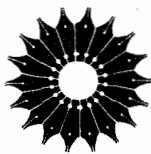


# آنالیز عددی مقدماتی

به شیوہ الگوریتمی

سموئل د. کونت، کارل دوپور

ترجمہ سراج الدین کاتبی



# آنالیز عددی مقدماتی

به شیوه الگوریتمی

سموئل د. کونت، کارل دوبور

ترجمه سراج الدین کاتبی

## بسم الله الرحمن الرحيم

### فهرست

صفحه	عنوان
۱	پیشگفتار
۳	مقدمه
۵	<b>فصل اول دستگاه اعداد و خطاها</b>
۵	۱۰۱ نمایش اعداد صحیح
۱۰	۲۰۱ نمایش کسرها
۱۳	۳۰۱ حساب با ممیز شناور
۲۰	۴۰۱ از دست دادن ارقام با معنی و بخش خطا، شرط و ناپایداری
۲۷	۵۰۱ روشهای کامپیوتری برای برآورد خطا
۲۹	۶۰۱ توضیحاتی در مورد همگرایی دنباله‌ها
۳۷	۷۰۱ برخی مقدمات ریاضی
۲۲	<b>فصل دوم درونیایی به وسیله بسجمله‌ایها</b>
۲۲	۱۰۲ صورتهای بسجمله‌ای
۵۳	۲۰۲ وجود و یکتایی بسجمله‌ای درونیاب
۵۸	۳۰۲ جدول تفاضل منقسم
۶۴	۴۰۲ درونیایی در یک عدد رو به افزایش از نقاط درونیایی
۶۹	۵۰۲ خطای بسجمله‌ای درونیاب
۷۵	۶۰۲ درونیایی یک تابع جدولی بر اساس نقاط با فواصل مساوی
۸۴	۷۰۲ تفاضل منقسم به عنوان تابعی از شناسه‌ها و درونیایی بوسانی‌اش

## فصل سوم حل معادلات غیرخطی

۹۶	
۹۸	۱.۳ بررسی اجمالی روشهای بارستی
۱۰۷	۲.۳ برنامه‌های فورترن برای بعضی از روشهای بارستی
۱۱۴	۳.۳ بارست نقطه ثابت
۱۲۳	۴.۳ شتاب همگرایی برای بارست نقطه ثابت
۱۲۸	۵.۳ همگرایی نیوتن و روشهای خط قاطع
۱۴۱	۶.۳ معادلات بسجمله‌ای: ریشه‌های حقیقی
۱۵۴	۷.۳ ریشه‌های همتافت و روش مولر

## فصل چهارم ماتریسها و دستگاههای معادلات خطی

۱۶۴	
۱۶۴	۱.۴ ویژگیهای ماتریسها
۱۸۹	۲.۴ حل دستگاههای خطی از راه حذف
۲۰۲	۳.۴ تدبیر لولاگزینی
۲۰۶	۴.۴ تجزیه به عوامل مثلثی
۲۱۷	۵.۴ خطا و باقیمانده یک جواب تقریبی؛ نرم‌ها
۲۲۷	۶.۴ تحلیل خطای پسرو و اصلاح بارستی
۲۳۸	۷.۴ دترمینانها
۲۴۳	۸.۴ مسئله ویژه مقدارها

## فصل پنجم دستگاه معادلات و بهینه‌سازی نامقید

۲۶۹	
۲۷۰	۱.۵ بهینه‌سازی و تندترین کاهش
۲۷۹	۲.۵ روش نیوتن
۲۸۹	۳.۵ بارست نقطه ثابت و روشهای واهلش

## فصل ششم تقریب

۳۰۵	
۳۰۵	۱.۶ تقریب یکنواخت به وسیله بسجمله‌ایها
۳۱۹	۲.۶ برازاندن داده‌ها
۳۲۶	۳.۶ بسجمله‌ایهای متعامد
۳۳۷	۴.۶ تقریب با روش کوچکترین مربعات به وسیله بسجمله‌ایها
۳۴۸	۵.۶ تقریب به وسیله بسجمله‌ایهای مثلثاتی
۳۶۰	۶.۶ تبدیلهای سریع فوریه



۷-۶ تقریب به وسیله بسجمله‌ای-نکته‌ای

۳۶۹

### فصل هفتم مشتقگیری و انتگرالگیری

۳۸۰

۱-۷ مشتقگیری عددی

۳۸۲

۲-۷ انتگرالگیری عددی: برخی قواعد اساسی

۳۹۲

۳-۷ انتگرالگیری عددی: قواعد گاوسی

۴۰۲

۴-۷ انتگرالگیری عددی: قاعده‌های مرکب

۴۱۳

۵-۷ انتگرال تطبیقی

۴۲۵

۶-۷ برونیابی به سمت حد

۴۳۳

۷-۷ انتگرالگیری رامبرگ

۴۴۱

### فصل هشتم حل معادلات دیفرانسیل

۴۴۹

۱-۸ مقدمات ریاضی

۴۴۹

۲-۸ معادلات تفاضلی ساده

۴۵۳

۳-۸ انتگرالگیری عددی به وسیله سری تیلر

۴۶۰

۴-۸ برآوردهای خطا و همگرایی روش اویلر

۴۶۷

۵-۸ روش رونگه-کوتا

۴۷۱

۶-۸ کنترل اندازه  $h$  در روشهای رونگه-کوتا

۴۷۶

۷-۸ فرمولهای چندمرحله‌ای

۴۸۴

۸-۸ روشهای پیشگو-تصحیحی

۴۹۲

۹-۸ روش ادمز-مولتن

۴۹۶

۱۰-۸ پایداری روشهای عددی

۵۰۵

۱۱-۸ کنترل و بخش خطای گرد کردن

۵۱۲

۱۲-۸ دستگاههای معادلات دیفرانسیل

۵۱۶

۱۳-۸ معادلات دیفرانسیل سخت

۵۲۱

### فصل نهم مسائل با مقدار مرزی در معادلات دیفرانسیل معمولی

۵۲۶

۱-۹ روشهای تفاضل متناهی

۵۲۷

۲-۹ روشهای نشانه‌گیری

۵۳۳

۳-۹ روشهای پهلوی هم‌گذاری

۵۳۹

ضمیمه: کتابهای زیر برنامه‌ای

۵۴۵

منابع

۵۴۷

فهرست اسامی خاص

۵۴۹

واژه‌نامه

۵۵۰

## پیشگفتار

این سومین چاپ کتاب آنالیز عددی مقدماتی است که مخصوصاً برای نیاز دانشجویان سالهای دوره کارشناسی در رشته‌های مهندسی، ریاضی و علوم و به‌ویژه رشته علوم کامپیوتر طرح‌ریزی شده است. به‌طور کلی دانشجویانی که دروس حساب دیفرانسیل و انتگرال را به‌طور مرتب با پایه محکمی گذرانیده باشند، در دنبال نمودن مطالب این کتاب مشکلی نخواهند داشت. هنگام به‌کار بردن مفاهیم پیشرفته ریاضی مانند نرمها و تعامد<sup>۲</sup> دقت شده است تا عرضه مطالب در سطحی مناسب برای دانشجویان دوره کارشناسی، با فرض اینکه دانشجوی هیچ‌گونه آشنایی با آنها ندارد، صورت گیرد. در فصل مربوط به معادلات، آشنایی مختصری با ماتریسها و در فصلهای ۸ و ۹ آشنایی با معادلات دیفرانسیل مسلم انگاشته شده است. در این چاپ بخشهایی وجود دارند که نسبت به چاپهای قبلی به کمال ریاضی بیشتری نیاز دارند. ولی همه این بخشها با «ستاره نماد» مشخص شده‌اند و استاد می‌تواند بی‌آنکه پیوستگی مطالب را از دست دهد، آنها را حذف کند.

این چاپ جدید متضمن مطالبی تازه و تغییراتی قابل ملاحظه نسبت به چاپ قبلی است. ترتیب فصول به‌گونه‌ای که فکر می‌کنیم طبیعیتر است درآمده است. درونبایی بسجمله‌ای اکنون حتی مقدم بر فصل مربوط به حل دستگاههای غیرخطی (فصل ۳) عرضه شده است و بعداً برای برخی مطالب مذکور در تمام فصول به‌کار رفته است.

روش حذف گاوس (فصل ۴) ساده شده است. بعلاوه در فصل ۴، اکنون استفاده دامنه‌داری از تحلیل خطای پسر و ویلکینسن<sup>۳</sup> به عمل آمده، و شرح مختصری از روشهای بسیار معروف برای مسئله ویژه-مقدارها و ویژه-بردارها گنجانده شده است. فصل ۵، فصل تازه‌ای است درباره دستگاه معادلات و بهینه‌سازی نامقید و متضمن مقدمه‌ای است بر روشهای تندترین کاهش، روش نیوتن برای دستگاه معادلات غیرخطی و روشهای واهلش<sup>۴</sup> حل دستگاههای بزرگ خطی باروش بارستی. فصل مربوط به تقریب (فصل ۶) بسط داده شده

- 
1. norms
  2. orthogonality
  3. Wilkinson
  4. relaxation

است. در این فصل بهترین تقریب و تقریب خوب به توسط بسجمله ایها، و نیز تقریب به کمک توابع مثلثاتی از جمله تبدیلیهای سریع فوریه، برآزاندن دادهها با روش کمترین مربعات، بسجمله ایهای متعامد و برآزاندن منحنی با روش قلمی<sup>۲</sup> مورد بررسی قرار گرفته است. مشتقگیری و انتگرالگیری در فصل ۷، که شامل بخش جدیدی در باب انتگرالگیری تطبیقی است، آمده است. فصل ۸ در زمینه معادلات دیفرانسیل معمولی است که مقدار زیادی مطالب تازه و چندین بخش جدید به آن افزوده شده است. بخش جدیدی در زمینه کنترل اندازه- مرحله در روشهای رونکه- کوتا، و بخش جدیدی درباره معادلات دیفرانسیل سخت، و نیز بخشی در زمینه ناپایداری عددی که به مقدار زیادی تجدید نظر شده، در این فصل آمده است. فصل ۹ حاوی مقدمه ای است کوتاه در باب پهلوی هم گذاری<sup>۳</sup> به عنوان روشی برای حل مسائل مقدار مرزی.

در این چاپ، همانند چاپ قبلی، فرض بر این است که دانشجویان به یک کامپیوتر دسترسی، و با برنامه سازی در برخی از زبانهای شیوه- پرداخته<sup>۴</sup> آشنایی دارند. تعداد زیادی الگوریتم در متن گنجانیده شده و برای بسیاری از آنها برنامه های فورترن ارائه شده است. در این چاپ برنامه های کامل تاحدی کمتر آورده شده است. همه برنامه ها به زبان فورترن ۷۷ که در آن از مفهوم برنامه نویسی ساختیافته جدید استفاده می کنند، بازنویسی شده است. تمامی برنامه ها روی یک یا چند کامپیوتر امتحان، و در بسیاری از این حالات بروندادهای حاصل از کامپیوتر عرضه شده است. زمانی که بروندهای عددی داده می شوند، خود متن مشخص خواهد ساخت که برای به دست آوردن این نتایج از کدام نوع کامپیوتر (یونیوک، سی. دی. سی، آی. بی. ام<sup>۵</sup>) استفاده شده است.

مطالب این کتاب بیش از آنی است که معمولاً در یک ترم تحصیلی عادی دوره کارشناسی برای رشته های علوم تدریس می شود. این امر در طرح و چگونگی ارائه درس، آزادی بیشتری به معلم خواهد داد. به همین دلیل توجه به این موارد مهم است که مطالبی که در زمینه درونیایی بسجمله ای در فصل ۲، دستگاههای خطی در فصل ۴ و مشتقگیری و انتگرالگیری در فصل ۷ بیان می شوند مسلماً برای فصول بعدی ضروری هستند. مطالبی که در هفت فصل اول کتاب (به غیر از قسمتهای ستاره دار) آمده اند ساختار مطلوبی برای اولین دوره این درس را تشکیل می دهند.

در اینجا فرصت را مغتنم شمرده از کلیه کسانی که از راه مکاتبه خطاهای چاپ و اشتباهات موجود در چاپ دوم را متذکر شده اند، و همچنین پیشنهادهای برای بهبود کتاب به ما داده اند تشکر می کنیم.

س. د. کونت

کارل. دو بور

1. Fast Fourier Transforms
2. splines
3. collocation
4. procedure-oriented
5. UNIVAC, CDC, IBM

## مقدمه

در این کتاب حل عملی مسائل با کامپیوتر مطرح می‌شود. در روند حل يك مسئله امکان دارد چندین مرحله کم و بیش متمایز وجود داشته باشد. نخستین مرحله صورتبندی<sup>۱</sup> است. در صورتبندی يك الگوریتم ریاضی از يك وضعیت فیزیکی، دانشمندان می‌باید از قبل واقعیت امید به حل يك مسئله با کامپیوتر را در نظر داشته باشند. لذا برای هدفهای خاص و نوع و مقدار پرونداد، درون داده‌های مناسب، آزمونهای کافی فراهم سازند.

بعد از آنکه مسئله‌ای صورتبندی شد، برای حل آن روشهای عددی همراه با تحلیل خطای مقدماتی می‌باید طرح شود. هر روش عددی که بتواند برای حل مسئله‌ای به کار رود، الگوریتم نامیده می‌شود. يك الگوریتم مجموعه کامل و روشنی از روندهاست که به حل يك مسئله ریاضی منجر می‌شود. انتخاب ساختن الگوریتمهای مناسب به ویژه در محدوده آنالیز عددی صورت می‌گیرد. متخصصان آنالیز عددی، پس از اینکه درباره الگوریتم خاص یا مجموعه‌ای از الگوریتمها، برای حل مسئله‌ای تصمیم گرفتند، می‌باید همه منابع خطا را که ممکن است بر نتایج اثر گذارند در نظر بگیرند. باید درجه دقت مورد نیاز را رعایت کنند، دامنه خطاهای گرد کردن و گسسته‌سازی را بر آورد نمایند، نمو مناسبی با تعداد بارست‌های لازم معین و ملاکهای کافی برای دقت تهیه کنند و عمل تصحیحی را در صورت عدم همگرایی منظور نمایند.

