \circ & & & \ddots & \\ \vdots & & & & \beta_{k-1} \\ \circ & \dots & & \beta_{k-1} & \alpha_k - \lambda \end{bmatrix} \quad (۲.۴.۹)$$

را برای $1 \leq k \leq n$ و $f_0(\lambda) \equiv 1$ معرفی می‌کنیم. با محاسبه مستقیم،

$$f_1(\lambda) = \alpha_1 - \lambda$$

$$f_2(\lambda) = (\alpha_2 - \lambda)(\alpha_1 - \lambda) - \beta_1^2$$

$$= (\alpha_2 - \lambda)f_1(\lambda) - \beta_1^2 f_0(\lambda)$$

فرمول $f_2(\lambda)$ یک رابطه بازگشتی کلی سه‌گانه را، که دنباله $\{f_k(\lambda)\}$ در آن صدق می‌کند نشان می‌دهد:

$$f_k(\lambda) = (\alpha_k - \lambda)f_{k-1}(\lambda) - \beta_{k-1}^2 f_{k-2}(\lambda) \quad 2 \leq k \leq n \quad (۳.۴.۹)$$

برای اثبات این رابطه، درمیان (۲.۴.۹) را برحسب آخرین سطر آن با استفاده از درمیانهای فرعی بسط می‌دهیم و نتیجه به آسانی به دست می‌آید. روش محاسبه $f_n(\lambda)$ پس از آنکه ضرایب $\{\beta_k^2\}$ محاسبه شده باشند، به $2n - 3$ ضرب نیاز دارد.

مثال گیریم

$$T = \begin{bmatrix} 2 & 1 & \circ & \circ & \circ & \circ \\ 1 & 2 & 1 & \circ & \circ & \circ \\ \circ & 1 & 2 & 1 & \circ & \circ \\ \circ & \circ & 1 & 2 & 1 & \circ \\ \circ & \circ & \circ & 1 & 2 & 1 \\ \circ & \circ & \circ & \circ & 1 & 2 \end{bmatrix} \quad (۴.۴.۹)$$

پس

$$f_0(\lambda) = 1 \quad f_1(\lambda) = 2 - \lambda$$

$$f_j(\lambda) = (2 - \lambda)f_{j-1}(\lambda) - f_{j-2}(\lambda) \quad j = 2, 3, 4, 5, 6 \quad (۵.۴.۹)$$

بدون استفاده از رابطه بازگشتی سه‌گانه (۵.۴.۹)، محاسبه $f_6(\lambda)$ بسیار پیچیده‌تر خواهد بود.

در اینجا، می‌توانیم مسأله را حل شده تلقی نماییم زیرا $f_n(\lambda)$ یک چندجمله‌بی است و روشهای زیادی برای به دست آوردن ریشه‌های یک چندجمله‌بی وجود دارند. یا ممکن است روش کلیتری، چون روش خط قاطع یا روش برنت را که هر دو در فصل ۲ تشریح شدند به کار بریم. ولی دنباله $\{f_k(\lambda) \mid 0 \leq k \leq n\}$ دارای ویژگیهای خاصی است که از آن یک دنباله استورم می‌سازد و این ویژگیها جدا کردن ویژه‌مقدارهای T را، نسبتاً آسان می‌کند. به محض اینکه ویژه‌مقدارها جدا شدند، یک روش مانند روش برنت [بخش ۸.۲ را ببینید] را می‌توان برای محاسبه سریع ریشه‌ها به کار برد. نظریه دنباله‌های استورم در هنریچی (۱۹۷۴، ص ۴۴۴) توضیح داده شده است، ولی ما فقط حالت خاص $\{f_k(\lambda)\}$ را مطالعه می‌کنیم.

قبل از بیان نتایج نظریه استورم برای $\{f_k(\lambda)\}$ ، می‌خواهیم ببینیم که اگر $\beta_l = 0$ ، چه اتفاقی خواهد افتاد. در این صورت مسأله ویژه‌مقدار را می‌توان به دو مسأله ویژه‌مقدار مجزای کوچکتر از مرتبه‌های l و $n-l$ تجزیه کرد. به عنوان یک مثال

$$T = \begin{bmatrix} \alpha_1 & \beta_1 & 0 & 0 & 0 \\ \beta_1 & \alpha_2 & 0 & 0 & 0 \\ 0 & 0 & \alpha_3 & \beta_3 & 0 \\ 0 & 0 & \beta_3 & \alpha_4 & \beta_4 \\ 0 & 0 & 0 & \beta_4 & \alpha_5 \end{bmatrix}$$

را در نظر می‌گیریم. آن را به دو بلوک T_1 و T_2 ، به ترتیب از مرتبه‌های ۲ و ۳، در امتداد قطر می‌شکنیم، و سپس می‌نویسیم

$$T = \begin{bmatrix} T_1 & 0 \\ 0 & T_2 \end{bmatrix}$$

از این رابطه داریم

$$\det[T - \lambda I_5] = \det[T_1 - \lambda I_2] \det[T_2 - \lambda I_3]$$

و می‌توانیم ویژه‌مقدارهای T را با محاسبه ویژه‌مقدارهای T_1 و T_2 به دست آوریم. مسأله ویژه‌بردار را نیز می‌توان به شکل مشابهی حل کرد. مثلاً، اگر $T_1 \hat{x} = \lambda \hat{x}$ با $\hat{x} \neq 0$ در \mathbf{R}^2 باشد، تعریف می‌کنیم

$$x = [\hat{x}^T, 0, 0, 0]^T$$

در این صورت $Tx = \lambda x$. این ساختمان را می‌توان برای محاسبه مجموعه کامل ویژه‌بردارهای T

با استفاده از ویژه‌بردارهای T_1 و T_2 به‌کار برد. در بقیه این بخش فرض می‌کنیم که در ماتریس T ، β_i ها مخالف صفرند. با این فرض، تمام ویژه‌مقدارهای T ریشه‌های ساده $f_n(\lambda)$ خواهند بود.

ویژگی دنباله استورم برای $\{f_k(\lambda)\}$ دنباله‌های $\{f_k(a)\}$ و $\{f_k(b)\}$ را می‌توان برای تعیین تعداد ریشه‌های $f_n(\lambda)$ که در بازه $[a, b]$ واقعند به‌کار برد. برای این کار، تابع صحیح‌المقدار $s(\lambda)$ را به شرح زیر معرفی می‌کنیم. $s(\lambda)$ را تعداد تطابق علامتهای اعضای متوالی دنباله $\{f_k(\lambda)\}$ تعریف می‌کنیم، و اگر مقدار عضوی نظیر $f_j(\lambda)$ مساوی صفر باشد، علامت آن را مخالف علامت $f_{j-1}(\lambda)$ می‌گیریم. می‌توان نشان داد که $f_j(\lambda) = 0$ ایجاب می‌کند که $f_{j-1}(\lambda) \neq 0$.

مثال دنباله $(f_0(\lambda), \dots, f_6(\lambda))$ در (۵.۴.۹) در آخرین مثال را در نظر می‌گیریم. به‌ازای $\lambda = 3$ ،

$$(f_0(\lambda), \dots, f_6(\lambda)) = (1, -1, 0, 1, -1, 0, 1)$$

دنباله علامتهای متناظر آن عبارت است از

$$(+, -, +, +, -, +, +)$$

$$s(3) = 2 \text{ و}$$

اکنون قضیه اساسی را که در محاسبه ریشه‌های $f_n(\lambda)$ و بنابراین در محاسبه ویژه‌مقدارهای T به‌کار رفته است بیان می‌کنیم. برهان از نظریه کلی که در هنریچی (۱۹۷۴) داده شده است نتیجه می‌شود.

قضیه ۵.۹. گریم T یک ماتریس متقارن حقیقی سه‌قطری از مرتبه n ، مانند (۲۲.۳.۹) باشد. فرض می‌کنیم دنباله $\{f_k(\lambda) \mid 0 \leq k \leq n\}$ مطابق با (۲.۴.۹) تعریف شده، و به‌ازای $i = 1, \dots, n-1$ ، و کلیه β_i ها مخالف صفر باشند. در این صورت تعداد ریشه‌های $f_n(\lambda)$ که از $a = \lambda$ بزرگترند برابر است با $s(a)$ که در پاراگراف قبل تعریف شده است. برای $a < b$ ، تعداد ریشه‌ها در بازه $a < \lambda \leq b$ برابر با $s(a) - s(b)$ خواهد بود.

محاسبه ویژه‌مقدارها قضیه ۵.۹ ابزار اساسی در تعیین جا و جداکردن ریشه‌های $f_n(\lambda)$ است. ابتدا، بازه‌ای را که شامل تمام ریشه‌ها باشد حساب می‌کنیم. با استفاده از قضیه ۱.۹ دایره‌گرشگورین،

همه ویژه مقادارها در بازه $[a, b]$ ، با

$$a = \min_{1 \leq i \leq n} \{ \alpha_i - |\beta_i| - |\beta_{i-1}| \}$$

$$b = \max_{1 \leq i \leq n} \{ \alpha_i + |\beta_i| + |\beta_{i-1}| \}$$

که $\beta_n = \beta_0 = 0$ قرار دارند.

از روش نیمسازي در بازه $[a, b]$ استفاده می‌کنیم تا آن را به زیر بازه‌های کوچکتری تقسیم کنیم. با استفاده از قضیه ۵.۹ معلوم می‌کنیم که چند ریشه در هر زیر بازه واقع است و زیر بازه‌هایی را جستجو می‌کنیم که هر کدام شامل یک ریشه باشد. اگر بعضی از ویژه مقادارها تقریباً مساوی باشند، تقسیم زیر بازه‌ها را آنقدر ادامه می‌دهیم تا ریشه با دقت کافی به دست آید. وقتی معلوم شد که زیر بازه تنها شامل یک ریشه است، می‌توانیم از یک روش همگرایی سریعتر استفاده کنیم.

مثال باز هم مثال (۴.۴.۹) را در نظر می‌گیریم. به موجب قضیه گرشگورین، ۱.۹، تمام ویژه مقادارها در $[0, 4]$ واقع‌اند. به آسانی می‌توان تحقیق کرد که نه $\lambda = 0$ یک ویژه مقدار است و نه $\lambda = 4$ یک روند منظم نیمسازي در $[0, 4]$ انجام گرفته است تا شش ویژه مقدار $f_6(\lambda)$ در زیر بازه‌های مجزا قرار گیرند. نتایج در جدول ۲.۹ به ترتیب به دست آمده، نشان داده شده‌اند. ریشه‌ها مانند زیر اندیس گذاری شده‌اند:

$$0 \leq \lambda_6 \leq \lambda_5 \leq \dots \leq \lambda_1 \leq 4$$

ریشه‌ها را می‌توان با ادامه روش نیمسازي پیدا کرد، که در این صورت به قضیه ۵.۹ دیگر نیازی نیست. ولی بهتر است از بعضی روشهای دیگر ریشه‌یابی استفاده کرد.

اگرچه تمام ریشه‌های یک ماتریس سه قطری را می‌توان با این راهکار پیدا کرد، ولی معمولاً در چنین حالتی استفاده از الگوریتم QR که در بخش بعد می‌آید، سریعتر جوابها را به ما خواهد

جدول ۲.۹ مثال در استفاده از قضیه ۵.۹

λ	$f_6(\lambda)$	$s(\lambda)$	شرح
0°	7_0°	۶	$\lambda_6 > 0$
4_0°	7_0°	۰	$\lambda_1 < 4$
2_0°	-1_0°	۳	$\lambda_2 < 2 < \lambda_3$
1_0°	1_0°	۴	$\lambda_5 < 1 < \lambda_4 < 2$
0_5°	$-1,421875$	۵	$0 < \lambda_6 < 0_5^{\circ} < \lambda_5 < 1$
3_0°	1_0°	۲	$2 < \lambda_3 < 3 < \lambda_2$
3_5°	$1,421875$	۱	$3 < \lambda_2 < 3_5^{\circ} < \lambda_1 < 4$

داد. برای ماتریسهای بزرگ، معمولاً تمام ریشه‌ها موردنظر نیستند، در چنین حالتی، روشی که در این بخش توضیح داده شده ترجیح داده می‌شود. اگر فقط بعضی از ریشه‌های مشخص، مثلاً پنج ریشه از همه بزرگتر یا همه ریشه‌ها در یک بازه داده شده، یا تمام ریشه‌های واقع در [۳، ۱] را بخواهیم، آنگاه با استفاده از قضیه ۵.۹، به آسانی می‌توان آنها را جابایی کرد.

۵.۹ روش QR

امروزه، روش QR کاراترین و متداولترین روش برای به دست آوردن همه ویژه مقدرهای یک ماتریس است. این روش ابتدا توسط فرانسیس^۱ در سال ۱۹۶۱ انتشار یافت و تاکنون به شدت موضوعی تحت بررسی قرار گرفته است. روش QR هم از لحاظ نظری و هم از لحاظ کاربردی، کاملاً پیچیده است، و ما فقط می‌توانیم مقدمه‌ای از نظریه را بیاوریم. برای الگوریتمهای اصلی، هم برای ماتریسهای متقارن و هم نامتقارن، خوانندگان را به آنچه که در EISPACK و ولکینسن و راینس (۱۹۷۱) آمده است ارجاع می‌دهیم.

برای ماتریس داده شده A ، یک تجزیه

$$A = QR$$

وجود دارد که R بالامثلثی و Q متعامد است. برای A ی حقیقی، هم Q و هم R را می‌توان حقیقی انتخاب کرد. ساختمان آنها در بخش ۳.۹ داده شده است. ما در سرتاسر این بخش فرض می‌کنیم A حقیقی است. گیریم $A_1 = A$ ، و دنباله ماتریسهای A_m ، Q_m و R_m را چنین تعریف می‌کنیم

$$A_m = Q_m R_m \quad A_{m+1} = R_m Q_m \quad m = 1, 2, \dots \quad (1.5.9)$$

چون $R_m = Q_m^T A_m$ داریم

$$A_{m+1} = Q_m^T A_m Q_m \quad (2.5.9)$$

ماتریس A_{m+1} از نظر تعامد مشابه با A_m ، و بنابراین با استقرا، مشابه با A_1 است. دنباله $\{A_m\}$ یا به یک ماتریس مثلثی می‌گراید که دارای ویژه مقدرهای A بر قطر اصلی خود است، یا به یک ماتریس تقریباً مثلثی که ویژه مقدرهای آن به سادگی قابل محاسبه است، می‌گراید. در این حالت، معمولاً همگرایی کند است و یک تکنیکی برای تسریع همگرایی به نام انتقال به‌کار برده می‌شود. روش انتقال در بخش بعد معرفی و توضیح داده می‌شود.

مثال گیریم

$$A_1 = \begin{bmatrix} 2 & 1 & 0 \\ 1 & 3 & 1 \\ 0 & 1 & 4 \end{bmatrix} \quad (3.5.9)$$

ویژه مقدارها عبارت اند از:

$$\lambda_1 = 3 + \sqrt{3} \doteq 4,7321 \quad \lambda_2 = 3,0 \quad \lambda_3 = 3 - \sqrt{3} \doteq 1,2679$$

بارستهای A_m به سرعت همگرا نمی شوند، و فقط چند تا از آنها برای نشان دادن رفتار کیفی همگرایی آنها داده شده اند

$$A_7 = \begin{bmatrix} 3,0000 & 1,0954 & 0 \\ 1,0954 & 3,0000 & -1,3416 \\ 0 & -1,3416 & 3,0000 \end{bmatrix} \quad A_7 = \begin{bmatrix} 3,7059 & 0,9558 & 0 \\ 0,9558 & 3,5214 & 0,9738 \\ 0 & 0,9738 & 1,7727 \end{bmatrix}$$

$$A_7 = \begin{bmatrix} 4,6792 & 0,2979 & 0 \\ 0,2979 & 3,0524 & 0,0274 \\ 0 & 0,0274 & 1,2684 \end{bmatrix} \quad A_8 = \begin{bmatrix} 4,7104 & 0,1924 & 0 \\ 0,1924 & 3,0216 & -0,0115 \\ 0 & -0,0115 & 1,2680 \end{bmatrix}$$

$$A_9 = \begin{bmatrix} 4,7233 & 0,1229 & 0 \\ 0,1229 & 3,0087 & 0,0048 \\ 0 & 0,0048 & 1,2680 \end{bmatrix} \quad A_{10} = \begin{bmatrix} 4,7285 & 0,0781 & 0 \\ 0,0781 & 3,0035 & -0,0020 \\ 0 & -0,0020 & 1,2680 \end{bmatrix}$$

عناصر واقع در موضع (۱، ۲) از لحاظ هندسی با نسبتی در حدود $0,64$ در هر بارست، کاهش می یابند و عناصری که در موضع (۲، ۳) هستند با نسبتی در حدود $0,42$ در هر بارست کاهش می یابند. مقدار موضع (۳، ۳) در ماتریس A_{10} برابر $1,2679$ خواهد بود که تا 5 رقم درست است.

تبدیل مقدماتی A به شکل ساده تر روش QR نسبتاً گران است، زیرا تجزیه QR وقتی چند بار تکرار شود وقتگیر است. برای کاهش هزینه، ماتریسی را که برای روش QR آماده می کنند، نخست به شکل ساده تری، که برای تجزیه QR کمتر هزینه دارد، تبدیل می کنند.

اگر A متقارن باشد، آن را به یک ماتریس سه قطری متقارن مشابه، درست همانگونه که در بخش ۳.۹ تشریح شد، بدل می کنند. اگر A نامتقارن باشد، آن را به یک ماتریس هسنبرگ مشابه،

بدل می‌کنند. ماتریس B هسنبرگ است اگر

$$b_{ij} = 0 \quad i > j + 1 \quad \text{به‌ازای جمیع مقادیر } 1 \quad (۴.۵.۹)$$

این ماتریس بالامثلثی است جز اینکه فقط یک زیرقطری ناصفر دارد. ماتریس A با استفاده از همان الگوریتمی که برای تبدیل ماتریسهای متقارن به شکل سه‌قطری به‌کار می‌رفت، به ماتریس هسنبرگ تبدیل می‌شود.

وقتی A سه‌قطری یا هسنبرگ باشد، ماتریسهای هاؤسهولدر بخش ۳.۹، در به‌دست آوردن تجزیه QR ، شکل ساده‌ای به خود می‌گیرند. ولی معمولاً دورانهای (۲۸.۳.۹) صفحه به‌جای ماتریسهای هاؤسهولدر به‌کار برده می‌شوند زیرا برای محاسبه و کاربرد در چنین وضعیتی، کارایی بیشتری دارند. وقتی $A_1 = Q_1 R_1$ و $A_2 = R_1 Q_1$ به‌دست آمدند، نیاز داریم بدانیم که شکل A_2 همانند شکل A_1 است تا با استفاده از شکل ارزاتر تجزیه QR ، کار را ادامه دهیم. فرض می‌کنیم A_1 به شکل ماتریس هسنبرگ است. با توجه به بخش ۳.۹، Q_1 در تجزیه $A_1 = Q_1 R_1$ دارای مقدار زیر است:

$$Q_1 = H_1 \dots H_{n-1} \quad (۵.۵.۹)$$

که هر یک از H_k ها، یک ماتریس هاؤسهولدر (۱۲.۳.۹) است:

$$H_k = I - 2w^{(k)}w^{(k)T} \quad 1 \leq k \leq n-1 \quad (۶.۵.۹)$$

چون ماتریس A_1 به شکل ماتریس هسنبرگ است، می‌توان نشان داد که بردارهای $w^{(k)}$ شکل خاص زیر را دارند.

$$w_i^{(k)} = 0 \quad i < k \quad \text{و} \quad i > k + 1 \quad \text{برای} \quad (۷.۵.۹)$$

این مطلب را می‌توان از روی معادلات مؤلفه‌های $w^{(k)}$ ، به‌ویژه (۱۰.۳.۹)، ثابت کرد. با توجه به (۷.۵.۹)، ماتریس H_k فقط در چهار عنصر در مواضع (k, k) ، $(k, k+1)$ ، $(k+1, k)$ و $(k+1, k+1)$ با ماتریس یکه، تفاوت دارد. و از این روی، محاسبه نسبتاً ساده نشان می‌دهد که Q_1 باید شکل هسنبرگ داشته باشد. یک لم لازم دیگر این است که ضرب یک ماتریس بالامثلثی و یک ماتریس هسنبرگ باز هم یک ماتریس هسنبرگ است. فقط دو شکل ماتریسها را در هم ضرب و طرز قرار گرفتن صفرهای متناظر را مشاهده کنید تا بتوانید لم را اثبات کنید. از ترکیب این

نتایج، و توجه به این نکته که R_1 یک ماتریس بالامثلثی است، ملاحظه می‌کنیم که $A_2 = R_1 Q_1$ به شکل هسبرگ است.

اگر A_1 متقارن و سه قطری باشد، به طور نمایان دیده می‌شود که هسبرگ است. با توجه به قضیه قبلی، A_2 نیز هسبرگ است. ولی A_2 متقارن است، زیرا

$$A_2^T = (Q_1^T A_1 Q_1)^T = Q_1^T A_1^T Q_1 = Q_1^T A_1 Q_1 = A_2$$

چون هر ماتریس هسبرگ متقارن سه قطری است، نشان داده‌ایم که A_2 سه قطری است. ملاحظه می‌کنیم که بارستهای مثال (۳.۵.۹) این نتیجه را نشان می‌دهند.

همگرایی روش QR قضایای همگرایی روش QR را می‌توان در گلوب و وُن لون (۱۹۸۳)، بخشهای ۵.۷ و ۲.۸، پارلت (۱۹۶۸) و (۱۹۸۰، فصل ۸) و ویلکینسن (۱۹۶۵، فصل ۸) ملاحظه کرد. قضیه زیر از مرجع آخری گرفته شده است.

قضیه ۶.۹ گیریم A یک ماتریس حقیقی از مرتبه n باشد، و ویژه‌مقدارهای $\{\lambda_i\}$ ی آن، در

$$|\lambda_1| > |\lambda_2| > \dots > |\lambda_n| > 0 \quad (۸.۵.۹)$$

صدق کنند. در این صورت بارستهای R_m از روش QR، که در (۱.۵.۹) تعریف شدند، به یک ماتریس بالامثلثی D ، که ویژه‌مقدارهای $\{\lambda_i\}$ عناصر قطر آن هستند، همگرا می‌شوند. اگر A متقارن باشد، دنباله $\{A_m\}$ به یک ماتریس قطری همگرا می‌شود. برای سرعت همگرایی، داریم

$$\|D - A_m\| \leq c \max_i \left| \frac{\lambda_{i+1}}{\lambda_i} \right| \quad (۹.۵.۹)$$

به عنوان یک مثال از این کران خطا، مثال (۳.۵.۹) را در نظر می‌گیریم. نسبت‌های ویژه‌مقدارهای متوالی در آن، عبارت‌اند از

$$\frac{\lambda_2}{\lambda_1} = ۰.۶۳ \quad \frac{\lambda_3}{\lambda_2} = ۰.۴۲ \quad (۱۰.۵.۹)$$

اگر هر یک از عناصر ناقطر A_m را در این مثال امتحان کنیم، ملاحظه می‌کنیم که آنها با یکی از دو نسبت (۱۰.۵.۹)، کاهش می‌یابند.

برای ماتریسهایی که ویژه‌مقدارهای آنها، در (۸.۵.۹) صدق نمی‌کنند، بارستهای A_m ممکن است به یک ماتریس مثلثی همگرا نشوند. اگر A متقارن باشد، دنباله $\{A_m\}$ به یک ماتریس

قطری بلوکی زیر همگرا می‌شود

$$A_m \rightarrow D = \begin{bmatrix} B_1 & & & \\ & B_2 & & \\ & & \ddots & \\ & & & B_r \end{bmatrix} \quad (۱۱.۵.۹)$$

که در آن همه بلوکهای B_i از مرتبه یک یا دو هستند. بنابراین ویژه‌مقدارهای A را، به آسانی از روی ویژه‌مقدارهای D می‌توان محاسبه کرد، اگر A حقیقی و نامتقارن باشد، وضعیت پیچیده‌تر ولی قابل قبول است. برای بحث در این مورد، ویلکینسن (۱۹۶۵، فصل ۸) و پارلت (۱۹۶۸) را ملاحظه کنید. برای اینکه بینیم $\{A_m\}$ همیشه به یک ماتریس قطری همگرا نمی‌شود، مثال متقارن ساده زیر را در نظر می‌گیریم

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

ویژه‌مقدارهای آن، $\lambda = \pm 1$ اند. چون A متعامد است داریم

$$Q_1 = A \quad R_1 = I \quad \text{که} \quad A = Q_1 R_1$$

و بنابراین

$$A_2 = R_1 Q_1 = A$$

و همه بارستها $A_m = A$. دنباله $\{A_m\}$ به یک ماتریس قطری میل نخواهد کرد.

روش QR با انتقال الگوریتم QR معمولاً با یک انتقال مبدأ برای ویژه‌مقدارها به کار برده می‌شود تا سرعت همگرایی افزایش یابد. برای یک دنباله از مقادیر ثابت $\{c_m\}$ ، تعریف می‌کنیم $A_1 = A$ و

$$A_m - c_m I = Q_m R_m \quad (۱۲.۵.۹)$$

$$A_{m+1} = c_m I + R_m Q_m \quad m = 1, 2, \dots$$

ماتریسهای A_m با A_1 متشابه‌اند، زیرا

$$R_m = Q_m^T (A_m - c_m I)$$

$$A_{m+1} = c_m I + Q_m^T (A_m - c_m I) Q_m$$

$$= c_m I + Q_m^T A_m Q_m - c_m I$$

$$A_{m+1} = Q_m^T A_m Q_m \quad m \geq 1 \quad (۱۳.۵.۹)$$

ویژه مقدارهای A_{m+1} همان ویژه مقدارهای A_m اند، و بنابراین با ویژه مقدارهای A یکی هستند. برای اینکه در انتخاب انتقال‌های $\{c_m\}$ دقیقتر باشیم، فقط یک ماتریس متقارن سه قطری A را در نظر می‌گیریم. برای A_m ، گیریم

$$A_m = \begin{bmatrix} \alpha_1^{(m)} & \beta_1^{(m)} & \circ & \dots & \circ \\ \beta_1^{(m)} & \alpha_2^{(m)} & \beta_2^{(m)} & & \vdots \\ \circ & & \ddots & & \\ \vdots & & & & \beta_{n-1}^{(m)} \\ \circ & & \dots & \beta_{n-1}^{(m)} & \alpha_n^{(m)} \end{bmatrix} \quad (۱۴.۵.۹)$$

دوروش برای انتخاب $\{c_m\}$ وجود دارد: (۱) می‌گیریم $c_m = \alpha_n^{(m)}$ و (۲) c_m را آن ویژه مقدار ماتریس

$$\begin{bmatrix} \alpha_{n-1}^{(m)} & \beta_{n-1}^{(m)} \\ \beta_{n-1}^{(m)} & \alpha_n^{(m)} \end{bmatrix} \quad (۱۵.۵.۹)$$

می‌گیریم که به $\alpha_n^{(m)}$ نزدیکتر است. روش دومی بهتر است، ولی در هر دو حالت ماتریسهای A_m به یک ماتریس قطری بلوکی میل می‌کنند که بلوکها، مانند (۱۱.۵.۹)، از مراتب یک یا دو هستند. می‌توان نشان داد که هریک از دو انتخاب $\{c_m\}$ تضمین می‌کند که

$$\beta_{n-1}^{(m)} \beta_{n-2}^{(m)} \rightarrow \circ \quad \text{وقتی } m \rightarrow \infty \quad (۱۶.۵.۹)$$

که معمولاً همگرایی با نرخ سریعتر از روش اولیة (۱۵.۹) QR صورت می‌گیرد. از (۱۳.۵.۹) و با استفاده از نرم ماتریس عملگر (۱۹.۳.۷) و مسألهٔ ۲۷ (ج) فصل ۷، داریم:

$$\|A_{m+1}\|_2 = \|Q_m^T A_m Q_m\|_2 = \|A_m\|_2$$

ماتریسهای $\{A_m\}$ یکنواخت-کراندارند و نتیجتاً همین ویژگی برای هریک از عناصر آنها نیز برقرار است. از (۱۶.۵.۹) و یکنواخت-کراندار بودن $\{\beta_{n-1}^{(m)}\}$ و $\{\beta_{n-2}^{(m)}\}$ ، وقتی $m \rightarrow \infty$ ، یا $\beta_{n-1}^{(m)} \rightarrow \circ$ یا $\beta_{n-2}^{(m)} \rightarrow \circ$. در حالت اول، $\alpha_n^{(m)}$ به یک ویژه مقدار A همگرا می‌شود. و در حالت دوم، دو ویژه مقدار را می‌توان با استفاده از حد زیرماتریس (۱۵.۵.۹) به آسانی به دست آورد. وقتی یک یا دو ویژه مقدار به این علت که $\beta_{n-1}^{(m)}$ یا $\beta_{n-2}^{(m)}$ در واقع صفر است، به دست آمدند، مرتبهٔ ماتریس A_m ممکن است به ترتیب یک یا دو سطر تقلیل یابد. در اینجا، روش QR با انتقال را می‌توان برای ماتریس کاهش یافته به کار گرفت. انتخاب انتقالها به گونه‌ای طراحی شده است که

همگرایی به صفر $\beta_{n-1}^{(m)}\beta_{n-2}^{(m)}$ ، نسبت به بقیه عناصر ناقطر ماتریس، سریعتر انجام پذیرد. بدین طریق، روش QR به یک روش همه منظوره سریع تبدیل می شود که از هر روش دیگر در حال حاضر سریعتر است. برای اثبات همگرایی روش QR با انتقال، ویلکینسن (۱۹۶۸) را ملاحظه کنید. برای بحث بسیار کاملتر روش QR، از جمله انتخاب انتقال، پارلت (۱۹۸۰، فصل ۸) را ببینید.