مرحله سوم حل مسئله برنامه‌نویسی است. برنامه‌نویس باید الگوریتمی را که به آن عقیده مند شده است، به مجموعه دستورالعملهای غیر مبهم گام به گامی برای کامپیوتر تبدیل نماید. نخستین گام این شیوه عمل «دند نما سازی»<sup>۲</sup> نامیده می‌شود. يك روند نما اصولا مجموعه‌ای است از شیوه‌های عمل، معمولا به شکل قالبی منطقی، که کامپیوتر دنبال می‌کند و می‌تواند به صورت دستورالعمل نموداری یا عبارتی باشد. پیچیدگی روند نما به پیچیدگی مسئله و میزان جزئیات مندرج در آن بستگی دارد. ولی باید غیر از خود برنامه‌نویس شخص

دیگری هم بتواند جریان اطلاعات را از روی روندنما دنبال کند. روندنما کمک مؤثری است برای برنامه نویسی که باید توابع اصلی را به برنامه تبدیل کند، و درعین حال وسیله‌ای مؤثر برای ارتباط دیگرانی است که می‌خواهند کاری را که برنامه انجام می‌دهد بفهمند. در این کتاب، روندنما را گاهی به شکل نموداری ولی اغلب به صورت دستورالعمل عبارتی عرضه می‌کنیم. هنگامی که روندنما به صورت نموداری به کار برده شود، قراردادهای استاندارد دنبال می‌شوند. ولی وقتی به صورت دستورالعمل عبارتی باشد، از زبان بی‌نیاز از توضیح الگول-مانندی استفاده خواهد شد. پس از تهیه یک روندنما، برنامه نویسی شیوه‌های مذکور را به مجموعه‌ای از دستورالعمل‌های ماشینی تبدیل می‌کند. این کار را می‌توان مستقیماً به زبان ماشین یا به زبان هم‌گذاری<sup>۱</sup> یا یک زبان شیوه-پرداخته انجام داد. در این کتاب منحصرأ از یک گویش<sup>۲</sup> خاص فورترین، به نام فورترین ۷۷ استفاده شده است. فورترین ۷۷ گویش جدیدی است از زبان فورترین که عبارتهای ممیزی ترازه‌ای را وارد می‌کند و بر مفاهیم برنامه نویسی ساختاریافته<sup>۳</sup> جدیدی تکیه دارد. درحالی که همگردان‌های<sup>۴</sup> فورترین ۴ تقریباً روی همه کامپیوترها وجود دارند. دستیابی به همگردان فورترین ۷۷ به آسانی میسر نیست ولی تبدیل برنامه از فورترین ۷۷ به فورترین ۴ نسبتاً به سادگی انجام می‌گیرد.

زبانهای شیوه-پرداخته مسانند فورترین یا الگول، گاهی زبان الگوریتمی نامیده می‌شوند. این زبانها امکان بیان یک الگوریتم ریاضی را به شکل مناسبتری برای ارتباط با کامپیوترها فراهم می‌سازند. یک شیوه عمل فورترینی، که یک الگوریتم را پیاده می‌کند، در حالت کلی از خود الگوریتم دقیقتر است. مثلاً اگر این الگوریتم ریاضی مشخص کننده یک شیوه عمل بارسستی برای پیدا کردن یک معادله باشد، برنامه فورترین مسی باید سه نکته: (۱) دقت لازم، (۲) تعداد بارستهایی که باید انجام گیرد و (۳) در صورت ناهمگرایی چه باید کرد، را مشخص کند. بیشتر الگوریتمهای این کتاب به شکل عادی ریاضی و به شکلی دقیقتر از یک شیوه عمل فورترینی داده شده‌اند.

در بسیاری از مراکز کامپیوتری هر یک از این مراحل حل مسئله را شخص جداگانه‌ای انجام می‌دهد. در برخی دیگر ممکن است شخص واحدی مسئول هر سه کار باشد. واضح است که بین این سه مرحله فعل و انفعالیهای وجود دارند. هر چه اجرای برنامه پیشرفت می‌کند، اطلاعات بیشتری به دست می‌آید و این اطلاعات ممکن است تغییراتی را در صورتبندی، در الگوریتمهای به کار رفته و در خود برنامه موجب شوند.

- 
- |                           |              |
|---------------------------|--------------|
| 1. assembly language      | 2. dialect   |
| 3. structured-programming | 4. compilers |



## دستگاه اعداد و خطاها

در این فصل روشهای نمایش اعداد در کامپیوتر و خطاهای ناشی از این نمایشها را بررسی می‌کنیم. بعلاوه، منشأ انواع خطاهای گوناگون محاسباتی و انتشار بعدی آنها را مورد بررسی قرار می‌دهیم، و اندکی هم به بحث در مقدمات ریاضی می‌پردازیم.

### ۱.۱ نمایش اعداد صحیح

در زندگی روزمره اعداد را در پایهٔ دستگاه دهدهی به کار می‌بریم. از این رو مثلاً عدد ۲۵۷ به صورت زیر قابل بیان است

$$\begin{aligned} 257 &= 2 \times 100 + 5 \times 10 + 7 \times 1 \\ &= 2 \times 10^2 + 5 \times 10^1 + 7 \times 10^0 \end{aligned}$$

عدد ۱۰ را پایهٔ این دستگاه می‌نامیم. هر عدد صحیح در دستگاه دهدهی را می‌توان به صورت یک بسجمله‌ای در پایهٔ ۱۰ بیان کرد که ضرایب آن صحیح و بین ۰ و ۹ باشند. برای مشخص کردن هر عدد صحیح مثبت در پایهٔ ۱۰، نشانگذاری زیر را به کار می‌بریم

$$\begin{aligned} N &= (a_n a_{n-1} \dots a_0)_{10} \\ &= a_n 10^n + a_{n-1} 10^{n-1} + \dots + a_0 10^0 \end{aligned} \quad (1.1)$$

در به کارگیری پایه ۱۰ هیچ دلیل غریزی وجود ندارد. تمدنهای دیگر، پایه‌های دیگری مانند ۱۲، ۲۰ یا ۶۰ را به کار برده‌اند. کامپیوترهای جدید تپهای<sup>۱</sup> ارسالی ابزارهای الکتریکی را دریافت می‌کنند. حالت یک ضربه الکتریکی به یکی از دو صورت وصل یا قطع است. بنابراین نمایش اعداد در کامپیوتر در دستگاه دودویی<sup>۲</sup> مناسبتر است. پایه این دستگاه ۲ است و ضرایب صحیح می‌توانند مقادیر صفر یا یک باشند. یک عدد صحیح نامنفی  $N$  در دستگاه دودویی به صورت

$$\begin{aligned} N &= (a_n a_{n-1} \dots a_1 a_0)_2 \\ &= a_n 2^n + a_{n-1} 2^{n-1} + \dots + a_1 2^1 + a_0 2^0 \end{aligned} \quad (2.1)$$

نمایش داده می‌شود، که در آن ضرایب  $a_k$  صفر یا یک هستند. ملاحظه می‌کنید که  $N$  باز به وسیله یک بسجمله‌ای، اما اکنون در پایه ۲، نمایش داده شده است. بسیاری از کامپیوترهایی که در کارهای علمی مورد استفاده قرار می‌گیرند ذاتاً در دستگاه دودویی عمل می‌کنند. لیکن استفاده کنندگان از این کامپیوترها ترجیح می‌دهند در دستگاه دهدهی، که برای آنها مأنوستر است، کار کنند. بنابراین، هنگام دادن اطلاعات بسد کامپیوتر روشهایی برای تبدیل اعداد دهدهی به دودویی، و هنگام دریافت اطلاعات از آن، روشهایی برای تبدیل اعداد دودویی به دهدهی، مورد نیاز است. تبدیل اعداد دودویی به دهدهی ممکن است مستقیماً طبق تعریف (۲.۱) انجام پذیرد. برای مثال داریم:

$$(11)_2 = 1 \times 2^1 + 1 \times 2^0 = 3$$

$$(1101)_2 = 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 = 13$$

تبدیل اعداد صحیح از یک پایه  $\beta$  به پایه ۱۰ به وسیله الگوریتم زیر، که روش به دست آوردن آن در فصل دوم آمده، نیز انجام پذیر است.

**الگوریتم ۱۰۱** ضرایب بسجمله‌ای

$$p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0 \quad (3.1)$$

یعنی  $a_0, a_1, \dots, a_n$  و عدد  $\beta$  داده شده‌اند. اعداد  $b_n, b_{n-1}, \dots, b_0$  را باید بدین‌طور بازگشتی<sup>۳</sup> به شکل زیر محاسبه کرد

1. pulses

2. binary

3. recursively

$$b_n := a_n$$

$$b_{n-1} := a_{n-1} + b_n \beta$$

$$b_{n-2} := a_{n-2} + b_{n-1} \beta$$

$$b_{n-3} := a_{n-3} + b_{n-2} \beta$$

.....

$$b_0 := a_0 + b_1 \beta$$

بنابراین

$$b_0 := p(\beta)$$

از آنجا که طبق تعریف (۲۰۱)، عدد صحیح دودویی  $(a_n a_{n-1} \dots a_0)_2$  نمایانگر مقدار بسجمله‌ای (۳۰۱) به ازای  $x=2$  است، برای پیدا کردن معادل دهدهی یک عدد صحیح دودویی می‌توانیم الگوریتم (۱۰۱) را با  $\beta=2$  به کار ببریم. بنابراین معادل دهدهی عدد  $(1101)_2$  که با استفاده از الگوریتم (۱۰۱) محاسبه شده به صورت زیر است:

$$b_4 = 1$$

$$b_3 = 1 + 1 \times 2 = 3$$

$$b_2 = 0 + 3 \times 2 = 6$$

$$b_0 = 1 + 6 \times 2 = 13$$

و عدد دهدهی معادل  $(10000)_2$  چنین است:

$$b_4 = 1$$

$$b_3 = 0 + 1 \times 2 = 2$$

$$b_2 = 0 + 2 \times 2 = 4$$

$$b_1 = 0 + 4 \times 2 = 8$$

$$b_0 = 0 + 8 \times 2 = 16$$

اگر کسی بخواهد حساب دودویی را به کار برد، می‌تواند تبدیل عدد صحیح دهدهی  $N$  به معادل دودویی آن را نیز به وسیله الگوریتم (۱۰۱) انجام دهد. زیرا در صورتی که

$N = (a_n a_{n-1} \dots a_0)_{10}$ ، آنگاه طبق تعریف (۱.۱) داریم  $N = p(10)$ ، و در اینجا همان بسجمله‌ای (۳.۱) است. لذا می‌توانیم نمایش دودویی  $N$  را چنین محاسبه کنیم که ابتدا ضرایب  $a_n, \dots, a_0$  را به اعداد صحیح دودویی تبدیل و سپس با استفاده از الگوریتم (۱.۱)، مقدار  $p(x)$  را به ازای  $x = 10 = (1010)_2$  در حساب دودویی محاسبه کنیم. برای مثال اگر  $N = 187$ ، آنگاه:

$$187 = (187)_{10} = 1 \times 10^2 + 8 \times 10^1 + 7 \times 10^0$$

$$= (1)_2(1010)_2^2 + (1000)_2(1010)_2^1 + (111)_2(1010)_2^0$$

اکنون با استفاده از الگوریتم (۱.۱) و حساب دودویی داریم،

$$b_2 = (1)_2$$

$$b_1 = (1000)_2 + (1)_2(1010)_2 = (1000)_2 + (1010)_2 = (10010)_2$$

$$b_0 = (111)_2 + (10010)_2(1010)_2 = (111)_2 + (10110100)_2$$

$$= (10111011)_2$$

بنابراین

$$187 = (10111011)_2$$

گرچه اعداد دودویی و حساب دودویی برای کامپیوترهای امروزی به‌طور آرمانی مناسب‌اند، استفاده از آنها به‌سبب تعداد ارقامی که حتی برای نمایش اعداد متوسط لازم است، برای شخص تاحدی خسته‌کننده است. چنانچه ملاحظه شد برای نمایش عدد دهدهی سه‌رقمی ۱۸۷، هشت‌رقم دودویی لازم است. دستگاه هشت‌هشتی (اکتال) که از پایهٔ ۸ بهره می‌گیرد، تقریباً واسطه‌ای بین دستگاه دودویی مطلوب کامپیوتر و دستگاه دهدهی مطلوب افراد است. از آنجا که سه‌رقم دودویی معرف یک‌رقم هشت‌هشتی است، به‌سادگی می‌توان دستگاه هشت‌هشتی را به‌دستگاه دودویی، و به‌عکس اعداد دودویی را به‌هشت‌هشتی تبدیل کرد. برای تبدیل یک عدد هشت‌هشتی به‌دودویی فقط باید به‌جای تمامی ارقام عدد هشت‌هشتی، اعداد دودویی هم ارزش آنها را گذاشت، بنابراین

$$(347)_8 = (011 \ 100 \ 111)_2 = (11100111)_2$$

به‌عکس، برای تبدیل اعداد دودویی به‌هشت‌هشتی، رقم‌های عدد دودویی را به‌گروه‌های سه‌تایی (از طرف راست) افراز می‌کنیم و سپس به‌جای هر گروه سه‌تایی معادل هشت‌هشتی آن را قرار می‌دهیم، بنابراین:

$$(10111011)_2 = (010111011)_2 = (273)_8$$

در صورتی که بخواهیم يك عدد صحیح دهدهی را با محاسبهٔ دستی به دودویی تبدیل کنیم، معمولاً سریعترین راه این است که ابتدا با استفاده از الگوریتم ۱۰۱، این عدد را به هشت هشتی و سپس از هشت هشتی به دودویی تبدیل کنیم. مثال قبلی را در نظر می‌گیریم

$$187 = (187)_{10} = (1)_8(12)_8^2 + (10)_8(12)_8 + (7)_8(12)_8^0$$

با استفاده از الگوریتم ۱۰۱ و با قراردادن  $(12)_8 = 10$  به جای ۲ و به کار گرفتن حساب هشت هشتی داریم

$$b_2 = (1)_8$$

$$b_1 = (10)_8 + (1)_8(12)_8 = (22)_8$$

$$b_0 = (7)_8 + (22)_8(12)_8 = (7)_8 + (264)_8 = (273)_8$$

بنابراین سرانجام داریم:

$$187 = (273)_8 = (010111011)_2$$

## تمرین

۱-۱۰۱ اعداد دودویی زیر را به دهدهی تبدیل کنید:

$$(1010)_2$$

$$(100101)_2$$

$$(10000001)_2$$

۲-۱۰۱ اعداد دهدهی زیر را به دودویی تبدیل کنید:

$$82$$

$$109$$

$$3433$$

۳-۱۰۱ اعداد تمرینهای ۱-۱۰۱ و ۲-۱۰۱ را ابتدا به صورت هشت هشتی تبدیل کنید و سپس اعداد به دست آمده را در دستگاه خواسته شده بنویسید.

۴-۱۰۱ زیر بر نامه‌ای<sup>۱</sup> به زبان فورترن بنویسید که عددی با NIN رقم را در پایهٔ مفروض BETIN واقع در يك آرایهٔ<sup>۲</sup> يك بعدی NUMIN بپذیرد و عدد NOUT رقمی هم‌ارز با آن در پایهٔ مفروض BETOUT را در يك آرایهٔ يك بعدی NUMOUT ذخیره کند. برای سادگی، هم BETIN و هم BETOUT را به اعداد ۲، ۴، ۸، ۱۰ محدود کنید.



## ۲.۱ نمایش کسرها

اگر  $x$  يك عدد مثبت حقیقی باشد، قسمت صحیح  $x$  را که بزرگترین عدد صحیح نابزرگتر از  $x$  است، با  $x_I$  و قسمت کسری  $x$  را با  $x_F$  نمایش می‌دهیم، بنا بر این

$$x_F = x - x_I$$

قسمت کسری را همیشه می‌توان به صورت کسر اعشاری زیر نوشت:

$$x_F = \sum_{k=1}^{\infty} b_k 10^{-k} \quad (۲.۱)$$

که در آن هر  $b_k$  يك عدد صحیح نامنفی کمتر از ۱۰ است. اگر به ازای جمیع  $k$  های بزرگتر از يك عدد صحیح معین، داشته باشیم  $b_k = 0$ ، آنگاه کسر بالا را مختوم گویند. از این رو کسر

$$\frac{1}{4} = 0.25 = 2 \times 10^{-1} + 5 \times 10^{-2}$$

مختوم است، درحالی که کسر

$$\frac{1}{3} = 0.333\dots = 3 \times 10^{-1} + 3 \times 10^{-2} + 3 \times 10^{-3} + \dots$$

مختوم نیست.

اگر قسمت صحیح  $x$  به عنوان يك عدد صحیح در دستگاه دهدهی به صورت

$$x_I = (a_n a_{n-1} \dots a_0)_{10}$$

و قسمت کسری آن طبق فرمول (۲.۱) داده شده باشد، معمولاً دو قسمت را پشت سرهم می‌نویسند و با يك «مميز» آنها را از هم جدا می‌کنند.

$$x = (a_n a_{n-1} \dots a_0 . b_1 b_2 b_3 \dots)_{10}$$

عیناً به همین طریق می‌توان قسمت کسری  $x$  را به صورت يك کسر دودویی نوشت:

$$x_F = \sum_{k=1}^{\infty} b_k 2^{-k}$$

که در آن هر  $b_k$  عدد صحیحی است نامنفی کوچکتر از ۲، یعنی یا صفر یا يك. اگر قسمت صحیح  $x$  به وسیله عدد صحیح دودویی

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

داده شده باشد، می توان عدد  $x$  را با استفاده از «ممیز»، به صورت زیر نوشت:

$$x = (a_n a_{n-1} \dots a_0 . b_1 b_2 b_3 \dots)_2$$

کسر دودویی  $(\dots b_3 b_2 b_1)_2$  را برای عدد مفروض  $x_F$ ، که بین صفر و یک است، می توان به صورت زیر محاسبه کرد: اگر

$$x_F = \sum_{k=1}^{\infty} b_k 2^{-k}$$

آنگاه

$$2x_F = \sum_{k=1}^{\infty} b_k 2^{-k+1} = b_1 + \sum_{k=1}^{\infty} b_{k+1} 2^{-k}$$

بنابراین  $b_1$  قسمت صحیح  $2x_F$  است، در حالی که

$$(2x_F)_F = 2x_F - b_1 = \sum_{k=1}^{\infty} b_{k+1} 2^{-k}$$

بنابراین اگر این روش را تکرار کنیم خواهیم دید که  $b_2$  قسمت صحیح  $2(2x_F)_F$ ،  $b_3$  قسمت صحیح  $2(2(2x_F)_F)_F$  و هکذا برای بقیه می باشند.

مثلا اگر  $x = 0.625 = x_F$ ، آنگاه

$$b_1 = 1 \quad \text{بنابراین} \quad 2(0.625) = 1.25$$

$$b_2 = 0 \quad \text{بنابراین} \quad 2(0.25) = 0.5$$

$$b_3 = 1 \quad \text{بنابراین} \quad 2(0.5) = 1.0$$

و کلیه  $b_k$  های بعدی صفرند، در نتیجه

$$0.625 = (0.101)_2$$

این مثال به گونه ای انتخاب شده بود که کسر دودویی آن مختوم باشد. متأسفانه تمامی کسرهای دهدهی مختوم به کسرهای دودویی مختوم منجر نمی شوند. این امر ناشی از این حقیقت است که کسر دودویی معادل  $0.1 = 10^{-1} = x_F$  مختوم نیست زیرا داریم

$$b_1 = 0 \quad \text{بنابراین} \quad 2(0.1) = 0.2$$

$$b_2 = 0 \quad \text{بنابراین} \quad 2(0.2) = 0.4$$

$$\begin{array}{lll}
 b_4 = 0 & \text{بنا بر این} & 2(004) = 008 \\
 b_4 = 1 & \text{بنا بر این} & 2(008) = 106 \\
 b_5 = 1 & \text{بنا بر این} & 2(006) = 102
 \end{array}$$

و اکنون به کسر ۲/۵ باز می‌گردیم، و ارقام دوباره به‌طور دوره‌ای تکرار می‌شوند. در نتیجه

$$\frac{2}{5} = (001100110011\dots)_2$$

شیوه‌ای که به اجمال شرح داده شد در الگوریتم زیر به‌صورت فرمول درمی‌آید.

**الگوریتم ۲۰۱** عدد  $x$  بین  $0$  و  $1$  و یک عدد صحیح  $\beta$  بزرگتر از  $1$ ، داده شده‌اند. اعداد  $b_1, b_2, b_3, \dots$  را به‌طور بازگشتی به‌طریق زیر تولید می‌کنیم:

$$\begin{aligned}
 c_0 &:= x \\
 b_1 &:= (\beta c_0)_I, \quad c_1 := (\beta c_0)_F \\
 b_2 &:= (\beta c_1)_I, \quad c_2 := (\beta c_1)_F \\
 &\dots
 \end{aligned}$$

سپس

$$x = (0b_1b_2b_3\dots)_\beta = \sum_{k=1}^{\infty} b_k \beta^{-k}$$

بنا بر دو دلیل این الگوریتم را برای یک پایهٔ کلی  $\beta$  بیان کرده‌ایم، و نه برای یک پایهٔ بخصوص دودویی  $\beta = 2$ . اگر این تبدیل به دودویی با مداد و کاغذ انجام پذیرد، معمولا اگر عدد فوق اول به هشت‌هشتی،  $\beta = 8$ ، و سپس از هشت‌هشتی به دودویی تبدیل شود، سرعت عمل بیشتر خواهد بود. همچنین می‌توان با انتخاب  $\beta = 10$  و استفاده از محاسبات دودویی (یا هشت‌هشتی)، این الگوریتم را برای تبدیل یک کسر دودویی (یا هشت‌هشتی) به دهدهی به‌کار برد.

برای مثال، اگر  $x = (0101)_2$ ، آنگاه به‌ازای  $\beta = 10 = (1010)_2$ ، و به‌کمک حساب دودویی و با استفاده از الگوریتم ۲۰۱، داریم

$$\begin{array}{lll}
 c_1 = (001)_2, b_1 = (110)_2 = 6 & \text{بنا بر این} & 10(0101)_2 = (110010)_2^* \\
 c_2 = (001)_2, b_2 = (10)_2 = 2 & \text{بنا بر این} & 10(001)_2 = (10010)_2 \\
 c_3 = 0, b_3 = (101)_2 = 5 & \text{بنا بر این} & 10(001)_2 = (10100)_2
 \end{array}$$

---

\* در عبارت  $(10010101)_2$  عدد  $10$  در پایهٔ دهدهی است. از آنجا که  $(1010)_2 = (10)_10$ ، بنا بر این  $(10010101)_2 = (1010)(101)_2 = (100010)_2$ .

بنابر این  $b$ های بعدی صفر خواهند بود. این نشان می‌دهد که

$$(0.101)_2 = 0.625$$

که مؤید محاسبات قبلی است. باید توجه داشت که هر گاه  $x$  يك كسر دودویی مختوم  $n$  رقمی باشد، در دستگاه دهدهی نیز کسری مختوم است و تعداد ارقام آن برابر با  $n$  زیر ا

$$(0.1)_2 = 0.5$$

### تمرین

۲-۲۰۱ کسرهای دودویی زیر را به کسرهای دهدهی تبدیل کنید.

$$(0.11111111)_2 \quad (0.1100011)_2$$

۲-۲۰۱ پنجم رقم اول  $0.1$  را در صورتی که به صورت يك كسر هشت هشتی نوشته شود پیدا، سپس به وسیله این کسر، ۱۵ رقم اول  $0.1$  را در دستگاه دودویی محاسبه کنید.

۳-۲۰۱ کسرهای هشت هشتی زیر را به کسرهای دهدهی تبدیل کنید:

$$(0.776)_8 \quad (0.614)_8$$

جوابهای این مسئله را با جوابهای مسئله ۲-۲۰۱ مقایسه کنید.

۴-۲۰۱ يك عدد دودویی پیدا کنید که  $n$  را با تقریب  $2-10$  نشان دهد.

۵-۲۰۱ اگر بخواهیم با استفاده از الگوریتم ۱-۱، يك عدد صحیح دهدهی  $N$  را به دودویی تبدیل کنیم، بایستی حساب دودویی را به کار گیریم. نشان دهید که چگونه این تبدیل با استفاده از الگوریتم ۲-۱ و حساب دهدهی انجام پذیر است.

(راهنمایی:  $N$  را بر توان مناسبی از ۲ تقسیم، و سپس نتیجه را به دودویی تبدیل کنید و «ممیز دودویی» را به طور مناسبی تغییر مکان دهید).

۶-۲۰۱ اگر بخواهیم يك كسر مختوم دودویی  $x$  را با استفاده از الگوریتم ۲-۱ به يك كسر دهدهی تبدیل کنیم، بایستی حساب دودویی را به کار گیریم. نشان دهید که چگونه این تبدیل با به کار بستن الگوریتم ۱-۱ و حساب دهدهی انجام پذیر است.

### ۳-۱ حساب با ممیز شناور

محاسبات علمی معمولاً در حساب با ممیز شناور انجام می‌گیرد. يك عدد با  $n$  رقم ممیز شناور در پایه  $\beta$  به شکل زیر است:

$$x = \pm (0.d_1 d_2 \dots d_n)_\beta \beta^e \quad (5.1)$$

که در آن  $(d_0 d_1 \dots d_n)_\beta$  کسری است در پایه  $\beta$  که جزء کسری یا مانتیس<sup>۱</sup> نامیده می‌شود و  $e$  عدد صحیحی است که آن را نما<sup>۲</sup> می‌نامند. هر گاه  $d_1 \neq 0$  و یا  $d_n = 0$  و  $d_1 = d_2 = \dots = d_n = 0$  چنین عدد باممیز شناور را نرمال شده<sup>۳</sup> نامند.

اگر چه در برخی از کامپیوترها  $\beta = ۱۶$  و در محاسبات دستی و حسابگرهای جیبی و رومیزی  $\beta = ۱۰$ ، لیکن در اکثر کامپیوترها  $\beta = ۲$  گرفته می‌شود.

معمولاً در هر کامپیوتری، دقت یا طول  $n$  در اعداد با ممیز شناور به طول کلمه در آن کامپیوتر بستگی دارد، از این رو ممکن است دقت این اعداد در بسیاری جاها متفاوت باشند (به شکل ۱۰۱ مراجعه شود). در دستگا‌های محاسباتی که برنامه‌های فورترن را می‌پذیرند قرار بر این است که اعداد با ممیز شناور را با دو طول مختلف به کار گیرند که یکی تقریباً دو برابر دیگری است. جز در شرایط خاص که به اصطلاح به دقت مضاعف<sup>۴</sup> نیاز است، در حالت عادی از اعداد با ممیز شناور با طول کوتاهتر به اصطلاح با دقت معمولی<sup>۵</sup> استفاده می‌شود. معمولاً انجام محاسبات با دقت مضاعف، در مقایسه با دقت ساده، نیاز به حافظه‌ای دو برابر دارد و زمان محاسبه را به بیش از دو برابر افزایش می‌دهد. نمای  $e$  معمولاً در محدوده<sup>۶</sup> زیر

$$m < e < M \quad (۶.۱)$$

تغییر می‌کند که  $m$  و  $M$  اعداد صحیح معینی هستند.

معمولاً  $m = -M$ ، لیکن ممکن است این محدوده در بسیاری از جاها تغییر یابد، به شکل ۱۰۱ نگاه کنید.

برای تبدیل یک عدد حقیقی  $x$  به یک عدد  $f(x)$  رقمی باممیز شناور در پایه  $\beta$ ، دو روش متداول است یکی گرد کردن<sup>۶</sup> و دیگری قطع کردن<sup>۷</sup>. در روش گرد کردن، عدد

Computer	$\beta$	$n$	$M = -m$
IBM 7094	2	27	$2^7$
Burroughs 5000 Series	8	13	$2^6$
IBM 360/370	16	6	$2^6$
CDC 6000 and Cyber Series	2	48	$2^{10}$
DEC 11/780 VAX	2	24	$2^7$
Hewlett Packard 67	10	10	99

شکل ۱۰۱ مشخصه‌های ممیز شناور.

1. mantissa
2. exponent
3. normalized
4. double precision
5. single precision
6. rounding
7. chopping



$fl(x)$  به عنوان نزدیکترین عدد نرمال شده با ممیز شناور به  $x$ ، انتخاب می شود. در حالتی که محدودیتی وجود داشته باشد، از بعضی قواعد خاص مانند گرد کردن متقارن (گرد کردن به یک عدد زوج) استفاده می شود. در روش قطع کردن عدد  $fl(x)$  به عنوان نزدیکترین عدد نرمال شده با ممیز شناور بین  $x$  و ۰ انتخاب می شود. برای مثال اگر اعداد دهدهی دورقمی با ممیز شناور به کار برده شوند، داریم

$$fl\left(\frac{7}{3}\right) = \begin{cases} (0.067) \times 10^0 & \text{گرد شده} \\ (0.066) \times 10^0 & \text{قطع شده} \end{cases}$$

$$fl(-138) = \begin{cases} -(0.084) \times 10^3 & \text{گرد شده} \\ -(0.083) \times 10^3 & \text{قطع شده} \end{cases}$$

و

در برخی از کامپیوترها این تعریف عدد  $fl(x)$  در حالت  $|x| \geq \beta^M$  (حالت سرریز) یا  $0 < |x| < \beta^{m-n}$  (حالت پی ریز)، که  $M$  و  $m$  کرانه های نما می باشند، اصلاح شده است؛ در این حالت یا  $fl(x)$  تعریف نشده است که این امر موجب قطع برنامه می شود و یا  $fl(x)$  به وسیله عدد خاصی نمایش داده شده است که خود این عدد هنگام ترکیب با اعداد با ممیز شناور معمولی از قواعد متداول حساب پیروی نمی کند.

تفاضل بین  $x$  و  $fl(x)$  خطای گرد کردن نامیده می شود. این خطا به مقدار  $x$  بستگی دارد و از این رو بهترین روش برای اندازه گیری آن سنجش آن نسبت به  $x$  است. زیرا اگر بنویسیم

$$fl(x) = x(1 + \delta) \quad (7.1)$$

که  $\delta = \delta(x)$  عددی است که به  $x$  بستگی دارد، آنگاه می توانیم، حداقل مادامی که  $x$  موجب بروز سرریز و پی ریز نشود،  $\delta$  را مستقل از  $x$  بدانیم. برای چنین  $x$  (ی) که موجب سرریز و پی ریز نشود) دشوار نیست که نشان دهیم

$$|\delta| < \frac{1}{\beta} \beta^{1-n} \quad \text{در گرد کردن} \quad (8.1)$$

در حالی که

$$-\beta^{1-n} < \delta \leq 0 \quad \text{در قطع کردن} \quad (9.1)$$

به تمرین ۳-۳۰۱ نگاه کنید. حداکثر مقدار ممکن برای  $|\delta|$  اغلب «واحد گرد کردن» نامیده شده و با  $u$  نشان داده می شود.