مثال مثال قبلی (۳.۵.۹) را به کار برده از روش اول انتخاب انتقال $c_m = \alpha_n^{(m)}$ استفاده می کنیم. بارستها چنین اند:

$$A_1 = \begin{bmatrix} 2 & 1 & 0 \\ 1 & 3 & 1 \\ 0 & 1 & 4 \end{bmatrix} \quad A_2 = \begin{bmatrix} 1,4000 & 0,4899 & 0 \\ 0,4899 & 3,2667 & 0,7454 \\ 0 & 0,7454 & 4,3333 \end{bmatrix}$$

$$A_3 = \begin{bmatrix} 1,2915 & 0,2017 & 0 \\ 0,2017 & 3,0202 & 0,2724 \\ 0 & 0,2724 & 4,6884 \end{bmatrix} \quad A_4 = \begin{bmatrix} 1,2737 & 0,0993 & 0 \\ 0,0993 & 2,9943 & 0,0072 \\ 0 & 0,0072 & 4,7320 \end{bmatrix}$$

$$A_5 = \begin{bmatrix} 1,2694 & 0,0498 & 0 \\ 0,0498 & 2,9986 & 0 \\ 0 & 0 & 4,7321 \end{bmatrix}$$

عنصر $\beta_2^{(m)}$ خیلی سریع به صفر همگرا می شود، ولی عنصر $\beta_1^{(m)}$ به طور هندسی و فقط با نسبتی در حدود ۵٪ به صفر می گراید.

باید به روش پیش از روش QR که خیلی زیاد انگیزه پیدایش آن بوده است اشاره ای کنیم. در ۱۹۵۸، روتیس هاوژر^۱ یک روش LR بر اساس تجزیه حذف گاوسی ماتریس به یک ماتریس پایین مثلثی ضرب در یک ماتریس بالامثلثی، معرفی کرد. طبق تعریف

$$A_m = L_m R_m \quad A_{m+1} = R_m L_m = L_m^{-1} A_m L_m$$

که L_m یک ماتریس پایین مثلثی و R_m یک ماتریس بالامثلثی است. این روش معمولاً از روش QR کاراتر است، اگر بتوانیم آن را به کار ببریم. ولی تبدیلات تشابهی غیرمتعامد می توانند یک انحراف در وضع ویژه مقدرهای ماتریسها ایجاد نمایند. و معمولاً پیاده کردن آن در برنامه نویسی خودکار مشکلتر است. یک بحث کامل از آن در ویلکینسن (۱۹۶۵، فصل ۸) داده شده است.

۶.۹ محاسبه ویژه بردارها و بارست معکوس

تواناترین ابزار در به دست آوردن ویژه بردارهای یک ماتریس روش بارست معکوس، یعنی روشی است که توسط ه. ویلنت^۱ (۱۹۴۴) ابداع شده است. ابتدا، بارست معکوس را تعریف کرده توضیح می دهیم، و سپس در محاسبه ویژه بردارها نظر کلیتری خواهیم داد. برای ساده کردن تحلیل، گیریم A ماتریسی باشد که شکل متعارف ژوردان آن قطری است.

$$P^{-1}AP = \text{diag}[\lambda_1, \dots, \lambda_n] \quad (۱.۶.۹)$$

گیریم ستونهای P با x_1, \dots, x_n نشان داده شده اند. پس

$$Ax_i = \lambda_i x_i \quad i = 1, \dots, n \quad (۲.۶.۹)$$

بی آنکه از کلیت کاسته شود، می توان فرض کرد که به ازای جميع مقادیر i ، $\|x_i\|_\infty = 1$. گیریم λ تقریبی برای یک ویژه مقدار ساده λ_k ی A باشد. اگر یک مقدار آغازی $z^{(0)}$ داده شده باشد، $\{w^{(m)}\}$ و $\{z^{(m)}\}$ را با روابط زیر تعریف می کنیم

$$(A - \lambda I)w^{(m+1)} = z^{(m)}, \quad z^{(m+1)} = \frac{w^{(m+1)}}{\|w^{(m+1)}\|_\infty} \quad m \geq 0 \quad (۳.۶.۹)$$

در اصل این روش، روش توانی است، که در (۲.۲.۹) - (۳.۲.۹) به جای A ، $(A - \lambda I)^{-1}$ آمده است. از دیدگاه مطالب بخش ۴.۸ فصل ۸، ماتریس $A - \lambda I$ بد - وضع است. ولی هرگونه اختلال بزرگ حاصله در جواب، بر ویژه بردار x_k از ویژه مقدار $\lambda_k - \lambda$ در $A - \lambda I$ اثر زیاد می گذارد، و این همان برداری است که ما دنبال آن هستیم. برای بحث بیشتر در این منشأ ناپایداری در حل دستگاه خطی، مطالبی را که به دنبال (۸.۴.۸) در بخش ۴.۸ آمده است ملاحظه نمایید. برای آنکه روش (۳.۶.۹) نتیجه بخش باشد، ما نمی خواهیم $A - \lambda I$ تکین باشد. بنابراین λ نباید دقیقاً λ_k باشد، اگرچه می تواند کاملاً به آن نزدیک باشد، همان گونه که مثال بعدی نشان خواهد داد. برای یک تحلیل دقیقتر، گیریم $z^{(0)}$ بر حسب پایه های ویژه بردار (۲.۶.۹) بیان شده باشد:

$$z^{(0)} = \sum_{i=1}^n \alpha_i x_i \quad (۴.۶.۹)$$

و فرض می کنیم $\alpha_k \neq 0$. در قیاس با فرمول (۴.۲.۹) برای روش توانی، می توانیم نشان دهیم

$$z^{(m)} = \frac{\sigma_m (A - \lambda I)^{-m} z^{(0)}}{\|(A - \lambda I)^{-m} z^{(0)}\|_\infty} \quad |\sigma_m| = 1 \quad (۵.۶.۹)$$

با استفاده از (۴.۶.۹)،

$$(A - \lambda I)^{-m} z^{(0)} = \sum_{i=1}^n \alpha_i \left[\frac{1}{\lambda_i - \lambda} \right]^m x_i \quad (۶.۶.۹)$$

گیریم $\epsilon = \lambda_k - \lambda$ و فرض می‌کنیم

$$|\lambda_i - \lambda| \geq c > 0 \quad i = 1, \dots, n \quad i \neq k \quad (۷.۶.۹)$$

با استفاده از (۶.۶.۹) و (۵.۶.۹)،

$$z^{(m)} = \sigma_m \frac{x_k + \epsilon^m \sum_{i \neq k} \frac{\alpha_i}{\alpha_k} \left[\frac{1}{\lambda_i - \lambda} \right]^m x_i}{\left\| x_k + \epsilon^m \sum_{i \neq k} \frac{\alpha_i}{\alpha_k} \left[\frac{1}{\lambda_i - \lambda} \right]^m x_i \right\|_{\infty}} \quad (۸.۶.۹)$$

که $|\sigma_m| = 1$ اگر $|\epsilon| < c$ ، آنگاه

$$\left\| \epsilon^m \sum_{i \neq k} \frac{\alpha_i}{\alpha_k} \left[\frac{1}{\lambda_i - \lambda} \right]^m x_i \right\|_{\infty} \leq \left[\frac{\epsilon}{c} \right]^m \sum_{i \neq k} \left| \frac{\alpha_i}{\alpha_k} \right| \quad (۹.۶.۹)$$

این کمیت وقتی $m \rightarrow \infty$ ، به صفر میل می‌کند. ترکیب آن با (۸.۶.۹) نشان می‌دهد که وقتی $m \rightarrow \infty$ ، $z^{(m)}$ به مضربی از $x_{(k)}$ میل می‌کند. این همگرایی خطی است، با یک نسبت کاهش خطای $|\epsilon/c|$ در هر بارست. در عمل، $|\epsilon|$ کاملاً کوچک است که همگرایی سریع را تضمین می‌کند. در پیاده‌کردن (۳.۶.۹)، از تجزیه $A - \lambda I$ با استفاده از تجزیه LU در بخش ۱.۸ از فصل ۸،

آغاز می‌کنیم. برای ساده‌کردن نمادگذاری، می‌نویسیم

$$A - \lambda I = LU$$

که در آن از محورگیری استفاده نشده است. در عمل، محورگیری به‌کار خواهد رفت. برای هر بارست $z^{(m+1)}$ مسأله را به شکل زیر حل می‌کنیم:

$$Ly^{(m+1)} = z^{(m)} \quad Uw^{(m+1)} = y^{(m+1)}$$

$$z^{(m+1)} = \frac{w^{(m+1)}}{\|w^{(m+1)}\|_{\infty}} \quad (۱۰.۶.۹)$$

چون $A - \lambda I$ تقریباً تکین است، آخرین عنصر قطری U تقریباً صفر خواهد شد. اگر دقیقاً صفر باشد، آن را به یک عدد کوچک تبدیل می‌کنیم یا λ را اندکی تغییر داده L و U را مجدداً به دست می‌آوریم. برای حدس آغازی $z^{(0)}$ ، ویلکینسن (۱۹۶۳، ص ۱۴۷) انتخاب زیر را توصیه می‌کند:

$$z^{(0)} = Le \quad e = [1, 1, \dots, 1]^T$$

پس در (۱۰.۶.۹)

$$y^{(1)} = e \quad Uw^{(m+1)} = e \quad (11.6.9)$$

این انتخاب بدین منظور صورت گرفته است که تضمین کند α_k در (۴.۶.۹)، نه صفر است و نه کوچک. ولی اگر هم کوچک باشد، این روش معمولاً با سرعت همگراست. برای مثال، فرض می‌کنیم، بعضی یا همه مقادیر α_i/α_k در (۹.۶.۹) در حدود 10^4 باشند. و فرض می‌کنیم $|\epsilon/c| = 10^{-5}$ ، که یک مقدار واقع‌بینانه در خیلی از حالات است. در این صورت کران (۹.۶.۹) چنین می‌شود

$$(10^{-5})^m \cdot n \cdot 10^4$$

و این کمیت، وقتی m بزرگ می‌شود، خیلی به سرعت کاهش می‌یابد.

مثال ماتریس قبلی (۳.۵.۹)

$$A = \begin{bmatrix} 2 & 1 & 0 \\ 1 & 3 & 1 \\ 0 & 1 & 4 \end{bmatrix} \quad (12.6.9)$$

را در نظر می‌گیریم. گیریم $\sqrt{3} - 3 = \lambda_2 = 1.2679 = \lambda$ ، که تا ۵ رقم دقیق است. لذا، به علت خطاهای گرد کردن داریم:

$$L = \begin{bmatrix} 1.0 & 0 & 0 \\ 1.3659 & 1.0 & 0 \\ 0 & 2.7310 & 1.0 \end{bmatrix} \quad U = \begin{bmatrix} 0.7321 & 1.0 & 0 \\ 0 & 0.3662 & 1.0 \\ 0 & 0 & 0.0011 \end{bmatrix}$$

با استفاده از $y^{(1)} = [1, 1, 1]^T$ داریم

$$\begin{aligned} w^{(1)} &= [3385, 2, -2477, 3, 908, 20]^T \\ z^{(1)} &= [1, 0, 0, 0, -0, 73180, 0, 26828]^T \\ w^{(2)} &= [20345, -14894, 5451, 9]^T \\ z^{(2)} &= [1, 0, 0, 0, -0, 73207, 0, 26797]^T \end{aligned} \quad (13.6.9)$$

و بردار $z^{(3)} = z^{(2)}$. جواب درست چنین است

$$\begin{aligned} x_2 &= [1, 1 - \sqrt{3}, 2 - \sqrt{3}]^T \\ &\doteq [1, 0, 0, 0, 0, -0, 73205, 0, 26795]^T \end{aligned} \quad (14.6.9)$$

و $z^{(2)}$ در حدود انباشتگیهای خطای گردکردن، با x_2 برابر است.

ویژه بردارهای ماتریسهای سه قطری متقارن گیریم A یک ماتریس سه قطری متقارن حقیقی از مرتبه n باشد. مانند قبل، فرض می‌کنیم بعضی یا همه ویژه مقدارهای آن با دقت محاسبه شده‌اند. روش بارست معکوس روش ارجح در محاسبه ویژه بردارهاست، و انجام آن خیلی آسان است. برای λ به عنوان یک ویژه مقدار تقریبی A ، محاسبه تجزیه LU ، حتی وقتی که محورگیری به کار رفته باشد، از نظر زمان و از نظر حافظه گران نخواهد بود. برای مثال، مطالب بخش ۳.۸ از فصل ۸ را در باب دستگاههای سه قطری ملاحظه کنید. مثال عددی قبلی نیز این روش را برای ماتریسهای سه قطری نشان می‌دهد. بعضی نتایج درباره خطا در توجیه بیشتر برای استفاده از بارست معکوس را می‌آوریم. فرض می‌کنیم که حساب رایانه‌ای ممیز شناور دودویی با گردکردن و با t رقم در قسمت اعشاری باشد. در ویلکینسن (۶۳، صص ۱۴۳۱۹-۱۴۷) نشان داده شده است که جواب محاسبه شده \hat{w} ی

$$(A - \lambda I)w^{(m+1)} = z^{(m)}$$

همان جواب دقیق

$$(A - \lambda I + E)\hat{w} = z^{(m)} \quad (15.6.9)$$

است که در آن

$$\|E\|_2 \leq K\sqrt{n} \cdot 2^{-t} \quad (16.6.9)$$

که K مقداری ثابت از مرتبه واحد است. اندازه این کران به همان اندازه‌ای است که از خطاهایی از مرتبه خطای گردکردن انتظار می‌رود.

اگر جواب \hat{w} ی (۱۵.۶.۹) کاملاً بزرگ باشد، این کران یک تقریب خوب برای ویژه بردار A خواهد بود. برای اثبات این امر، مطلب را با معرفی نماد زیر آغاز می‌کنیم

$$\hat{z} = \frac{\hat{w}}{\|\hat{w}\|_2}$$

پس

$$(A - \lambda I + E)\hat{z} = \frac{z^{(m)}}{\|\hat{w}\|_2}$$

$$\eta \equiv (A - \lambda I)\hat{z} = -E\hat{z} + \frac{\hat{z}^m}{\|\hat{w}\|_2}$$

برای این مانده η ,

$$\|\eta\|_2 \leq \|E\|_2 + \frac{\sqrt{n}}{\|\hat{w}\|_2}$$

$$\leq K\sqrt{n} \cdot 2^{-t} + \frac{\sqrt{n}}{\|\hat{w}\|_2} \quad (17.6.9)$$

که اگر $\|\hat{w}\|_2$ بزرگ باشد، کمیتی کوچک است. برای اثبات اینکه (۱۷.۶.۹) ایجاب می‌کند که \hat{z} به یک ویژه بردار A نزدیک باشد، $\{x_i \mid i = 1, \dots, n\}$ را یک مجموعه از ویژه بردارهای یکا متعامد می‌گیریم. و فرض می‌کنیم

$$\lambda \doteq \lambda_k \quad \hat{z} = \sum_{i=1}^n \alpha_i x_i \quad \|\hat{z}\|_2^2 = \sum_{i=1}^n \alpha_i^2 = 1$$

که λ_k یک ویژه مقدار تنهای A است. همچنین، فرض می‌کنیم

$$|\lambda_i - \lambda_k| \geq c > 0 \quad \text{برای هر } i \neq k$$

با

$$c \gg |\lambda_k - \lambda|$$

با این مفروضات، اکنون می‌توانیم یک کران برای خطا در \hat{z} پیدا کنیم. ویژه بردارها را به عنوان

پایه به کار برده η را بر حسب آنها بیان می‌کنیم:

$$\begin{aligned} \eta &= A\hat{z} - \lambda\hat{z} = \sum_1^n \alpha_i(\lambda_i - \lambda)x_i \\ \|\eta\|_2^2 &= \sum_1^n \alpha_i^2(\lambda_i - \lambda)^2 \\ &\geq \sum_{i \neq k} \alpha_i^2(\lambda_i - \lambda)^2 \geq c^2 \sum_{i \neq k} \alpha_i^2 \\ \sum_{i \neq k} \alpha_i^2 &\leq \frac{1}{c^2} \|\eta\|_2^2 \end{aligned}$$

که با استفاده از (۱۷.۶.۹)، کاملاً کوچک است. از $\|\hat{z}\| = 1$ نتیجه می‌شود که $\alpha_k = 1$ و

$$\|\hat{z} - \alpha_k x_k\|_2 = \sqrt{\sum_{i \neq k} \alpha_i^2} \leq \frac{1}{c} \|\eta\|_2 \quad (18.6.9)$$

که نتیجه مطلوب را نشان می‌دهد. برای بحث بیشتر در مورد خطا، ویلکینسن (۱۹۶۳، صص ۱۴۲ تا ۱۴۶) و (۱۹۶۵، صص ۳۲۱ تا ۳۳۰) را ببینید.

روش دیگر محاسبه ویژه بردارها، ظاهراً حل مستقیم

$$(A - \lambda I)x = 0$$

خواهد بود، پس از حذف یک معادله و گرفتن یکی از مؤلفه‌های مجهول برابر یک ثابت ناصفر، مثلاً $x_1 = 1$. این روندی است که اغلب در کتابهای جبرخطی دوره‌های کارشناسی به کار می‌رود. ولی به‌عنوان یک روش عددی کلی، ممکن است ناموفق باشد. یک بحث کامل از این مسأله در ویلکینسن (۱۹۶۵، صص ۳۱۵-۳۲۱)، همراه با یک مثال عالی، داده شده است. ما فقط همان مثال قبلی را به کار می‌بریم که نشان دهیم این جوابها ممکن است به خوبی جوابهای روش بارست معکوس نباشند.

مثال مثال قبلی (۱۲.۶.۹) را با $\lambda = 1.2679$ در نظر می‌گیریم. معادله $(A - \lambda I)x = 0$ را انتخاب و آخرین معادله را حذف می‌کنیم تا چنین به دست آید:

$$0.7321x_1 + x_2 = 0$$

$$x_1 + 1.7321x_2 + x_3 = 0$$

با گرفتن $x_1 = 1$ ، ویژه بردار تقریبی زیر را به دست می آوریم

$$x = [1.00000, -0.73210, 0.26807]^T$$

که در مقایسه با جواب درست (۱۴.۶.۹)، این یک جواب کمی ضعیفتر از جواب (۱۳.۶.۹) است که از بارست معکوس حاصل شده است. به طور کلی، نتایج به دست آمده از این راه ممکن است خیلی ضعیف باشند، و هنگام کاربرد این روش باید خیلی مواظب بود.

روش بارست معکوس در اجرا نیاز به توجه بسیار زیادی دارد. در پرداختن به یک ماتریس خاص، هرگونه مشکلات پیش‌بینی نشده‌ای ممکن است رخ دهد. ولی در یک برنامه رایانه‌ای عمومی زیادی است، با ویژه‌مقدارهایی سروکار خواهیم داشت که یا چندگانه‌اند یا بهم نزدیک‌اند، که اگر به آنها توجه نشود، ممکن است مشکلاتی به بار آورند. برای ماتریسهای نامتقارن، که شکل متعارف ژوردان آنها قطری نیست، مشکلاتی هم درگزینش یک پایه صحیح از ویژه‌مقدارها، وجود دارند. بهترین مرجع برای این موضوع، ویلکینسن (۱۹۶۵) است و همچنین می‌توان گلوب و ون لون (۱۹۸۳، صص ۲۳۸-۲۴۰) و پارلت (۱۹۸۰، ۶۲-۶۹) را ملاحظه کرد. برای چندین برنامه عالی، ویلکینسن و راینش (۱۹۷۱، صص ۴۱۸-۴۳۹) و گاربو^۱ و همکاران (۱۹۷۷) را ببینید.

۷.۹ حل دستگاههای خطی به روش کمترین مربعات

اکنون حل دستگاههای خطی فرا معین به شکل

$$\sum_{j=1}^n a_{ij}x_j = b_i \quad i = 1, \dots, m \quad (1.7.9)$$

را مطالعه می‌کنیم، که در آنها $m > n$. این دستگاهها در کاربردهای متنوعی ظاهر می‌شوند، که معروفترین آنها برازش توابع با یک مجموعه از داده‌های $\{(t_i, b_i) \mid i = 1, \dots, m\}$ است، که درباره آنها بعداً بیشتر صحبت خواهیم کرد. شاید به نظر آید که منطقیترین جا برای بررسی این دستگاهها، فصل هشتم باشد، ولی بعضی ابزارهای حل (۱.۷.۹) شامل تبدیلات متعامد هستند که در این فصل بررسی شده‌اند. حل عددی (۱.۷.۹) هم از لحاظ نظری و هم از جنبه عملی ممکن است کاملاً پیچیده باشد ولی ما فقط بعضی از قسمتهای مهم آن را بیان می‌کنیم.

حل دستگاههای خطی به روش کمترین مربعات ۷۲۱

یک دستگاه فرامعین (۱.۷.۹) در حالت کلی، جوابی ندارد. بدین دلیل یک بردار $x = (x_1, \dots, x_n)$ را جستجو می‌کنیم که به تعبیری جواب تقریبی (۱.۷.۹) باشد. نمادهای

$$A = [a_{ij}] \quad x = [x_1, \dots, x_n]^T \quad b = [b_1, \dots, b_m]^T$$

را که در آن A یک ماتریس $m \times n$ است وارد می‌کنیم. در این صورت معادله (۱.۷.۹) را می‌توانیم به صورت زیر بنویسیم

$$Ax = b \quad (۲.۷.۹)$$

برای سادگی، A و b حقیقی فرض می‌کنیم. در بین راههای ممکن پیدا کردن یک جواب تقریبی، می‌توانیم بردار x را جویا شویم که

$$\|Ax - b\|_p \quad (۳.۷.۹)$$

را مینیمم سازد، که در آن $1 \leq p \leq \infty$. در این بخش، فقط حالت کلاسیک $p = 2$ را بررسی می‌کنیم، گرچه در سالهای اخیر، کارهای زیادی برای حالت‌های $p = 1$ و $p = \infty$ شده است.

جواب x^* ی

$$\text{Minimize } \|Ax - b\|_2 \quad (۴.۷.۹)$$

$x \in \mathbb{R}^n$

را جواب کمترین مربعات دستگاه خطی $Ax = b$ می‌نامند. دلایل متعددی برای حل دستگاه $Ax = b$ بدین طریق وجود دارد. اولاً، تئوری و روشهای عملی مینیمم‌سازی $\|Ax - b\|_2$ آسانتر است، تا حدی به این علت که $\|Ax - b\|_2$ یک تابع پیوسته مشتق‌پذیر بر حسب x_1, \dots, x_n است و ثانیاً، مسائل برازش منحنی که به دستگاههای (۱.۷.۹) می‌انجامند، اغلب یک چارچوب آماری دارند که به (۴.۷.۹) منجر می‌شوند، که بر مینیمم‌سازی $\|Ax - b\|_p$ ، $p \neq 2$ ، رجحان دارند.

برای درک بهتر ماهیت حل (۴.۷.۹)، ساختمان نظری زیر را عرضه می‌کنیم. از این ساختمان نیز می‌توان به عنوان یک روش عددی عملی نیز استفاده کرد، اگرچه معمولاً روشهای کارا تر دیگری

وجود دارند. نکته مهم در این نظریه، تجزیه تکین - مقدار زیر است،

$$V^T A U = F = \begin{bmatrix} \mu_1 & & \dots & & & & 0 \\ & \circ & & \ddots & & & \\ & & & & \mu_r & & \\ & & & & & \circ & \\ \vdots & & & & & & \ddots \\ & \circ & & \dots & & & 0 \end{bmatrix} \quad (5.7.9)$$

ماتریسهای U و V متعامدند، و مقدارهای تکین μ_i در روابط زیر صدق می‌کنند

$$\mu_1 \geq \mu_2 \geq \dots \geq \mu_r > 0$$

برای آگاهی بیشتر، قضیه ۵.۷ در فصل ۷ را ببینید، بعداً در این بخش روش ساختن تجزیه تکین - مقدار A را ذکر می‌کنیم.

قضیه ۷.۹ گیریم A ماتریس حقیقی $m \times n$ است، $m \geq n$. تعریف می‌کنیم $z = U^T x$ $c = V^T b$. در این صورت $x^* = U z^*$ جواب (۴.۷.۹) با رابطه

$$z_i^* = \frac{c_i}{\mu_i} \quad i = 1, \dots, r \quad (6.7.9)$$

داده می‌شود که در آن z_n, \dots, z_{r+1} دلخواه‌اند. وقتی $r = n$ x^* یکتاست. وقتی $r < n$ جواب نرم اقلیدسی مینیمال (۴.۷.۹) با قراردادن

$$z_i^* = 0 \quad i = r + 1, \dots, n \quad (7.7.9)$$

به دست می‌آید. [این نیز جواب کمترین مربعات (۴.۷.۹) خوانده می‌شود، ولو اینکه این یک مینیم‌ساز یکتا برای $\|Ax - b\|_2$ نیست.] مینیمم (۴.۷.۹) با رابطه زیر داده می‌شود

$$\|Ax^* - b\|_2 = \left[\sum_{j=r+1}^m c_j^2 \right]^{1/2} \quad (8.7.9)$$

برهان مسأله ۱۳ (الف) از فصل ۷ را به یاد می‌آوریم. برای هر $x \in \mathbf{R}^n$ و هر ماتریس متعامد P

$$\|Px\|_2 = \|x\|_2$$

ابزار مهمی در مطالعه حل عددی $Ax = b$ است. بعضی از ویژگیهای دیگر A^+ در مسائل ۲۷ و ۲۸ مطرح شده‌اند.

برای ساده‌کردن مطالعه بقیه روشها در پیدا کردن x^* و تحلیل پایداری آن، A را به حالت رده کامل $r = n$ محدود می‌سازیم. این حالت مهمترین حالت در کاربردهاست. یادآوری می‌کنیم: برای تکین مقدارهای A ,

$$\mu_1 \geq \mu_2 \geq \dots \geq \mu_n > 0 \quad (۱۳.۷.۹)$$

مفهوم نرم ماتریسی را، که در بخش ۳.۷ برای ماتریسهای مربعی داده شده بود، می‌توان برای A تعمیم داد. تعریف می‌کنیم

$$\|A\| = \text{Supremum}_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{\|Ax\|_2}{\|x\|_2} \quad (۱۴.۷.۹)$$

با استفاده از تجزیه تکین مقدار ماتریس A ، می‌توان نشان داد که،

$$\|A\| = \sqrt{r_\sigma(A^T A)} = \mu_1 \quad (۱۵.۷.۹)$$

در قیاس با تحلیل خطا در بخش ۴.۸، ضریب وضعیت برای $Ax = b$ را با رابطه زیر تعریف می‌کنیم

$$\text{cond}(A)_2 = \|A\| \|A^+\| = \frac{\mu_1}{\mu_n} \quad (۱۶.۷.۹)$$

با استفاده از این نماد، یک قضیه پایداری از گلوب و وِن لون (۱۹۸۳، ص ۱۴۱) را در اینجا می‌آوریم. این قضیه مشابه قضیه ۴.۸، برای تحلیل اختلال در دستگاههای خطی ناتکین است. گیریم $b + \delta b$ و $A + \delta A$ به ترتیب شکلهای اختلال یافته b و A باشند. تعریف می‌کنیم

$$\begin{aligned} x^* &= A^+ b & \hat{x}^* &= (A + \delta A)^+(b + \delta b) \\ r &= b - Ax^* & \hat{r} &= (b + \delta b) - (A + \delta A)\hat{x}^* \end{aligned} \quad (۱۷.۷.۹)$$

فرض می‌کنیم

$$\epsilon \equiv \text{Max} \left[\frac{\|\delta A\|}{\|A\|}, \frac{\|\delta b\|_2}{\|b\|_2} \right] < \frac{1}{\text{cond}(A)_2} \quad (۱۸.۷.۹)$$

و

$$\sin(\theta) \equiv \frac{\|r\|_2}{\|b\|_2} < 1 \quad (۱۹.۷.۹)$$

که به طور ضمنی $0 \leq \theta < \pi/2$ را تعریف می‌کند. در این صورت

$$\frac{\|\hat{x}^* - x^*\|_2}{\|x^*\|_2} \leq \epsilon \left[\frac{2 \operatorname{cond}(A)_2}{\cos \theta} + \tan \theta [\operatorname{cond}(A)_2]^2 \right] + O(\epsilon^2) \quad (20.7.9)$$

$$\frac{\|\hat{r} - r\|_2}{\|b\|_2} \leq \epsilon [1 + 2 \operatorname{cond}(A)_2] \operatorname{Min}\{1, m - n\} + O(\epsilon^2) \quad (21.7.9)$$

برای حالت $m = n$ مرتبه (A) است، مانده r صفر می‌شود و بنابراین (۲۰.۷.۸) به قضیه قبلی ۴.۸ بَدَل می‌شود.

نتایج قبلی بیان می‌کنند که تغییر در r ممکن است کاملاً کوچک و تغییر در x^* کاملاً بزرگ باشد. توجه کنید که کران در (۲۰.۷.۹)، در مقایسه با وابستگی خطی $\operatorname{cond}(A)$ در حالت ناتکین $m = n$ به مربع $\operatorname{cond}(A)_2$ بستگی دارد [(۱۸.۴.۸) را ببینید]. اگر ستونهای A تقریباً وابسته باشند، آنگاه $\operatorname{cond}(A)_2$ ممکن است خیلی بزرگ باشد، که کران بزرگتری را در (۲۰.۷.۹) در مقایسه با (۲۱.۷.۹)، ایجاد خواهد کرد [مسئله ۳۴ (الف) را ببینید]. قابل قبول بودن یا نبودن، به مسئله بستگی دارد، به اینکه آیا مقادیر کوچک r را می‌خواهند یا مقادیر دقیق x^* را.