هرگاه عملی حسابی روی دو عدد یا ممیز شناور انجام گیرد، نتیجه معمولاً یک عدد یا

ممیز شناور با همان طول نیست. اگر برای مثال با اعداد دهدهی دورقمی کار کنیم و بگیریم

$$x = (0.20)10^1 = 2 \quad y = (0.77)10^{-6} \quad z = (0.30)10^1 = 3$$

آنگاه

$$x + y = (0.2000000077)10^1 \quad x \cdot y = (0.154)10^{-5}$$

$$\frac{x}{z} = (0.666\dots)10^0$$

بنابراین، اگر  $\omega$  معرف یکی از اعمال حسابی (جمع، تفریق، ضرب و یا تقسیم) و  $\omega^*$  همان عمل با ممیز شناور باشد که به توسط کامپیوتر انجام شده است، آنگاه گرچه به ازای دو عدد با ممیز شناور  $x$ ،  $y$  نتیجهٔ محاسبهٔ کامپیوتر  $x\omega^*y$  می باشد، اما می توان مطمئن بود که معمولاً

$$x\omega^*y \neq x\omega y$$

گرچه ممکن است عمل با ممیز شناور  $\omega^*$  متناظر با  $\omega$ ، در ماشینهای مختلف در بعضی از جزئیات کمی متفاوت باشد، اما معمولاً  $\omega^*$  طوری اختیار می شود که

$$x\omega^*y = fl(x\omega y) \quad (10.1)$$

به عبارت دیگر، مجموع (تفاضل، حاصلضرب یا خارج قسمت) با ممیز شناور دو عدد با ممیز شناور معمولاً مساوی عددی است با ممیز شناور که مجموع (تفاضل، حاصلضرب یا خارج قسمت) دقیق این دو عدد را نمایش دهد. بنابراین (در صورتی که سرریز و پی ریزی پیش نیاید) داریم:

$$x\omega^*y = (x\omega y)(1 + \delta) \quad |\delta| \leq u \quad (11.1 \text{ الف})$$

که  $u$  واحد گرد شده در این فرمول است. در برخی موارد مناسبتر آن است که فرمول معادل زیر به کار گرفته شود

$$x\omega^*y = (x\omega y)/(1 + \delta) \quad |\delta| \leq u \quad (11.1 \text{ ب})$$

معادلهٔ (11.1) بیانگر موضوع اساسی تحلیل خطاها به طریق «پسرو»<sup>۱</sup> می باشد (به \* [۲۴] J. H. Wilkinson نگاه کنید). توضیحاً اضافه می کنیم که معادلهٔ (11.1) امکان می دهد که یک نتیجه با ممیز شناور به عنوان نتیجه ای متناظر با حساب معمولی، اما با داده هایی که اندکی تغییر یافته است، تلقی شود. از این رو، تحلیل نتیجه با حساب ممیز شناور را می توان بر حسب حساب معمولی انجام داد.

برای مثال، مقدار تابع  $f(x) = x^2$  در نقطهٔ  $x_0$  را می توان با  $n$  بار مربع کردن

## 1. backward

\* عدد درون کروشه هر بوط به مراجعی است که در آخر کتاب داده شده است.

یعنی با انجام مراحل متوالی زیر انجام داد.

$$x_1 := x_0^y, x_2 := x_1^y, \dots, x_n := x_{n-1}^y$$

که در رابطه فوق  $f(x_0) = x_n$  در حساب باممیز شناور، به جای آن، طبق معادله (۱۱.۱ الف) دنباله اعداد زیر را محاسبه می‌کنیم.

$$\hat{x}_1 = x_0^y(1 + \delta_1), \quad \hat{x}_2 = (\hat{x}_1)^y(1 + \delta_2), \dots, \hat{x}_n = (\hat{x}_{n-1})^y(1 + \delta_n)$$

که در آن به ازای جميع مقادیر  $i$ ،  $|\delta_i| \leq u$ . بنا بر این جواب محاسبه شده برابری است با

$$\hat{x}_n = x_0^{y^n} (1 + \delta_1)^{y^{n-1}} \dots (1 + \delta_{n-1})^y (1 + \delta_n)$$

برای ساده کردن عبارت فوق توجه کنید که اگر  $|\delta_1|, \dots, |\delta_r| \leq u$ ، آنگاه به ازای بعضی از مقادیر  $\delta$  با  $|\delta| \leq u$  داریم:

$$(1 + \delta_1) \dots (1 + \delta_r) = (1 + \delta)^r$$

(به تمرین ۳۰۱-۶ نگاه کنید). همچنین به ازای بعضی از مقادیر  $\eta$  با شرط  $|\eta| \leq u$ ، داریم

$$(1 + \delta)^r = (1 + \eta)^{r+1}$$

در نتیجه به ازای بعضی از مقادیر  $\delta$  با  $|\delta| \leq u$

$$\hat{x}_n = x_0^{y^n} (1 + \delta)^{y^n} = f(x_0(1 + \delta))$$

یا اگر با عبارت بیان کنیم، مقدار محاسبه شده  $\hat{x}_n$  برای  $f(x_0)$  برابر است با مقدار دقیق  $f(x)$  در شناسه آشفته  $x = x_0(1 + \delta)$ .

اکنون با بررسی این امر که چگونه مقدار تابع (مقدار دقیق محاسبه شده)  $f(x)$  با آشفته شدن (تغییر جزئی) شناسه  $x$  تغییر می‌کند، می‌توانیم اثری را که محاسبه ممیز شناور بر دقت عمل مقدار محاسبه شده برای  $f(x_0)$  برجا گذاشته است اندازه گیری کنیم، همان گونه که در بند بعدی انجام داده ایم. علاوه توجه داریم که در مثال مورد بحث این خطا اساساً با خطای ناشی از تبدیل داده اولیه  $x_0$  به یک عدد باممیز شناور، قابل مقایسه است.

به عنوان مثال دوم، که مخصوصاً در فصل ۴ مورد توجه قرار خواهد گرفت، محاسبه عدد  $S$  در معادله

$$a_1 b_1 + \dots + a_r b_r + a_{r+1} s = c \quad (12.1)$$

را به وسیله فرمول

$$s = \left( c - \sum_{k=1}^r a_k b_k \right) / a_{r+1}$$

مورد بررسی قرار می‌دهیم. اگر  $s$  از مراحل

$$s_0 := c$$

$$s_i := s_{i-1} - a_i b_i \quad i = 1, \dots, r$$

$$s := s_r / a_{r+1}$$

به دست آید، آنگاه اعداد متناظر محاسبه شده در حساب با ممیز شناور، در روابط زیر صدق می‌کنند:

$$\hat{s}_0 = c$$

$$\hat{s}_i = [\hat{s}_{i-1} - a_i b_i (1 + \delta)] (1 + \delta), \quad i = 1, \dots, r$$

$$\hat{s} = \hat{s}_r / [a_{r+1} (1 + \delta)]$$

در اینجا معادله (۱۱.۱ الف) و (۱۱.۱ ب) را به کار برده‌ایم و به مشخص کردن  $\delta$  های مختلف به وسیله زیر نمایه اعتنا نکرده‌ایم. در نتیجه داریم:

$$\begin{aligned} a_{r+1} (1 + \delta) \hat{s} &= \hat{s}_r \\ &= \hat{s}_{r-1} (1 + \delta) - a_r b_r (1 + \delta)^2 \\ &= \hat{s}_{r-2} (1 + \delta)^2 - a_{r-1} b_{r-1} (1 + \delta)^3 - a_r b_r (1 + \delta)^2 \\ &\vdots \\ &= \hat{s}_0 (1 + \delta)^r - a_1 b_1 (1 + \delta)^{r+1} - \dots - a_r b_r (1 + \delta)^2 \end{aligned}$$

این رابطه نشان می‌دهد که مقدار محاسبه شده  $s$  برای  $\hat{s}$ ، در معادله آشفته زیر صدق می‌کند.

$$a_1 b_1 (1 + \delta)^{r+1} + \dots + a_r b_r (1 + \delta)^2 + a_{r+1} (1 + \delta) \hat{s} = c (1 + \delta)^r \quad (13.1)$$

توجه داریم که در حالت  $a_{r+1} = 1$ ، یعنی در صورتی که نیازی به انجام آخرین تقسیم نباشد، می‌توانیم کلیه نماها را به اندازه ۱ کاهش دهیم.

## تمرین

۱-۳۰۱ دریک کامپیوتر دهدهی\* با مانیتس نرمال شده چهار رقمی، اعداد زیر داده شده‌اند

$$\text{الف) } 0.4523 \times 10^4 \quad \text{ب) } 0.2115 \times 10^{-3} \quad \text{پ) } 0.2583 \times 10^1$$

## 1. subscripts

\* منظور کامپیوتری است که با اعداد در دستگاه دهدهی کار می‌کند.م.

عملیات زیر را با اتخاذ روش گرد کردن متقارن انجام دهید و خطای موجود در نتیجه را تعیین کنید.

$$(الف) (a) + (b) + (c) \quad (ب) (a)/(c)$$

$$(پ) (a) - (b) \quad (ت) (a) - (b) - (c)$$

$$(ث) (a)(b)/(c) \quad (ج) (b)/(c) \cdot (a)$$

۳-۳۰۱ فرض می‌کنیم عدد  $fl(x)$  از راه قطع کردن داده شده باشد. نشان دهید که (جز در حالت سرریز و بی‌ریز) روابط زیر برقرارند

$$fl(-x) = -fl(x)$$

$$fl(\beta' x) = \beta' fl(x)$$

۳-۳۰۲ فرض می‌کنیم  $fl(x)$  به وسیله قطع کردن  $x$  داده شده باشد.  $\delta = \delta(x)$  را طوری انتخاب می‌کنیم که  $fl(x) = x(1 + \delta)$ . (در حالت  $x = 0$ ,  $\delta = 0$  انتخاب شود) نشان دهید که در این صورت کرانه‌های  $\delta$  از فرمول (۹-۱) به دست می‌آید.

۴-۳۰۱ با ذکر چند مثال نشان دهید که در حساب با ممیز شناور، اکثر قواعد حساب صادق نمی‌کنند. (دانه‌مایی: به قواعدی که مستلزم سه عملوند هستند توجه کنید).

۵-۳۰۱ یک تابع فورترن  $FL(X)$  بنویسید که نتیجه بازگشت آن یک عدد  $n$  رقمی اعشاری با ممیز شناور، حاصل از گرد کردن  $X$  باشد.  $n$  را ۴ فرض کنید و محاسبات خود در تمرین ۳-۳۰۱ را کنترل کنید. برای تعیین  $e$  با شرط  $10^e < |x| \leq 10^{e-1}$ ، از  $ALOG(ABS(X))$  استفاده کنید.

۶-۳۰۱ گیریم  $U = \{u : |\delta| \leq u, 1 + \delta\}$ . نشان دهید که به ازای جمیع مقادیر  $\alpha_1, \alpha_2, \dots, \alpha_r \in U$  مقداری مانند  $\alpha \in U$  وجود دارد به طوری که  $\alpha^r = \alpha_1 \alpha_2 \dots \alpha_r$ . همچنین نشان دهید که اگر  $a_1, a_2, \dots, a_r$  همگی دارای یک علامت باشند، به ازای مقداری چون  $\alpha \in U$  رابطه

$$a_1 \alpha_1 + a_2 \alpha_2 + \dots + a_r \alpha_r = (a_1 + a_2 + \dots + a_r) \alpha$$

برقرار است.

۷-۳۰۱ برای محاسبه حاصلضرب داخلی  $s = a_1 b_1 + a_2 b_2 + \dots + a_n b_n$  یک تحلیل خطای پسرو انجام دهید. تحلیل را با این فرض که انباشتگی با دقت مضاعف به کار گرفته شده است تکرار کنید. این فرض بدین معنی است که نتایج هر ضرب با دقت مضاعف نگهداری



و با دقت مضاعف به حاصلجمع افزوده می‌شود و فقط در پایان، حاصلجمع نهایی با دقت معمولی گرد می‌گردد.

### ۴.۱ ازدست دادن ارقام بامعنی و پخش خطا، شرط و پایداری

اگر عدد  $x^*$  تقریبی برای جواب دقیق  $x$  باشد، آنگاه تفاضل  $x - x^*$  را خطای  $x^*$  می‌نامیم. لذا

$$\text{خطا} + \text{مقدار تقریبی} = \text{مقدار دقیق} \quad (14.1)$$

بنا بر تعریف، عدد  $(x - x^*)/x$ ، خطای نسبی<sup>۱</sup> در  $x^*$ ، تقریبی برای  $x$  است. توجه کنید که این عدد (اگر اصلاً بسیار کوچک باشد) به عدد  $(x - x^*)/x^*$  نزدیک است [تحقیقاً، اگر  $\alpha = (x - x^*)/x$ ، آنگاه  $\alpha = (x - x^*)/x^* = \alpha/(1 - \alpha)$  در یک روند محاسباتی از هر عمل بامعنی شناور ممکن است خطایی حاصل شود که میزان آن در عملیات بعدی احتمالاً وسعت یا کاهش می‌یابد.

یکی از عادیترین راههای افزایش اهمیت یک خطا (که غالباً اجتناب پذیر است) به اصطلاح ازدست دادن ارقام بامعنی است. اگر  $x^*$  تقریبی از  $x$  باشد،  $x^*$  را تقریب  $x$  تا  $r$  رقم بامعنی در پایهٔ  $\beta$  گوئیم، مشروط بر آنکه خطای مطلق  $|x - x^*|$  حداکثر برابر با  $1/2$  رقم  $r$ ام  $x$  در پایهٔ  $\beta$  باشد. این مطلب را می‌توان با رابطهٔ

$$|x - x^*| \leq \frac{1}{2} \beta^{s-r+1} \quad (15.1)$$

بیان کرد، که در آن  $s$  بزرگترین عدد صحیحی است که  $|x| \leq \beta^s$ . برای مثال،  $x^* = 3$  تا یک رقم بامعنی (اعشاری) تقریبی برای  $x = \pi$  است، در حالی که  $x^* = 3.14159 \dots$  تا ۷ رقم بامعنی (اعشاری) تقریبی برای  $\pi$  تا سه رقم بامعنی صحیح است. اکنون فرض کنید می‌خواهیم عدد

$$z = x - y$$

را محاسبه کنیم و تقریبهای  $x^*$  و  $y^*$  را که هر کدام تا  $r$  رقم اعشاری صحیح است به ترتیب برای  $x$  و  $y$  در دست داریم. در این صورت

$$z^* = x^* - y^*$$

تقریبی است از  $z$  که تا  $r$  رقم معتبر است، مگر آنکه یک یا چند رقم از  $x^*$  و  $y^*$  با هم یکی باشند. در حالت اخیر، هنگام عمل تفریق برخی از ارقام، صفر و در نتیجه تعداد ارقام صحیح  $z^*$  کمتر از  $r$  خواهد شد.

برای مثال، فرض می‌کنیم

$$x^* = (0.765425421)10^1 \quad y^* = (0.76544200)10^1$$

به ترتیب تقریبی از  $x$  و  $y$  تا هفت رقم با معنی صحیح باشند. در این صورت تفاضل دقیق بین  $x^*$  و  $y^*$  در حساب با ممیز شناور هشت رقمی برابر است با

$$z^* = x^* - y^* = (0.12210000)10^{-3}$$

اما به عنوان تقریبی برای  $z = x - y$ ،  $z^*$  فقط تاسه رقم صحیح است زیرا رقم چهارم در  $z^*$  از هشتمین ارقام  $x^*$  و  $y^*$  به دست آمده است که هر دو اینها احتمالا متضمن خطا هستند. بنابراین در عین اینکه حداکثر مقدار خطا در  $z^*$  (به عنوان تقریبی برای  $z = x - y$ ) برابر با مجموع خطاها در  $x^*$  و  $y^*$  است، خطای نسبی در  $z^*$  احتمالا ۱۰۰۰۰ برابر خطای نسبی در  $x^*$  یا  $y^*$  است. بنابراین از دست دادن ارقام با معنی زمانی خطرناک است که بخواهیم خطای نسبی را کوچک نگاهداریم.

اغلب می‌توان با پیشبینی، جلوی از دست رفتن ارقام با معنی را گرفت. برای مثال، محاسبه مقدار تابع

$$f(x) = 1 - \cos x$$

را در نظر می‌گیریم. عملیات را با حساب با شش رقم اعشاری انجام می‌دهیم. از آنجا که برای  $x$ های نزدیک به صفر داریم  $\cos x \approx 1$ ، لذا اگر  $f(x)$  را بدین طریق محاسبه کنیم که اول مقدار  $\cos x$  را پیدا و سپس این مقدار را از یک کم کنیم، احتمال از دست دادن ارقام با معنی بسیار خواهد بود. از آنجا که نمی‌توانیم  $\cos x$  را تا بیش از شش رقم محاسبه کنیم، پس خطای محاسبه برای  $x$ های نزدیک به صفر ممکن است به اندازه  $5 \times 10^{-7}$  و بنابراین به اندازه  $f(x)$  و یا بزرگتر از آن باشد. اگر بخواهیم مقدار تابع  $f(x)$  را در نزدیکی  $x = 0$  تا شش رقم با معنی در حساب شش رقمی محاسبه کنیم، باید از فرمول دیگری مانند

$$f(x) = 1 - \cos x = \frac{1 - \cos^2 x}{1 + \cos x} = \frac{\sin^2 x}{1 + \cos x}$$

استفاده کنیم، که با این فرمول می‌توان  $f(x)$  را به طور کاملا دقیق به ازای مقادیر کوچک  $x$  محاسبه کرد؛ و یا از بسط  $f(x)$  به سری تیلر (به بخش ۷.۱ نگاه کنید)

$$f(x) = \frac{x^2}{2} - \frac{x^4}{24} + \dots$$

استفاده کنیم. رابطه بالا نشان می‌دهد که برای مثال، به ازای  $|x| \leq 10^{-3}$ ،  $|x^2/2|$  حداقل تا شش رقم با معنی  $f(x)$  تطبیق می‌کند.

تعیین ریشه‌های معادله درجه دوم

$$ax^2 + bx + c = 0 \quad (16.1)$$

مثال دیگری برای این مطلب است. در درس جبر دیده‌ایم که جواب معادله فوق از دستور درجه دوم

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} \quad (17.1)$$

به دست می‌آید. فرض کنید  $b^2 - 4ac > 0$  و  $b > 0$  باشد و بخواهیم با به کار بستن فرمول (۱۷.۱)، ریشه‌ای را که قدر مطلقش کمتر است یعنی

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \quad (18.1)$$

را پیدا کنیم. اگر  $4ac$  در مقایسه با  $b^2$  کوچک باشد، آنگاه  $\sqrt{b^2 - 4ac}$  در تعداد زیادی از ارقام با  $b$  یکی خواهد بود. بنابراین با فرض آنکه  $\sqrt{b^2 - 4ac}$  تنها با همان تعداد ارقام صحیحی که در محاسبات به کار می‌رود محاسبه شود، تعداد ارقام صحیح صورت کسر (۱۸.۱) و در نتیجه تعداد ارقام صحیح ریشه محاسبه شده نیز کمتر از تعداد ارقام به کار برده شده در محاسبه خواهد بود. برای اینکه مورد خاصی بررسی شود، معادله

$$x^2 + 11111x + 12121 = 0 \quad (19.1)$$

را در نظر می‌گیریم. با به کار گرفتن معادله (۱۸.۱) و حساب قطع شده با ممیز شناور با پنج رقم دهدهی داریم:

$$b^2 = 12345$$

$$b^2 - 4ac = 12340$$

$$\sqrt{b^2 - 4ac} = 111209$$

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} = -0.01000$$

در حالی که ریشه دقیق معادله تا ۵ رقم دهدهی برابر با  $x_1 = -0.010910$  است. در اینجا نیز از دست دادن ارقام بامعنی قابل اجتناب است و این کار با به کار گرفتن فرمول دیگری برای محاسبه ریشه‌ای که قدر مطلقش کوچکتر است، یعنی

$$x_1 = \frac{-2c}{b + \sqrt{b^2 - 4ac}} \quad (20.1)$$

انجام می‌پذیرد. با استفاده از این فرمول و محاسبه تا پنج رقم دهدهی، ریشهٔ  $x$  را تا پنج رقم صحیح به دست می‌آوریم

$$x_1 = -0.010910$$

هنگامی که خطایی صورت گرفته باشد، این خطا نهایتاً به نتایج بعدی سرایت می‌کند. انتشار خطا در محاسبات بعدی را می‌توان به آسانی بر حسب دو مفهوم مرتبط به هم، "شرط" و "ناپایداری" مورد مطالعه قرار داد.

واژهٔ شرط برای بیان حساسیت مقدار تابع  $f(x)$  در رابطه با تغییرات شناسهٔ  $x$  به کار برده می‌شود. شرط معمولاً به وسیلهٔ حداکثر تغییر نسبی در مقدار تابع  $f(x)$  که به توسط یک واحد تغییر نسبی در شناسهٔ  $x$  صورت گرفته اندازه‌گیری می‌شود. بسا یک فرمول نسبتاً غیررسمی شرط یک تابع را چنین تعریف می‌کنند.

$$\text{شرط } f \text{ در } x = \max \left\{ \left| \frac{f(x) - f(x^*)}{f(x)} \right| / \left| \frac{x - x^*}{x} \right| : |x - x^*| \text{ "کوچک"} \right\}$$

$$\approx \left| \frac{f'(x)x}{f(x)} \right| \quad (21.1)$$

هر قدر که این شرط بزرگتر باشد، به همان اندازه تابع  $f$ ، به اصطلاح بد شرط‌تر<sup>۲</sup> است. در اینجا از این حقیقت استفاده شده است که

$$f(x) - f(x^*) \approx f'(x)(x - x^*)$$

یعنی با تبدیل شناسهٔ  $x$  به  $x^*$ ، مقدار تابع تقریباً به اندازهٔ  $f'(x)(x - x^*)$  تغییر می‌یابد. (به بخش ۷.۱ نگاه کنید.) اگر، برای مثال، داشته باشیم

$$f(x) = \sqrt{x}$$

در این صورت  $f'(x) = \frac{1}{2\sqrt{x}}$  و بنابراین شرط  $f$  تقریباً چنین است

$$\left| \frac{f'(x)x}{f(x)} \right| = \frac{\left[ \frac{1}{2\sqrt{x}} \right] x}{\sqrt{x}} = \frac{1}{2}$$

رابطهٔ فوق نشان می‌دهد که جذرگیری یک روند "بهبود" است زیرا خطای نسبی را واقعاً کاهش می‌دهد. بعکس اگر داشته باشیم:

1. error propagation

2. ill-conditioned

3. well-conditioned

$$f(x) = \frac{10}{1-x^2}$$

در این صورت  $f'(x) = 20x/(1-x^2)^2$  و بنا بر این

$$\left| \frac{f'(x)x}{f(x)} \right| = \left| \frac{[20x/(1-x^2)^2]x}{10/(1-x^2)} \right| = \frac{2x^2}{|1-x^2|}$$

هر گاه  $|x|$  به ۱ نزدیک شود عدد فوق می تواند کاملاً بزرگ گردد. بنا بر این به ازای  $x$  نزدیک به ۱ یا -۱، تابع فوق کاملاً بد شرط است. این فرمول، خطاهای نسبی حاصل به توسط شناسه در  $f(x)$  را خیلی بزرگ می کند.\*

**ناپایداری، مفهوم وابسته به شرط، مبین حساسیت يك فرایند عددی در محاسبه  $f(x)$  به ازای  $x$  است، در خطاهای اجتناب ناپذیری که هنگام اجرای عمل گرد کردن در محاسبات با دقت محدود مرتکب می شویم. تأثیر دقیق این خطاها را بر مقدار محاسبه شده  $f(x)$  به دشواری می توانیم تعیین کنیم، مگر اینکه محاسبات را عملاً برای حساب ویژه بسا دقت محدود انجام دهیم و سپس نتیجه را با جواب دقیق مقایسه کنیم. اما این تأثیرها را می توانیم تقریباً از بررسی گرد کردن خطاها يك به يك تخمین بزنیم. بدین معنی که هر يك از مراحل محاسباتی را که کل فرایند را تشکیل می دهند تک تک مورد بررسی قرار دهیم. فرض کنید که  $n$  چنین مرحله ای وجود داشته باشد. برونداد مرحله  $i$ ام را با  $x_i$  نشان می دهیم و می گیریم  $x_i = x_{i-1}$ . سپس این  $x_i$ ها به عنوان درونداد برای يك یا چند مرحله بعدی به کار گرفته می شوند، و بدین طریق روی جواب نهایی  $f(x) = x_n$  اثر می گذارند. حال تابعی را که وابستگی جواب نهایی را به نتیجه میانی  $x_i$  بیان می کند با  $f_i$  نشان می دهیم. به ویژه  $f_i$  همان تابع  $f$  خواهد بود. در این حال کل فرایند را ناپایدار گوئیم هر گاه يك یا چند تا از این توابع  $f_i$  بد شرط باشند. دقیقتر بگوئیم، فرایند ناپایدار است هر گاه يك یا چند تا از توابع  $f_i$  شرطی بسیار بزرگتر از شرط تابع  $f = f_n$  داشته باشند. زیرا این شرط  $f_i$  است که اثر نسبی خطای حاصل از گرد کردن اجتناب ناپذیر در مرحله  $i$ ام را بر جواب نهایی می رساند. به عنوان يك مثال ساده، تابع**

$$f(x) = \sqrt{x+1} - \sqrt{x}$$

را به ازای  $x$ های «بزرگ» مثلاً،  $x \approx 10^4$  بررسی می کنیم. شرط این تابع به صورت زیر است

$$\left| \frac{f'(x)x}{f(x)} \right| = \frac{1}{2} \frac{|1/\sqrt{x+1} - 1/\sqrt{x}|x}{\sqrt{x+1} - \sqrt{x}} = \frac{1}{2} \frac{x}{\sqrt{x+1}\sqrt{x}} \approx \frac{1}{2}$$

\* یعنی هر خطای نسبی در شناسه  $x$  موجب ایجاد خطایی به مراتب بزرگتر در مقدار  $f(x)$  می شود. —

که شرط کاملاً خوبی است. اما اگر  $f(12345)$  را با حساب شش رقم اعشاری محاسبه کنیم، مشاهده خواهیم کرد که

$$\begin{aligned} f(12345) &= \sqrt{12346} - \sqrt{12345} \\ &= 111.113 - 111.108 = 0.005 \end{aligned}$$

درحالی که مقدار دقیق  $f(12345)$  عبارت است از

$$f(12345) = 0.0004500003262627751 \dots$$

بنابراین جواب محاسبه شده دارای ۱۰ درصد خطاست. اکنون روند محاسباتی را مورد تحلیل قرار می‌دهیم. این روند شامل چهار مرحله محاسباتی به شرح زیر است

$$x_0 := 12345$$

$$x_1 := x_0 + 1$$

$$x_2 := \sqrt{x_1} \quad (22.1)$$

$$x_3 := \sqrt{x_0}$$

$$x_4 := x_2 - x_3$$

اکنون برای مثال تابع  $f_3$ ، یعنی تابعی را که چگونگی وابستگی جواب نهایی  $x_4$  به  $x_4$  را بیان می‌کند در نظر می‌گیریم. داریم

$$f_3(t) = x_2 - t$$

بنابراین شرط تابع  $f_3$  تقریباً چنین است

$$\left| \frac{f_3'(t)t}{f_3(t)} \right| = \left| \frac{t}{x_2 - t} \right|$$

این عدد معمولاً نزدیک به ۱ است، یعنی  $f_3$  معمولاً خوش شرط است، مگر اینکه  $t$  نزدیک به  $x_4$  باشد. درحالت اخیر  $f_3$  می‌تواند کاملاً بد شرط باشد. مثلاً در مورد خاص مسئله ما  $t \approx 111.11$ ، و حال آنکه  $0.005 \approx x_4 - t$ ، در نتیجه شرط تابع  $f_3$  تقریباً برابر با  $22222$  یا  $40000$  مرتبه بزرگتر از شرط خود  $f$  است.

از اینجا نتیجه می‌گیریم که روند مذکور در (22.1) راه ناپایداری برای محاسبه  $f$  است. البته اگر ابتدای این قسمت را با دقت مطالعه کرده باشید می‌دانید که یک راه پایدار برای محاسبه این تابع همان استفاده از فرمول هم ارز آن، یعنی فرمول زیر است

$$f(x) = \frac{1}{\sqrt{x+1} + \sqrt{x}}$$

در حساب با شش رقم اعشاری داریم

$$f(12345) = \frac{1}{\sqrt{12346} + \sqrt{12345}} = \frac{1}{2229221} = 0.000450002$$

که فقط دارای خطایی برابر ۰۰۰۰۰۳ درصد است. روند محاسباتی به شرح زیر است

$$x_0 = 12345$$

$$x_1 = x_0 + 1$$

$$x_2 = \sqrt{x_1}$$

$$x_3 = \sqrt{x_0} \quad (23.1)$$

$$x_4 = x_2 + x_3$$

$$x_5 = 1/x_4$$

برای مثال، در اینجا  $f_3(t) = 1/(x_2 + t)$ ، و شرط این تابع تقریباً به ازای  $t = x_3$  عبارت است از

$$\left| \frac{f_3'(t)t}{f_3(t)} \right| = \left| \frac{t}{x_2 + t} \right| \approx \frac{1}{2}$$

بنابراین شرط  $f_3$  کاملاً خوب، و به خوبی شرط خود  $f$  است. مثالهایی از ناپایداریهای بزرگ را، بخصوص در مبحث حل معادلات دیفرانسیل، مشاهده خواهید کرد.

### تمرین

۱-۴۰۱ با استفاده از فرمول (۱۸.۱) و (۲۰.۱) کوچکترین ریشه معادله

$$x^2 + 0.4002 \times 10^6 x + 0.8 \times 10^{-4} = 0$$

را به دست آورید. محاسبات را با ممیز شناور و مانتیس چهاررقمی اعشاری انجام دهید.

۲-۴۰۱ در صورتی که خطای مطلق  $x$  برابر با  $10^{-6}$  باشد، خطای ناشی از محاسبه  $f(x) = (\cos x) \exp(10x^2)$  را در حول نقطه  $x = 2$  برآورد کنید.

۳-۴۰۱ برای محاسبه عبارت‌های

$$f(x) = \frac{x - \sin x}{\tan x} \quad (\text{الف})$$

$$f(x) = (\alpha + x)^n - \alpha^n \quad (\text{ب})$$

$$f(x) = \sin(\alpha + x) - \sin \alpha \quad (\text{پ})$$

$$f(x) = x - \sqrt{x^2 - \alpha} \quad (\text{ت})$$

راهی پیدا کنید که  $f(x)$  را درست با همان تعداد ارقام  $x$  به دست دهد، وقتی که مقدار  $x$  در (الف) تا (پ) نزدیک به صفر و در (ت) به مراتب بزرگتر از  $\alpha$  باشد.

۴-۴۰۱ فرض کنید کامپیوتری با مانیتیس چهاررقمی اعشاری در دست باشد. اعداد زیر را اول به ترتیب صعودی (از کوچکتر تا بزرگتر) و بعد به ترتیب نزولی با هم جمع کنید. در حین عمل حاصلجمعهای جزئی را گسرد کنید. نتایج خود را با جواب درست  $105 \times 10^{23} \times 1010710105$ ، مقایسه نمایید.