مسئله برازش منحنی به روش کمترین مربعات منشأ بسیاری از دستگاههای خطی فرامعین مسئله برازاندن داده‌ها با یک تابع از یک خانواده از توابع معین است. گیریم $\{(t_i, b_i) \mid i = 1, \dots, m\}$ یک مجموعه از داده‌ها باشد، احتمالاً معرف یک تابع $b = g(t)$. گیریم $\varphi_1(t), \dots, \varphi_n(t)$ توابع داده شده باشند، و گیریم \mathcal{F} خانواده تمام ترکیبات خطی $\varphi_1, \dots, \varphi_n$ باشد:

$$\mathcal{F} = \left\{ \sum_{j=1}^n x_j \varphi_j(t) \mid x_j \in \mathbf{R} \right\} \quad (22.7.9)$$

می‌خواهیم عنصری از \mathcal{F} بیابیم به گونه‌ای که تقریباً در

$$\sum_{j=1}^n x_j \varphi_j(t_i) = b_i \quad i = 1, \dots, m \quad (23.7.9)$$

صدق کند. این همان دستگاه (۱.۷.۹)، با $a_{ij} = \varphi_j(t_i)$ است.

به دلایل مدل‌سازی آماری، در پی آن هستیم که

$$E(x) = \left[\frac{1}{m} \sum_{i=1}^m \left[b_i - \sum_{j=1}^n x_j \varphi_j(t_i) \right]^2 \right]^{1/2} \quad (24.7.9)$$

را مینیمم سازیم و بنابراین می‌خواهیم برازش داده‌ها را به مفهوم کمترین مربعات بیان کنیم. کمیت $E(x^*)$ ، که برای آن $E(x)$ مینیمم می‌شود، میانگین مربع خطاها، در تقریب داده‌ها با تابع

$$g^*(t) = \sum_{j=1}^n x_j^* \varphi_j(t) \quad (25.7.9)$$

نامیده می‌شود. با استفاده از نماد قبلی

$$E(x) = \frac{1}{\sqrt{m}} \|b - Ax\|_2$$

و مینیمم کردن $E(x)$ هم‌ارز با پیدا کردن جواب کمترین مربعات (۲۳.۷.۹) است. از گرفتن مشتقات جزئی (۲۴.۷.۹) نسبت به هر x_i ، و برابر صفر قرار دادن آنها، دستگاه معادلات زیر را به دست می‌آوریم

$$A^T A x = A^T b \quad (26.7.9)$$

این دستگاه شرط لازم برای هر مینیمم‌ساز $E(x)$ است، و می‌توان نشان داد که شرط کافی هم هست. دستگاه (۲۶.۷.۹) معادلهٔ نرمال برای مسألهٔ کمترین مربعات خوانده می‌شود. اگر A دارای رتبهٔ n باشد، آنگاه $A^T A$ یک ماتریس $n \times n$ و ناتکین است و (۲۶.۷.۹) دارای یک جواب یکتاست.

برای اثبات هم‌ارزی (۲۶.۷.۹) با جواب قبلی مسألهٔ کمترین مربعات، تجزیهٔ تکین - مقدار A را به‌کار می‌بریم تا (۲۶.۷.۹) را به یک شکل ساده‌تر برگردانیم. از گذاردن $A = V F U^T$ در (۲۶.۷.۹)

$$U F^T F U^T x = U F^T V^T b$$

از ضرب در U و استفاده از نماد قبلی $z = U^T x$ ، $c = V^T b$ ، داریم

$$F^T F z = F^T c$$

این رابطه یک هم‌ارزی ریاضی کامل بین معادلهٔ نرمال و مینیمم‌سازی قبلی $\|Ax - b\|_2$ در قضیهٔ ۷.۹ برقرار می‌سازد.

فرض می‌کنیم رتبهٔ A برابر n باشد، جواب x^* را می‌توان از حل معادلهٔ نرمال به دست آورد. چون $A^T A$ متقارن و معین مثبت است، تجزیهٔ چولسکی را برای پیدا کردن جواب می‌توان مورد استفاده قرار داد [(۸.۳.۸) - (۱۷.۳.۸)] را ببینید. اثر خطاهای گرد کردن روی x^* هم با واحد

خطای گرد کردن در رایانه و هم با ضریب وضعیت $A^T A$ ، متناسب است. از تجزیهٔ تکین - مقدار مربوط به A به سادگی می‌توان دید که

$$\text{cond}(A^T A)_r = \frac{\mu_1^2}{\mu_n^2} = [\text{cond}(A)_r]^2 \quad (27.7.9)$$

بنابراین حساسیت x^* نسبت به خطاها متناسب است با $[\text{cond}(A)_r]^2$ ، که با کران خطای اختلال (۲۰.۷.۹) سازگار است.

نتیجهٔ (۲۷.۷.۹) معمولاً به‌عنوان دلیل اصلی احتراز از استفاده از معادلهٔ نرمال در حل مسألهٔ کمترین مربعات ذکر شده است. این گوشزد خوبی است، ولی دلایل خیلی ظریفتر از اینها هستند. از (۲۰.۷.۹) دیده می‌شود که اگر $\|r\|_2 \ll r$ نزدیک صفر باشد، آنگاه $\sin \theta \cong 0$ ، و کران متناسب با $\text{cond}(A)_r$ می‌شود. در مقابل، کران خطا در روش چولسکی کمیت $[\text{cond}(A)_r]^2$ را نشان می‌دهد که اگر $\text{cond}(A)_r$ بزرگ باشد، این کمیت بزرگتر می‌شود. یک دلیل دیگر وقتی به‌وجود می‌آید که ستونهای A تقریباً وابسته باشند. در این صورت استفاده از حساب رایانه‌ی محدود ممکن است به یک معادلهٔ نرمال تقریبی منجر شود که اطلاعات ضروری موجود در A را ندارد. در چنین حالتی، $A^T A$ تقریباً تکین است، و جواب معادلهٔ نرمال دقت بسیار کمتری در x^* را در مقایسه با سایر روشها، که مستقیماً با $Ax = b$ کار می‌کنند، به‌دست خواهد داد. برای یک بحث دقیقتر از این موضوع لاونس و هسنس^۱ (۱۹۷۴، صص ۱۲۶-۱۲۹) را ببینید.

مثال داده‌های جدول ۳.۹ و نمودار مربوط به آن را در شکل (۲.۹) در نظر می‌گیریم. یک چندجمله‌ی درجهٔ ۳ برای برازش این داده‌ها به‌کار می‌بریم، و به مینیمم کردن عبارت زیر می‌رسیم

$$E(x) = \left[\frac{1}{m} \sum_{i=1}^m \left[b_i - \sum_{j=1}^4 x_j t_i^{j-1} \right]^2 \right]^{1/2}$$

از اینجا، دستگاه خطی فرامعین

$$\sum_{j=1}^4 x_j t_i^{j-1} = b_i \quad i = 1, \dots, m \quad (28.7.9)$$

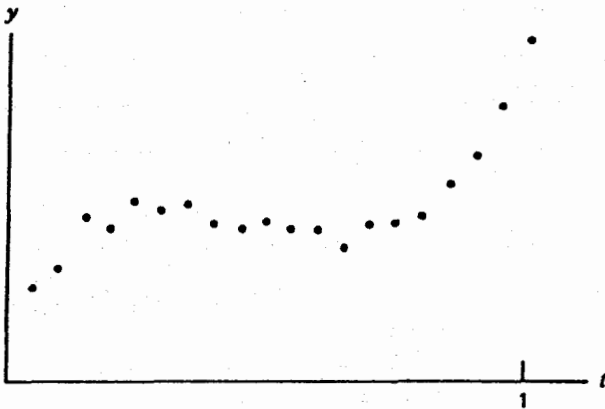
نتیجه می‌شود.

و معادلات نرمال عبارت‌اند از:

$$\sum_{j=1}^n x_j \left[\sum_{i=1}^m t_i^{j+k-2} \right] = \sum_{i=1}^m t_i^{k-1} b_i \quad k = 1, 2, 3, 4$$

جدول ۳.۹ داده‌ها برای یک برازش کمترین مربعات درجه ۳

x_i	y_i	x_i	y_i
۰٫۰۰	۰٫۴۸۶	۰٫۵۵	۱٫۱۰۲
۰٫۰۵	۰٫۸۶۶	۰٫۶۰	۱٫۰۹۹
۰٫۱۰	۰٫۹۴۴	۰٫۶۵	۱٫۰۱۷
۰٫۱۵	۱٫۱۴۴	۰٫۷۰	۱٫۱۱۱
۰٫۲۰	۱٫۱۰۳	۰٫۷۵	۱٫۱۱۷
۰٫۲۵	۱٫۲۰۲	۰٫۸۰	۱٫۱۵۲
۰٫۳۰	۱٫۱۶۶	۰٫۸۵	۱٫۲۶۵
۰٫۳۵	۱٫۱۹۱	۰٫۹۰	۱٫۳۸۰
۰٫۴۰	۱٫۱۲۴	۰٫۹۵	۱٫۵۷۵
۰٫۴۵	۱٫۰۹۵	۱٫۰۰	۱٫۸۵۷
۰٫۵۰	۱٫۱۲۲		



شکل ۲.۹ نمودار داده‌های جدول ۳.۹

اگر آن را به شکل (۲۶.۷.۹) بنویسیم، داریم

$$A^T A = \begin{bmatrix} 21 & 10.5 & 7.175 & 5.5125 \\ 10.5 & 7.175 & 5.5125 & 4.51666 \\ 7.175 & 5.5125 & 4.51666 & 3.85416 \\ 5.5125 & 4.51666 & 3.85416 & 3.38212 \end{bmatrix}$$

$$A^T b = [24.1180, 13.2345, 9.468365, 7.5594405]^T \quad (29.7.9)$$

جواب چنین است

$$x^* = [0.05747, 4.7259, -11.1282, 7.6687]^T \quad (30.7.9)$$

این جواب نسبت به تغییرات در b خیلی حساس است. به این حساسیت می‌توان از ضریب وضعیت

$$\text{cond}(A^T A) = 1210.5 \quad (31.7.9)$$

پی برد. به‌عنوان یک گواه دیگر، بردار سمت راست $A^T b$ را با بردار

$$[0.01, -0.01, 0.01, -0.01]^T$$

جمع کرده اختلال ایجاد می‌کنیم. این اختلال با اندازه خطاهای موجود در مقادیر داده‌های b_i سازگار است. با این مقدار جدید سمت راست، معادله نرمال جواب اختلال‌یافته زیر را دارد

$$\hat{x}^* = [0.7408, 2.6825, -6.1538, 4.4550]^T$$

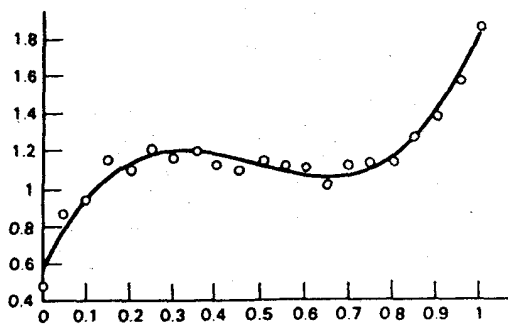
که تفاوت زیادی با x^* دارد. نمودار برازش کمترین مربعات

$$g^*(t) = x_0^* + x_1^* t + x_2^* t^2 + x_3^* t^3$$

همراه با داده‌ها در شکل ۳.۹ نشان داده شده است. خطای ریشه دوم میانگین نمودار چنین است

$$E(x^*) = 0.421$$

ماتریس A در این مثال، دارای ستونهایی است که تقریباً خطی - وابسته‌اند، و $A^T A$ ضریب وضعیت بزرگی دارد. برای بهبود آن، می‌توانیم یک مجموعه توابع پایه $\{\varphi_i(t)\}$ ی بهتری از



شکل ۳.۹ برازش کمترین مربعات $g^*(t)$

چند جمله‌بیهی درجه نایزگتر از ۳، برای خانواده \mathcal{F} انتخاب کنیم. از بررسی ضرایب $A^T A$ ، داریم

$$[A^T A]_{jk} = \sum_{i=1}^m \varphi_j(t_i) \varphi_k(t_i) \quad 1 \leq j, k \leq n \quad (۳۲.۷.۹)$$

اگر نقاط $\{t_i\}$ در سرتاسر بازه $[a, b]$ خوب توزیع شده باشند، آنگاه مجموع، وقتی در $(b-a)/m$ ضرب شود، تقریبی برای

$$\int_a^b \varphi_k(t) \varphi_j(t) dt$$

خواهد بود. برای به دست آوردن یک ماتریس $A^T A$ که ضریب وضعیت کوچکتری داشته باشد، توابع $\varphi_j(t)$ را یکا متعامد انتخاب می‌کنیم. در این صورت $A^T A$ ماتریس تقریباً یکا متعامد است و A ستونهای تقریباً یکا متعامدی خواهد داشت، که به ضریب وضعیت نزدیک به واحد منتهی خواهد شد. در واقع آنچه از همه مهمتر است متعامد بودن خانواده $\{\varphi_j(t)\}$ است، زیرا در آن صورت ماتریس $A^T A$ تقریباً قطری است، که یک ماتریس خوش-وضع است.

مثال مثال قبلی را تکرار می‌کنیم، و از چند جمله‌بیهی لژاندار که بر بازه $[0, 1]$ یکا متعامدند استفاده می‌کنیم. نخستین چهار چند جمله‌بی لژاندار که در بازه $[0, 1]$ یکا متعامدند عبارت‌اند از:

$$\begin{aligned} \varphi_0(t) &= 1 & \varphi_1(t) &= \sqrt{3}s & \varphi_2(t) &= \frac{\sqrt{5}}{2}(3s^2 - 1) \\ \varphi_3(t) & & & & &= \frac{\sqrt{7}}{4}(5s^3 - 3s) \end{aligned} \quad (۳۳.۷.۹)$$

که در آنها $s = 2t - 1$ ، $0 \leq t \leq 1$ ، برای معادله نرمال (۲۶.۷.۹)،

$$A^T A = \begin{bmatrix} 21,0000 & 0 & 2,3479 & 0 \\ 0 & 23,1000 & 0 & 5,1164 \\ 2,3479 & 0 & 25,4993 & 0 \\ 0 & 5,1164 & 0 & 28,3889 \end{bmatrix}$$

$$A^T b = [24,1118, 4,0721, 3,4015, 4,8519]^T$$

$$x^* = [1,1454, 0,1442, 0,0279, 0,1449]^T \quad (۳۴.۷.۹)$$

ضریب وضعیت $A^T A$ اکنون برابر است با

$$\text{cond}(A^T A) = 1,58 \quad (۳۵.۷.۹)$$

که بسیار کوچکتر از مقدار قبلی در (۳۱.۷.۹) است.

حل دستگاههای خطی به روش کمترین مربعات ۷۳۱

حل با روش QR تجزیه به QR در بخش ۳.۹ و به دنبال آن (۱۱.۳.۹) را به یاد می‌آوریم. مانند آنجا، ماتریسهای هاوسهولدر از مرتبه $m \times m$ را در نظر می‌گیریم

$$P_j = I - 2w^{(j)}w^{(j)T} \quad j = 1, \dots, n$$

تا عناصر زیر قطر در A را به صفر بدل کنیم. ماتریسهای متعامد P_j را متوالیاً به‌کار می‌بریم تا عناصر زیر قطر در ستونهای اول تا n به صفر بدل شوند. بردار $w^{(j)}$ در موضع j ام تا m ام، عناصر ناصفر دارد. این فرایند به ماتریس زیر می‌انجامد

$$R = P_n \dots P_1 A = Q^T A \quad (۳۶.۷.۹)$$

اگر این عملیات را در سمت راست دستگاه $Ax = b$ هم اعمال کنیم، دستگاه هم‌ارز زیر را به‌دست می‌آوریم

$$Rx = Q^T b \quad (۳۷.۷.۹)$$

ماتریس R به شکل زیر است

$$R = \begin{bmatrix} R_1 \\ \circ \end{bmatrix} \quad (۳۸.۷.۹)$$

که R_1 یک ماتریس مربعی بالامتثلی از مرتبه $n \times n$ است. ماتریس R_1 ناکتین است، زیرا A و $R = Q^T A$ یک رتبه دارند که n است. برطبق (۳۸.۷.۹) می‌نویسیم

$$Q^T b = \begin{bmatrix} g_1 \\ g_2 \end{bmatrix} \quad g_1 \in \mathbf{R}^n \quad g_2 \in \mathbf{R}^{m-n}$$

پس

$$\begin{aligned} \|Ax - b\|_2 &= \|Q^T Ax - Q^T b\|_2 = \|Rx - Q^T b\|_2 \\ &= \left[\|R_1 x - g_1\|_2^2 + \|g_2\|_2^2 \right]^{1/2} \end{aligned}$$

جواب کمترین مربعات $Ax = b$ از حل دستگاه بالامتثلی ناکتین زیر به‌دست می‌آید

$$R_1 x = g_1 \quad (۳۹.۷.۹)$$

پس مینیمم چنین است

$$\|Ax^* - b\|_2 = \|g_2\|_2 \quad (۴۰.۷.۹)$$

روش QR برای محاسبه x^* از لحاظ عملیات، کمی از روش چولسکی گرانتر است. تعداد عملیات (ضرب و جمع) روش چولسکی، از جمله تشکیل $A^T A$ ، در حدود

$$\frac{1}{2}mn^2 + \frac{n^3}{6}$$

است و تعداد عملیات روش هاوسهولدر QR در حدود

$$mn^2 - \frac{n^3}{3}$$

است. معذالک معمولاً روش QR برای حل کمترین مربعات توصیه می‌شود. این روش مستقیماً روی A انجام می‌شود و به همین علت و به علت استفاده از تبدیلات متعامد، اثر خطاهای گرد کردن با این روش از نتیجه آن در استفاده از تجزیه به عوامل چولسکی در حل معادله نرمال بهتر است. برای بحث عمیقتر گلوب و ون لون (۱۹۸۳، صص ۱۴۷-۱۴۹) و لاونسن و هنسن (۱۹۷۴، فصل ۱۶) را ملاحظه کنید.

مثال مثال قبلی دستگاه خطی (۲۸.۷.۹)

$$A_{ij} = [t_i^{j-1}] \quad 1 \leq i \leq 21, \quad 1 \leq j \leq 4$$

را همراه با داده‌های جدول ۳.۹ در نظر می‌گیریم. در اینجا $\text{cond}(A) = 110^{\circ}01$.

$$R_1 = \begin{bmatrix} -4,5826 & -2,2913 & -1,5657 & -1,2029 \\ \circ & 1,3874 & 1,3874 & 1,2688 \\ \circ & \circ & -0,3744 & -0,5617 \\ \circ & \circ & \circ & -0,9887 \end{bmatrix}$$

$$g_1 = [-5,2630, 0,8472, -0,1403, -0,7566]^T$$

این جواب x^* با جواب x^* در (۳۰.۷.۹) یکی است، همان‌طور که خطای ریشه دوم میانگین یکی است.

تجزیه تکین - مقدار تجزیه تکین - مقدار ابزار خیلی باارزشی برای تحلیل و حل مسائل کمترین مربعات و مسائل دیگر در جبر خطی است. برای مسائل کمترین مربعات در حالتی که رتبه کامل نیست، روش QR که هم‌اکنون توضیح داده شد احتمالاً به ماتریس مثلثی R_1 می‌انجامد که نانکین است، ولی دارای عناصر قطری بسیار کوچک است. در این صورت تجزیه تکین - مقدار A ممکن است در روش‌ترکردن ساختار A بسیار مفید باشد. اگر بعضی از تکین - مقدارهای μ_i

نزدیک به صفر باشند، آنگاه تعیین اثر صفر گذاردن آنها برای یافتن جواب x^* ممکن است ساده تر از تعیین نتیجه در سایر روشها برای حل کردن x^* باشد. از این رو توجه فراوانی برای پیدا کردن روشهای کارا برای محاسبه تجزیه تکین - مقدار A وجود دارد.

یکی از بهترین راههای شناخته شده برای محاسبه تجزیه تکین - مقدار A ، راه منسوب به گلوب^۱ و راینش^۲ و کاهان^۳ است و بحث کامل آن در گلوب و ون لون (۱۹۸۳، بخش ۵.۶) داده شده است. ما فقط نشان می دهیم که چگونه تجزیه تکین - مقدار در (۵.۷.۹) را می توان از حل مسأله ویژه مقدار یک ماتریس متقارن همراه با تجزیه QR به دست آورد.

اگر A حقیقی و $m \times n$ با $m \geq n$ باشد، می دانیم که $A^T A$ یک ماتریس حقیقی $n \times n$ است. به علاوه به سادگی می توان نشان داد که $A^T A$ متقارن و نیمه معین مثبت است [به ازای جمیع مقادیر x ، $x^T A x \geq 0$]. با استفاده از یک برنامه برای حل مسأله ویژه مقدار متقارن، یک ماتریس قطری D و یک ماتریس متعامد U پیدا می کنیم که

$$U^T(A^T A)U = D \quad (41.7.9)$$

گیریم $D = \text{diag}[\lambda_1, \dots, \lambda_n]$ ، با ویژه مقدارهایی که به ترتیب نزولی مرتب شده اند. اگر یکی از λ_i ها یک عدد کوچک منفی باشد آن را صفر می گیریم، زیرا همه ویژه مقدارهای $A^T A$ باید نامنفی باشند مگر آنکه اختلالاتی بر اثر خطاهای گرد کردن حاصل شده باشد.

در (41.7.9)، B از مرتبه $m \times n$ را چنین تعریف می کنیم $B = AU$. در این صورت از (41.7.9) نتیجه می شود

$$B^T B = D$$

پس ستونهای B متعامدند. به علاوه اگر یکی از λ_i ها صفر باشد، ستون متناظر با آن در B باید متحد با صفر باشد، زیرا نرم آن صفر است. با به کار بردن روش QR ، یک ماتریس متعامد V پیدا می کنیم که برای آن

$$V^T B = R \quad (42.7.9)$$

زیر قطر اصلی، همه ستونها صفر باشند. ماتریس R در رابطه زیر صدق می کند

$$R^T R = B^T V^T V B = B^T B = D$$

باز هم، ستونهای R باید متعامد باشند، و اگر یکی از λ_i ها صفر باشد، ستون متناظر آن در R نیز باید صفر باشد. چون R یک بالامثلثی است، می‌توانیم از تعامد استفاده کرده نشان دهیم که ستونهای R در تمام مواضع بالای قطر دارای عناصر صفرند. بنابراین R به شکل ماتریس F در (۵.۷.۹) است، به جز آنکه بعضی عناصر قطر R ممکن است منفی باشند. در این حالت، علامت این عنصر قطر و ستون متناظر آن در U را عوض می‌کنیم. در این صورت داریم $R = F$ ، $\mu_i = \sqrt{\lambda_i}$ در (۴۲.۷.۹) قرار می‌دهیم $B = AU$ ، شکل تجزیه تکین - مقدار مورد نظر را خواهیم داشت:

$$V^T AU = R$$

یکی از معایب ممکن این شیوه آن است که $A^T A$ را باید تشکیل داد و این کار ممکن است به از دست دادن بعضی اطلاعات ناشی از استفاده از حساب با طول متناهی رایانه منجر شود. ولی اگر مسأله ویژه مقدار متقارن حلپذیر باشد، اجرای این روش ساده است.

مثال مجدداً ماتریس A در (۲۸.۷.۹) را بر پایه داده‌های جدول ۳.۹ در نظر می‌گیریم. ماتریس $A^T A$ در (۲۹.۷.۹) داده شده است. با به‌کار بردن برنامه‌های EISPACK و LINPACK خواهیم داشت

$$U = \begin{bmatrix} ۰٫۷۸۲۷ & ۰٫۵۹۶۳ & -۰٫۱۷۶۴ & ۰٫۲۵۶ \\ ۰٫۴۵۳۳ & -۰٫۳۵۹۶ & ۰٫۷۴۸۹ & -۰٫۳۲۳۱ \\ ۰٫۳۳۲۶ & -۰٫۴۹۹۸ & -۰٫۰۹۸۹ & ۰٫۷۹۳۶ \\ ۰٫۲۶۷۰ & -۰٫۵۱۵۰ & -۰٫۶۳۱۱ & -۰٫۵۱۵۰ \end{bmatrix}$$

مقدارهای تکین چنین‌اند

$$۳۲٫۰۱۰۲، ۳٫۸۹۳۵، ۰٫۱۶۷۴، ۰٫۰۰۲۶$$

ماتریس V متعامد و از مرتبه ۲۱×۲۱ است و ما به دلایلی که روشن است آن را نمی‌نویسیم. در عمل، این ماتریس محاسبه نمی‌شود، زیرا حاصلضرب چهار ماتریس هاؤسهولدر است که می‌توان آنها را به شکل ساده‌تری ذخیره کرد.

برای یک بحث مفصلتر در حل مسائل کمترین مربعات گلوب و وِن‌لون (۱۹۸۳، فصل ۶) و کتاب لائوسن و هِنسن (۱۹۷۴) را ملاحظه کنید. مسائل عملی بیشتری وجود دارند که باید مورد بحث قرار گیرند، از جمله معین‌کردن رتبه ماتریس وقتی که خطای گردکردن موجب شود که

ماتریس به غلط رتبه کامل پیدا کند. برای برنامه‌ها، پیوست کتاب (۱۹۷۴) و LINPACK را ببینید. برای تجزیه تکین - مقدار به LINPACK یا EISPACK مراجعه کنید.

بحث در آثار خواندنی

منبع اصلی اطلاعات این فصل، کتاب دایرةالمعارف گونه معروف ویلکینسن (۱۹۶۵) بوده است. منابع دیگر عبارت بودند از گلوب و ون لون (۱۹۸۳)، گولی و واتسن (۱۹۷۶) استوارت (۱۹۶۴)، نوبل (۱۹۶۹، فصلهای ۹-۱۲)، پارلت (۱۹۸۰)، استوارت (۱۹۷۳) و ویلکینسن (۱۹۶۳). برای ماتریسهای با اندازه‌های معمولی، حل عددی مسأله ویژه مقدار نسبتاً خوب فهمیده شده است. برای جنبه دیگر روش QR ، واتکینز (۱۹۸۲) را ببینید و برای نگاه دقیق به بارست معکوس، پیترز و ویلکینسن (۱۹۷۹) را ملاحظه کنید. الگوریتمهای بسیار عالی برای اغلب مسائل ویژه مقدار در ویلکینسن و و راینش (۱۹۷۱) و راهنماهای EISPACK توسط اسمیت و همکاران (۱۹۷۶) و گاربو و همکاران (۱۹۷۷) داده شده‌اند. برای تاریخ طرح EISPACK به دونگارا و مولر (۱۹۸۴) مراجعه کنید. یک گزارش کلی عالی از مسائل توسعه نرم‌افزار ریاضی برای مسائل ویژه مقداری و سایر مسائل ماتریسی، در راس (۱۹۸۱) داده شده است. بسته نرم‌افزاری EISPACK پایه بسیاری از برنامه‌های ویژه مقدار است که در کتابخانه‌های نرم‌افزاری IMSL و NAG وجود دارند. بیشتر به علت محدودیت جا، شماری از مسائل و روشهای عددی در این فصل مورد بحث قرار نگرفته‌اند. برای مسأله ویژه مقدار متقارن، روش ژاکوبی ذکر نشده است. این روش یک روش خوب و همگرای سریع برای محاسبه تمام ویژه مقدارهای یک ماتریس متقارن است، و برنامه‌ریزی آن نسبتاً ساده است. برای توضیح روش ژاکوبی، گلوب و ون لون (۱۹۸۳، بخش ۴.۸) و پارلت (۱۹۸۰، فصل ۹) و ویلکینسن (۱۹۶۵، صص ۲۶۶-۲۸۲) را ببینید. یک برنامه به زبان ALGOL ویلکینسن و راینش (۱۹۷۱، صص ۲۰۲-۲۱۱) داده شده است. مسأله تعمیم‌یافته ویژه مقدار $Ax = \lambda Bx$ نیز حذف شده است. در سالهای اخیر این مسأله یک مسأله مهم شده است. متداولترین روش حل آن از مولر و استوارت (۱۹۷۳) است و توضیحات دیگر مسأله و حل آن در گلوب و ون لون (۱۹۸۳، بخشهای ۷.۷ و ۶.۸) و پارلت (۱۹۸۰، فصل ۱۵) داده شده است. برنامه‌های EISPACK برای مسأله ویژه مقدار تعمیم‌یافته در گاربو و همکاران (۱۹۷۷) داده شده است.

مسأله پیدا کردن ویژه مقدارها و ویژه بردارهای ماتریسهای بزرگ تنک یک زمینه فعال پژوهشی است. هنگامی که ماتریسها دارای مرتبه بسیار بزرگ باشند (مثلاً $n \geq 300$)، استفاده از اغلب روشهای این فصل به دلیل ملاحظات حافظه‌ای، دشوارتر است. به علاوه، این روشها اغلب توجه خاصی به این موضوع ندارند که در عمل بیشتر این ماتریسها تنک هستند. یک شکل معمول

مسأله شامل ماتریسهای متقارن کراندار است. برنامه‌های مربوط به این مسأله در ویلکینسن و راینش (۱۹۷۱، صص ۲۶۶-۲۸۳) و گاریو و همکاران (۱۹۷۷) داده شده‌اند. برای بحثهای کلیتر در مسأله ویژه مقدار برای ماتریسهای تنک، جنینگز (۱۹۸۵) و پیسانتسکی^۱ (۱۹۸۴، فصل ۶) را ببینید. برای بحث در مورد نرم‌افزار مربوط به مسأله ویژه مقدار ماتریسهای تنک، داف^۲ (۱۹۸۴، صص ۱۷۹-۱۸۲) و هیث^۳ (۱۹۸۲) را ملاحظه کنید. یک روش مهم در حل مسأله ویژه مقدار برای ماتریسهای تنک متقارن، روش لانتسوس^۴ است. برای بحث در این مورد، اسکات^۵ (۱۹۸۱) و کتابها و برنامه‌های بسیار جامع کالوم^۶ و ویلی^۷ (۱۹۸۴ و ۱۹۸۵) را ببینید. حل کمترین مربعات دستگاههای خطی فرامعین یک وسیله بسیار باارزش است، وسیله‌ای که در علوم پزشکی، بیولوژیکی و اجتماعی کاربرد گسترده‌ای پیدا کرده است. ما فقط به معرفی برخی از جنبه‌های این موضوع که نقش قاطع تجزیه تکین مقدار را نشان می‌دهد، پرداخته‌ایم. یک مقدمه خیلی جامع برای حل کمترین مربعات دستگاههای خطی در لاؤسن و هینسن (۱۹۷۴) داده شده است. این کتاب نحوه به‌کار بردن کلی این نظریه و اجرای عملی روشها و طرق استفاده از مجموعه‌های داده‌های بزرگ را به نحو کارایی ارائه می‌دهد. برای مراجع دیگر در حل کمترین مربعات دستگاههای خطی، گلوب و وِن‌لون (۱۹۸۳، فصل ۶) و رایس (۱۹۸۱، فصل ۱۱) را ببینید. برنامه‌هایی برای بعضی از مسائل کمترین مربعات در LINPACK نیز داده شده‌اند.