$0.1580 \times 10^0$	$0.6266 \times 10^2$	$0.8999 \times 10^4$
$0.2653 \times 10^0$	$0.7555 \times 10^2$	
$0.2581 \times 10^1$	$0.7889 \times 10^3$	
$0.4288 \times 10^1$	$0.7767 \times 10^3$	

۵-۴۰۱ يك روش آشكارا ناپایدار برای محاسبه  $f(x) = e^x$ ، به ازای  $x$  منفی، بسط آن به سری تیلر (۳۶.۱) است. مقدار  $e^{-12}$  را از راه بسط به سری تیلر (۳۶.۱) و به ازای  $x = -12$  محاسبه کنید و مقدار محاسبه شده را با مقدار دقیق

$$e^{-12} = 0.0000006144212354 \dots$$

مقایسه کنید. [داهنمایی: به موجب (۳۶.۱) اختلاف بین  $e^x$  و حاصلجمع جزئی  $S_n = \sum_{j=0}^n x^j / j!$  وقتی  $x$  منفی باشد، از قدر مطلق جمله بعدی یعنی  $|x^{n+1} / (n+1)!|$  کوچکتر است. بنا بر این بجا خواهد بود که سری را تا  $S_n = S_{n+1}$  جمع کنیم.]

۶-۴۰۱ نتیجه تمرین ۵-۴۰۱ را از راه مقایسه شرط  $f(x) = e^x$  در مجاورت  $x = -12$  با شرط برخی از توابع  $f$  موجود در روند محاسباتی فوق، بیان کنید. سپس يك روش پایدار برای محاسبه  $e^{-12}$  از راه سری تیلر (۳۶.۱) پیدا کنید. (داهنمایی:  $e^{-x} = 1/e^x$ .)

## ۵.۱ روشهای کامپیوتری برای برآورد خطا

هدف ما در این فصل آگاه ساختن دانشجویان به منشأ ممکن خطا و ارائه تکنیکهایی است



که می توان برای اجتناب از آنها به کار گرفت. در ارزیابی نتایج حاصله از کامپیوتر این خطاها را باید منظور کرد. بر آورد واقعی خطای کلی در یک مسئله عملی کاری است دشوار و هنوز یک نظریه ریاضی مناسب برای آن وجود ندارد. یک فکر جالب استفاده از خود کامپیوتر برای تهیه این گونه برآوردهاست. روشهای گوناگونی از این نوع پیشنهاد شده اند که پنج تای آنها را مختصراً شرح می دهیم. ساده ترین روش استفاده از دقت مضاعف است. در این روش یک مسئله اساساً دوبار حل می شود - یکبار با دقت معمولی و یکبار با دقت مضاعف. از تفاضل این دو نتیجه می توان کل خطای ناشی از گرد کردن را (با فرض اینکه همه خطاهای دیگر ناچیز باشند) تخمین زد. در این صورت می توان فرض کرد که در مسائل دیگری که با همین زیر برنامه حل می شوند همان اتباشنگی خطای ناشی از گرد کردن پدید خواهد آمد. این روش در رابطه با وقت مصرفی کامپیوتر فوق العاده پرهزینه است زیرا که در برخی از کامپیوترها حساب با دقت مضاعف وقت مورد نیاز را هشت برابر افزایش می دهد و بعلاوه جدا کردن خطاهای دیگر همواره ممکن نیست.

یک روش دیگر حساب بازه ای<sup>۱</sup> است. در اینجا هر عدد به زبان ماشین به وسیله دو عدد بیان می شود که ماکسیمم و مینیمم مقادیری هستند که آن عدد ممکن است اختیار کند. هر وقت عملی انجام گیرد، این مقادیر ماکسیمم و مینیمم محاسبه می شوند. بنابراین اصلاً در هر مرحله دو جواب به دست می آید که جواب واقعی ازوماً در محدوده این دو مقدار ماکسیمم و مینیمم جواب قرار دارد. این روش مستلزم بیش از دو برابر وقت کامپیوتر و تقریباً دو برابر ذخیره سازی استانده است. بعلاوه این فرض متداول، که جواب واقعی در وسط محدوده ما قرار دارد، در حالت کلی صادق نیست. از این رو ممکن است این محدوده به اندازه ای وسیع باشد که هر برآوردی از خطای گرد کردن بر پایه این روش بسیار اغراق آمیز باشد.

نگرش سوم حساب رقم - بامعنی است. همان طوری که قبلاً اشاره شد هر گاه دو عدد تقریباً نزدیک به هم در زبان ماشین، از هم تفریق شوند این خطر وجود دارد که برخی از ارقام بامعنی از دست بروند. در حساب رقم - بامعنی سعی برای این است که رد اعداد بامعنی که بدین طریق از دست می روند حفظ شود. در یکی از این شیوه ها فقط ارقام با معنی هر عدد نگهداشته و کلیه ارقام دیگر حذف می شود. بدین ترتیب در انتهای محاسبه مطمئنیم که تمام ارقام حفظ شده بامعنی هستند. اشکال اصلی این روش این است که برخی اطلاعات با حذف ارقام از بین می روند و بانه نتیجه به دست آمده احتمالاً باید بسیار محتاطانه عمل کرد. گرچه تجربیات به دست آمده با این تکنیک نوید بخش نبوده است، اما آزمایش با این تکنیک هنوز ادامه دارد.

چهارمین روش که نوید بخش فراهم کردن یک نظریه ریاضی مناسب برای انتشار خطای گرد کردن است، بر نگرش آماری استوار شده است. این نگرش با فرض مستقل بودن خطاهای گرد کردن آغاز می شود. البته این فرض معتبری نیست، زیرا اگر مسئله ای چندین بار

روی يك ماشين حل شود جوابها همواره يكسان خواهند بود. اما می توانیم يك الكوی تصادفی<sup>۱</sup> از انتشار خطای گرد کردن، که در آن به خطاهای موضعی به صورت متغیرهای تصادفی نگریسته می شود، در نظر بگیریم. بنا بر این می توانیم فرض کنیم که مقادیر خطاهای موضعی به طور یکنواخت و یا به طور نرمال بین مقادیر فرینه<sup>۲</sup> توزیع شده اند. با استفاده از روشهای آماری می توان انحراف معیار<sup>۳</sup> واریانس<sup>۴</sup> توزیع را به دست آورد و در نتیجه خطای گرد کردن انباشته را بر آورد کرد. این نگرش آماری به طور نسبتاً مشروح به توسط همینگ<sup>۵</sup>[۱] و هنریتسی<sup>۶</sup>[۲] بررسی شده است. این روش مستلزم تحلیل اساسی و وقت اضافی برای کامپیوتر است، اما در آزمایشهایی که تا این تاریخ انجام شده نتایجی از بر آورد خطاها به دست آمده است که تا حد زیادی با مدارك تجربی موجود توافق دارد.

روش پنجم، روش تحلیل خطای پسر و است که در بخش ۳.۱ معرفی شد. همان طوری که مشاهده کردیم، این روش تحلیل اثرات خطای گرد کردن را به بررسی آشفته گیها در حساب دقیق و در نهایت به يك مسئله مربوط به شرط، بدل می کند. در فصل ۴ این روش به نحو مؤثری مورد استفاده قرار خواهد گرفت.

## ۶.۱ توضیحاتی در مورد همگرایی دنباله ها

حساب دیفرانسیل و انتگرال و به طور کلی آنالیز، بر پایه مفهوم همگرایی استوار است. مفاهیم پایه ای مانند مشتق و انتگرال و پیوستگی، بر حسب دنباله های همگرا و توابع مقدماتی مانند  $\ln x$  یا  $\sin x$  تعریف می شوند. از طرفی برای مسائل مهندسی و علمی هرگز جوابهای عددی کاملاً دقیق مورد نیاز نیستند، بلکه جوابهایی تقریبی، لازم اند که «تاعده معینی ارقام اعشاری» یا در حد قابل تحملی مانند  $\epsilon$  دقیق باشند.

از این روش شگفت آور نیست اگر در پیدا کردن  $\alpha$ ، جواب يك مسئله مفروض، بسیاری از روشهای عددی، تنها (چند جمله اول)  $\alpha_1$  و  $\alpha_2$  و  $\alpha_3$  و... از يك دنباله را که همگرا بودنش به سمت جواب مطلوب نشان داده شده است، به دست دهند. یادآوری می کنیم که طبق تعریف:

يك دنباله اعداد (حقیقی یا مختلط)  $\alpha_1, \alpha_2, \alpha_3, \dots$  به سمت  $\alpha$  همگراست اگر فقط اگر، به ازای جميع مقادیر  $\epsilon > 0$  عدد صحیحی مانند  $n_0(\epsilon)$  وجود داشته باشد به طوری که به ازای هر  $n \geq n_0$ ، نامساوی  $|\alpha - \alpha_n| < \epsilon$  برقرار باشد.

بنابر این اگر يك روش عددی داشته باشیم که دنباله  $\alpha_1, \alpha_2, \alpha_3, \dots$  را تولید کند و این دنباله به سمت جواب مطلوب  $\alpha$  همگرا باشد، آنگاه می توان  $\alpha$  را تا هر دقت دلخواهی فقط با محاسبه  $\alpha_n$  به ازای مقدار «به اندازه کافی بزرگ»  $n$  محاسبه کرد.

- |               |            |                       |
|---------------|------------|-----------------------|
| 1. stochastic | 2. extreme | 3. standard deviation |
| 4. variance   | 5. Hamming | 6. Henrici            |

از دیدگاه محاسباتی این تعریف به دلایل زیر رضایتبخش نیست:

(۱) غالباً (بدون دانستن جواب  $\alpha$ ) امکان ندارد که درایم چه موقع  $n$  "به اندازهٔ کافی بزرگ" است. به عبارت دیگر پیدا کردن تابع  $n_0(\varepsilon)$  که در تعریف همگرایی ذکر شد، بسیار مشکل است.

(۲) حتی وقتی هم اطلاعاتی دربارهٔ  $n_0(\varepsilon)$  در دست باشد، ممکن است  $n$  مورد نیاز به اندازه‌های بزرگ شود که محاسبهٔ  $\alpha_n$  را غیرممکن سازد.

□ مثال: عدد  $\pi/4$  برابر با مقدار رشتهٔ نامتناهی<sup>۱</sup> زیر است

$$\sum_{i=0}^{\infty} \frac{(-1)^i}{2i+1} = 1 - \sum_{j=1}^{\infty} \frac{2}{16j^2-1}$$

از این رو با

$$\alpha_n = 1 - \sum_{j=1}^n \frac{2}{16j^2-1} \quad n=1, 2, \dots$$

دنبالهٔ  $\alpha_1, \alpha_2, \dots$  در حالی که یکنواکم می‌شود به حد خود یعنی  $\pi/4$  میل می‌کند. علاوه داریم

$$0 \leq \alpha_n - \frac{\pi}{4} \leq \frac{1}{4n+3} \quad n=1, 2, \dots$$

برای اینکه از دنبالهٔ بالا برای محاسبهٔ  $\pi/4$  تا  $10^{-6}$  رقم صحیح استفاده کنیم، لازم است داشته باشیم  $4n+3 \leq 10^6$  یا تقریباً  $n = 250000$ . اما در کامپیوتری که با ممیز شناور هشت رقم اعشاری کار می‌کند خطای گرد کردن در محاسبهٔ  $\alpha_{250000}$  احتمالاً خیلی بزرگتر از  $10^{-6}$  خواهد شد. بنابراین  $\pi/4$  را نمی‌توان با استفاده از دنبالهٔ فوق تا  $10^{-6}$  رقم صحیح محاسبه نمود (مگر، احتمالاً با اضافه کردن جملات به یکدیگر از کوچکترین تا بزرگترین). □

برای پرداختن به این گونه مسائل، برخی قراردادهای نمادگذاری مفیدند. بخصوص آنکه می‌خواهیم سرعت همگرایی دنباله‌ها را اندازه‌گیری کنیم. مانند تمام اندازه‌گیریهای دیگر، این عمل را نیز به وسیلهٔ مقایسه با دنباله‌های استاندارد<sup>۲</sup>، مانند دنباله‌های

$$\left. \begin{array}{l} 1/n \\ 1/n^r \\ r^n \\ 1/(\ln n) \end{array} \right\} n=1, 2, 3, \dots$$

انجام می‌دهیم. مقایسه به گونه‌ی زیر انجام می‌گیرد: گوئیم  $\alpha_n$  از مرتبه  $\beta_n$  (یا  $\alpha_n$  مهضراً  $\beta_n$ ) است و می‌نویسیم

$$\alpha_n = O(\beta_n) \quad (24.1)$$

وقتی به ازای مقدار ثابتی مانند  $K$  و جمیع مقادیر به اندازه کافی بزرگ  $n$  داشته باشیم

$$|\alpha_n| \leq K |\beta_n| \quad (25.1)$$

بنابراین

$$\left. \begin{array}{l} 1/n \\ 10,000/n \\ 10/n - 40/n^2 + e^{-n} \\ 1/n^2 \end{array} \right\} = O(1/n)$$

و انگهی هر گاه به محض اینکه  $n$  به قدر کافی بزرگ شد، بتوان ثابت  $K$  را در (25.1) به دلخواه کوچک انتخاب کرد، یعنی اگر حالت

$$\lim_{n \rightarrow \infty} \alpha_n / \beta_n = 0$$

پیش آید، آنگاه گفته می‌شود که مرتبه  $\alpha_n$  از مرتبه  $\beta_n$  بیشتر (یا  $\alpha_n$  کهضراً  $\beta_n$ ) است و می‌نویسیم

$$\alpha_n = o(\beta_n) \quad (26.1)$$

بنابراین وقتی  $\sin(1/n) \neq o(1/n)$ ، داریم

$$\left. \begin{array}{l} 1/n^2 \\ 1/(n \ln n) \end{array} \right\} = o(1/n)$$

معمولاً علامت مرتبه فقط در سمت راست يك معادله ظاهر می‌شود و هدف آن تشریح جنبه اصلی يك جمله خطا، بدون توجه به ضرب ثابتها یا سایر جزئیات است. برای مثال، می‌توانیم حالت نامطلوب موضوع را در مثال قبلی، احتمالاً به وسیله روابط زیر بیان کنیم

$$1 - \sum_{j=2}^n 1/(j^2 - 1) = \pi/4 + O(1/n)$$

اما همچنین

$$1 - \sum_{j=2}^n 1/(j^2 - 1) \neq \pi/4 + o(1/n)$$

یعنی، این دنباله با همان سرعت  $1/n$  (سرعتی که  $1/n$  به سمت صفر میل می‌کند) و نه سریعتر، به سمت  $\pi/4$  همگرا می‌شود. مرتبه یا میزان همگرایی  $1/n$  به اندازه‌ای کم است که در محاسبات نمی‌تواند مورد استفاده قرار گیرد.

□ مثال: اگر  $\alpha_n = \alpha + o(1)$ ، آنگاه بنا بر تعریف داریم

$$\lim_{n \rightarrow \infty} \frac{\alpha_n - \alpha}{1} = 0$$

بنا بر این  $\alpha_n = \alpha + o(1)$  فقط نوعی سلیقه برای بیان همگرایی دنبالهٔ  $\alpha_1, \alpha_2, \dots$  به سمت  $\alpha$  است.