در بحث حل کمترین مربعات دستگاههای معادله‌های خطی فرامعین، از بحث در جنبه آماری موضوع خودداری کرده‌ایم. بخشی به علت کمبود جا، و بخشی به دلیل بی‌اعتمادی در استفاده از توجیه آماری بوده است. زیرا این مسأله اغلب به فرضیهایی درباره توزیع خطا ارتباط پیدا می‌کند که اعتبار بخشیدن به آن مشکل است. ما خواننده را به مطالعه یکی از کتابهای درسی آماری زیادی که یک چارچوب آماری برای روش کمترین مربعات در برازش منحنی داده‌ها ارائه می‌دهند ارجاع می‌دهیم.

مراجع

- Chatelin, F. (1987). *Eigenvalues of Matrices*. Wiley, London.
- Conte, S., and C. de Boor (1980). *Elementary Numerical Analysis*, 3rd ed. McGraw-Hill, New York.
- Cullum, J., and R. Willoughby (1984, 1985). *Lanczos Algorithms for Large Symmetric Eigenvalue Computations*, Vol. 1, Theory; Vol. 2, Programs. Birkhäuser, Basel.

1. Pissanetzky

2. Duff

3. Heath

4. Lanczos

5. Scott

6. Culum

7. Willoughby

- Dongarra, J., and C. Moler (1984). EISPACK— A package for solving matrix eigenvalue problems. In *Sources and Development of Mathematical Software*, W. Cowell (Ed.), pp. 68–87. Prentice-Hall, Englewood Cliffs, N.J.
- Dongarra, J., J. Bunch, C. Moler, and G. Stewart (1979). *LINPACK User's Guide*. SIAM Pub., Philadelphia.
- Duff, I. (1984). A survey of sparse matrix software, In *Sources and Development of Mathematical Software*, W. Cowell (Ed.). Prentice-Hall, Englewood Cliffs, N.J.
- Garbow, B., J. Boyle, J. Dongarra, and C. Moler (1977). *Matrix Eigensystems Routines—EISPACK Guide Extension, Lecture Notes in Computer Science*, Vol. 51. Springer-Verlag, New York.
- Golub, G., and C. Van Loan (1983). *Matrix Computations*. Johns Hopkins Press, Baltimore.
- Gourlay, A., and G. Watson (1976). *Computational Methods for Matrix Eigenproblems*. Wiley, New York.
- Gregory, R., and D. Karney (1969). *A Collection of Matrices for Testing Computational Algorithms*. Wiley, New York.
- Heath, M., Ed. (1982). *Sparse Matrix Software Catalog*. Oak Ridge National Laboratory, Mathematics and Statistics Dept., Tech. Rep. Oak Ridge, Tenn.
- Henrici, P. (1974). *Applied and Computational Complex Analysis*, Vol. I. Wiley, New York.
- Householder, A. (1964). *The Theory of Matrices in Numerical Analysis*. Ginn (Blaisdell), Boston.
- Jennings, A. (1985). Solutions of sparse eigenvalue problems. In *Sparsity and Its Applications*, D. Evans (Ed.), pp. 153–184. Cambridge Univ. Press, Cambridge, England.
- Lawson, C., and R. Hanson (1974). *Solving Least Squares Problems*. Prentice-Hall, Englewood Cliffs, N.J.
- Moler, C., and G. Stewart (1973). An algorithm for generalized matrix eigenvalue problems, *SIAM J. Numer. Anal.* 10, 241–256.
- Noble, B. (1969). *Applied Linear Algebra*. Prentice-Hall, Englewood Cliffs, N.J.
- Parlett, B. (1968). Global convergence of the basic QR algorithm on Hessenberg matrices, *Math. Comput.* 22, 803–817.
- Parlett, B. (1980). *The Symmetric Eigenvalue Problem*. Prentice-Hall, Englewood Cliffs, N.J.
- Peters, G., and J. Wilkinson (1979). Inverse iteration, ill-conditioned equations and Newton's method, *SIAM Rev.* 21, 339–360.
- Pissanetzky, S. (1984). *Sparse Matrix Technology*. Academic Press, New York.

- Rice, J. (1981). *Matrix Computations and Mathematical Software*. McGraw-Hill, New York.
- Scott, I. (1981). The Lanczos algorithm. In *Sparse Matrices and Their Uses*, I. Duff (Ed.), pp. 139-160. Academic Press, London.
- Smith, B. T., J. Boyle, B. Garbow, Y. Ikebe, V. Klema, and C. Moler (1976). *Matrix Eigensystem Routines—EISPACK Guide*, 2nd ed., *Lecture Notes in Computer Science*, Vol. 6. Springer-Verlag, New York.
- Stewart, G. (1973). *Introduction to Matrix Computations*. Academic Press, New York.
- Watkins, D. (1982). Understanding the QR algorithm, *SIAM Rev.* **24**, 427-440.
- Wilkinson, J. (1963). *Rounding Errors in Algebraic Processes*. Prentice-Hall, Englewood Cliffs, N.J.
- Wilkinson, J. (1965). *The Algebraic Eigenvalue Problem*. Oxford Univ. Press, Oxford, England.
- Wilkinson, J. (1968). Global convergence of the tridiagonal QR algorithm with origin shifts. *Linear Algebra Its Appl.* **1**, 409-420.
- Wilkinson, J., and C. Reinsch, Eds. (1971). *Linear Algebra*. Springer-Verlag, New York.

مسائل

۱. از قضیه ۱.۹، گرشگورین، برای تعیین جای تقریبی ویژه مقدارهای ماتریسهای زیر، استفاده کنید.

$$\begin{bmatrix} -2 & 1 & 1 \\ 1 & 3 & 1 \\ 1 & -1 & 3 \end{bmatrix} \quad (\text{ب}) \quad \begin{bmatrix} 1 & -1 & 0 \\ 1 & 5 & 1 \\ -2 & -1 & 9 \end{bmatrix} \quad (\text{الف})$$

هرجا که ممکن است از این نتایج برای تعیین حقیقی یا مختلط بودن ویژه مقدارها استفاده کنید. برای بررسی این نتایج، ویژه مقدارها را مستقیماً از طریق یافتن ریشه‌های چندجمله‌یی مشخصه پیدا کنید.

۲. (الف) چندجمله‌یی زیر داده شده است:

$$p(\lambda) = \lambda^n + a_{n-1}\lambda^{n-1} + \dots + a_0.$$

نشان دهید که برای ماتریس زیر $p(\lambda) = \det[\lambda I - A]$

$$A = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & & \vdots \\ \vdots & & & \ddots & \\ 0 & & & & 1 \\ -a_0 & -a_1 & \dots & & -a_{n-1} \end{bmatrix}$$

حل دستگاههای خطی به روش کمترین مربعات ۷۳۹

ریشه‌های $p(\lambda)$ ویژه‌مقدارهای A هستند. ماتریس A را ماتریس همراه چندجمله‌یی $p(\lambda)$ می‌نامند.
(ب) با استفاده از قضیه ۱.۹ گرشگورین، نشان دهید که:

$$|r| \leq 1 \quad \text{یا} \quad |r + a_{n-1}| \leq |a_0| + \dots + |a_{n-2}|$$

کرانه‌های ریشه‌های $p(\lambda)$ هستند. اگر این کرانه‌ها در صفحه مختلط نواحی مجزا باشند درباره تعداد ریشه‌ها در هر ناحیه چه می‌توان گفت؟

(ج) قضیه گرشگورین را برای ستونهای A به‌کار برید تا کرانه‌های دیگری برای ریشه‌های $p(\lambda)$ به‌دست آورید.

(د) از نتایج قسمتهای (ب) و (ج) برای یافتن حدود ریشه‌های معادله‌های چندجمله‌یی زیر استفاده کنید.

$$\lambda^{10} + 8\lambda^9 + 1 = 0 \quad (\text{i})$$

$$\lambda^6 - 4\lambda^5 + \lambda^4 - \lambda^3 + \lambda^2 - \lambda + 1 = 0 \quad (\text{ii})$$

۳. دستگاه خطی (۵.۸.۸) از فصل ۸ را که هنگام حل عددی معادله پواسون به‌وجود آمده به‌خاطر آورید. اگر معادلات به‌صورت مذکور در ۱۲.۸.۸ و پس از آن مرتب شوند، آنگاه این دستگاه خطی، متقارن با عناصر قطری مثبت است. براساس قضیه ۷.۸ شرط لازم و کافی برای همگرایی روش بارستی گاوس - زایدل این است که A معین مثبت باشد. از قضیه ۱.۹ گرشگورین برای اثبات معین مثبت بودن A استفاده کنید. همچنین لازم است قضیه ۸.۸ را برای اینکه $\lambda = 0$ ویژه‌مقدار ماتریس A نیست ذکر کنید.

۴. مقادیر $\lambda = -8.02861$ و

$$x = [1.0, 2.50146, -0.75773, -2.56421]^T$$

یک تقریب ویژه‌مقدار و یک تقریب ویژه بردار ماتریس زیرند:

$$A = \begin{bmatrix} 2 & 1 & 3 & 4 \\ 1 & -3 & 1 & 5 \\ 3 & 1 & 6 & -2 \\ 4 & 5 & -2 & -1 \end{bmatrix}$$

نتیجه ۲۲.۱.۹ را برای محاسبه یک کران خطا برای λ به‌کار برید.

۵. برای ماتریس مثال ۱۷.۱.۹ با $\epsilon = 0.001$ و $\lambda = 2$ کران خطای اختلال ۳۶.۱.۹ را محاسبه کنید. همین کران در ۳۸.۱.۹ برای ویژه‌مقدار دیگر، $\lambda = 1$ ، داده شده بود.

۶. قضیه اختلال ویژه بردار (۴۱.۱.۹) را ثابت کنید.

راهنمایی: فرض کنید $\lambda_k(\epsilon)$ و $u_k(\epsilon)$ توابع پیوسته مشتق پذیر ϵ باشند. با استفاده از ۳۲.۱.۹ داریم: $\lambda'_k(0) = u_k^* B u_k / s_k$. بنویسید:

$$u_k(\epsilon) = u_k(0) + \epsilon u'_k(0) + O(\epsilon^2)$$

و آن را نسبت به $u'_k(0)$ حل کنید. چون $\{u_1, \dots, u_n\}$ یک پایه است، بنویسید

$$u'_k(0) = \sum_{j=1}^n a_j u_j$$

برای یافتن a_j ابتدا از ۴۰.۱.۹ نسبت به ϵ مشتق بگیرید و سپس فرض کنید $\epsilon = 0$. به جای $u'_k(0)$ نمایش قبلی را بگذارید. از ۲۹.۱.۹ و از دوه دو متعامد بودن ۲۸.۱.۹، یعنی $u_i^* u_j = 0$ استفاده کنید.

۷. برای ماتریسهای $A(\epsilon)$ در زیر، ویژه مقادیر و ویژه بردارها را به ازای $\epsilon = 0$ و $\epsilon > 0$ تعیین کنید. رفتار آنها را هنگامی که $\epsilon \rightarrow 0$ بررسی کنید.

$$\begin{matrix} \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & \epsilon \\ 0 & \epsilon & 1 \end{bmatrix} & \text{(د)} & \begin{bmatrix} 1 & \epsilon \\ 0 & 1 \end{bmatrix} & \text{(ج)} & \begin{bmatrix} 1 & 1 \\ 0 & 1 + \epsilon \end{bmatrix} & \text{(ب)} & \begin{bmatrix} 1 & 1 \\ \epsilon & 1 \end{bmatrix} & \text{(الف)} \end{matrix}$$

این مثالها چه چیزی از پایداری زیرفضاهای ویژه بردارها را بیان می کنند؟

۸. روش توانی را برای محاسبه مقدار غالب و ویژه بردار متناظر آن، برای ماتریسهای زیر به کار برید.

$$\begin{matrix} \begin{bmatrix} 2 & 1 & 3 & 4 \\ 1 & -3 & 1 & 5 \\ 3 & 1 & 6 & -2 \\ 4 & 5 & -2 & -1 \end{bmatrix} & \text{(ب)} & \begin{bmatrix} 6 & 4 & 4 & 1 \\ 4 & 6 & 1 & 4 \\ 4 & 1 & 6 & 4 \\ 1 & 4 & 4 & 6 \end{bmatrix} & \text{(الف)} \\ & & \begin{bmatrix} 1 & 3 & -2 \\ -1 & -2 & 3 \\ 1 & 1 & 2 \end{bmatrix} & \text{(ج)} \end{matrix}$$

با محاسبه نسبتهای R_m در ۱۴.۲.۹ سرعت همگرایی را بررسی کنید.

هنگامی که نسبتهای R_m تقریباً ثابت‌اند از برون‌یابی ایتکن برای بهبود سرعت همگرایی ویژه‌مقدار و ویژه‌بردار استفاده کنید. با استفاده از نسبتهای ویژه‌مقدار R_m همگرایی ویژه‌بردارهای $\{z^{(m)}\}$ را تسریع کنید.

۹. روش توانی را برای یافتن ویژه‌مقدار غالب ماتریس

$$A = \begin{bmatrix} 7 & 13 & -16 \\ 13 & -10 & 13 \\ -16 & 13 & 7 \end{bmatrix}$$

به‌کار برید. از حدس اولیه $z^{(0)} = [1, 0, 1]^T$ استفاده کنید. هریک از بارستهای $z^{(m)}$ و $\lambda_1^{(m)}$ را بنویسید. در باب نتایج توضیح دهید. اگر $\lambda_1^{(m)}$ با $\alpha_m = \lambda_1^{(m)}$ تعریف شود چه پیش می‌آید؟
۱۰. برای یک ماتریس A از مرتبه n ، فرض کنید شکل متعارف ژوردان آن قطری است و ویژه‌مقدارها با $\lambda_1, \dots, \lambda_n$ نشان داده شده‌اند. فرض کنید $\lambda_1 = \lambda_2 = \dots = \lambda_r$ ، برای $r > 1$ و

$$|\lambda_r| > |\lambda_{r+1}| \geq \dots \geq |\lambda_n| \geq 0$$

نشان دهید که روش توانی (۲.۲.۹) - (۳.۲.۹) به λ_1 و برای اغلب بردارهای اولیه $z^{(0)}$ به یک ویژه‌بردار متناظر آن همگرا می‌شود.

۱۱. فرض کنید A یک ماتریس متقارن از مرتبه n باشد که ویژه‌مقدارهایش به‌صورت زیر مرتب شده‌اند:

$$\lambda \leq \lambda_2 \leq \dots \leq \lambda_n$$

با استفاده از ضرب داخلی استاندارد، تعریف می‌کنیم:

$$\mathcal{R}(x) = \frac{(Ax, x)}{(x, x)} \quad x \neq 0 \quad x \in \mathbb{R}^n$$

نشان دهید هنگامی که $x \neq 0$ در \mathbb{R}^n است داریم:

$$\text{Max} \mathcal{R}(x) = \lambda_1 \quad \text{Min} \mathcal{R}(x) = \lambda_n$$

تابع $\mathcal{R}(x)$ خارج‌قسمت ریلی نامیده می‌شود و می‌توان آن را برای مشخص‌کردن ویژه‌مقدارهای باقی‌مانده A ، نیز به‌کار برد. این ماکسیمم‌سازی و مینیمم‌سازی $\mathcal{R}(x)$ پایه بعضی از روشهای عددی کلاسیک برای محاسبه ویژه‌مقدارهای A است.

۱۲. برای به دست آوردن یک تعبیر هندسی برای ماتریس $n \times n$ هاوسهولدر، $P = I - 2ww^T$ فرض کنید $u^{(2)}, \dots, u^{(n)}$ یک پایهٔ یکا متعامد برای زیرفضای $(n-1)$ بعدی عمود بر w باشد. تعریف می‌کنیم

$$T(x) = (I - 2ww^T)x \quad x \in \mathbf{R}^n$$

با استفاده از پایهٔ $\{w, u^{(2)}, \dots, u^{(n)}\}$ برای \mathbf{R}^n می‌نویسیم:

$$x = a_1 w + a_2 u^{(2)} + \dots + a_n u^{(n)}$$

T را برای این نمایش به کار برید و نتیجه‌ها را تعبیر کنید.

۱۳. الف) فرض کنید A یک ماتریس متقارن باشد و λ و x یک جفت ویژه مقدار و ویژه بردار A با ویژگی $\|x\|_2 = 1$ باشند. فرض می‌کنیم P ماتریس یکا متعامدی باشد که

$$Px = e_1 = [1, 0, \dots, 0]^T$$

ماتریس متشابه $B = PAP^T$ را در نظر بگیرید و نشان دهید که اولین سطر و ستون آن صفرند به جز عنصر روی قطر که برابر λ است.

راهنمایی: Be_1 را محاسبه و از آن استفاده کنید.

ب) برای ماتریس

$$A = \begin{bmatrix} 2 & 10 & 2 \\ 10 & 5 & -8 \\ 2 & -8 & 11 \end{bmatrix}$$

$\lambda = 9$ یک ویژه مقدار متناظر با ویژه بردار $x = [2/3, 1/3, 2/3]^T$ است. یک ماتریس هاوسهولدر P بسازید، به طوری که $Px = e_1$ و سپس $B = PAP^T$ را بسازید. مسألهٔ ویژه مقدار ماتریس B ، به آسانی به یک مسألهٔ با ماتریس 2×2 بدل می‌شود. از این شیوه برای محاسبهٔ بقیهٔ ویژه مقدارها و ویژه بردارهای A استفاده کنید. فرایند تبدیل A به B و سپس حل مسألهٔ ویژه مقدار یک ماتریس که مرتبه‌اش از مرتبهٔ A یک واحد کمتر است به فرایند کاهش معروف است. این فرایند برای تعمیم قابلیت اعمال روش توانی برای ویژه مقدارهای دیگر، غیر از ویژه مقدار غالب، به کار برده می‌شود. برای بحث بیشتر ویلکینسن (۱۹۶۵، صص ۵۸۴-۵۹۸) و پارلت (۱۹۸۰، فصل ۵) را ببینید.

۱۴. ماتریسهای هاوسهولدر را برای تجزیه QR در ماتریسهای زیر به کار برید.

$$\begin{bmatrix} 1 & 3 & -2 \\ -1 & -2 & 3 \\ 1 & 1 & 2 \end{bmatrix} \quad (\text{ب}) \quad \begin{bmatrix} 1 & 1 & 1 \\ 2 & -1 & -1 \\ 2 & -4 & 5 \end{bmatrix} \quad (\text{الف})$$

۱۵. ماتریس دوران مرتبه n در زیر را در نظر بگیرید،

$$R^{(k,l)} = \begin{bmatrix} 1 & 0 & 0 & & \dots & & 0 \\ & 0 & 1 & 0 & & & \vdots \\ & \vdots & & \ddots & & & \\ & 0 & & \alpha & 0 & \dots & \beta & 0 & \dots & 0 \\ & & & 0 & 1 & & 0 & & & \\ \vdots & & & \vdots & & \ddots & & & & \\ & & & -\beta & 0 & & \alpha & 0 & & 0 \\ & & & & & & & 1 & & \\ & & & & & & & & \ddots & \vdots \\ & 0 & & & & \dots & & & & 1 \end{bmatrix}$$

سطر k

سطر l

با $\alpha^2 + \beta^2 = 1$. اگر Rb را، برای مقدار داده شده $b \in \mathbf{R}^n$ ، محاسبه کنیم، تنها عناصری که در مواضع k و l هستند تغییر می کنند. با انتخاب مناسب α و β می توانیم کاری کنیم که یک صفر در موضع l داشته باشیم. α و β را چنان انتخاب کنید که به ازای مقداری از γ

$$\begin{bmatrix} \alpha & \beta \\ -\beta & \alpha \end{bmatrix} \begin{bmatrix} b_k \\ b_l \end{bmatrix} = \begin{bmatrix} \gamma \\ 0 \end{bmatrix}$$

(الف) فرمولهایی برای α و β به دست آورید و نشان دهید که $\gamma = \sqrt{b_k^2 + b_l^2}$
 (ب) با استفاده از دنباله ای از ضربهای ماتریسهای دوران:

$$\hat{b} = R^{(1,2)} R^{(1,3)} R^{(1,4)} b$$

$b = [1, 1, 1, 1]^T$ را به شکل $\hat{b} = [c, 0, 0, 0]$ بدل کنید.

۱۶. نشان دهید که چگونه یک ماتریس دوران $R^{(k,l)}$ را برای تجزیه QR یک ماتریس می‌توان مورد استفاده قرار داد.

۱۷. (الف) مانند بخش ۳.۹، عملها را در به دست آوردن تجزیه QR یک ماتریس به استفاده از ماتریسهای هاؤسهولدر شمارش کنید. مطابق معمول ضربها و تقسیمها را یکجا و جذرها را جداگانه بشمارید.

(ب) قسمت (الف) را تکرار کنید، اما این بار برای تبدیل از ماتریسهای دوران $R^{(k,l)}$ استفاده کنید.

۱۸. فرمولهای صریحی برای محاسبه تجزیه QR در یک ماتریس سه قطری متقارن به دست آورید. تعداد عملها را بشمارید و نتیجه را با آنچه در مسأله ۱۷ به دست آمد مقایسه کنید.

۱۹. قضیه ۵.۹ را برای جدا کردن ریشه‌های

$$\begin{bmatrix} 1 & 2 & 0 & 0 & 0 \\ 2 & 2 & 3 & 0 & 0 \\ 0 & 3 & 3 & 4 & 0 \\ 0 & 0 & 4 & 4 & 5 \\ 0 & 0 & 0 & 5 & 5 \end{bmatrix} \quad (\text{ب}) \quad \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 2 \end{bmatrix} \quad (\text{الف})$$

به کار برید. سپس تقریبهایی دقیقی از ریشه‌ها با استفاده از روش تنصیف یا یک تکنیک دیگر ریشه‌یابی به دست آورید.

۲۰. (الف) با استفاده از ماتریسهای هاؤسهولدر برای تبدیلات تشابهی، برنامه‌ای برای تبدیل یک ماتریس متقارن به شکل سه قطری بنویسید. برای کارایی بیشتر در ضربهای ماتریسی، همانند ضربی که در بخش ۱۸.۳.۹ نشان داده شده عمل کنید.

(ب) این برنامه را برای تبدیل ماتریسهای زیر به شکل سه قطری به کار برید:

$$\begin{bmatrix} 4 & 6 & 242 & 12 \\ 6 & 225 & 3 & 18 \\ 242 & 3 & 25 & 6 \\ 12 & 18 & 6 & 0 \end{bmatrix} \quad (\text{iii}) \quad \begin{bmatrix} 5 & 4 & 1 & 1 \\ 4 & 5 & 1 & 1 \\ 1 & 1 & 4 & 2 \\ 1 & 1 & 2 & 4 \end{bmatrix} \quad (\text{ii}) \quad \begin{bmatrix} 1 & 2 & 3 \\ 2 & 3 & 5 \\ 3 & 5 & 8 \end{bmatrix} \quad (\text{i})$$

(ج) ویژه مقدارهای ماتریس سه قطری تبدیل یافته را با دقت ممکن محاسبه کنید.

۲۱. گیریم $\{p_n(x) | n \geq 0\}$ یک خانواده از چند جمله‌یهای متعام نسبت به یک تابع وزن $\omega(x)$ بر یک بازه $a < x < b$ باشد. به علاوه فرض کنید در این چند جمله‌یها ضریب جمله

پیشرو ۱ باشد:

$$p_n(x) = x^n + \sum_{j=0}^{n-1} a_{n,j} x^j$$

یک ماتریس متقارن سه قطری پیدا کنید که $p_n(\lambda)$ چندجمله‌یی مشخصه آن باشد بنابراین محاسبه ریشه‌های یک چندجمله‌یی متعامد (وگره‌های فرمول تربیع گاوس) به حل یک مسأله ویژه مقدار برای ماتریسهای سه قطری متقارن بدل می‌شود.

راهنمایی: رابطه بازگشتی سه تایی برای $\{p_n(x)\}$ را به خاطر بیاورید و آن را با فرمول (۳.۴.۹) مقایسه کنید.

۲۲. روش QR را (الف) بدون انتقال، (ب) با انتقال، برای محاسبه ویژه مقدارهای ماتریسهای زیر به کار برید

$$\begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 2 \end{bmatrix} \quad (\text{ج}) \quad \begin{bmatrix} 2 & 1 & 0 & 0 \\ 1 & 2 & 1 & 0 \\ 0 & 1 & 2 & 1 \\ 0 & 0 & 1 & 2 \end{bmatrix} \quad (\text{ب}) \quad \begin{bmatrix} 3 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 1 \end{bmatrix} \quad (\text{الف})$$

۲۳. فرض کنید A یک ماتریس هسنبرگ باشد و فرض کنید در تجزیه $A = QR$ متعامد Q و R بالامتثلی باشد.

(الف) با یادآوری بحثی که در پی (۵.۵.۹) آمده؛ نشان دهید که (۷.۵.۹) برقرار است.
(ب) نشان دهید که نتیجه (۷.۵.۹) مستلزم شکلی برای H_k در (۶.۵.۹) است که Q در آن یک ماتریس هسنبرگ است.

(ج) نشان دهید که حاصلضرب یک ماتریس هسنبرگ و یک ماتریس بالامتثلی از هر مرتبه، باز هم یک ماتریس هسنبرگ است. با ترکیب کردن این نتایج نشان دهید که RQ باز هم یک ماتریس هسنبرگ است، همان گونه که در بند پس از (۷.۵.۹) ادعا شده بود.

۲۴. در ماتریس A در مسأله ۴، $\lambda = 7,9329$ و $\lambda = 5,6689$ دو ویژه مقدار تقریبی دیگر هستند. بارست معکوس را برای محاسبه ویژه بردارهای متناظر آنها به کار برید.

۲۵. برنامه‌های محاسبه ویژه مقدارهای یک ماتریس متقارن حقیقی در مرکز کامپیوتران را بررسی کنید. از چنین برنامه‌ای برای محاسبه ویژه مقدارهای ماتریسهای H_n هیلبرت برای $n = 3, 4, 5, 6, 7$ استفاده کنید. برای بررسی صحت جوابهای خود گزگوری^۱ و کارنی^۲ (صص ۶۶-۷۳) را ببینید.

۲۶. محاسبه ویژه مقدارها و ویژه تابعهای $x(t)$ ی متناظر با آنها را که در

$$\int_0^1 \frac{x(t)dt}{1+(s-t)^2} = \lambda x(s) \quad 0 \leq s \leq 1$$

صدق می کنند در نظر بگیرید. یک راه یافتن ویژه مقدارهای تقریبی گسسته سازی معادله با استفاده از انتگرالگیری عددی است. فرض کنید $h = 1/n$ به ازای مقداری از $n \geq 1$ و تعریف کنید $t_j = (j - 1/2)h$, $j = 1, \dots, n$. با گذاردن t_i به جای s در این معادله و تقریب زدن انتگرال با استفاده از روش انتگرالگیری عددی میانگامی دستگاه زیر به دست می آید:

$$h \sum_{j=1}^n \frac{\hat{x}(t_j)}{1+(t_i-t_j)^2} = \lambda \hat{x}(t_i) \quad i = 1, \dots, n$$

که در آن $\hat{x}(s)$ تابعی است که انتظار داریم $x(s)$ را تقریب بزند. این دستگاه مسأله ویژه مقدار برای یک ماتریس متقارن از مرتبه n است. بزرگترین دو ویژه مقدار این ماتریس را به ازای $n = 2, 4, 8, 16, 32$ پیدا کنید. همگرایی این ویژه مقدارها را هنگامی که n افزایش می یابد امتحان و سعی کنید خطا را در دقیقترین حالت ($n = 32$)، با مقایسه با ویژه مقدارهای واقعی مجهول در معادله انتگرال، پیش بینی کنید.

۲۷. نشان دهید که معکوس تعمیم یافته A^+ در (۱۱.۷.۹) در شرایط مور-پنروزا زیر صدق می کند.

$$(AA^+)^T = AA^+ \quad ۳ \quad AA^+A = A \quad ۱$$

$$(A^+A)^T = A^+A \quad ۴ \quad A^+AA^+ = A^+ \quad ۲$$

همچنین نشان دهید:

$$(AA^+)^2 = AA^+ \quad ۶ \quad (A^+A)^2 = A^+A \quad ۵$$

شرایط ۳ تا ۶ نشان می دهند که AA^+ و A^+A به ترتیب معرف تصویرهای قائم بر \mathbf{R}^m و \mathbf{R}^n هستند.

۲۸. برای ماتریس $m \times n$ دلخواه A ، نشان دهید که به ازای $\alpha > 0$

$$\lim_{\alpha \rightarrow 0^+} (\alpha I + A^T A)^{-1} A = A^+$$

راهنمایی: از تجزیه تکین - مقدار ماتریس A استفاده کنید.

۲۹. برخلاف وضعیت ماتریسهای مربع ناتکین، نیاز نیست وارون تعمیم یافته A^+ ، با تغییرات A ، به طور پیوسته تغییر کند. برای تأیید این مطلب یک خانواده $\{A(\epsilon)\}$ بیابید به طوری که ماتریس $A(\epsilon)$ به $A(0)$ همگرا باشد ولی $A(\epsilon)^+$ به $A(0)^+$ همگرا نباشد.

حل دستگاههای خطی به روش کمترین مربعات ۷۴۷

۳۰. برازش کمترین مربعات چندجمله‌یی خطی را برای داده‌های زیر حساب کنید. نمودار داده‌ها و برازش کمترین مربعات را رسم کنید. همچنین خطای ریشه دوم میانگین مربع را در برازش کمترین مربعات بیابید.

t_i	b_i	t_i	b_i	t_i	b_i
-۱٫۰	۱٫۰۳۲	-۰٫۳	۱٫۱۳۹	۰٫۴	-۰٫۴۱۵
-۰٫۹	۱٫۵۶۳	-۰٫۲	۰٫۶۴۶	۰٫۵	-۰٫۱۱۲
-۰٫۸	۱٫۶۱۴	-۰٫۱	۰٫۴۷۴	۰٫۶	-۰٫۸۱۷
-۰٫۷	۱٫۳۷۷	۰٫۰	۰٫۴۱۸	۰٫۷	-۰٫۲۳۴
-۰٫۶	۱٫۱۷۹	۰٫۱	۰٫۰۶۷	۰٫۸	-۰٫۶۲۳
-۰٫۵	۱٫۱۸۹	۰٫۲	۰٫۳۷۱	۰٫۹	-۰٫۵۳۶
-۰٫۴	۰٫۹۱۰	۰٫۳	۰٫۱۸۳	۱٫۰	-۱٫۱۷۳

۳۱. یک برازش مربعات درجه دوم برای داده‌های زیر درست کنید، از شکل استانده

$$g(t) = x_1 + x_2 t + x_3 t^2$$

و معادله نرمال ۲۶.۷.۹ استفاده کنید. ضریب وضعیت $A^T A$ چیست؟

t_i	b_i	t_i	b_i	t_i	b_i
-۱٫۰	۷٫۹۰۴	-۰٫۳	۰٫۳۳۵	۰٫۴	-۰٫۷۱۱
-۰٫۹	۷٫۴۵۲	-۰٫۲	-۰٫۲۷۱	۰٫۵	۰٫۲۲۴
-۰٫۸	۵٫۸۲۷	-۰٫۱	-۰٫۹۶۳	۰٫۶	۰٫۶۸۹
-۰٫۷	۴٫۴۰۰	۰٫۰	-۰٫۸۴۷	۰٫۷	۰٫۸۶۱
-۰٫۶	۲٫۹۰۸	۰٫۱	-۱٫۲۷۸	۰٫۸	۱٫۳۵۸
-۰٫۵	۲٫۱۴۴	۰٫۲	-۱٫۳۳۵	۰٫۹	۲٫۶۱۳
-۰٫۴	۰٫۵۸۱	۰٫۳	-۰٫۶۵۶	۱٫۰	۴٫۵۹۹

۳۲. برای ماتریس A که در برازش منحنی کمترین مربعات مسأله ۳۱ به دست آمد تجزیه آن به QR و تجزیه تکین - مقدار و همچنین وارون تعمیم‌یافته آن را به دست آورید. و از آنها برای حل مسأله کمترین مربعات استفاده کنید.

۳۳. تجزیه QR ، تجزیه تکین مقدار و وارون تعمیم یافته ماتریسهای زیر را به دست آورید. همچنین $\text{cond}(A)_2$ را نیز پیدا کنید.

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 3 & 4 \\ 3 & 4 & 5 \\ 4 & 5 & 6 \end{bmatrix} \quad (\text{ب}) \quad A = \begin{bmatrix} 0.9 & 1.1 \\ -1.0 & -1.0 \\ 1.1 & 0.9 \end{bmatrix} \quad (\text{الف})$$

۳۴. (الف) فرض می‌کنیم A ماتریسی $m \times n$ ، $m \geq n$ باشد و ستونهای b تقریباً وابسته باشند. دقیقتر بگوییم فرض کنید $A = [u_1, \dots, u_n]$ ، $u_j \in \mathbf{R}^m$ و فرض کنید بردار

$$v = \alpha_1 u_1 + \dots + \alpha_n u_n$$

در مقایسه با $\|\alpha\|_2$ ، $\alpha = [\alpha_1, \dots, \alpha_n]^T$ بسیار کوچک باشد. نشان دهید که A ضریب وضعیت بزرگی خواهد داشت.

(ب) برخلاف قسمت (الف) فرض می‌کنیم ستونهای A یکا متعامد باشند نشان دهید که $\text{cond}(A)_2 = 1$.

پیوستها

نرم افزار ریاضی

در آغاز دهه ۱۹۷۰، چند برنامه پیشرفته رایانه‌یی در موضوعهای اصلی آنالیز عددی نوشته شده‌اند. این برنامه‌ها، برای هر نوع رایانه کارا، قابل اعتماد و قابل اجرا هستند. این موضوعهای آنالیز عددی و علوم رایانه‌یی اکنون نرم افزار ریاضی نامیده می‌شوند، اصطلاحی که از زبان رایس^۱ (۱۹۷۱) متداول شده است. اکنون بسته‌هایی برنامه‌یی با کیفیت عالی برای بیشتر موضوعهای آنالیز عددی وجود دارند. اغلب این برنامه‌ها در دسترس عموم هستند و اغلب جزئی از کتابخانه‌های آنالیز عددی تجاری شده‌اند. برای یک مطالعه کلی و گزارش تاریخی این بسته‌های آنالیز عددی کاول^۲ (۱۹۸۴) را ببینید و برای سایر کتبی که بحث در این مورد را به طریق قابل توجهی انجام داده‌اند رایس (۱۹۷۷)، (۱۹۸۱) و (۱۹۸۳) را ملاحظه کنید. در کتاب فورسایت و دیگران (۱۹۷۷) چند برنامه در مسائل مختلف آنالیز عددی دقیقاً انتخاب شده‌اند که کارایی زیادی دارند و نسبتاً به سادگی قابل فهم‌اند.

دو دسته کتابخانه‌های تجاری بزرگی در آنالیز عددی برای رایانه‌های بزرگ، کتابخانه‌های IMSL و کتابخانه NAG وجود دارند که هر دو در همه خطوط رایانه‌یی، از جمله ریزرایانه‌ها، در دسترس هستند. در این کتابخانه‌ها خیلی از بسته‌های خاصی که در پاراگراف بعدی به آنها اشاره خواهد شد وجود دارند. دو شرکت IMSL و INAG همچنین متنهای خاص دیگری برای استفاده در خط لوله رایانه‌های برداری به وجود آورده‌اند. کتابخانه‌های تجاری آنالیز عددی دیگری نیز وجود دارند، ولی معمولاً در انواع کمتری از رایانه‌ها به خدمت گرفته می‌شوند.

ما قویاً به دانشجویان و پژوهشگران توصیه می‌کنیم که ابتدا به این کتابخانه‌های تجاری برای برنامه مورد نیاز خود مراجعه کنند. برنامه‌ها خوب نوشته شده‌اند، توجه خاصی برای سهولت استفاده، دقت و کنترل خطا و کارایی آنها مبذول شده است. این برنامه‌ها بطور مداوم بهتر می‌شوند تا آخرین پژوهشها را در روشهای عددی و تغییرات سخت افزاری رایانه‌ها منعکس کنند. در مقایسه با بسته‌های نرم‌افزاری خاصی که در ذیل به آنها اشاره می‌کنیم، این کتابخانه‌های تجاری بزرگ از نظر کاربرد به مراتب ساده‌ترند و پیاده‌سازی آنها را کارمندان مرکز محاسبات شما می‌توانند انجام دهند بی‌آنکه نیاز به پرداختن به زمینه‌های تخصصی داشته باشند. بسته‌های خاصی که بعداً توضیح داده می‌شوند فقط باید در مواردی به کار روند که کتابخانه‌های تجاری نیاز کاربر را برآورده نمایند. اغلب این بسته‌های خاص، ثابت‌اند و برای زمان خاصی نوشته شده‌اند و به روز در نیامده‌اند که تغییرات در الگوریتمها و رایانه‌ها را منعکس کنند.

در صورت امکان و در صورت دسترسی به بسته‌های استاندارد، ما نباید برای مسائل خود

برنامه‌نویسی کنیم. نوشتن یک برنامه خوب رایانه‌یی بسیار وقت‌گیر و مشکل است، و غیر محتمل است که یک کاربر معمولی بتواند برنامه‌نویسی خود را به پای برنامه کتابخانه‌های تجاری استاندارد برساند. کتابخانه IMSL به سه مؤلفه اصلی تقسیم شده است. (۱) کتابخانه آنالیز عددی، (۲) کتابخانه آمار و (۳) کتابخانه توابع خاص. هر کتابخانه از زیربرنامه‌هایی به زبان فورترن ۷۷ تشکیل یافته است که از برنامه ارائه شده کاربر فراخوانی می‌شوند. به علاوه، IMSL، بسته‌های اضافی نیز دارد که استفاده از کتابخانه اصلی را در سطح بسیار بالاتری به شکل محاوره‌یی، برای کاربر، ممکن می‌سازد. IMSL همچنین به عنوان یک عامل توزیع چند برنامه برای عموم نیز عمل می‌کند. نشانی IMSL چنین است:

IMSL , Inc.
NBC Building
7500 Bellaire Boulevard
Houston, Texas 77036-5085

از آنجا که کتابخانه‌های آنها به طور مداوم بهنگام می‌شوند، برای اطلاع از موضوعات روز کتابخانه‌ها، می‌توانید با IMSL مکاتبه کنید.

کتابخانه NAG تقریباً همان موضوعات کتابخانه IMSL را دارد و بسته‌های کمکی نیز برای ساده کردن استفاده از کتابخانه اصلی را دارا می‌باشد. برنامه‌های این کتابخانه هم به زبان فورترن و هم به زبان Algol در دسترس اند^۱. برنامه‌های کتابخانه را دانشمندان دانشگاه‌های مختلف انگلستان نوشته‌اند و بسیاری از بسته‌های نرم‌افزاری خاص که بعداً توضیح داده خواهند شد جزئی از کتابخانه NAG شده‌اند. برای اطلاعات بیشتر با آدرس زیر مکاتبه کنید:

Numerical Algorithm Group, Ltd.
NAG Central office
Mayfield House,
256 Bandburg Road
Oxford ox2 7DE, United Kingdom.

بسته‌های رایانه‌یی زیادی با کیفیت بالا برای مسائل زمینه‌های خاص از ۱۹۷۲ نوشته شده‌اند. اولین، و شاید معروفترین آنها، بسته EISPACK برای حل مسائل ویژه مقادیر ماتریسها است. پس از آن بسته‌هایی در انواع گوناگون، برای بیشتر زمینه‌های اصلی آنالیز عددی، تهیه شده‌اند،

۱. از آنجا که این کتاب در دهه ۱۹۸۰ نوشته شده است به زبانهای پاسکال و C اشاره نشده است.

و بسیاری از آنها تاکنون در کتابهای درسی گنجانده شده‌اند. بیشتر این نرم‌افزارها در کتاب کاول (۱۹۸۲) به نحو مطلوبی وارد شده‌اند. ما اکنون بسته‌هایی را که از آنها اطلاع داریم ذیلاً می‌آوریم و از اینکه بعضی را از قلم انداخته‌ایم پوزش می‌خواهیم. در بعضی موارد منبع برنامه‌ها را نیز ذکر می‌کنیم. به علاوه، تقریباً تمام این برنامه‌ها، را (بویژه آنهایی را که در آزمایشگاههای ملی ایالات متحده، تولید شده‌اند)، با محدودیتهایی، می‌توان از مؤسسه زیر دریافت کرد:

National Energy Software Center
Argonne National Laboratory
9700 South Cass Avenue
Argonne, Illinois 60439

۱. **EISPACK** این بسته نرم‌افزاری مسأله ویژه مقدار ماتریس را برای انواع گوناگون ماتریسها حل می‌کند، همچنین شامل برنامه‌هایی برای مسأله ویژه مقدار تعمیم‌یافته $Ax = \lambda Bx$ و برای محاسبه تجزیه ماتریس است. راهنمای کاربران برای این بسته در اسمیت و همکاران و (۱۹۷۶) و در گاریو و همکاران (۱۹۷۷) داده شده است. برای شرح توسعه این بسته، کاول (۱۹۸۴)، فصل ۴) را ببینید. IMSL یکی از توزیع‌کنندگان این بسته نرم‌افزاری است.

۲. **LINPACK** این بسته نرم‌افزاری برای حل دستگاههای معادلات خطی است. این برنامه چهار متن دارد: حساب حقیقی و مختلط با دقت ساده و مضاعف. علاوه بر برنامه‌های معمولی برای انواع مختلف دستگاههای ناتکین مربعی، LINPACK حاوی برنامه‌هایی برای حل دستگاههای خطی کمترین مربعات نیز می‌باشد. برای راهنمای کاربران دونگارا و همکاران (۱۹۷۹) را ببینید، و برای شرح توسعه LINPACK، کاول (۱۹۸۴)، فصل ۲) را ملاحظه کنید. IMSL یکی از توزیع‌کنندگان این بسته نرم‌افزاری است.

۳. **MINPACK** این بسته نرم‌افزاری برای حل دستگاههای معادلات غیرخطی و مسائل بهینه‌سازی است. در حال حاضر فقط متن اول این بسته وجود دارد که برای (الف) حل n معادله غیرخطی n مجهولی و (ب) برای حل کمترین مربعات دستگاههای معادلات فراتامعین غیرخطی، به کار می‌رود. متن دیگر قرار است شکلهای مختلف دیگر مسائل بهینه‌سازی مقید و نامقید را در بر گیرد. برای راهنمای کاربران، موره^۱ و همکاران (۱۹۸۰) و برای شرح پروژه کاول (۱۹۸۴)، فصل ۵) را ببینید. IMSL یکی از توزیع‌کنندگان این بسته نرم‌افزاری است.

۴. **LLSQ** این یک بسته نرم‌افزاری برای حل کمترین مربعات دستگاه‌های معادلات خطی است. این برنامه همراه کتاب لاوسن و هسن آمده و از **IMSL** می‌توان آن را تهیه نمود.

۵. **QUADPACK** این بسته نرم‌افزاری برای انتگرالگیری عددی توابع یک متغیره است. هم برنامه‌های کلی و هم برنامه‌هایی برای انتگرالهایی با شکل خاص در آن وجود دارند. برای راهنمای کاربران، پیسنز و همکاران (۱۹۸۳) را ببینید.

۶. **DEPAC** این بسته برنامه‌هایی است برای حل دستگاه‌های معادلات دیفرانسیل معمولی. در حال حاضر شامل سه برنامه اصلی است: یک برنامه متغیر-مرتبه ادامه، یک برنامه ثابت مرتبه رونگه-فلبرک، و یک برنامه متغیر مرتبه برای مسائل سرسخت بر پایه فرمولهای مشتقگیری پسرو. این برنامه‌ها در آزمایشگاه ملی سندیا^۱ در شهر آلبوکرک^۲ در نیومکزیکو^۳ موجودند. برای بحث در مورد تولید نرم‌افزاری در معادلات دیفرانسیل معمولی کاول (۱۹۸۴، فصل ۶) را ملاحظه کنید.

۷. **GRDI** این یک بسته برنامه‌هایی برای حل معادلات دیفرانسیل معمولی و جزئی است که این آخری براساس روش خطوط پایه‌گذاری شده است. برنامه‌هایی برای معادلات دیفرانسیل معمولی بر پایه روشهایی است که هایندمارش و همکاران او در آزمایشگاه لاورنس لیورمور^۴ ایجاد کرده‌اند. برای کسب اطلاعات یا دریافت این بسته نرم‌افزاری با آدرس زیر مکاتبه کنید:

Dr. W. E. Schiesser

Whitaker Laboartory No.5

Lehigh University

Bethlehen, Pennsylvania 18015

۸. **FISHPAK** این یک بسته نرم‌افزاری برای حل تقریبهای تفاضلات متناهی برای بعضی از دسته‌های معادلات دیفرانسیل جزئی تفکیک‌پذیر معروف (مثلاً معادله پواسن و معادله هلمهولتز) در نواحی خاص، مانند مستطیله‌ها یا دوایر، است. این بسته برای محاسبات مکانیک سیالات در **NCAR** (مرکز ملی برای تحقیقات جوی) تولید شده است. برای بحث در مورد این بسته اسوارتس تراوبر^۵ و سویت^۶ (۱۹۷۹) را ببینید. برای یک آشنایی کلیتر با نرم‌افزار حل معادلات دیفرانسیل جزئی کاول (۱۹۸۴، فصل ۹) را ملاحظه کنید. این کتاب شامل بحثی است در **ELLPACK** [برای معادلات دیفرانسیل جزئی بیضوی، رایس (۱۹۷۷، صص ۳۱۹-۳۴۱) را ببینید] و **ITPACK**

1. Sandia

2. Albuquerque

3. New Mexico

4. Lawronce Livormore

5. Swarztrauber

6. Sweet

[برای روشهای بارستی حل معادلات تفاضلات متناهی، هاگمن و یانگ (۱۹۸۱) را ببینید]، می باشد.

علاوه بر بسته‌های نرم‌افزاری فوق‌الذکر، بسیاری از برنامه‌ها و الگوریتمها در ACM Transactions on Mathematical Software انتشار یافته‌اند. برای لیست بعضی از الگوریتمهای موجود به یک شماره این مجله مراجعه کنید. برنامه‌های اصلی را می‌توان از IMSL تهیه نمود.

یک مبادله نرم‌افزاری الکترونیکی یک وسیله به دست آوردن بسیاری از نرم‌افزارهای فوق‌الذکر توسط پست الکترونیکی مستقر در آزمایشگاه ملی آرگون در شیکاگو، ایلینویز، است. این سیستم NETLIB خوانده شده و می‌توان از طریق آدرس پست الکترونیکی زیر با آن تماس برقرار نمود

NETLIB@AML-MCS on ARPANET

پیام زیر را به این آدرس بفرستید تا راههای استفاده از NETLIB را دریافت دارید

send index

این سیستم در زمان نوشتن این کتاب وجود داشته است ولی هیچ تضمینی برای ادامه موجودیت آن وجود ندارد. یک ابزار بسیار مفیدی است برای به دست آوردن نرم‌افزار ریاضی با کیفیت بالا.

نرم‌افزار برای ریزرایانه‌ها. کتابخانه‌های IMSL و NAG، زیر مجموعه‌های خاصی از کتابخانه‌های خود را برای تعدادی از ریزرایانه‌ها از جمله ماشینهای سازگار با IBM شخصی، در دسترس دارند. به علاوه بسته‌های دیگری نیز وجود دارند که پیوسته توسعه می‌یابند. توجه خاص ما به بسته‌هایی است که دسترسی آسانتری به ابزارهای پیچیده آنالیز عددی را به وجود می‌آورند. دو بسته متداول از این نوع، PC-MATLAB و GAUSS هستند، که هر دو دسترسی به تعدادی از برنامه‌های موجود در بسته‌های EISPACK و LINPACK را همراه با تسهیلات آنالیز عددی و گرافیک دیگر مهیا می‌سازند. این نرم‌افزارها در اغلب ریزرایانه‌ها و ایستگاههای علمی در دسترس‌اند. برای اطلاعات در مورد این بسته‌ها با آدرس زیر مکاتبه کنید.

برای GAUSS:

APTECH System, Inc.

P. O. Box 6487

Kent, Washington 98064

The Math Works, Inc.
 24 Prime Park Way
 Natick, MA01760-1500
 Fax: (508) 653-2997
 E-mail: info@Mathworks.com

مراجع

- Cowell, W., Ed. (1984). *Sources and Development of Mathematical Software*. Prentice-Hall, Englewood Cliffs, N.J.
- Dongarra, J., J. Bunch, C. Moler, and G. Stewart (1979). *LINPACK User's Guide*. SIAM Pub., Philadelphia.
- Forsythe, G., M. Malcolm, and C. Moler (1977). *Computer Methods for Mathematical Computations*. Prentice-Hall, Englewood Cliffs, N.J.
- Garbow, B., J. Boyle, J. Dongarra, and C. Moler (1977). *Matrix Eigensystem Routines—EISPACK Guide Extension. Lecture Notes in Computer Science 51*. Springer-Verlag, New York.
- Hageman, L., and D. Young (1981). *Applied Iterative Methods*. Academic Press, New York.
- Lawson, C., and R. Hanson (1974). *Solving Least Squares Problems*. Prentice-Hall, Englewood Cliffs, N.J.
- Moré, J., B. Garbow, and K. Hillstom (1980). *User Guide for MINPACK-1*. Argonne National Laboratory Rep. ANL-80-74, Chicago, Ill.
- Piessens, R., E. deDoncker-Kapenga, C. Überhuber, and D. Kahaner (1983). *QUADPACK: A Subroutine Package for Automatic Integration*. Springer-Verlag, New York.
- Rice, J., Ed. (1971). *Mathematical Software*. Academic Press, New York.
- Rice, J., Ed. (1977). *Mathematical Software III*. Academic Press, New York.
- Rice, J. (1981). *Matrix Computations and Mathematical Software*. McGraw-Hill, New York.
- Rice, J. (1983). *Numerical Methods, Software, and Analysis*. McGraw-Hill, New York.
- Smith, B., J. Boyle, B. Garbow, Y. Ikebe, V. Klema, and C. Moler (1976). *Matrix Eigensystem Routines—EISPACK Guide*, 2nd ed., *Lecture Notes in Computer Science 6*. Springer-Verlag, New York.
- Swarztrauber, P., and R. Sweet (1979). Algorithm 541: Efficient Fortran subprograms for the solution of separable elliptic partial differential equations. *ACM Trans. Math. Softw.* 5, 352–364.

پاسخهای تمرینهای انتخابی

فصل ۱

۵. (الف)

$$p_{2n-1}(x) = \sum_{j=0}^{n-1} (-1)^j \frac{x^{2j}}{(2j+1)(j!)}$$

$$| \text{خطا} | \leq \frac{|x|^{2n}}{(2n+1)(n!)}$$

$$p_2(x, y) = p_1(x+y) - \frac{1}{4}x^2 + \frac{1}{4}xy - \frac{1}{8}y^2, p_1(x, y) = 1 + x - \frac{1}{4}y \quad \text{الف) ۷}$$

$$۱۰. \text{الف) } ۲۱,۶۲۵ \quad \text{ج) } \frac{2}{3} \quad \text{د) } \frac{2}{3}$$

$$۱۱. \text{ب) } ۱۱۱۱۱۱$$

$$۱۲. \text{الف) } ۱۱۰۱۰ \quad \text{ج) } ۰۰۰۰۱۱۰۰۱۱۰۰ \dots$$

$$۲۱. \text{الف) } ۴ \quad \text{ب) } ۲ \quad \text{ج) } ۳$$

$$۲۲. \text{الف) } [۲,۰۵۲۶۵, ۲,۰۵۳۷۵]$$

$$\left[\frac{۸,۴۷۲۵}{۰,۰۶۴۵}, \frac{۸,۴۷۳۵}{۰,۰۶۳۵} \right] \doteq [۱۳۱,۳۵۶۶, ۱۳۳,۴۴۰۹] \quad \text{د) ۱۳}$$

۲۶. ب) برای $x \doteq 0$ تقریب سری درجه دوم تیلر عبارت است از

$$f(x) \doteq \frac{-5}{24} - \frac{11}{48}x - \frac{379}{1920}x^2$$

همچنین $f(0)$ دقیقاً برابر است با $\frac{-5}{24}$.

۲۸. (الف) $0.0030 -$ (ج) 0.00068

۳۰. (الف) با $x_A = 3.14$ و $y_A = 2.685$ خطا تقریباً برابر است با $y_T - y_A = -1.71(x_T - x_A) - 1.46(x_T - x_A) - 1.71(y_T - y_A)$ و x_T اعداد واقعی گردننده مربوط به x_A و y_A هستند همچنین $0.0098 \leq$ خطای نسبی $|$

فصل ۲

۱. $\prod_{j=0}^{\infty} [1 + r^{2^j}] = 1/(1 - r)$ و حاصلضرب نامتناهی همگراست اگر و تنها اگر $|r| < 1$.

۳. (الف) ریشه برابر است با 4.493409458

۴. (ب) 1.8392868 (د) 1.1284251

۵. (الف) 4.493409458 (ب) 98.95006282

۶.

B	ریشه
۱	-0.5884017765
۵	-0.4049115482
۱۰	-0.3265020101
۲۵	-0.2374362439
۵۰	-0.1832913333

۱۲. (ه) $2^{-15} 10^{-16} \approx 3.05 \times 10^{-21}$ $|\text{Rel}(x_r)| \leq$

۱۳. بارست x_r به اندازه کافی دقیق خواهد بود.

۱۵. (الف) همگرا نیست. (ب) همگراست.

۱۹. $[a, b] = [1, 1 + \frac{\pi}{4}]$, $a = 2.1322677$

۲۴. (الف) واگراست. (ج) همگراست.

۲۵. $\text{Limit}_{n \rightarrow \infty} \frac{\sqrt{a} - x_{n+1}}{(\sqrt{a} - x_n)^3} = \frac{1}{4a}$

۲۹. (الف) نرخ ≈ 0.625

۳۹. (ب) $2.470638970 + 4.640533162i$

۴۰. (الف) این چندجمله‌یی، چندجمله‌یی لژاندار از درجه ۱۲ بر بازه $1 \leq x \leq -1$ است.

ریشه‌های آن عبارت‌اند از:

± 0.1252334085

± 0.7699026742

$$\begin{aligned} \pm ۰٫۳۶۷۸۳۱۴۹۹۰ & \quad \pm ۰٫۹۰۴۱۱۷۲۵۶۴ \\ \pm ۰٫۵۸۷۳۱۷۹۵۴۳ & \quad \pm ۰٫۹۸۱۵۶۰۶۳۴۲ \end{aligned}$$

.۵۱

(x_0, y_0)	روش همگراست به
(۱٫۲, ۲٫۵)	(۱٫۳۳۶۳۵۵۳۷۷, ۱٫۷۵۴۲۳۵۱۹۸)
(-۲, ۲٫۵)	(-۰٫۹۰۱۲۶۶۱۹۰۸, -۲٫۰۸۶۵۸۷۵۹۵)
(-۱٫۲, -۲٫۵)	(-۰٫۹۰۱۲۶۶۱۹۰۸, -۲٫۰۸۶۵۸۷۵۹۵)
(۲٫۰, -۲٫۵)	(-۳٫۰۰۱۶۲۴۸۸۷, ۰٫۱۴۸۱۰۷۹۹۵۰)

در آخرین حالت، این روش ظاهراً به طور تصادفی بالا و پایین پرش دارد. و بعد از ۱۶ بارست به ریشه‌ای که قبلاً داده شده همگرا می‌شود. معادلات حتماً ریشه دیگری دارند.

۵۲. (الف) دو ریشه وجود دارد (a, b) و (b, a) با $a = ۰٫۶۷۳۰۰۷۱۷۰$ و $b = ۱٫۹۴۵۰۲۶۸۲$

فصل ۳

.۴

$$\begin{aligned} | \text{خطای درونیابی} | & \leq \left(\frac{h^2}{8}\right)e^2 \doteq ۹٫۲ \times ۱۰^{-۵} \\ | \text{خطای گردکردن} | & \leq ۵ \times ۱۰^{-۵} \\ | \text{خطای کل} | & \leq ۱٫۴۲ \times ۱۰^{-۴} \end{aligned}$$

۷. h را برابر $۱۰^{-۷}$ انتخاب و فرض کنید درایه‌های جدول هفت رقم بامعنی دارند. خطای کل با $۱۰^{-۷} \times ۸٫۴$ کراندار می‌شود.

۱۰. درایه‌های جدول را با ۶ رقم بامعنی انتخاب کنید. یک افزاز ممکن $۱ \leq x \leq ۱۰$ با اندازه شبکه پیشنهادی و خطای درونیابی کلی حاصل در جدول زیر نشان داده شده است.

بازه	h	خطای کلی
$1 \leq x \leq 2$	$۰٫۱$	$۶٫۸۱ \times ۱۰^{-۷}$
$2 \leq x \leq 3$	$۰٫۲۵$	$۷٫۳۵ \times ۱۰^{-۷}$
$3 \leq x \leq 6$	$۰٫۵$	$۸٫۸۵ \times ۱۰^{-۷}$
$6 \leq x \leq 10$	۱	$۸٫۸۵ \times ۱۰^{-۷}$

.۲۰. درجه ۳

۲۲. ششمین ورودی ۴۱۹۳۲۷ باید به ۴۱۹۲۳۷ یا ۴۱۹۲۳۸ تغییر یابد. دو خطای دیگر وجود دارند و اثر آنها در تفاضلات مرتبه بالاتر روی یکدیگر می‌افتند.

۲۹. شرط لازم و کافی برای اینکه چندجمله‌یی درونیاب با درجه نایبتر از ۲ یکتا باشد

$$x_1 \neq \frac{1}{2}(x_0 + x_2)$$

۳۴. $s(x)$ بر بازه $[x_0, x_2]$ با فرمول زیر داده شده است:

$$s(x) = \frac{(x_2 - x)^2 M_0 + (x - x_0)^2 M_2}{6(h_0 + h_2)} + \frac{(x_2 - x)y_0 + (x - x_0)y_2}{h_0 + h_2} - \frac{h_0 + h_2}{6} [(x_2 - x)M_0 + (x - x_0)M_2]$$

$s(x_1) = y_1$ را تشکیل دهید تا یک شرط درونیاب در x_1 برقرار کنید، که یک معادله برای M_2 و M_0 خواهد داد. شکل مشابهی برای $s(x)$ بر بازه $[x_{n-2}, x_n]$ با رابطه $s(x_{n-1}) = y_{n-1}$ به‌کار برده می‌شود.