□ مثال: اگر  $|r| < 1$ ، آنگاه مجموع سری هندسی  $\sum_{i=0}^{\infty} r^i$  برابر است با  $s_n = \sum_{i=0}^n r^i \cdot 1/(1-r)$  داریم

$$s_n = (1 - r^{n+1}) / (1 - r) = \frac{1}{1-r} - \frac{r^{n+1}}{1-r}$$

بنا بر این

$$s_n = \frac{1}{1-r} + \theta(r^n)$$

بعلاوه اگر  $|r| > r$ ؛ آنگاه

$$s_n = \frac{1}{1-r} + o(r^n)$$

بنا بر این هر گاه تساوی  $\alpha_n = \alpha + \theta(r^n)$  به ازای يك  $|r| < 1$  برقرار باشد، گوییم همگرایی (حد اقل) هندسی است، زیرا در این صورت مرتبهٔ همگرایی (حد اقل) مساوی مرتبهٔ همگرایی سری هندسی خواهد بود.

□

با وجودی که دانستن رابطهٔ  $\alpha_n = \alpha + \theta(\beta_n)$  از ندانستن آن بهتر است، آگاهی از مرتبهٔ همگرایی زمانی کاملاً مفید است که به طور دقیق‌تری رابطهٔ  $\alpha_n = \alpha + \beta_n + o(\beta_n)$  را بشناسیم. رابطهٔ بالا بدین معنی است که به ازای مقدار «به اندازهٔ کافی بزرگ»  $n$ ،  $\alpha_n \approx \alpha + \beta_n$  به عبارت دیگر

$$\begin{aligned}\alpha_n &= \alpha + \beta_n + o(\beta_n) \\ &= \alpha + \beta_n + \beta_n o(1) \\ &= \alpha + \beta_n(1 + \varepsilon_n)\end{aligned}$$

در رابطه بالا،  $\varepsilon_p$ ،  $\varepsilon_q$ ، ... دنباله‌ای است که به صفر همگراست. گرچه «به اندازه کافی بزرگ» بودن مقداری از  $n$  را نمی‌توانیم ثابت کنیم، اما می‌توانیم از مقایسه  $\alpha_k - \alpha_{k+1}$  با  $\beta_k - \beta_{k+1}$ ، این فرض را که  $n$  «به اندازه کافی بزرگ است» آزمایش کنیم. اگر برای  $k$  نزدیک به  $n$  مثلاً  $n-1$ ،  $n-2$ ،  $k=n-2$  داشته باشیم

$$\frac{|\alpha_{k+1} - \alpha_k|}{|\beta_{k+1} - \beta_k|} \approx 1$$

آنگاه می‌توانیم درستی این فرض را، که  $n$  «به قدر کافی بزرگ» است، برای

$$\alpha_n \approx \alpha + \beta_n$$

پذیریم و بنابراین  $|\beta_n|$  را به عنوان یک برآورد مطلوب از خطای  $|\alpha - \alpha_n|$  قبول کنیم.

□ مثال: گیریم  $p > 1$ . پس سری  $\sum_j 1/(p^j + 1)$  مانند سری هندسی  $\sum_j 1/p^j$  به سمت حد خود  $\alpha$  همگراست، یعنی

$$\alpha_n = \sum_{j=1}^n 1/(p^j + 1) = \alpha + o(1/p^{n+1})$$

برای اینکه حکم دقیقتری به دست آوریم، معادله زیر را در نظر می‌گیریم

$$\beta_n = \sum_{j=n+1}^{\infty} 1/p^j = (1/p^{n+1}) / (1 - 1/p) = 1/[p^n(p-1)]$$

پس

$$\begin{aligned}\alpha_n &= \alpha - \sum_{j=n+1}^{\infty} 1/(p^j + 1) = \alpha - \beta_n + \sum_{j=n+1}^{\infty} [1/p^j - 1/(p^j + 1)] \\ &= \alpha - \beta_n + o(\beta_n)\end{aligned}$$

زیرا که

$$\begin{aligned}0 \leq \sum_{j=n+1}^{\infty} [1/p^j - 1/(p^j + 1)] &= \sum_{j=n+1}^{\infty} 1/(p^j(p^j + 1)) \leq \sum_{j=n+1}^{\infty} (1/p^2)^j \\ &= o((1/p^2)^{n+1})\end{aligned}$$

برای نسبت فوق داریم:

$$\left| \frac{\alpha_{n+1} - \alpha_n}{\beta_{n+1} - \beta_n} \right| = p^{n+1} / (p^{n+1} + 1)$$

که مثلاً برای  $n=3$  و  $p=2$  در حدود  $1/10$  از  $1$  است. بنابراین  $\beta_4 = 1/8 = 0.125$ ،  
شاخص خوبی است از خطا در  $0.0644444$  زیرا  $\alpha_4 = 0.0764449$ ،  
که در این صورت خطا در  $\alpha_4$  برابر است با  $0.012005$ . □

این نمادگذاری می‌تواند برای توابعی از یک متغیر حقیقی بسط داده شود. اگر

$$\lim_{h \rightarrow 0} T(h) = A$$

همگرایی را  $\theta(f(h))$  گوئیم به شرطی که برای مقدار ثابت منتهای  $K$  و کلیه  $h$ ‌های به اندازهٔ کافی کوچک، داشته باشیم

$$\frac{|T(h) - A|}{|f(h)|} \leq K$$

اگر رابطهٔ فوق به‌ازای جميع مقادیر  $0 < K$  صادق باشد، یعنی اگر

$$\lim_{h \rightarrow 0} \frac{T(h) - A}{f(h)} = 0$$

آنگاه همگرایی را  $o(f(h))$  می‌نامیم.

□ مثال: به‌ازای  $h$  «نزدیک» به صفر داریم

$$\frac{\sin h}{h} = 1 - \left(\frac{1}{3!}\right) h^2 + \left(\frac{1}{5!}\right) h^4 - \dots = 1 + \theta(h^2)$$

$$= 1 - \frac{1}{6} h^2 + o(h^2)$$

بنابراین، به‌ازای جميع مقادیر  $\gamma < 2$  داریم

$$\frac{\sin h}{h} = 1 + o(h^\gamma)$$

□

□ مثال: اگر تابع  $f(x)$  ریشه‌ای از مرتبهٔ  $\gamma$  در  $x = \xi$  داشته باشد، آنگاه:

$$\square \quad f(\xi + h) \neq o(h^\gamma) \quad \text{اما} \quad f(\xi + h) = \theta(h^\gamma)$$

قواعد محاسبه با نمادهای مرتبه‌ای در لم ۲ زیر گردآوری شده است.

لم ۱۰۱ اگر  $\alpha_n = \alpha + \theta(f(n))$  و  $\lim_{n \rightarrow \infty} f(n) = 0$  و  $c$  عدد ثابتی باشد، آنگاه:

$$c\alpha_n = c\alpha + \theta(f(n))$$

همچنین اگر  $\beta_n = \beta + \theta(g(n))$  و  $g(n) = \theta(f(n))$ ، در این صورت داریم:

$$\alpha_n + \beta_n = \alpha + \beta + \theta(f(n)) \quad \text{و} \quad \alpha_n \beta_n = \alpha\beta + \theta(f(n)) \quad (27.1)$$

بعلاوه اگر  $\beta \neq 0$ ، همچنین داریم

$$\frac{\alpha_n}{\beta_n} = \frac{\alpha}{\beta} + \theta(f(n))$$

در صورتی که اگر  $\alpha = \beta = 0$ ، آنگاه

$$\alpha_n \beta_n = \theta(f(n)g(n))$$

بالاخره کلیه احکام فوق با قراردادن  $0$  به جای  $\theta$  در همه جا، صحیح خواهد ماند.

محاسبه تقریبی يك عدد  $\alpha$  به وسیله دنباله  $\alpha_1, \alpha_2, \dots$  که به سمت  $\alpha$  همگراست، خواه مرتبه همگرایی معلوم باشد و خواه نباشد، همیشه امری است توأم با اطمینان. اگر همگرایی این دنباله به سمت  $\alpha$  معلوم باشد، يك تحلیلگر مجرب روشهای عددی، پس از اطمینان به اینکه به ازای مقادیر کوچک  $n$ ، تفاوت  $\alpha_n$  با  $\alpha$  «به اندازه کافی کوچک است»، «به اندازه کافی بزرگ» بودن عدد  $n$  را تحقیق خواهد کرد. همچنین در صورتی که بداند همگرایی به صورت  $(\beta_n + o(\beta_n))$  است، بررسی می کند که آیا وضعیت این دنباله در نزدیکی  $n$  به همان ترتیب است یا خیر. همچنین اگر بداند که  $\alpha$  در معادلات و یا نامعادلاتی صدق می کند  $\alpha$  می تواند جواب يك معادله باشد. بررسی می کند که آیا  $\alpha_n$  در آن معادلات و یا نامعادلات «به اندازه کافی خوب» صدق می کند یا نه. کوتاه سخن آنکه يك تحلیلگر مجرب روشهای عددی، مطمئن می شود که  $n$  در تمام شرایط متصور لازم برای «به اندازه کافی بزرگ» بودن صدق می کند. در صورتی که تمامی این شرایط برقرار باشند و شرط کافی برای اینکه  $n$  «به اندازه کافی بزرگ» باشد در دست نباشد، با اطمینان  $\alpha_n$  را تقریب «به اندازه کافی خوبی» برای  $\alpha$  قبول می کند. به عبارت دیگر، يك تحلیلگر تمامی روشهایی را که در اختیار دارد برای تمیز يك تقریب «به اندازه کافی خوب» از يك تقریب بد به کار می گیرد. بیش از این نمی تواند کرد. انجام دهد (و کمتر هم نباید انجام دهد).

از اینجا معلوم می شود که نتایج عددی به دست آمده از این طریق را نباید با جواب نهایی اشتباه گرفت، بلکه چنانچه بررسیهای بعدی صحت آنها را مورد تردید قرار دهند، بایستی با خیال راحت در درستی آنها تأمل کرد.

دانشجویان باید این مطلب را به عنوان مثال دیگری از تفاوت اساسی بین آنالیز



عددی و آنالیز تلقی کنند. آنالیز زمانی به صورت یک نظام دقیق درآمد که محدودیتهای محاسبات عملی را کلاً به عهدهٔ مسائلی گذاشت که بر حسب الگوی مجردی از دستگاه عددی، به نام اعداد حقیقی، طرح شده بودند. این الگوی مجرد چنان طرح شده است که دادن یک تعریف دقیق و مفید از حد را ممکن می‌سازد و در صورتی که این گونه مسائل به مصطلحات چنین الگویی ترجمه شوند، این تعریف راه را برای حل تجریدی یا نمادی تعداد قابل توجهی از مسائل عملی می‌گشاید. و این امر به نوبهٔ خود کار ترجمهٔ حل تجریدی یا نمادی را به حل عملی بازمی‌گرداند. آنالیز عددی این وظیفه را به عهده می‌گیرد و همراه با آن، محدودیتهای موجود در محاسبات عملی را که آنالیز با ظرافت از آنها اجتناب می‌ورزد، تقبل می‌کند. از این رو آنالیز عددی الزاماً غیر دقیق است و جوابهای عددی آن معمولاً آزمایشی هستند و در نهایت فقط تا حدود معینی صحیح‌اند.

بنابراین، کار آنالیز عددی فقط ساختن روشهای عددی نیست، بلکه قسمت اعظم آن استخراج کرانه‌های خطای مفید یا برآوردهای خطا برای جوابهای حاصل از الگوریتمهای عددی است. در سراسر این کتاب دانشجویان قبل از هر چیز با این کرانه‌های خطا که در آنالیز عددی فراوان هستند مواجه خواهند شد.

## نهمین

۱-۶۰۱  $\ln 2$  را می‌توان از سری

$$\ln 2 = 1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \dots$$

محاسبه کرد. در درس آنالیز دیده‌ایم که این سری همگرا و اندازهٔ خطا در هر حاصلجمع جزئی آن کمتر از اندازهٔ اولین جملهٔ صرف نظر شده است. تعداد جملات لازم برای محاسبهٔ  $\ln 2$  را تا ده رقم اعشاری برآورد کنید.

۲-۶۰۱ به ازای  $h$  نزدیک به صفر می‌توان نوشت

$$\frac{\tan h}{h} = 1 + \mathcal{O}(h^2)$$

و

$$\frac{\tan h}{h} = 1 + o(h^2)$$

مقادیر  $\gamma$  و  $\delta$  را طوری پیدا کنید که تساویهای فوق برقرار باشند.

۳-۶۰۱ سعی کنید حد دنباله زیر را به وسیله يك کامپیوتر محاسبه نمایید

$$\alpha_n = (\tan \lambda^{-n} - \sin \lambda^{-n}) \cdot \lambda^{2n} \quad n = 0, 1, 2, \dots$$

از لحاظ نظری دنباله  $\alpha_n$  را  $\alpha = \lim_{n \rightarrow \infty} \alpha_n$  و مرتبه همگرایی آن را پیدا کنید.

## ۷.۱ برخی مقدمات ریاضی

فرض بر این است که دانشجویان با مباحث آشنای دروس حساب دیفرانسیل و انتگرال و هندسه تحلیلی در دوره لیسانس آشنایی دارند. این درس شامل مفاهیم مقدماتی دستگاههای اعداد حقیقی و مختلط، پیوستگی، مفهوم حد، دنباله‌ها، رشته‌ها و مشتقگیری و انتگرالگیری است. برای فصل چهارم آشنایی کمی با دترمینانها مسلم فرض شده است. برای فصل هشتم و نهم نیز داشتن آشنایی اندکی با حل معادلات دیفرانسیل معمولی مسلم انگاشته شده است، گرچه این فصول ممکن است تدریس نشوند. و به ویژه، کراراً از قضیه‌های زیر استفاده خواهیم کرد.

**قضیه ۱۰۱:** قضیه مقدار میانه برای توابع پیوسته. فرض می‌کنیم تابع  $f(x)$  در بازه  $[a, b]$  پیوسته باشد. اگر به ازای اعداد  $\alpha$  و  $\bar{x} \in [a, b]$  داشته باشیم  $f(\bar{x}) \leq \alpha \leq f(\bar{x})$ ، آنگاه حداقل يك عدد  $\xi \in [a, b]$  وجود دارد به طوری که  $\alpha = f(\xi)$ . این قضیه غالباً به شکل زیر به کار برده می‌شود:

**قضیه ۲۰۱:** گیریم  $f(x)$  تابعی پیوسته در بازه  $[a, b]$ ، و  $x_1, \dots, x_n$  نقاطی بر  $[a, b]$  باشند. همچنین فرض می‌کنیم  $g_1, \dots, g_n$  اعداد حقیقی همعلامت باشند، در این صورت عددی مانند  $\xi \in [a, b]$  وجود دارد به طوری که

$$\sum_{i=1}^n f(x_i)g_i = f(\xi) \sum_{i=1}^n g_i$$

به منظور اشاره به اثبات قضیه، بی آنکه خطایی در کلیت حاصل آید، فرض می‌کنیم  $i = 1, \dots, n$ ،  $g_i \geq 0$ . اگر  $f(x) = \min_i f(x_i)$  و  $f(\bar{x}) = \max_i f(x_i)$ ، آنگاه

$$f(x) \sum_{i=1}^n g_i = \sum_{i=1}^n f(x)g_i \leq \sum_{i=1}^n f(x_i)g_i \leq \sum_{i=1}^n f(\bar{x})g_i = f(\bar{x}) \sum_{i=1}^n g_i$$

بنابراین  $\alpha = \sum_{i=1}^n f(x_i)g_i$  عددی است بین دو مقدار  $f(x) \sum_{i=1}^n g_i$  و  $f(\bar{x}) \sum_{i=1}^n g_i$  از تابع پیوسته  $f(x) \sum_{i=1}^n g_i$ ، با توجه به قضیه ۱۰۱ حکم قضیه ۲۰۱ نتیجه می‌شود.