۳۶. الف) برای $n = 96$ زیربازه، خطای $s(x)$ برای شرایط مرزی مختلف عبارت‌اند از:

$$\text{الف) } 1.0 \times 10^{-35} \quad \text{ب) } 8.32 \times 10^{-11} \quad \text{ج) } 1.0 \times 10^{-25}$$

$$B_i^{(m)}(x) = \frac{x_{i+m} - x}{x_{i+m} - x_{i+1}} B_{i+1}^{(m-1)}(x) - \frac{x - x_i}{x_{i+m-1} - x_i} B_i^{(m-1)}(x) \quad ۳۹$$

$$۴۳. \text{ الف) } d_0 = 1; d_k = 0 \text{ به‌ازای } 1 \leq k \leq m - 1$$

فصل ۴

۱. چندجمله‌یی درجه ۴ برنشتاین عبارت است از:

$$p_4(x) = 2\sqrt{2}x(1-x)^2 + 6x^2(1-x)^2 + 2\sqrt{2}x^2(1-x)$$

چندجمله‌یی درجه ۴ تیلر عبارت است از:

$$f_4(x) = 1 - \frac{\pi^2}{2} \left(x - \frac{1}{2}\right)^2 + \frac{\pi^4}{24} \left(x - \frac{1}{2}\right)^4$$

جدول زیر مقادیر $p_4(x)$ و $f_4(x)$ را در چندین نقطه همراه با خطاها هنگامی که با $\sin(\pi x)$

مقایسه می‌شوند نشان می‌دهد.

x	$p_f(x)$	خطا	$f_f(x)$	خطا
۰٫۰	۰٫۰	۰٫۰	۰٫۲۰۰	$-۲٫۰E-۲$
۰٫۱	۰٫۲۵۷۳	۰٫۰۵۱۷	۰٫۳۱۴۳	$-۵٫۳E-۳$
۰٫۲	۰٫۴۶۱۳	۰٫۱۲۶۵	۰٫۵۸۸۷	$-۹٫۶E-۴$
۰٫۳	۰٫۶۰۹۱	۰٫۱۹۹۹	۰٫۸۰۹۱	$-۸٫۵E-۵$
۰٫۴	۰٫۶۹۸۶	۰٫۲۵۰۵	۰٫۹۵۱۱	$-۱٫۳E-۶$
۰٫۵	۰٫۷۲۸۶	۰٫۲۷۱۴	۱٫۰۰۰۰	۰٫۰

۳. (ب) فرض کنید

$$S = \sum_{j=1}^{\infty} \frac{(-1)^j x^{2j}}{j^2}$$

بازه همگرایی این سری $|x| \leq 1$ برای خطا

$$\left| S - \sum_{j=1}^n (-1)^j \cdot \frac{x^{2j}}{j^2} \right| \leq \frac{|x|^{2n+2}}{(n+1)^2}$$

اگر $n \geq 316$ ، خطا برای همه مقادیر $|x| \leq 1$ از 10^{-5} کمتر است. اما برای $|x| < 1$ احتمالاً بسیار بزرگ است، زیرا عامل $|x|^{2n+2}$ برای تخمین خطا به حساب نیامده است.

۶. تقریب پادا عبارت است از $R(x) = (1 + \frac{1}{4}x)/(1 - \frac{1}{4}x)$ در خطای

$$e^x - R(x) = -\left[\frac{1}{12}x^2 + \frac{1}{12}x^4 + \frac{13}{240}x^6 + \dots \right]$$

و خطای سری تیلر $e^x - p_2(x) = \frac{1}{6}x^3 + \frac{1}{24}x^4 + \dots$ برای مقدار کوچک $|x|$ خطای $R(x)$ در حدود نصف خطای $p_2(x)$ است. اما برای $|x| = 1$ خطای $R(x)$ بزرگتر است.

$$10. \text{ (الف) } \ln(x) - p_1 \|_{\infty} = 0.072; \quad \ln(x) - p_1 \|_{\infty} = 0.072$$

$$\text{(ب) } q_1^*(x) = -0.6633 + 0.6931x; \quad \ln(x) - q_1^* \|_{\infty} = 0.030$$

۱۳. مقدار مینیمم ساز α عبارت است از $\alpha = \sqrt{e}$. مقدار مینیمم برابر است با

$$e + 1 - 2\sqrt{e} \approx 0.42$$

$$20. \text{ (الف) } \psi_2(x) = x^2 - \frac{5}{4}x + \frac{11}{102}, \quad \psi_1(x) = x - \frac{1}{4}, \quad \psi_0(x) \equiv 1$$

$$28. \text{ (الف) } \rho_n(\sin x) \leq \frac{\pi^{n+1}}{(n+1)! 2^{n+2}}$$

$$31. \text{ (ج) } I_2(x) = 0.968706x - 0.187130x^2; \quad \|\tan^{-1} x - I_2\|_{\infty} = 0.00590$$

(د)

n	$\ f - I_n\ _\infty$	n	$\ f - I_n\ _\infty$
۱	۳٫۷۲E-۲	۵	۱٫۱۴E-۵
۲	۴٫۳۷E-۳	۶	۱٫۶۹E-۶
۳	۵٫۷۲E-۴	۷	۲٫۵۵E-۷
۴	۷٫۹۴E-۵	۸	۳٫۹۱E-۸

(ج).۳۲

$$F_T(x) = ۰٫۹۷۲۴۲۰x - ۰٫۱۹۱۸۹۸x^2; \|\tan^{-1} x - F_T\|_\infty = ۰٫۰۰۵۰۰$$

$$q_{n-1}^*(x) = p(x) - \frac{a_n}{2^{n-1}} T_n(x), \rho_{n-1}(p) = \frac{|a_n|}{2^{n-1}} \quad .۳۶$$

(د).۴۰

n	$\rho_n(f)$	n	$\rho_n(f)$
۱	۲٫۹۸E-۲	۵	۸٫۶۹E-۶
۲	۳٫۴۲E-۳	۶	۱٫۲۸E-۶
۳	۴٫۴۲E-۴	۷	۱٫۹۲E-۷
۴	۶٫۰۷E-۵	۸	۲٫۹۳E-۸

فصل ۵

۱.

$\int_0^1 e^{-x^2} dx$ (الف)			$\int_{-1}^1 \frac{dx}{1+x^2}$ (ج)		
n	I_n	R_n	n	I_n	R_n
۲	۰٫۷۳۱۳۷۰۲۵		۲	۲٫۲۳۵۲۹۴	
۴	۰٫۷۴۲۹۸۴۱۰		۴	۲٫۹۱۷۶۴۷	
۸	۰٫۷۴۸۶۶۵۶۱	۴٫۰۳	۸	۲٫۶۵۸۸۲۴	۵٫۰۹
۱۶	۰٫۷۴۶۵۸۴۶۰	۴٫۰۱	۱۶	۲٫۶۵۰۵۰۷	۳٫۱۱
۳۲	۰٫۷۴۶۷۶۴۲۵	۴٫۰۰	۳۲	۲٫۶۵۱۳۴۷	-۹٫۹۰
۶۴	۰٫۷۴۶۸۰۹۱۶	۴٫۰۰	۶۴	۲٫۶۵۱۵۶۳	۳٫۸۹
۱۲۸	۰٫۷۴۶۸۲۰۳۹	۴٫۰۰	۱۲۸	۲٫۶۵۱۶۱۷	۴٫۰۰
۲۵۶	۰٫۷۴۶۸۲۳۲۰	۴٫۰۰	۲۵۶	۲٫۶۵۱۶۳۱	۴٫۰۰

$\int_0^1 e^{-x^2} dx$ (الف)			$\int_{-2}^2 \frac{dx}{1+x^2}$ (ج)		
n	I_n	R_n	n	I_n	R_n
۲	۰٫۷۴۷۱۸۰۴۲۸۹۰۹۵		۲	۵٫۴۹۰۱۹۶۰۷۸۴	
۴	۰٫۷۴۶۸۵۵۳۷۹۷۹۱۰		۴	۲٫۴۷۸۴۳۱۳۷۲۵	
۸	۰٫۷۴۶۸۲۶۱۲۰۵۲۷۵	۱۱٫۱	۸	۲٫۵۷۲۵۴۹۰۱۹۶	-۳۲٫۰
۱۶	۰٫۷۴۶۸۲۴۲۵۷۴۳۵۷	۱۵٫۷	۱۶	۲٫۶۴۷۷۳۴۵۶۳۵	۱٫۲۵
۳۲	۰٫۷۴۶۸۲۴۱۴۰۶۰۷۰	۱۵٫۹	۳۲	۲٫۶۵۱۶۲۷۲۸۳۰	۱۹٫۳
۶۴	۰٫۷۴۶۸۲۴۱۳۳۲۹۹۷	۱۶٫۰	۶۴	۲٫۶۵۱۶۳۵۲۸۰۷	۴۸۷
۱۲۸	۰٫۷۴۶۸۲۴۱۳۲۸۴۲۹	۱۶٫۰	۱۲۸	۲٫۶۵۱۶۳۵۳۲۴۴	۱۸۳
۲۵۶	۰٫۷۴۶۸۲۴۱۳۸۱۴۳	۱۶٫۰	۲۵۶	۲٫۶۵۱۶۳۵۳۲۷۲	۱۶٫۰

۶. با افزایش n ، نسبتهای R_n به مقادیر ثابت میل می‌کند. مقادیر برای $n = ۱۶$ به صورت

زیرند:

α	۰٫۲۵	۰٫۵	۰٫۷۵	۱٫۰
R_n	۲٫۱۱	۲٫۵۲	۳٫۰۳	۳٫۶۴

$$I_h(f) = \frac{2}{3}h[f(0) + 3f(2h)] \quad .9$$

.۱۵

n	I_n	خطا
۲	۱٫۲۶۳۱۶	۱٫۳۸۸
۴	۲٫۰۴۷۲۹	۰٫۶۰۴۴
۶	۲٫۴۱۱۶۹	۰٫۲۳۹۹
۸	۲٫۵۶۰۰۸	۰٫۰۹۱۵۶
۱۰	۲٫۶۱۷۲۵	۰٫۰۳۴۲۸

$$\int_0^1 f(x) \ln\left(\frac{1}{x}\right) dx \doteq w_1 f(x_1) + w_2 f(x_2) \quad .۱۷$$

$$x_1 = \frac{15 - \sqrt{106}}{42} \doteq ۰٫۱۱۲۰۰۸۸۰۶۲ \quad x_2 = \frac{15 + \sqrt{106}}{42} \doteq ۰٫۶۰۲۲۷۶۹۰۸۱$$

$$w_1 = \frac{21}{\sqrt{106}} \left[x_2 - \frac{1}{4} \right] \doteq ۰٫۷۱۸۵۳۹۳۱۹ \quad w_2 = 1 - w_1 = ۰٫۲۸۱۴۶۰۶۸۱$$

۲۰. برای انتگرالگیری گاوسی، $\|f^{(r)}\|_{\infty} \leq (0.8^n/n^r)$ | $I - I_n$ | برای قاعده

$$|I - I_n| \leq (0.18/n^r) \|f^{(r)}\|_{\infty}, [-1, 1] \text{ سیمپسون بر بازه}$$

$$\sum_{n=1}^{\infty} \frac{1}{n^{0.75}} = 4.959511254 + E \quad 0 < -E < 1.34 \times 10^{-6} \quad 24$$

۲۶. با فرض $I - I_n \doteq c/n^p$ و با استفاده از I_{16} و I_{32} و I_{64} به دست

می‌آوریم $p = 3.44$ ، $c = 0.154$. برای برونیاب ایتکن، $I_{64} = 0.28571428586$ و

$$I - I_{64} \doteq \tilde{I}_{64} - I_{64} = 9.43 \times 10^{-9} \quad \text{برای اینکه داشته باشیم } |I - I_n| \leq 10^{-11} \text{ به}$$

$n \geq 469$ احتیاج خواهیم داشت. انتخاب منطقی استفاده از $n = 512$ است.

$$\tilde{I}_n = \frac{1}{2} [I_n^{(T)} + 2I_n^{(M)}] \quad 28 \text{ که قاعده سیمپسون خواهد بود.}$$

۳۸. (الف) انتگرال $I = \int_0^1 (f_1 + f_1^{r_n}) \cos(x) \ln(x) dx$ را بشکنید. با استفاده از سری تیلر

$$\int_0^1 \cos(x) \ln(x) dx \doteq - \left(1 - \frac{1}{2!9} + \frac{1}{4!25} - \frac{1}{6!49} + \frac{1}{8!81} \right)$$

$$\doteq -0.9460830727$$

که در آن خطا کمتر از $10^{-9} \times 2.28$ است. برای محاسبه انتگرال باقی مانده می‌توان از یک

روش استاندارد استفاده کرد. به عنوان مثال، انتگرالگیری رامبرگ با $n = 129$ گره به مقدار

$$0.5460781580 \text{ منجر می‌شود. که در آن خطا کمتر از } 10^{-9} \times 2.8 \text{ است. بنابراین انتگرال}$$

برابر است با 1.4921612307 که در آن خطا از $10^{-9} \times 5.1$ کمتر است.

(ب) از آنجایی که $|x^2 \sin(\frac{1}{x})| \leq x^2$ ، $\varepsilon > 0$ را چنان انتخاب کنید که

$$\left| \int_0^{\varepsilon} x^2 \sin\left(\frac{1}{x}\right) dx \right| \leq \int_0^{\varepsilon} x^2 dx \leq 0.0005$$

یک انتخاب مناسب $\varepsilon = 0.1$ است. سپس $\int_{0.1}^{2/\pi} (x^2) \sin(1/x) dx$ را با دقت حداقل 0.0005

محاسبه کنید.

$$D_h^{(r)} f(x) = \frac{f(x) - 2f(x+h) + f(x+2h)}{h^r} \quad \text{(الف) 41}$$

$$f''(x) - D_h^{(r)} f(x) \doteq h f^{(r)}(x)$$

فصل ۶

$$20 \text{ (د)} \quad 0 \text{ (ج)} \quad 1 \text{ (ب)} \quad 2 \text{ (الف)}$$

$$3 \text{ (ب)}$$

$$y_1(0) = 1 \quad y_1' = y_2$$

$$y_2(0) = 0 \quad y_2' = 0.1(1 - y_1^2)y_2 - y_1$$

۵. (ج) جواب واقعی $Y(x) = 2e^x + \sin(x) - \cos(x)$. بارستهای پیکار عبارتند از:

$$Y_0(x) = 1, Y_1(x) = 1 + x + 2\sin(x)$$

$$Y_2(x) = 3 + x + \frac{x^2}{2} + 2\sin(x) - 2\cos(x)$$

$$Y_3(x) = 3 + 3x + \frac{1}{2}x^2 + \frac{1}{6}x^3 - 2\cos(x)$$

برای $x = 0$ آنها را با استفاده از بسط سری تیلر و $Y(x)$ با هم مقایسه کنید.

۸. فرض کنید خطای اولیه $e_0 = 0$. برای کران از (۱۳.۲.۶)

$$\text{Max}_{0 \leq x \leq b} |Y(x_n) - y_n| \leq \frac{h}{4}(e^{2b} - 1)$$

خطای مجانبی به وسیله فرمول زیر داده شده است

$$Y(x) - y_h(x) = h \frac{\ln(1+x)}{(1+x)^2} + O(h^2) \quad x \geq 0$$

این عبارت نشان می‌دهد که با افزایش x خطا کاهش می‌یابد، در حالی که کران قبلی پیش‌بینی می‌کند که وقتی بازه $[0, b]$ افزایش می‌یابد خطا نیز افزایش می‌یابد.

$$y_{1,n+1} = y_{1,n} + h y_{2,n} \quad (\text{ب}) \quad 10$$

$$y_{2,n+1} = y_{2,n} + h [0.1(1 - y_{1,n}^2)y_{2,n} - y_{1,n}] \quad n \geq 0$$

$$T_n(Y) = -\frac{5}{8}h^2 Y^{(3)}(x_n) + O(h^2) \quad 14$$

۲۶. برای $c = \frac{1}{4}$ ، $u_n = c_1 + c_2(\frac{1}{4})^n + c_3(-\frac{1}{4})^n$ ، c_2, c_1 و c_3 دلخواه. هنگامی که

$u_n \rightarrow c_1, m \rightarrow \infty$ که c_1 یک عدد ثابت است.

$$31. \quad (الف) \quad 0 \leq a_0 < 2 \quad (\text{ب}) \quad a_0 = 0$$

(ج) در $a_0 = 0$ هیچ ناحیه‌ای برای پایداری مطلق وجود ندارد. هنگامی که a_0 از ۰ به ۲ می‌رود

ناحیه پایداری مطلق در حد به $0 < h\lambda < 2$ افزایش می‌یابد. این رابطه فقط برای مقادیر حقیقی در ناحیه پایداری مطلق برقرار است.

$$33. \quad (الف) \quad a_0 = -9 - 3b_2 \quad a_1 = 9 \quad a_2 = 1 + 3b_2$$

$$b_1 = 6 + 4b_2 \quad b_0 = 6 + b_2$$

(ب) معادله مشخصه برای $h = 0$ چنین است

$$\rho(r) \equiv r^3 - a_0 r^2 - 9r + a_0 + 8 = 0$$

با تقسیم بر $(r - 1)$ دو ریشه باقیمانده در معادله زیر صادق اند.

$$r^2 + (1 - a_0)r - (a_0 + \lambda) = 0 \quad (*)$$

به ازای همه مقادیر حقیقی a_0 ، عبارت مبین

$$(1 - a_0)^2 + 4(a_0 + \lambda) = (1 + a_0)^2 + 32 \geq 32 > 0$$

مثبت است. بنابراین برای هر مقدار a_0 هر دو ریشه r_1 و r_2 معادله $(*)$ حقیقی اند. یک امتحان مستقیم این ریشه‌ها نشان خواهد داد که این ریشه‌ها همزمان به ازای یک مقدار a_0 در $1 < r < -1$ صدق نمی‌کنند.

۳۶. $y_{n+1} = y_{n-1} + 2hf(x_{n-1}, y_{n-1})$ از مرتبه اول است و پایدار نسبی است اما در شرط ریشه قوی صدق نمی‌کند.

۴۰. از $Y(x) = \frac{1}{4}x^2 = 1$ استفاده کنید بنابراین $Y''(x) = 1$. در این صورت خطای $e_n = Y(x_n) - y_n$ در رابطه زیر صدق می‌کند.

$$e_n = \frac{h}{2\lambda} [(1 + h\lambda)^n - 1] \quad n \geq 0$$

.۴۸

$$\gamma_1 + \gamma_2 + \gamma_3 = 1$$

$$\gamma_2\alpha_2 + \gamma_3\alpha_3 = \frac{1}{2}$$

$$\gamma_1\beta_{21} + \gamma_3(\beta_{21} + \beta_{22}) = \frac{1}{2}$$

$$\gamma_2\alpha_2^2 + \gamma_3\alpha_3^2 = \frac{1}{3}$$

$$\gamma_2\alpha_2\beta_{21} + \gamma_3\alpha_3(\beta_{21} + \beta_{22}) = \frac{1}{3}$$

$$\gamma_2\beta_{21}^2 + \gamma_3(\beta_{21} + \beta_{22})^2 = \frac{1}{3}$$

$$\gamma_2\alpha_2\beta_{22} = \frac{1}{6}$$

$$\gamma_3\beta_{22}\beta_{21} = \frac{1}{6}$$

این معادلات وابسته‌اند و می‌توان آنها را به ۶ معادلهٔ مستقل بدل کرد. یک جواب خاص عبارت است از

$$Y_{n+1} = Y_n + \frac{h}{\rho}(V_1 + 4V_2 + V_3)$$

$$V_1 = f(x_n, y_n) \quad V_2 = f\left(x_n + \frac{1}{3}h, y_n + \frac{1}{3}hV_1\right)$$

$$V_3 = f\left[x_n + h, y_n + h(2V_2 - V_1)\right]$$

اگر $f(x, y)$ به y وابسته نباشد این قاعده همان قاعدهٔ سیمپسن است.

۵۱. (الف) ریشه‌های $h\lambda$ می‌باشند $1 + h\lambda + \frac{1}{4}(h\lambda)^2 < 1$. برای λ حقیقی، این باره $0 < h\lambda < 2$ است.

فصل ۷

۱. (الف) وابسته (ب) مستقل

۷. (الف) $(Ax, x) = x^T(Ax)$. اما همچنین

$$(Ax, x) = (x, Ax) = (Ax)^T x = x^T A^T x = -x^T Ax$$

از آنجایی که $x^T Ax = -x^T Ax$ ، باید داشته باشیم $x^T Ax = 0$ که نتیجهٔ مطلوب است.

۸. $u^{(2)} = (-7, 4, 1)$. برای نرمال کردن، هر بردار را بر طولش تقسیم می‌کنیم.

$$\|u^{(1)}\|_2 = \sqrt{6} \quad \|u^{(2)}\|_2 = \sqrt{11} \quad \|u^{(3)}\|_2 = \sqrt{66}$$

$$A = \frac{1}{4} \begin{bmatrix} 7 & -4 & -4 \\ -4 & 1 & -8 \\ -4 & -8 & 1 \end{bmatrix} \quad (الف) \quad ۹.$$

$$x = [2, -1]^T, \quad \lambda_1 = -1 \quad (الف) \quad ۱۱.$$

$$x = [2, 1]^T, \quad \lambda_2 = 3$$

۱۳. $\|Ux\|_2^2 = (Ux, Ux) = (x, U^*Ux) = (x, Ix) = (x, x) = \|x\|_2^2$ (الف)

نشان می‌دهد که برای هر x, y ، $\|Ux\|_2 = \|x\|_2$. برای فاصلهٔ بین Ux و Uy ، از نتیجهٔ قبل،

$$\|Ux - Uy\|_2 = \|U(x - y)\|_2 = \|x - y\|_2$$

۱۴. چون A ارمیتی است، فرض کنید $x^{(1)}, \dots, x^{(n)}$ یک پایهٔ یکامتعامد برای C^n مربوط به ویژه مقادیرهای حقیقی $\lambda_1, \dots, \lambda_n$ باشد. برای هر x در C^n می‌توانیم بنویسیم

$$x = \sum_1^n \alpha_i x^{(i)} \quad \text{با} \quad \alpha_i = (x, x^{(i)}) \quad i = 1, \dots, n$$

از یکا متعامد بودن داریم $\|x\|_2^2 = \sum_1^n |\alpha_i|^2$. با استفاده از این شکل برای x نشان دهید

$$(Ax, x) = \sum_1^n \lambda_i |\alpha_i|^2$$

چون $\alpha_1, \dots, \alpha_n$ برای بدست آوردن عناصر مختلف x از C^n می‌توانند به طور دلخواه تغییر کنند این فرمول ثابت می‌کند که A معین مثبت است اگر و تنها اگر همهٔ λ_i ها مثبت باشند.

۲۳. ابتدا ثابت کنید

$$\|x\|_\infty \leq \|x\|_p \leq n^{1/p} \|x\|_\infty \quad x \in C^n \quad \text{به ازای همهٔ مقادیر}$$

از این رابطه به آسانی دیده می‌شود که هرگاه $p \rightarrow \infty$ ، $\|x\|_p \rightarrow \|x\|_\infty$.

۲۷. (الف) فرض کنید λ و x یک جفت ویژه مقدار- ویژه بردار ماتریس A باشند. آنگاه

جفت ویژه‌های مربوط به تغییرات مختلف A عبارت‌اند از: (۱) λ^m و x برای A^m ؛ (۲) $1/\lambda$ و x برای A^{-1} ؛ و (۳) $\lambda + c$ برای $A + cI$.

۲۸. با استفاده از ستونهای A_{*j} ماتریس A ، بنویسید $A = [A_{*1}, \dots, A_{*n}]$. بنابراین

$$F(A) = \sqrt{\|A_{*1}\|_2^2 + \dots + \|A_{*n}\|_2^2}$$

(الف) $F(UA) = \sqrt{\|UA_{*1}\|_2^2 + \dots + \|UA_{*n}\|_2^2}$. با استفاده از مسألهٔ ۱۳ (الف)

برای جمیع مقادیر j ، $\|UA_{*j}\|_2 = \|A_{*j}\|_2$. بنابراین $F(UA) = F(A)$.
(ب) فرض کنید $U^*AU = D = \text{diag}[\lambda_1, \dots, \lambda_n]$. برای به دست آوردن

$$F(A) = F(U^*AU) = F(D) = \sqrt{\lambda_1^2 + \dots + \lambda_n^2}$$

قسمت (الف) را بکار برید.

$$\|A\|_\infty \leq 1, \|A^{-1}\|_\infty < \frac{1}{\gamma} \quad ۳۲$$

فصل ۸

۱. (الف)

$$L = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ -2 & 3 & 1 \end{bmatrix} \quad U = \begin{bmatrix} 1 & 1 & -1 \\ 0 & 1 & -1 \\ 0 & 0 & 2 \end{bmatrix} \quad x = \begin{bmatrix} 2 \\ 2 \\ 3 \end{bmatrix}$$

۲. (الف) بدون محورگیری، ماتریس افزوده $[A | b]$ به

$$\left[\begin{array}{ccc|c} 6000 & 2000 & 2000 & -2000 \\ 0 & 00001000 & -03333 & 1667 \\ 0 & 0 & 5555 & -27790 \end{array} \right] \quad \begin{array}{l} m_{21} = 03333 \\ m_{21} = 01667 \\ m_{32} = 16670 \end{array}$$

بدل می‌شود. جواب با استفاده از جایگذاری پسرو عبارت است از

$$x_1 = 1335 \quad x_2 = 0 \quad x_3 = -5003$$

همه اعمال حساب انجام شده تا ۴ رقم اعشاری با ممیز شناور گردشده‌اند. این برای قسمت (ب) نیز درست است.

(ب) با محورگیری $[A | b]$ به

$$\left[\begin{array}{ccc|c} 6000 & 2000 & 2000 & -2000 \\ 0 & 1667 & -1333 & 03334 \\ 0 & 0 & 03332 & 1667 \end{array} \right] \quad \begin{array}{l} m_{21} = 03333 \\ m_{21} = 1667 \\ m_{32} = 000005999 \end{array}$$

بدل می‌شود. جواب با استفاده از جایگذاری پسرو عبارت است از

$$x_1 = 2602 \quad x_2 = -3801 \quad x_3 = -5003$$

۵. (الف) در شکل افرازشده،

$$\begin{bmatrix} A_1 & -A_2 \\ A_2 & A_1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

(ب) فرض کنید دستگاه ۱ معرف دستگاه حقیقی قسمت (الف) باشد، و فرض کنید دستگاه

۲ دستگاه مختلط اولیه $Ax = b$ را نشان دهد. برای ذخیره‌سازی ماتریس، دستگاه ۱ به $4n^2$

حافظه نیاز دارد و دستگاه ۲ به $2n^2$ حافظه. (هر عدد مختلط به دو حافظه ذخیره‌سازی نیاز دارد). برای حل، دستگاه ۱ تقریباً به $\frac{1}{3}(2n)^3 = \frac{8}{3}n^3$ عملیات ضرب و تقسیم نیاز دارد. دستگاه ۲؛ به $\frac{1}{3}n^3$ عملیات ضرب و تقسیم مختلط نیاز دارد. چون هر ضرب مختلط به ۴ ضرب حقیقی نیاز دارد تعداد عملیات واقعی $\frac{4}{3}n^3$ است. بنابراین دستگاه ۲ به نصف حافظه ذخیره‌سازی و در حدود نصف زمان عملیات نسبت به دستگاه ۱ نیاز دارد.

۹. چون L^T ناکین است $\langle L^T x, L^T x \rangle = (L L^T x, x) = (L^T x, L^T x) = \|L^T x\|^2 > 0$ برای هر $x \neq 0$ همچنین $\det(A) = \det(L)^2$

۱۰. (الف)

$$L = \begin{bmatrix} 1r^\circ & 0 & 0 \\ -2r^\circ & 1r^\circ & 0 \\ 3r^\circ & -4r^\circ & 3r^\circ \end{bmatrix}$$

۱۴

$$L = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ 1 & \frac{5}{3} & 0 & 0 & 0 \\ 0 & 1 & \frac{11}{5} & 0 & 0 \\ 0 & 0 & 1 & \frac{29}{11} & 0 \\ 0 & 0 & 0 & 1 & \frac{70}{29} \end{bmatrix} \quad U = \begin{bmatrix} 1 & -\frac{1}{3} & 0 & 0 & 0 \\ 0 & 1 & -\frac{2}{5} & 0 & 0 \\ 0 & 0 & 1 & -\frac{5}{11} & 0 \\ 0 & 0 & 0 & 1 & -\frac{11}{29} \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$Lz = b$ دارای جواب $z = [\frac{2}{3}, -\frac{7}{5}, \frac{11}{11}, \frac{-41}{29}, 1]^T$ است

$Ux = z$ دارای جواب $x = [1, -1, 1, -1, 1]^T$ است.

$$\text{cond}(A)_1 = \text{cond}(A)_\infty = 39601, \text{cond}(A)_2 = 39206.18$$

۲۶. برای نرخهای همگرایی μ از (۵-۶-۸) و η از (۱۷-۶-۸)، در این حالت داریم: $\mu = \frac{2}{3}$

و $\eta = \frac{2}{3}$. نرخهای واقعی کاهش مشاهده شده، ۰٫۶۱ برای روش گاوس-ژاکوبی، و ۰٫۳۷ برای روش گاوس-زایدل هستند.

۲۷. (الف) همگراست اگر و تنها اگر $r_\sigma(A^{-1}B) < 1$ ، نرخ همگرایی اساساً $r_\sigma(A^{-1}B)$

یا $\|A^{-1}B\|$ است بسته به نرمی که استفاده شده باشد.

(ب) باز هم همگراست اگر و تنها اگر $r_\sigma(A^{-1}B) < 1$ ، اما نرخ همگرایی در حدود $r_\sigma(A^{-1}B)^2$

یا $\|A^{-1}B\|^2$ است که سریعتر از قسمت (الف) است.

۳۲. (الف) $A^{-1} - C_m = A^{-1}R_m = A^{-1}R_0^{m \times m}$ جمله C_m به A^{-1} همگراست اگر و تنها اگر $r_\sigma(R_0) < 1$.
 (ب) با نوشتن $(I - R_0)^{-1}$ به صورت یک حاصلضرب نامتناهی،

$$A^{-1} = C_0(I + R_0)(I + R_0^1)(I + R_0^2)(I + R_0^3) \dots$$

$$C_m = C_0(I + R_0)(I + R_0^1)(I + R_0^2) \dots (I + R_0^{m-1})$$

فصل ۹

۱. (الف) با به کار بردن قضیه گرگورین برای سطرهای ماتریس A دایره‌های زیر به دست می‌آیند $|\lambda - 1| \leq 1$ ، $|\lambda - 5| \leq 2$ ، و $|\lambda - 9| \leq 3$. دایره‌های دوم و سوم متقاطع‌اند و بنابراین باید شامل دو ویژه مقدار در قسمت اشتراکشان باشند. اولین دایره متمایز از دو دایره دیگر است. چون چندجمله‌یی مشخصه A ضرایب حقیقی دارد، ویژه مقدارها اگر مختلط باشند، باید مزدوج همدیگر باشند. بنابراین $|\lambda - 1| \leq 1$ فقط شامل یک ویژه مقدار حقیقی در بازه $[2, 0]$ است. با به کار بردن قضیه برای ستونهای A دایره‌ها عبارت‌اند از $|\lambda - 1| \leq 3$ ، $|\lambda - 5| \leq 2$ ، و $|\lambda - 9| \leq 1$. با استدلالهای مشابه و استفاده از نتایج قبلی، ویژه مقدارها در هر یک از بازه‌های $[2, 0]$ ، $[3, 7]$ ، و $[8, 10]$ قرار دارند. ویژه مقدارهای واقعی عبارت‌اند از:

$$1.331927689$$

$$4.856933692$$

$$8.811138616$$

۲. (ج) n دایره عبارت‌اند از:

$$|r + a_{n-1}| \leq 1 \quad \text{و} \quad |r| \leq 1 + |a_j| \quad j = 1, \dots, n-2 \quad |r| \leq |a_0|$$

د. (i) ۹ ریشه وجود دارند که در $|\lambda - 1| \leq 1$ صدق می‌کنند و یک ریشه حقیقی در $-7 \leq r \leq -9$ صدق می‌کند.

$$\text{Min}_{1 \leq i \leq 2} |\lambda - \lambda_i| \leq 0.0000356. ۴$$

۸. (الف) $\lambda = 15$ ، $x = [1, 1, 1, 1]^T$. در هر بارست خطا با ضریب $\frac{1}{4}$ کاهش می‌یابد. بنابراین $\lambda = 5$ احتمالاً بزرگترین ویژه مقدار دوم است.

۱۳. (ب) با استفاده از $w = [1/\sqrt{6}, -1/\sqrt{6}, -2/\sqrt{6}]^T$ برای ساختن

بنابراین $Px = e_1$ می‌آوریم $P = I - 2ww^T$

$$B = PAP^T = \begin{bmatrix} 9 & 0 & 0 \\ 0 & 18 & 0 \\ 0 & 0 & -9 \end{bmatrix}$$

در این حالت ماتریس B قطری است اما معمولاً بدین سادگی نخواهد بود. با استفاده از $w = (1/\sqrt{30})[5, 1, 2]^T$ به یک P برای هر $Px = -e_1$ می‌رسیم. اما $B = PAP^T$ همان شکل را برای سطر و ستون ۱ خواهد داشت. در این حالت خاص

$$B = \begin{bmatrix} 9 & 0 & 0 \\ 0 & \frac{18}{25} & -\frac{224}{25} \\ 0 & -\frac{224}{25} & \frac{207}{25} \end{bmatrix}$$

اولین شکل P با علامت مخالف با $(9-3-9)$ ساخته شد. در حالی که دومین انتخاب P براساس $(9-3-9)$ ساخته می‌شود. در حالت عادی شکل P که با $(9-3-9)$ مطابقت می‌کند به دلیل دقت در ملاحظاتی که پس از $(9-3-9)$ آورده‌ایم، شکل ارجح خواهد بود.

۱۴

$$Q = \frac{1}{3} \begin{bmatrix} -1 & 2 & 2 \\ -2 & 1 & -2 \\ -2 & -2 & 1 \end{bmatrix} \quad R = 3 \begin{bmatrix} -1 & 1 & -1 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\beta = \frac{b_l}{\gamma}, \alpha = \frac{b_k}{\gamma} \quad \text{(الف) ۱۵}$$

(ب)

$$R^{(1,2)} = \begin{bmatrix} \frac{1}{\sqrt{2}} & 0 & 0 & \frac{1}{\sqrt{2}} \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ -\frac{1}{\sqrt{2}} & 0 & 0 & \frac{1}{\sqrt{2}} \end{bmatrix} \quad R^{(1,2)} = \begin{bmatrix} \frac{\sqrt{2}}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} & 0 \\ 0 & 1 & 0 & 0 \\ -\frac{1}{\sqrt{2}} & 0 & \frac{\sqrt{2}}{\sqrt{2}} & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$R^{(1,2)} = \begin{bmatrix} \frac{\sqrt{2}}{2} & \frac{1}{2} & 0 & 0 \\ -(\frac{1}{2}) & \frac{\sqrt{2}}{2} & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\hat{b} = Ub = [2, 0, 0, 0]^T$$

$$U = R^{(1,2)}R^{(1,3)}R^{(1,4)} = \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{-1}{2\sqrt{3}} & \frac{\sqrt{2}}{2} & \frac{-1}{2\sqrt{3}} & \frac{-1}{2\sqrt{3}} \\ \frac{-1}{\sqrt{6}} & 0 & \frac{\sqrt{2}}{\sqrt{3}} & \frac{-1}{\sqrt{6}} \\ \frac{-1}{\sqrt{2}} & 0 & 0 & \frac{1}{\sqrt{2}} \end{bmatrix}$$

۱۹. الف) $f_5(\lambda) = 1$ ، $f_4(\lambda) = -\lambda$ ، $f_3(\lambda) = (1 - \lambda)f_{j-1}(\lambda) - f_{j-2}(\lambda)$ ، $f_2(\lambda) = (2 - \lambda)f_1(\lambda) - f_0(\lambda)$ ، فرض کنید ازای $z = 2, 3, 4$ ، به موجب قضیه گرسگورین همه λ_i ها در بازه $[-1, 3]$ قرار دارند. بنابراین همان گونه که در جدول زیر نشان داده شده است، می توان ریشه ها را با استفاده از قضیه ۵-۹ جدا کرد. ریشه های واقعی را می توان توسط روش دیگری پیدا کرد. به عنوان مثال روش خط قاطع، سریعاً به جواب $\lambda_5 = 2.90211303$ می رسد.

λ	$f_5(\lambda)$	$s(\lambda)$	ملاحظات
-1.0	2.0	5	$\lambda_1 > -1$
3.0	-2.0	0	$\lambda_5 < 3$
1.0	0.0	2	$\lambda_3 = 1.0, \lambda_4 > 1$
0.0	1.0	3	$\lambda_2 < 0.0$
-0.5	-1.78	4	$-1 < \lambda_1 < 0.5 < \lambda_2 < 0$
2.0	-1.0	2	$\lambda_4 > 2$
2.5	1.78	1	$2 < \lambda_2 < 2.5 < \lambda_5 < 3$

۲۰. (ب) (۲) با استفاده از تبدیلات هاسهولدر مانند (۲۱.۳.۹) تا (۲۳.۳.۹) به دست

می آوریم $T = Q^T A Q$

$$T = \begin{bmatrix} 5.0 & -4.2426406871 & 0.0 & 0.0 \\ -4.2426406871 & 6.0 & 1.4142135624 & 0.0 \\ 0.0 & 1.4142135624 & 5.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 2.0 \end{bmatrix}$$

$$Q = \begin{bmatrix} 1.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & -0.9428090416 & 0.3333333333 & 0.0 \\ 0.0 & -0.2357022604 & -0.6666666667 & -0.7071067812 \\ 0.0 & -0.2357022604 & -0.6666666667 & 0.7071067812 \end{bmatrix}$$

در اینجا Q صریحاً محاسبه شده است، معمولاً در عمل ساختن Q به طور صریح اتلاف وقت خواهد بود.

۲۲. (۱) ویژه مقادیر عبارتند از

$$\lambda_1 = 2 - \sqrt{3} \quad \lambda_2 = 2 \quad \lambda_3 = 2 + \sqrt{3}$$

(الف) برای روش QR بدون انتقال عناصری که خارج از قطر قرار دارند بطور خطی به صفر همگرا می‌شوند. در هر بارست عنصرهای موضع $(2, 1)$ با ضریب 0.536 کاهش می‌یابند و عناصر موضع $(2, 3)$ با ضریب -0.134 در هر بارست کاهش می‌یابند.

(ب) برای روش QR با انتقال، با استفاده از انتخاب

$$c_m = a_{p,p}^{(m)}$$

همگرایی سریع خواهیم داشت. بعد از پنج بارست اندازه عنصر موضع $(2, 3)$ کمتر از 5×10^{-11} خواهد بود. با نشان دادن ماتریس اولیه A با A_0 نتیجه زیر را خواهیم داشت:

$$A_0 = \begin{bmatrix} 3.7316925974 & 0.249060210 & 0 \\ 0.249060210 & 2.0003582102 & \varepsilon \\ 0 & \varepsilon & 0.2679491924 \end{bmatrix}$$

با $5 \times 10^{-11} < |\varepsilon|$ وقتی روش QR با انتقال، برای ماتریس تبدیل شده حاصل از حذف سطر و ستون سوم A_0 به کار رود، ماتریس بسیار سریع همگرا می‌شود.

۲۴. برای $\lambda = 7.93229$ ویژه بردار عبارت است از

$$x = [0.7211774817, 0.2724739430, 1.0, 0.2515508351]^T$$

برای $\lambda = 5.6689$

$$x = [0.5736191640, 0.5489544105, -0.8148078321, 1.0]^T$$

۲۸. مقادیر به دست آمده در جدول زیر داده شده‌اند

n	$\lambda_i^{(n)}$	$\lambda_j^{(n)}$
۲	۰٫۹	۰٫۱
۴	۰٫۸۸۴۵۶۶	۰٫۱۰۳۹۵۷
۸	۰٫۸۸۰۹۷۱	۰٫۱۰۵۳۶۰
۱۶	۰٫۸۸۰۰۸۵	۰٫۱۰۵۷۱۵
۳۲	۰٫۸۷۹۸۶۵	۰٫۱۰۵۸۰۵

با بررسی نسبتهای تفاضلات پیاپی به طور تجربی می توان نشان داد که

$$\lambda_i - \lambda_i^{(n)} = O\left(\frac{1}{n^2}\right)$$

این نرخ همگرایی به صورت نظری قابل توجه است. برای برآورد خطاهای ایجاد شده برونمایی ریچاردسون را می توان به کار برد.

$$\lambda_1 - \lambda_1^{(32)} \doteq 7.3 \times 10^{-5} \quad \lambda_2 - \lambda_2^{(32)} \doteq 3.0 \times 10^{-5}$$

۲۷. A^+ , $A(1)$ را با استفاده از شکل تجزیه تکین - مقداری شان بیان کنید: $A = VFU^T$ و $A^+ = UF^+V^T$ بنابراین

$$\begin{aligned} AA^+A &= (VFU^T)(UF^+V^T)(VFU^T) \\ &= VFF^+FU^T \\ &= VFU^T \quad \text{چون با محاسبه مستقیم خواهیم داشت } FF^+F = F \\ &= A \end{aligned}$$

۳۰. خطای ریشه دوم میانگین برابر است با $p(x) = -1.269091x + 0.392952$.
۰.۲۴۳
۳۳. الف)

$$QR = \begin{bmatrix} -0.5179 & -0.7517 & 0.4082 \\ 0.5754 & 0.0470 & 0.8165 \\ -0.6330 & 0.6578 & 0.4082 \end{bmatrix} \begin{bmatrix} -1.7378 & -1.7148 \\ 0 & -0.2819 \\ 0 & 0 \end{bmatrix}$$

$$SVD = \begin{bmatrix} 0.5774 & 0.7071 & 0.4082 \\ -0.5774 & 0 & 0.8165 \\ 0.5774 & -0.7071 & 0.4082 \end{bmatrix} \begin{bmatrix} 2.4495 & 0 \\ 0 & 0.2000 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0.707 & 0.7071 \\ -0.7071 & 0.7071 \end{bmatrix}$$

$$A^+ = \begin{bmatrix} -\frac{7}{4} & -\frac{1}{6} & \frac{1}{3} \\ \frac{1}{4} & -\frac{1}{6} & -\frac{7}{3} \end{bmatrix} \quad \text{cond}(A) = 12.2474$$

نمایه

- اتحاد کریستوفل-درايو ۲۴۵
 ارقام بامعنى ۲۰
 اسكالرها ۵۲۳
 اعداد بزولى ۳۲۱
 اعداد دودويى ۱۳
 الحاقى ۵۲۶
 انتشار خطاها ۲۷، ۳۳
 انتگرالده نامتاهى ۳۴۶
 انتگرالگيرى (انتگرالگيرى عددى را ببينيد).
 انتگرالگيرى (انتگرالگيرى عددى را ببينيد).
 انتگرالگيرى تطبيقى ۳۳۹
 قاعده سيمپسون ۳۳۹
 انتگرالگيرى چندمتغيره ۳۶۲-۳۶۳
 انتگرالگيرى رامبرگ ۳۳۵
 انتگرالگيرى عددى ۲۸۰
 انتگرالگيرى گاوس-لژاندر ۳۱۰
 انتگرالگيرى گاوسى ۳۰۳
 انتگرالهاى تكين ۳۴۵
 برنامه خودكار ۳۴۱
 برونيابى ايتكن ۳۲۹
 برونيابى ريچارلسن ۳۳۲
 تطبيقى ۳۳۸
 درجه دقت ۲۹۸
 روش رامبرگ ۳۳۵
 فرمول اويلر-مكلورن ۳۲۲، ۳۲۷
 فرمول باز ۳۰۳
 فرمول كرونرود ۳۱۹
 فرمولهاى نيوتن-كوتس ۲۹۵
 قاعده بول ۲۹۸
 قاعده دوزنقه‌يى ۲۸۳
 قاعده دوزنقه‌يى تصحيح شده ۲۸۶
 قاعده سه هشتم ۲۹۶
 قاعده سيمپسون ۲۸۷
 قاعده مستطيلى ۳۸۵
 قاعده ميانگاهى ۳۰۲
 قاعده‌هاى پاترسن ۳۱۹

روش IMT ۳۴۷	مقایسه برنامه‌ها ۳۴۳
محاسبه تحلیلی ۳۵۱	همگرایی ۲۹۹
انتگرالهای حاصلضرب ۳۵۱	انتگرالگیری عددی خودکار ۳۳۷
اندازه گامها ۲۸۳	انتگرالگیری تطبیقی ۳۳۸
ایتکن ۹۷	قاعده سیمپسون ۳۳۹
بارست پیکار ۵۱۳	انتگرالگیری عددی میانگاهی ۳۰۲، ۳۶۷
بارست خطی (بارست تک نقطه‌ای را ببینید.)	انتگرالگیری گاوسی ۳۰۳ (انتگرالگیری
بارست گاوس-ژاکوبی ۶۱۸	گاوس-لژاندر را ببینید.)
بارست معکوس ۷۱۴	درجه دقت ۳۰۶
بارست نقطه ثابت، روشهای بارستی تک نقطه‌ای	فرمول ۳۰۶، ۳۰۹
۸۷	فرمول خطا ۳۰۶، ۳۰۹
بازه انتگرالگیری نامتناهی ۳۴۵	لاگر ۳۴۹
باؤتر-فیکه ۶۶۹	مثبت بودن وزنها ۳۰۹
برآورد خطا:	وزنها ۳۰۶
ریشه‌یابی ۷۳، ۸۱، ۱۴۰	همگرایی ۳۱۲
کلی (فراگیر) ۴۹۲، ۴۴۲	انتگرالگیری گاوسی ۳۰۳ (انتگرالگیری
برآورد خطای ریچاردسن ۳۳۴، ۴۱۹	گاوس-لژاندر را ببینید.)
برازش داده‌ها به روش کمترین مربعات (تقریب	انتگرالگیری نیوتن-کوش ۲۹۵
کمترین مربعات را ببینید.)	باز ۳۰۳
بردارها ۵۲۳	بسته ۳۰۳
دوبه‌دو متعامد ۶۷۵	فرمول خطا ۲۹۶
زاویه بین ۵۳۰	همگرایی ۲۹۹
مستقل ۵۲۴	انتگرالگیری-گاوس لژاندر ۳۱۰
نرم ۵۴۵ (نرم برداری متعامد را ببینید.)	بحث همگرایی ۳۱۳-۳۱۴
وابسته ۵۲۴	فرمول خطا ۳۰۹
همگرا ۵۴۸	قاعده دوزنقه‌ی مرکب ۲۸۵
برونیایی ایتکن:	کرونرود ۳۱۹
انتگرالگیری عددی ۳۲۹	گره‌ها و وزنها ۳۱۰
نرخ همگرایی ۱۴۰	انتگرالهای تکین ۳۴۵
ویژه‌مقدارها ۶۸۷	انتگرالگیری حاصلضرب ۳۵۱
برونیایی ریچاردسن ۳۳۲	انتگرالگیری گاوسی ۳۴۸
بسطهای چندجمله‌ی چیشف ۲۴۹	تبدیل متغیر ۳۴۵

تابع وزن ۲۳۳، ۲۸۲	بسطهای چندجمله‌یی لژاندر ۲۴۸
تبدیلات متعامد:	بعد ۵۲۵
پیدا کردن ماتریس متقارن ۶۹۷	بلوک ژوردان ۵۴۴
تبدیل بردار ۶۹۲	بهبود بارستی ۶۱۲
تجزیه QR ۶۹۴	بهترین تقریب (تقریب مینیماکس را ببینید.)
حفظ طول ۶۹۲	بهینه‌سازی ۱۲۹، ۱۳۴
دوران در صفحه ۷۰۱	روش امتداد مزدوج ۶۴۰
ماتریسهای هاوسهولدر ۶۹۰	روش تندترین کاهش ۱۳۲
تبدیل اعداد پایه ۵۵	روش گرادیان مزدوج ۱۳۲
تبدیل سریع فوریه ۲۰۷	روش نلدرو مد ۱۳۲
تبدیل متناهی فوریه ۲۰۵	روش نیوتن ۱۳۱
تجزیه LU	روشهای شبه‌نیوتنی ۱۳۲
بارست معکوس ۷۱۴	روشهای کاهش ۱۳۲
ذخیره ۵۷۷	ریاداری ۴۱
ماتریسهای سه‌قطری ۵۹۶	روش اویلر ۳۹۲
یکتایی ۶۹۶	روشهای عددی ۴۶
تجزیه LU ۵۷۷	ضعیف ۴۱۱
تجزیه QR ۶۹۴	مطلق ۴۶۰
حل دستگاه خطی ۶۹۷	معادلات دیفرانسیل ۳۷۹
یکتا ۶۹۵	روشهای عددی ۳۹۲، ۴۴۷
تجزیه تکین-مقدار ۵۴۱، ۵۶۸، ۷۲۲	نسبی ۴۵۷
محاسبه ۷۳۴	ویژه‌مقدارها ۶۶۹، ۶۷۷
تحلیل بازه‌یی ۴۸	A-پایداری (پایداری مطلق) ۴۱۸، ۴۶۳، ۴۶۷
تحلیل پسرو خطا ۶۰۷	پایداری مطلق ۴۶۰
ترازها ۵۲۶	پایداری نسبی ۴۵۷
ترازها در مزدوج ۵۲۶	پایه ۱۳، ۵۲۵
ترکیب خطی ۵۲۴	استانده ۵۲۵
تساوی پارسوال ۲۴۷	متعامد ۵۳۱
تصویر قائم ۵۳۲، ۵۶۶، ۶۶۲	پایه استاندارد ۵۲۵
تفاضلات پیشرو نیوتن ۱۶۹ (تفاضل پیشرو را ببینید.)	تابع دلتای کرونگر ۵۲۸
تفاضلات متناهی ۱۶۸	تابع زوج/فرد ۲۶۰
	تابع نفوذ ۴۳۲

- فرمول درونبایی ۱۶۹، ۱۷۲
- تفاضل پسر ۱۷۲
- تفاضل پسر نیون ۱۷۲ (تفاضل پسر را ببینید.)
- تفاضل پیشرو ۱۶۸
- خطی بودن ۱۷۰
- رابطه با تفاضلات منقسم ۱۷۰
- رابطه با مشتقات ۱۷۳
- فرمول درونبایی ۱۶۹
- کشف نوفه در داده‌ها ۱۷۴
- تفاضل منقسم ۱۵۸، ۹
- پیوستگی ۱۶۷
- چندجمله‌یها ۱۶۷
- درونبایی ۱۶۰
- رابطه با مشتقات ۱۶۴-۱۶۵
- فرمول ۱۵۹، ۱۶۴
- فرمول ارمیت-جنوکی ۱۶۵
- مشتگیری ۱۶۷
- نمودار محاسبه ۱۶۰
- تفاضل منقسم نیون ۱۵۸ (تفاضل منقسم را ببینید.)
- تقریب پاره ۲۶۹، ۲۷۳
- تقریب توابع ۲۲۳
- توابع زوج/فرد ۲۶۰
- درون‌بایی ۱۸۰
- سری چیبشف ۲۴۹، ۲۵۶
- صرفه‌جویی چندجمله‌یی تیلر ۲۷۹
- قضیه تیلر ۵، ۲۲۵
- قضیه جکسن ۲۰۶، ۲۵۵
- قضیه دولاواله پواسن ۲۵۳
- قضیه هموسانی ۲۵۴
- کمترین مربعات ۲۳۱، ۲۳۳، ۲۴۵
- مینیماکس ۲۲۸، ۲۵۲
- نزدیک مینیماکس ۲۵۶
- تقریب کمترین مربعات ۲۳۱
- مسأله گسسته
- پایداری ۲۲۴
- جواب تکن مقدار ۲۲۲
- حل با روش QR ۲۳۱
- مسأله برازاندن ۲۲۵
- معکوس تعمیم‌یافته ۷۲۳، ۷۴۶
- وزن ۲۴۵
- همگرایی ۲۴۷
- تقریب مینیماکس ۲۲۸
- خطا ۲۲۸
- قضیه هموسانی ۲۵۴
- تقریب یک‌نواخت (تقریب توابع را ببینید.)
- تقریب‌های نزدیک مینیماکس ۲۵۶
- سطح‌های چندجمله‌یی چیبشف ۲۴۹
- درونبایی ۲۵۹
- نوسان اجباری ۲۶۴
- تقلیل چندجمله‌یها ۱۱۱-۱۱۲
- ماتریس ۶۹۰، ۷۴۲
- تکنین-مقدار ۵۴۱
- تلسکوپی کردن چندجمله‌یی تیلر ۲۷۹
- توابع برازا ۱۹۰
- B-برازا ۱۹۸
- برازای درونیاب طبیعی ۲۲۱
- بهینگی ۱۹۴
- خطا ۱۹۳
- درون‌یاب برازای کامل ۱۹۳
- ساختن ۱۹۱
- شرط نبود یک‌گره ۱۹۵
- توابع خاص ۲۶۹
- توابع مثلثاتی، متعامد گسسته ۲۰۴، ۲۲۲

تجزیه LU ۵۷۷	جایگذاری پسر ۵۷۴
تحلیل پسر خطا ۶۰۷	جایگزینی پی در پی ۶۲۱ (روش گاوس-زایدل را ببینید.)
تحلیل خطا ۵۹۹	جایگزینی همزمان ۶۱۸ (روش گاوس-ژاکوبی را ببینید.)
تصحیح مانده ۶۱۲	جبرخطی ۵۲۳
دستگاه مختلط ۶۵۳	جزء کسری ۱۳
دستگاههای سه قطری ۵۹۶	جواب مزاحم ۴۱۱، ۴۵۵
روش چولسکی ۵۹۲	چندجمله‌یی برنشتاین ۲۲۵
روش گاوس-یوردان ۵۹۰	چندجمله‌یی مشخصه ۴۵۰، ۵۳۳
روشهای فشرده ۵۹۱	چندجمله‌یهای برنولی ۳۲۱، ۳۶۹
شمارش عملیات ۵۷۹	چندجمله‌یهای چیشف ۲۴۰
صورت‌های دیگر ۵۹۰	صفرها ۲۵۹
قضیه ویلکینسن ۶۰۷	ماکسیمم ۲۵۸
کرانه‌های خطا ۶۰۵، ۶۱۱	نوع دوم ۲۷۶
ماتریس معین مثبت ۵۹۲، ۶۵۴	چندجمله‌یهای لاگر ۲۴۰، ۲۴۴
محورگیری ۵۸۲	چندجمله‌یهای لژاندر ۲۳۹، ۲۴۵
مقیاس دهی ۵۸۵	چندجمله‌یهای متعامد ۲۳۵
حساب ممیز شناور ۱۲-۲۰، ۴۸	اتحاد کریستوفل-دارو ۲۴۵
حل دستگاههای خطی سه قطری ۵۹۶	چندجمله‌یهای چیشف ۲۴۰
حل موضعی ۴۱۴	رابطه بازگشتی سه‌گانه ۲۴۳
خارج قسمت ریلی ۶۸۸، ۷۴۱	صفرها ۲۴۲، ۲۷۵
خانواده چگال ۳۰۰	قضیه گرام-اشمیت ۲۳۷
خانواده متعامد ۲۴۰	لاگر ۲۴۰، ۲۴۴
خطا درگام ۴۸۷	لژاندر ۲۳۹، ۲۴۵
خطا درگام واحد ۴۸۷	چندجمله‌یهای مثلثاتی ۲۰۲
خطا واحد گرد کردن ۱۷	حاصلضرب داخلی ۳۹، ۲۳۶، ۵۲۹
خطا:	خطا ۳۸
انتشار ۲۷، ۳۳	حاصلضرب نامتناهی ۱۳۷
برشی ۲۳	حذف گاوسی ۵۷۴
تحلیل آماری ۳۷	برآورد خطا ۴۴۳
تحلیل پسر ۶۰۷	بهبود بارستی ۶۱۲
تعاریف ۲۰	
داده‌ها ۲۳، ۳۴، ۱۷۴	

- قطع کردن ۱۴
 کاهش ارقام با معنی ۳۰، ۳۲
 گرد کردن ۱۴
 ماشین ۲۳
 مجموعیابی ۳۴
 منشأ ۲۱
 نوفه ۲۵
 واحد گرد کردن ۱۷
 خطاهای ماشین ۲۳
 خطای برشی ۲۳، ۳۸۴
 خطای داده‌ها ۲۳، ۳۴
 خطای سرریز ۱۹، ۲۶
 خطای کاهش ارقام با معنی ۲۹
 خطای کلی ۳۸۶
 برآورد ۴۴۳
 خطای گرد کردن ۱۴
 انتگرالگیری عددی ۳۶۷
 حذف گاوسی ۶۰۷
 درونیابی ۱۵۶، ۲۱۴
 معادلات دیفرانسیل ۳۶۰
 خطای موضعی ۴۱۵
 خطای نسبی ۲۰
 خطی-مستقل ۵۲۴
 خطی-وابسته ۵۲۴
 خمهای تراز ۳۷۵
 دترمینان ۵۲۸، ۵۳۴
 محاسبه ۵۷۹
 درجه دقت ۲۹۸
 درونیابی ارمیت ۱۸۱
 درجه ۳ تکه‌ای ۱۸۹
 فرمول خطا ۱۸۴
 فرمولها ۱۸۴، ۲۱۷
 مسأله درونیابی کلی ۱۸۵
 درونیابی ارمیت-برکاف ۲۱۸
 درونیابی چندجمله‌یی ۱۵۰
 درونیابی چندجمله‌یی تکه‌یی ۱۸۶، ۲۱۰
 ارمیت ۱۸۹
 توابع برازا ۱۹۰
 لاگرانژ ۱۸۷
 محاسبه ۱۸۸-۱۸۹
 درونیابی چندمتغیره ۲۱۱
 درونیابی در صفرهای چبیشف ۲۵۹
 درونیابی گویا ۲۱۵
 درونیابی مثلثاتی ۲۰۲
 تبدیل سریع فوریه ۲۰۷
 همگرایی ۲۰۷
 درونیابی وارون ۱۶۲
 درونیابی:
 توابع برازا ۱۹۰
 چندجمله‌یی ۱۵۰
 ارمیت ۱۸۱، ۱۸۵
 ارمیت-برکاف ۲۱۸
 تفاضل پیشرو ۱۷۰
 خطا رفتاری ۱۷۹
 خطای گرد کردن ۱۵۶
 صفرهای چبیشف ۲۵۹
 فرمول تفاضل پسرو ۱۷۲
 فرمول تفاضل منقسم ۱۶۰
 فرمول خطا ۱۵۲، ۱۶۳
 فرمول گرانیگاهی ۲۱۴
 فرمول لاگرانژ ۱۵۲
 مثال رونگه ۱۸۰
 مثالی از \log_3^3 ۱۵۵
 مسأله تقریب ۱۸۰

- ۵۹۲ روش دولیتل
 ۵۹۲ روش کردت
 ۶۲۱ روش گاوس-زایدل
 ۶۱۸ روش گاوس-ژاکوبی
 ۶۴۲، ۶۳۸ روش گرایان مزدوج
 ۶۱۷ روشهای بارستی
 ۵۹۲ روشهای چولسکی
 ۵۹۱ روشهای فشرده
 ۵۹۶ سه قطری
 ۶۰۱، ۶۰۰ ضریب وضعیت
 ۶۰۵ کرانه‌های خطا
 ۵۷۶ ماتریس افزوده
 ۶۳۱ معادله پواسن
 ۵۸۵ مقیاس دهی
 ۷۰۴ دنباله استورم
 ۷۸ دنباله فیوناتچی
 ۷۰۹، ۷۰۱ دوران در صفحه
 ۶۴۸ رایانه‌های برداری
 ۶۴۹ رایانه‌های موازی
 ۵۲۸ رتبه
 ۴۸ روابط بازگشتی سه‌جمله‌یی
 ۷۰۷ روش QR
 ۷۰۸ آماده‌سازی مقدماتی
 ۷۱۱ با انتقال
 ۷۱۱ همگرایی
 ۱۴۲ روش استیفنس
 ۳۸۲ روش اویلر
 ۴۰۰، ۳۹۵ برآورد خطای مجانبی
 ۴۵۸، ۳۹۲ پایداری
 ۴۶۳ پسرو
 ۳۹۳ تحلیل خطای گرد کردن
 ۳۸۶ تحلیل همگرایی
- ۱۸۰ ناهمگرایی
 ۱۶۲ وارون
 ۲۱۰ چندجمله‌یی تکه‌ای
 ۲۱۱ چندمتغیره
 ۲۱۵ گویا
 ۲۰۲ مثلثاتی
 ۵۷۲ دستگاه خطی چگال
 ۴۴۹، ۳۹۹، ۳۸۱ دستگاه معادلات دیفرانسیل
 ۴۹۶
 دستگاه معادلات دیفرانسیل (دستگاه معادلات
 خطی را ببینید.)
 ۶۴۸، ۵۷۳ دستگاههای خطی تنگ
 ۷۳۵ مسأله ویژه مقدار
 ۶۳۱ معادله پواسن
 ۶۵۳ دستگاههای خطی مختلط
 ۱۱۹ دستگاههای غیرخطی
 ۱۲۸ روش نیوتن
 ۱۱۹ نظریه نقطه ثابت
 ۱۲۴، ۱۲۲ همگرایی
 دستگاههای غیرخطی (ریشه‌یابی را ببینید.)
 دستگاههای معادلات خطی ۵۷۲
 بحث در تجزیه QR ۶۹۵، ۷۳۱
 تجزیه LU ۵۷۷
 تحلیل خطا ۵۹۹
 تنگ ۵۷۳
 چگال ۵۷۲
 حذف گاوسی ۵۷۳
 صورتهای دیگر ۵۹۰
 حل عددی ۵۷۲
 حلپذیری ۵۲۸
 روش SQR ۶۳۰
 روش تصحیح مانده ۶۱۲

- خطای برشی ۳۸۴
 دستگاہ ۳۹۹
 کران خطا ۳۸۸، ۳۹۰
 مشتگیری ۳۸۴
 روش بارستی (ریشه‌یابی را ببینید).
 بارستی تک‌نقطه‌ای ۸۷
 دستگاہ خطی:
 حذف گاوسی ۶۲۹
 روش گاوس-زایدل ۶۲۱
 روش گاوس-ژاکوبی ۶۱۸
 روش گرادیان مزدوج ۶۳۸، ۶۴۲
 روشهای چندشبکه‌یی ۶۲۶
 معادلهٔ پواسن ۶۳۱
 دستگاہ معادلات غیرخطی ۱۱۹
 ریشه‌یابی چندجمله‌یها ۱۰۸-۱۱۹
 معادلهٔ دیفرانسیل ۴۱۳
 ویژه‌بردار ۶۸۲، ۷۱۴
 ویژه‌مقدار ۶۸۲
 روش برنت ۱۰۴
 محک همگرایی ۱۰۴
 مقایسه با روش نیم‌سازی ۱۰۶-۱۰۷
 روش پسرو اویلر ۴۶۳
 روش پیشگو-تصحیح‌کننده ۴۱۶-۴۲۰ (روشهای چندگامی را ببینید).
 روش تصحیح‌مانده ۶۱۲
 کران خطا ۶۱۶
 همگرا ۶۱۵
 روش توانی ۶۸۲
 برون‌یابی ایتکن ۶۸۷
 روشهای شتاب ۶۸۶
 کاهش ۶۸۹، ۷۴۲
 همگرایی ۶۸۴
- روش چندگامی صریح ۴۰۲
 روش چولسکی ۵۹۲، ۷۲۷
 روش خط قاطع ۷۵
 فرمول خطا ۷۶
 مقایسه با روش نیوتن ۸۱
 همگرایی ۷۹
 روش خطوط ۴۷۰
 روش صریح ۴۷۲
 معادلهٔ گرما ۴۷۰
 روش دولیتل ۵۹۲
 روش دوزنقه‌یی:
 انتگرالگیری عددی ۴۱۲
 انتگرالده دوره‌یی ۳۲۵
 برون‌یابی ریچاردسن ۳۳۲
 فرمول اویلر-مکلورن ۳۲۲
 فرمول خطا ۲۸۳
 فرمول خطای مجانبی ۲۸۶
 قاعده جاصلضرب ۳۵۲
 قاعدهٔ دوزنقه‌یی تصحیح‌شده ۲۸۶
 محاسبهٔ مرتبهٔ دوم گاوسی ۳۱۶
 مرکب ۲۸۵
 هسته‌یی پثانو ۲۹۱
 معادلات دیفرانسیل ۴۱۳
 A-پایداری ۴۱۸
 برآورد ریچاردسن ۴۱۹
 پایداری ۴۱۷
 خطای کلی ۴۲۷
 خطای مجانبی ۴۱۷
 خطای موضعی ۴۱۵
 راه‌حل بارستی ۴۱۴
 نواحی پایداری ۴۶۳
 همگرایی ۴۱۷

- روش رونگه-کوتا ۴۷۷
 برآورد خطا ۴۸۵
 برآورد خطای ریچاردسن ۴۸۱-۴۸۲
 برنامه‌های خودکار ۴۹۰
 پایداری ۴۸۵
 خطای برشی ۴۷۷
 خطای کلی ۴۹۲
 خطای مجانبی ۴۸۴
 روشهای ضمنی ۴۹۲
 روشهای فلبرک ۴۸۷
 سازگاری ۴۸۳
 فرمول کلاسیک مرتبه چهار ۴۸۱
 فرمولهای مراتب بالاتر ۴۸۰
 مشتقات ۴۷۹
 همگرایی ۴۸۴
 روش سری تیلر، معادلات دیفرانسیل ۴۷۶
 روش ضرایب نامعین ۳۵۹
 روش ضمنی چندگامی ۴۰۲
 راه حل بارستی ۴۱۴، ۴۳۰
 روش کروت ۵۹۲
 روش گاوس-یوردان ۵۹۰
 وارونیایی ماتریس ۵۹۱
 روش گرادیان مزدوج ۱۳۲، ۶۳۸، ۶۴۲
 بهینگی ۶۴۳
 شتابی ۶۳۰
 قضیه همگرایی ۶۴۳، ۶۴۴
 روش مولر ۸۳
 روش میانگامی، معادلات دیفرانسیل ۴۰۷
 خطای کران ۴۰۶
 ضعیف-پایدار ۴۱۱
 فرمول خطای مجانبی ۴۰۸
 معادله مشخصه ۴۱۰
 روش میلن ۴۳۴
 روش نلدرو مید ۱۳۲
 روش نیم‌ساز ۶۴-۶۵
 همگرا ۶۴-۶۵
 روش نیوتن ۶۷
 برآورد خطا ۷۳
 چندجمله‌ای ۱۱۱
 دستگاههای غیرخطی ۱۲۵
 روش نیوتن-فوریه ۷۱
 ریشه‌های چندگانه ۱۰۱
 ریشه‌های دوم ۱۳۸
 مسأله مقدار مرزی ۴۹۹
 مقایسه با روش خط قاطع ۸۱
 همگرایی ۶۹
 روش نیوتن-فوریه ۷۱
 روش همگرایی خطی ۶۴
 روش هم‌مکانی ۵۰۵
 روش هورنر ۱۱۱
 روش-گاوس-زایدل ۶۲۳
 شتابی ۶۳۰
 همگرایی ۶۲۰، ۶۲۴
 روشهای آدامز ۴۳۵
 پایداری ۴۵۸
 متغیر-مرتبه ۴۴۰
 نواحی پایداری ۴۶۲
 روشهای آدامز-بشفورت ۴۳۵
 روشهای آدامز-مولتن ۴۳۷
 روشهای امتداد مزدوج ۶۴۰
 روشهای بارستی تک‌نقطه‌ای ۸۷
 دستگاه غیرخطی ۱۱۹
 روش نیوتن ۶۷
 قضیه همگرایی ۸۸-۹۵

- مرتبه همگرایی ۴۰۴
 معادلات دیفرانسیل سرسخت ۴۶۴
 معادله نمونه ۴۴۹
 ناحیه‌های پایداری ۴۵۸
 نسبی-پایدار ۴۵۷
 هسته پتانو ۴۳۱
 همگرایی ۴۵۴، ۴۰۵
 روشهای حصارکشی (خصر) در ریشه‌یابی ۶۷
 روش برنت ۱۰۴
 روشهای رنگه-کوتا-فلبرک ۴۸۷
 روشهای شبه‌نیوتنی ۱۳۲، ۱۲۸
 روشهای شتاب:
 انتگرالگیری عددی ۲۸۶، ۳۳۲
 ویژه‌مقدارها ۶۸۶
 روشهای فشرده ۵۹۱
 ریاضیات نمادی ۵۰
 ریشه اصلی ۴۵۰
 ریشه‌های چندجمله‌یها ۱۰۸-۱۱۹
 بد وضع ۱۱۵
 پایداری ۱۱۲
 تقلیل ۱۱۱، ۱۱۷
 روش نیوتن ۱۱۲
 ماتریس همراه ۷۳۹
 محدود کردن ریشه‌ها ۱۰۹
 ریشه‌های چندگانه ۹۹
 اثر نوفه ۱۰۰
 بازه عدم قطعیت ۱۰۰
 تعیین چندگانگی ۱۰۳
 روش نیوتن ۱۰۱
 ناپایداری ۱۱۷
 ریشه‌های مشخصه ۴۵۰
 ریشه‌یابی ۶۰
- معادله دیفرانسیل ۴۶۹، ۴۱۳
 همگرایی خطی ۶۴
 همگرایی مرتبه بالاتر ۹۴-۹۵
 روشهای برونابی:
 انتگرالگیری عددی ۳۳۲
 معادلات دیفرانسیل ۵۰۶
 روشهای پرتابی ۴۹۷
 روشهای تک‌گامی ۴۷۵ (روشهای رنگه-کوتا را ببینید.)
 روشهای چندگامی ۴۰۱
 آدامز-بشفورت ۴۳۵
 آدامز-مولتن ۴۳۷
 انتگرالگیری عددی ۴۳۴
 پایداری ۴۰۷، ۴۴۷
 پیدا کردن ۴۲۹
 جواب مزاحم ۴۱۱، ۴۵۵
 خطای برشی ۴۰۲
 خطای کران ۴۰۶
 راه حل بارستی ۴۱۴، ۴۳۰
 روش دوزنقه‌ی ۴۱۲
 روش ضرایب نامعین ۴۳۰
 روش میانگامی ۴۰۷
 روش میلن ۴۳۴
 روشهای آدامز ۴۳۵
 شرط ریشه ۴۴۷
 شرط سازگاری ۴۰۴، ۴۴۷
 شرط قوی ریشه ۴۵۸
 شکل کلی ۴۰۱
 صریح ۴۰۲
 ضمنی ۴۰۲، ۴۳۰
 متغیر-مرتبه ۴۴۰، ۵۰۶
 مثال ناپایدار ۴۴۷

شکلهای متعارف ۵۳۶	الگوریتم برنت ۱۰۴
تجزیهٔ تکین-مقدار ۵۴۱	برآورد خطا ۷۳، ۸۰
ژوردان ۵۴۴	برونایی ایتکن ۹۶
شور ۵۳۷	بهینه‌سازی ۱۲۹
ماتریسهای مقارن ۵۴۱	چندجمله‌یها ۱۰۸-۱۱۹
شمارش عملیات ۶۵۴	دقیق ۹۷
حذف گاوسی ۵۷۹	روش استیفنس ۱۴۲
صرفه‌جویی سری تیلر ۲۷۹	روش حصارکشی (خصر) ۶۷
صفرها (ریشه‌یابی را ببینید.)	روش خط قاطع ۷۵
صورت درجهٔ دوم ۵۶۵	روش مولر ۸۳
ضرایب دوجمله‌ای ۱۷۰	روش نیم‌سازی ۶۴
ضرب تودرتو ۱۱۱	روش نیوتن ۶۷، ۱۲۵
ضریب وضعیت ۴۳	روشهای بارستی تک‌نقطه‌ای ۸۷
قضیهٔ گاستیل ۶۰۳	ریشه‌های چندگانه ۹۹
ماتریس ۶۰۰	معادلات غیرخطی ۱۱۹
ماتریس هیلبرت ۶۰۴	زاویهٔ بین بردارها ۵۳۱
محاسبه ۶۰۹	زیر رایانه‌ها ۴۹
ویژه‌مقدار ۶۷۲، ۶۷۷	زیر روالهای اساسی جبرخطی ۵۹۰، ۶۴۷
ضریبها ۵۷۵	سری فوریه ۲۰۶
عکس نابرابری مثلثی ۲۲۷	سری همگرا ۲۷۲
عنصر محور ۵۸۲	سری هندسی ۶
غالب قطری ۶۱۹	نرم ماتریسی ۵۵۸
فرمول استرلینگ ۳۱۴	شانزده‌شانزدهی ۱۳
فرمول اویلر-مک‌لورن ۳۲۲	شرایط مور-ینروز ۷۴۶
تعمیم ۳۲۷	شرط ریشه ۴۴۷
فرمول مجموعیابی ۳۲۶	شرط سازگاری ۴۰۲، ۴۴۷
فرمول تصحیح‌کننده ۴۱۶	روش رونگه‌کوتا ۴۸۳
فرمول خطای مجانبی:	شرط قوی ریشه ۴۵۷
انتگرالگیری عددی ۲۸۶، ۳۲۰	شرط لیپشیتس ۳۷۶، ۴۰۰، ۴۸۳
تعریف ۲۸۶	شعاع طیفی ۵۵۰
روش اویلر ۳۹۵	شکل متعارف ژوردان ۵۴۴
روش رونگه‌کوتا ۴۸۴	شکل نرمال شور ۵۳۷

- ۱۶۵ قضیه ارمیت-جنوکی
 ۲۲۵، ۵ قضیه تیلر
 ۵ بسطهای تیلر
 ۸ فضای دوبعدی
 ۶ قضیه سری هندسی
 ۲۵۵، ۲۰۶ قضیه جکسن
 ۲۵۳ قضیه دولاوله یواسن
 ۵۶۸ قضیه کیلی-همیلتن
 ۶۰۳ قضیه گاستینل
 ۲۷۵، ۲۳۷ قضیه گرام-اشمیت
 ۶۶۵ قضیه گرشگورین
 ۵۳۹ قضیه محورهای اصلی
 ۴ قضیه مقدار میانگین
 ۴۰ قضیه مقدار میانی
 ۲۲۴ قضیه وایرستراس
 ۶۷۳ قضیه ویلانت-هوفمن
 ۲۵۴ قضیه هموسانی چیشف
 ۱۴ قطع کردن
 ۱۷۴ کشف نوفه در داده‌ها
 ۵۲۵ ماتریس
 ۵۲۸ ارمیتی
 ۵۴۴ پوچتوان
 ۵۷۷ تجزیه LU
 ۵۴۱ تجزیه تکین-مقدار
 ۶۶۲، ۵۶۶ تصویر
 ۵۸۵ جایگشت
 ۵۲۸ رتبه
 ۵۹۶ سه قطری
 ۵۴۴ شکل متعارف زوردان
 ۵۳۷ شکل نرمال شور
 ۵۳۶ شکلهای متعارف
- ۳۲۶، ۳۲۲ فرمول اوپلر-مکلورن
 ۲۸۶ قاعده دوزنقه‌یی
 ۲۹۰ قاعده سیمپسون
 ۴۱۵، ۴۰۹ معادلات دیفرانسیل
 ۱۶۰ فرمول درونیابی تفاضل منقسم
 ۱۵۲ فرمول لاگرانژ درونیاب
 ۴۱۶ فرمولهای پیشگو
 ۳۱۹ فرمولهای کرونرود
 ۴۶۵ فرمولهای مشتقگیری پسر
 ۵۲۳ فضای برداری
 ۵۲۵ بعد
 ۵۲۵ پایه
 ۵۳۱ پایه متعامد
 ۵۲۹ حاصلضرب داخلی
 ۴۱۱ ضعیف-پایدار
 ۲۹۸ قاعده بول
 ۲۸۹ قاعده ترکیبی سیمپسون
 ۳۶۶، ۲۸۶ قاعده دوزنقه‌یی تصحیح شده
 ۲۸۵ قاعده دوزنقه‌یی مرکب
 ۲۸۷ قاعده سیمپسون
 ۳۲۹ برونیابی ایتکن
 ۳۳۲ برونیابی ریچاردسن
 ۲۸۹ ترکیبی
 ۳۳۸ تطبیقی
 ۲۹۰ فرمول خطای جانبی
 ۲۹۱، ۲۸۸ فرمولهای خطا
 ۳۵۳ قاعده حاصلضرب
 ۴۳۴ معادلات دیفرانسیل
 ۲۹۲ هسته پتانو
 ۵۸۱ قاعده کرامر
 ۳۹۶ قاعده مرتبه دوم سه هشتم
 ۳۱۹ قاعده‌های پاترسن

تشابه با ماتریس قطری ۵۳۹	صفر ۵۲۷
خارج قسمت ریلی-ریس ۶۸۸، ۷۴۱	ضریب وضعیت ۶۰۱، ۶۰۰
دنباله استورم ۷۰۴	عملیات ۵۲۶
روش QR ۷۰۷	غالب قطری ۶۱۹
روش زاکوبی ۷۳۵	قضیه اختلالها ۵۶۱
قضیه ویلانت-هوفمن ۶۷۳	قضیه سری هندسی ۵۵۸
کاهش ۷۴۲	قضیه محورهای اصلی ۵۳۹
کران خطا ویژه مقدارها ۶۷۳	کاهش ۶۸۹
محاسبه ویژه بردارها ۷۱۴	متعامد ۵۳۱ (تبدیلات متعامد را ببینید.)
محاسبه ویژه مقدارها ۷۰۲، ۷۰۷	مقارن ۵۲۸
معین مثبت ۵۶۷، ۶۵۴	مرتبه ۵۲۶
ویژه مقدارها ۵۳۹	مشابه ۵۳۶
ماتریس معین مثبت ۵۶۷، ۵۹۲	معادله مشخصه ۵۳۳
روش چولسکی ۵۹۲	معین مثبت ۵۶۷، ۵۹۲
ماتریس نواری ۵۹۶	نرم ۵۴۵ (نرم ماتریس، نرم برداری را ببینید.)
ماتریس وارون ۵۲۸	نواری ۵۹۶
روش بارستی ۶۶۰	وارون ۵۲۸
عمل کاهش ۵۸۱	هاؤسهولدر ۶۹۰
کران خطا ۶۱۰	همانی ۵۲۷
محاسبه ۵۹۲	هیلبرت ۴۵، ۲۳۵، ۶۰۴
ماتریس واندرموند ۱۵۰، ۲۱۳	یکانی ۵۳۱، ۵۶۶
ماتریس هرمیتی ۵۲۸ (ماتریس مقارن را ببینید.)	ماتریس افزوده ۵۷۶
ماتریس هسنبگ ۷۰۹	ماتریس پادمقارن ۵۲۸
ماتریس هیلبرت ۴۵، ۲۳۵، ۶۰۴	ماتریس پوچ توان ۵۴۴
ضرایب وضعیت ۶۰۴	ماتریس تبدیل پایه ۵۳۶
مقدار ویژه ۶۷۰	ماتریس تصویر ۵۶۶، ۶۶۲
ماتریس یکانی ۵۳۱، ۵۶۶ (تبدیلات متعامد را ببینید.)	ماتریس جایگشت ۵۸۵
ماتریسهای دوران در صفحه ۱، ۷۰۱، ۷۴۳	ماتریس چگال ۵۷۲
ماتریسهای مشابه ۵۳۵	ماتریس زاکوبی ۱۲۱، ۴۰۰
ماتریسهای هاؤسهولدر ۶۹۰، ۷۴۲	ماتریس مقارن ۵۲۸
پیدا کردن ماتریس مقارن ۶۹۷	پایداری ویژه مقدارها ۶۷۴
	پیدا کردن ماتریس سه قطری ۶۹۷

- تبدیل بردار ۶۹۲
تجزیه QR ۶۹۴
ماتریسهای همراه ۷۳۹
متعامد ۲۳۷ (تبدیلات متعامد را ببینید.)
پایه ۵۳۱
خانواده ۲۴۰
ماتریس ۵۳۱
متناهی-بعد ۵۲۵
مجموعهٔ دوه‌دو متعامد ۶۷۵
محاسبهٔ ریشه‌های دوم ۱۳۸
محورگیری ۵۸۲
محورگیری جزئی ۵۸۲
محورگیری کلی ۵۸۲
مدل‌سازی ریاضی ۲۱
مسئلهٔ دیریکله (معادلهٔ پواسن را ببینید.)
مسئلهٔ مقدار اولیه ۳۷۳
مسائل بدحالت ۴۲، ۳۷۲
مسائل بدوضع ۴۴
چندجمله‌ای ۱۱۳
دستگاه خطی ۶۰۲
مسائل معکوس ۴۹
معادلات دیفرانسیل ۳۸۱
مقادیر ویژه ۶۷۷
مسائل خوش حالت ۴۲ (پایداری را ببینید.)
مسائل مقدار مرزی ۴۹۳
روش هم‌مکانی ۵۰۵
روشهای پرتابی ۴۹۷
روشهای تفاضل متناهی ۵۰۱
روشهای معادلهٔ انتگرالی ۵۰۵
مسائل ناپایدار ۴۲ (مسائل بدحالت را ببینید.)
مشتقگیری عددی ۳۵۶
بدحالتی ۳۷۲
- پایهٔ خطای درونیایی ۳۵۶
خطای فرمول ۳۵۸
ضرایب نامعین ۳۵۹
معادلات دیفرانسیل سرسخت ۴۶۴
روش پسرو اویلر ۴۶۳
روش خطوط ۴۷۰
روشهای A-پایداری ۴۱۸، ۴۶۳، ۴۶۷
فرمول بارستی ۴۶۸
فرمولهای مشتقگیری پسرو ۴۶۵
قاعدهٔ دوزنقه‌ی ۴۶۷
معادلات دیفرانسیل مراتب بالا ۳۸۲
معادلات دیفرانسیل:
بارست بیکار ۵۱۳
بدوضع ۳۸۰
برنامه‌های خودکار ۴۹۰
برنامه‌های مقدار مرزی ۵۰۵، ۵۰۸
خطا در طول گام واحد ۴۲۱
خطای کلی ۴۴۲
روش آدامز ۴۴۱
کنترل خطا ۴۴۲
کنترل خطای موضعی ۴۲۱، ۴۴۲
متغیر-مرتب ۴۴۰
مقایسه ۵۰۸
پایداری ۳۷۸
حل عددی:
A-پایداری ۴۱۸، ۴۶۳، ۴۶۷
انتگرالگیری عددی ۴۳۴
اندازهٔ گامها ۳۸۳
برآورد خطای فراگیر (کلی) ۴۱۹، ۴۴۳
پایداری ۳۹۲، ۴۰۷، ۴۴۷
حل موضعی ۴۱۴
روش اویلر ۳۸۲

میدان سوهایی ۳۷۴	روش پرسرو اوایلر ۴۶۳
نظریه وجود ۳۷۶	روش خطوط ۴۷۰
معادله انتگرالی ۵۱۳، ۵۰۵	روش دوزنقه‌یی ۴۱۲
معادله پواسن ۶۳۱	روش سری تیلر ۴۷۵
تعمیم ۶۶۱	روش ضمنی ۴۰۲
تقریب تفاضل منتهای ۶۳۱-۶۳۲	روش میانگاهی ۴۰۷
روش SQR ۶۳۰، ۶۳۶	روشهای آدامز ۴۳۵
روش گاوس-زایدل ۶۲۷، ۶۳۵	روشهای برونایی ۵۰۷
معادله تفاضلی خطی ۴۰۹، ۴۵۰	روشهای تک‌گامی ۴۷۵
معادله دیفرانسیل ۳۷۳، ۴۵۱	روشهای چندگامی ۴۰۱
معادله دیفرانسیل جزئی ۴۷۰، ۶۳۱	روشهای رونگه-کوتا ۴۷۷
معادله گرما ۴۷۰	روشهای صریح ۴۰۲
معادله مشخصه:	روشهای متغیر-مرتب ۴۴۰، ۵۰۶
ماتریسها ۵۳۳	ضرایب نامعین ۴۲۹
معادلات دیفرانسیل ۴۰۹، ۴۵۰	فرمول پیشگو ۴۱۶
معادله نرمال ۷۲۶	فرمول تصحیح‌کننده ۴۱۶
معکوس تعمیم‌یافته ۷۲۳، ۷۴۶	مسائل سرسخت ۴۶۴
مقیاس دهی ۵۸۵	مسائل مقدار مرزی ۴۹۳
میانگین مربع خطاها ۲۳۲، ۷۲۶	معادله مشخصه ۴۵۰
میدان سوهایی ۳۷۴	معادله نمونه ۴۴۹
نابرابری کوشی-شوارتس ۲۳۶، ۵۳۰	ناحیه‌های پایداری ۴۵۸
نابرابری مثلثی ۱۱، ۲۲۷، ۵۳۰	نظریه همگرایی ۴۰۵، ۴۵۴
نابایداری ۴۲ (مسائل بدوضع را ببینید.)	نقاط گره‌یی ۳۸۳
ناحیه پایداری مطلق ۴۶۰	خطی ۳۷۴، ۳۸۲
ناحیه‌های پایداری ۴۵۸	دستگاه ۳۸۱
نامتناهی-بعد ۵۲۵	روشهای معادله انتگرالی ۵۰۵
نامساوی بسل ۲۴۷	سرسخت ۳۸۱، ۴۶۴
نرم ۵۴۵	مراتب بالاتر ۳۸۱
نرم بینهایت ۲۲۶	مرتبه اول خطی ۳۷۴
نرم چیشف ۲۲۶	مسأله مقدار اولیه ۳۷۳
نرم دو ۲۳۶	مسأله نمونه ۴۴۹
نرم ستونی ۵۵۳	مسائل مقدار مرزی ۴۹۳

نقاط گره‌یی ۳۸۳، ۲۸۱	نرم سطری ۵۵۴
نقطه ثابت ۸۷	نرم عملگر ۵۵۰
نما ۱۳	نرم فروبنیوس ۵۴۹
نمایش با ممیز شناور ۱۴	نرم ماکسیم ۵۴۵، ۲۲۶
پی‌ریز ۱۹	نرم یکنواخت ۲۲۶
تبدیل ۱۳	نرم افزار ریاضی ۷۵۰، ۵۰
جزء کسری ۱۳	نرم اقلیدسی ۲۳۶
دقت ۱۷	نرمها ۵۴۵
سرریز ۱۹	اقلیدسی ۵۲۹، ۲۳۶
قطع کردن ۱۴	بردار ۵۴۶
گرد کردن ۱۴	پیوسته ۵۴۶
ممیز اعشاری ۱۴	دو ۲۳۶
نما ۱۳	سازگار ۵۴۹
واحد گرد کردن ۱۷	شعاع طیفی ۵۵۰
نوفه در محاسبه تابع ۱۰۰، ۲۵	عملگر ۵۵۰
نیم‌ساز ۶۵	فروبنیوس ۵۴۹
نیوتن ۱۱۲	ماتریسی ۵۴۹
وزنها ۲۸۱	ماکسیم ۵۴۵، ۲۲۶
ویژه‌بردار ۵۳۳	هم‌ارزی ۵۴۷
تقریب عددی:	یکنواخت ۲۲۶
بارست معکوس ۷۱۴	نرمهای برداری ۵۴۵، ۱۲
پایداری ۶۷۹	D-نرم ۵۴۵
روش توانی ۶۸۲	پیوستگی ۵۴۶
کران خطا ۷۸۱	ماکسیم ۵۴۵، ۲۲۶
ویژه‌مقدارها ۵۳۳	هم‌ارزی ۵۴۷
بدوضع ۶۷۸	نرمهای ماتریسی ۵۴۹
پایداری ۶۶۹	سازگار ۵۵۰
تحت تبدیلات یکانی ۶۷۹	ستونی ۵۵۳
تقریب عددی:	سطری ۵۵۴
دنباله استورم ۷۰۴	عملگر ۵۵۰
روش QR ۷۰۷	فروبنیوس ۵۴۹
روش توانی ۶۸۲	نظریه اختلال چندجمله‌ای ۱۱۲-۱۱۳

- | | |
|--|---|
| هسته‌یی پتانو ۳۶۷، ۳۱۵، ۲۹۱ | روش ژاکوبی ۷۳۵ |
| معادلات دیفرانسیل ۴۳۱-۴۳۲ | ماتریسهای تنک ۷۳۵ |
| همگرایی خطی ۶۴ | جایابی ۶۶۴ |
| شتاب دادن ۹۶ | چندگانگی ۵۳۵ |
| نرخ ۶۴ | حل عددی (مقادیر ویژه تقریب عددی را ببینید.) |
| همگرایی مربعی ۶۴ | ضرایب وضعیت ۶۷۲، ۶۸۰ |
| همگرایی مرتبه ۶۴ | قضیهٔ باؤتر-فیکه ۶۶۹ |
| همگرایی: | قضیهٔ گرشگورین ۶۶۵ |
| بازه ۶۴ | کاهش ۶۸۹ |
| بردار ۵۴۸ | ماتریسهای دارای شکل ناقطری ژوردان ۶۸۰ |
| خطی ۶۴ | ماتریسهای سه‌قطری ۷۰۲ |
| مربعی ۶۴ | ماتریسهای متقارن ۷۰۲ |
| مرتبه ۶۴ | یک کران خطا برای ماتریسهای متقارن ۶۷۳ |
| نرخ همگرایی ۶۴ | ویژه‌مقدارهای یک ماتریس سه‌قطری ۷۰۲ |
| $\text{Cond}(A)_p, \text{Cond}(A)_*$ ۶۰۱ | دنبالهٔ استورم ۷۰۴ |
| $\text{Cond}A$ ۶۰۰ | رشو هاؤسهولدر ۷۰۲ |
| Trace ۵۳۴ | روش گیونز ۷۰۲ |