به نحو مشابه می‌توان حکم متناظری را برای حاصلجمعهای نامتناهی<sup>۱</sup> یا انتگرالها ثابت کرد.

**قضیه ۳.۱:** قضیه مقدار میانگین برای انتگرالها. گیریم  $g(x)$  تابعی نامنفی یا نامثبت و انتگرالپذیر در  $[a, b]$  باشد. هر گاه  $f(x)$  روی  $[a, b]$  پیوسته باشد، به ازای مقداری مانند  $\xi \in [a, b]$  داریم

$$\int_a^b f(x)g(x) dx = f(\xi) \int_a^b g(x) dx \quad (28.1)$$

**توجه:** چنانچه از مثال ساده

$$f(x) = g(x) = x, [a, b] = [-1, 1]$$

برمی آید، فرض یک علامت داشتن  $g(x)$  در قضیهٔ ۳.۱، فرضی است اساسی.

**قضیه ۴.۱:** گیریم  $f(x)$  تابعی پیوسته در بازهٔ بسته و کراندار  $[a, b]$  باشد. پس  $f(x)$  دارای «مقادیر ماکسیمم و مینیمم» در روی  $[a, b]$  است، یعنی نقاط  $\underline{x}$  و  $\bar{x}$  در بازهٔ  $[a, b]$  طوری وجود دارند که به ازای هر  $x \in [a, b]$  نامساویهای

$$f(\underline{x}) \leq f(x) \leq f(\bar{x})$$

برقرارند.

**قضیه ۵.۱:** قضیه رول. فرض می‌کنیم  $f(x)$  در بازهٔ (بسته و منتهای)  $[a, b]$  پیوسته و در بازهٔ  $(a, b)$  مشتقپذیر باشد. اگر  $f(a) = f(b) = 0$ ، آنگاه مقداری مانند  $\xi \in [a, b]$  وجود دارد که  $f'(\xi) = 0$ .

در اثبات این قضیه عمده‌تأ از قضیهٔ ۴.۱ استفاده می‌شود. زیرا طبق قضیهٔ ۴.۱ نقاطی مانند  $\underline{x}$  و  $\bar{x}$  در  $[a, b]$  وجود دارند به طوری که به ازای جمیع مقادیر  $x \in [a, b]$  نامساویهای  $f(\underline{x}) \leq f(x) \leq f(\bar{x})$  برقرارند. اما اگر نه  $\underline{x}$  و نه  $\bar{x}$  هیچکدام در  $(a, b)$  نباشند، آنگاه  $f(x) \equiv 0$ ، و هر  $\xi \in (a, b)$  در این نامساویها صدق می‌کند. در غیر این صورت یا  $\underline{x}$  و یا  $\bar{x}$  در  $(a, b)$  واقع است، مثلاً  $\bar{x} \in (a, b)$ . در این صورت  $f'(\bar{x}) = 0$ ، زیرا

$$0 \leq \lim_{h \rightarrow 0^-} \frac{f(\bar{x}) - f(\bar{x} + h)}{-h} = f'(\bar{x}) = \lim_{h \rightarrow 0^+} \frac{f(\bar{x} + h) - f(\bar{x})}{h} \leq 0$$

$f(\bar{x})$  بزرگترین مقداری است که  $f(x)$  روی  $[a, b]$  اختیار کرده است.

**قضیه ۶.۱:** قضیه مقدار میانگین برای مشتقها. اگر  $f(x)$  در بازهٔ (بسته و منتهای)  $[a, b]$  پیوسته و در بازهٔ  $(a, b)$  مشتقپذیر باشد، آنگاه مقداری مانند  $\xi \in (a, b)$  وجود دارد به طوری که

$$\frac{f(b)-f(a)}{b-a} = f'(\xi) \quad (۲۹.۱)$$

اگر در قضیه ۵.۱ به جای تابع  $f(x)$  تابع

$$F(x) = f(x) - f(a) - \frac{f(b)-f(a)}{b-a}(x-a)$$

را قرار دهیم، قضیه ۶.۱ به دست خواهد آمد. روشن است که  $F(x)$ ، هم در  $a$  و هم در  $b$  صفر می شود.

از قضیه ۶.۱ مستقیماً نتیجه می شود که اگر  $f(x)$  روی  $[a, b]$  پیوسته و روی  $(a, b)$  مشتق پذیر و  $c$  نقطه ای از  $[a, b]$  باشد، آنگاه مقداری مانند  $\theta \in (0, 1)$  وجود دارد به طوری که

$$f(x) = f(c) + (x-c)f'(c + \theta(x-c)) \quad (۳۰.۱)^*$$

قضیه اساسی حساب دیفرانسیل و انتگرال حکم دقیقتری را در اختیار ما قرار می دهد: اگر  $f(x)$  پیوسته و مشتق پذیر باشد، آنگاه به ازای هر  $x \in [a, b]$  داریم

$$f(x) = f(c) + \int_c^x f'(s) ds \quad (۳۱.۱)$$

از آنجا که  $f'(x)$  پیوسته است، می توانیم با استفاده از رابطه فوق و قضیه میانگین انتگرال (۲۸.۱)، رابطه (۳۰.۱) را به دست آوریم. \*\*

قضیه ۷.۱: فرمول تیلر با باقیمانده (انتگرالی). اگر  $f(x)$  دارای  $n+1$  مشتق پیوسته روی  $[a, b]$  و  $c$  يك نقطه در  $[a, b]$  باشد، آنگاه به ازای هر  $x \in [a, b]$  داریم:

\* چنانچه در قضیه ۶.۱ به جای بازه  $[a, b]$  بازه  $[c, x]$  را در نظر بگیریم، داریم

$$\frac{f(x)-f(c)}{x-c} = f'(\xi)$$

اگر  $\xi = c + \theta(x-c)$ ، رابطه ۳۰.۱ نتیجه می شود. م.

\*\* فرض کنیم  $g(x) = 1$ . طبق قضیه ۳.۱ داریم:

$$\int_c^x f'(s) ds = f'(\xi) \int_c^x ds = (x-c)f'(\xi) = (x-c)f'(c + \theta(x-c))$$

و در نتیجه رابطه ۳۰.۱ به دست می آید. م.

$$f(x) = f(c) + f'(c)(x-c) + \frac{f''(c)(x-c)^2}{2!} + \dots + \frac{f^{(n)}(c)(x-c)^n}{n!} + R_{n+1}(x) \quad (۳۲.۱)$$

که در فرمول بالا  $R_{n+1}(x)$  برابر است با

$$R_{n+1}(x) = \frac{1}{n!} \int_c^x (x-s)^n f^{(n+1)}(s) ds \quad (۳۳.۱)$$

رابطه (۳۲.۱) را می‌توان با در نظر گرفتن تابع زیر به‌جای  $f(x)$  از رابطه (۳۱.۱) به‌دست آورد

$$F(x) = f(x) + f'(x)(c-x) + \frac{f''(x)(c-x)^2}{2!} + \dots + \frac{f^{(n)}(x)(c-x)^n}{n!}$$

زیرا در این صورت  $F'(x) = f^{(n+1)}(x)(c-x)^n/n!$  و بنابراین با استفاده از (۳۱.۱) داریم

$$F(x) = F(c) + \frac{1}{n!} \int_c^x (c-s)^n f^{(n+1)}(s) ds$$

اما از آنجا که  $F(c) = f(c)$ ، داریم

$$f(c) = F(x) + \frac{1}{n!} \int_x^c (c-s)^n f^{(n+1)}(s) ds$$

که بعد از قرارداد  $x$  به‌جای  $c$  و  $c$  به‌جای  $x$ ، رابطهٔ فوق همان رابطهٔ (۳۲.۱) خواهد شد. در عمل لازم نیست که  $f^{(n+1)}(x)$  پیوسته باشد تا (۳۲.۱) برقرار گردد. گرچه اگر در (۳۲.۱) تابع  $f^{(n+1)}(x)$  پیوسته باشد با استفاده از قضیهٔ ۳.۱ شکلی آشنا تر، ولی با سودمندی کمتری برای باقیمانده یعنی:

$$R_{n+1}(x) = \frac{f^{(n+1)}(\xi)(x-c)^{n+1}}{(n+1)!}, \quad \xi = c + \theta(x-c) \quad (۳۴.۱)$$

به‌دست می‌آید. با قرارداد  $h = x - c$ ، روابط (۳۲.۱) و (۳۴.۱) شکل زیر را به‌خود می‌گیرند:

$$f(c+h) = f(c) + hf'(c) + \frac{h^2}{2!} f''(c) + \dots + \frac{h^n}{n!} f^{(n)}(c) + \frac{h^{n+1}}{(n+1)!} f^{(n+1)}(c + \theta h) \quad (۳۵.۱)$$

به ازای مقداری مانند  $\theta \in (0, 1)$

□ مثال: تابع  $f(x) = e^x$  در نزدیکی  $c = 0$  بسط تیلری به صورت زیر دارد

$$e^x = 1 + x + \frac{x^2}{2!} + \dots + \frac{x^n}{n!} + \frac{x^{n+1}e^{\xi}}{(n+1)!} \quad (۳۶.۱)$$

به ازای مقداری مانند  $\xi$  بین  $0$  و  $x$ .

بسط  $f(x) = \ln x = \log_e x$  پیرامون  $c = 1$  عبارت است از

$$\ln x = (x-1) - \frac{(x-1)^2}{2} + \frac{(x-1)^3}{3} - \frac{(x-1)^4}{4} + \dots$$

$$- \frac{(-1)^n (x-1)^n}{n} + \frac{(-1)^n (x-1)^{n+1} \xi^{-(n+1)}}{n+1}$$

□ در فرمول بالا  $0 < x \leq 2$  و  $\xi$  بین  $1$  و  $x$  است.

فرمول مشابهی برای توابع چند متغیره برقرار است. این فرمول از قضیه ۷.۱ به کمک قضیه زیر به دست می آید.

قضیه ۸.۱: قاعده مشتق تابع مضاف. اگر تابع  $f(x, y, \dots, z)$  دارای مشتق اول جزئی پیوسته نسبت به هر یک از متغیرهایش باشد و توابع  $x = x(t), y = y(t), \dots, z = z(t)$  نیز به طور پیوسته نسبت به  $t$  مشتقپذیر باشند، آنگاه

$$g(t) = f(x(t), y(t), \dots, z(t))$$

نیز به طور پیوسته مشتقپذیر خواهد بود و

$$g'(t) = \frac{\partial f}{\partial x} x'(t) + \frac{\partial f}{\partial y} y'(t) + \dots + \frac{\partial f}{\partial z} z'(t)$$

از این قضیه و با دخالت تابع

$$g(t) = f(a + t(x-a), b + t(y-b), \dots, c + t(z-c))$$

رابطه‌ای برای بیان  $f(x, y, \dots, z)$  بر حسب مقدار تابع و مشتقهای جزئی آن در  $(a, b, \dots, c)$  به دست می آید و سپس سری بسط تیلر این تابع پیرامون  $t = 0$  در  $t = 1$  ارزیابی می شود. برای مثال توجه به این نکته قضیه زیر را می دهد

**قضیه ۹۰۱:** اگر  $f(x, y)$  در همسایگی  $D$  از نقطهٔ  $(a, b)$  واقع در صفحهٔ  $(x, y)$  دارای مشتقهای جزئی اول و دوم پیوسته باشد، آنگاه به ازای کلیهٔ مقادیر  $(x, y)$  نزدیک به  $D$  و نقطه‌ای مانند  $(\xi, \eta) \in D$  وابسته به  $(x, y)$  داریم:

$$f(x, y) = f(a, b) + f_x(a, b)(x-a) + f_y(a, b)(y-b) + R_2(x, y) \quad (37.1)$$

و

$$R_2(x, y) = \frac{f_{xx}(\xi, \eta)(x-a)^2}{2} + f_{xy}(\xi, \eta)(x-a)(y-b) + \frac{f_{yy}(\xi, \eta)(y-b)^2}{2}$$

در فرمول بالا زیرنما بهای  $f$  معرف مشتقات جزئی هستند.

برای مثال بسط  $e^{x \sin y}$  پیرامون  $(a, b) = (0, 0)$  چنین است:

$$e^{x \sin y} = 1 + 0 \times x + 0 \times y + R_2(x, y) \quad (38.1)$$

بالاخره در بحث ویژه مقدارهای ماتریسها و در جاهای دیگر قضیهٔ زیر لازم می‌شود.

**قضیه ۱۰۰۱:** قضیهٔ بنیادی جبر. اگر  $p(x)$  یک بسجمله‌ای از درجهٔ  $n \geq 1$  باشد، یعنی اگر

$$p(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n$$

و  $a_0, a_1, \dots, a_n$  اعداد حقیقی یا مختلط باشند و  $a_n \neq 0$ ، آنگاه  $p(x)$  حداقل دارای یک ریشه است، یعنی حداقل یک عدد مختلط مانند  $\xi$  وجود دارد که  $p(\xi) = 0$ .

این قضیهٔ نسبتاً پر محتوا را با قضیهٔ سراسر «یک بسجمله‌ای از درجهٔ  $n$  حداکثر  $n$  ریشه (با احتساب ریشه‌های مکرر) دارد»، که در فصل دوم اثبات و مثلاً در مبحث درونیابی بسجمله‌ای از آن استفاده می‌کنیم، نباید اشتباه گرفت.

## تمرین

**۷۰۱-۱** در قضیهٔ مقدار میانگین برای انتگرالها، قضیهٔ ۳۰۱، می‌گیریم  $f(x) = e^x$  و  $g(x) = x$  و  $[a, b] = [0, 1]$ . نقطهٔ  $\xi$  را که به وسیلهٔ این قضیه مشخص شده است پیدا و تحقیق کنید که این نقطه در بازهٔ  $(0, 1)$  قرار دارد.

**۷۰۱-۲** در قضیهٔ مقدار میانگین برای مشتق، قضیهٔ ۶۰۱، می‌گیریم  $f(x) = x^2$ . نقطهٔ  $\xi$  را

کسه به وسیله این قضیه مشخص شده است پیدا و تحقیق کنید کسه این نقطه در بازه  $(a, b)$  قرار دارد.

۳-۷۰۱ در بسط (۳۶۰۱) برای  $e^x$ ،  $n$  را طوری پیدا کنید که به ازای جمیع مقادیر  $x$  روی  $[0, 1]$  حاصل جمع توانی نتیجه، تقریبی به دست دهد که تا پنج رقم با معنی دقیق باشد.

۴-۷۰۱ با استفاده از فرمول تیلر (۳۲۰۱) بسط سری توانی  $\sin(\pi x/2)$  را پیرامون  $c=0$  پیدا کنید. عبارتی برای باقیمانده تعیین نمایید و به وسیله آن تعداد جملات ضروری برای تضمین دقت عملی تا شش رقم با معنی در محاسبه  $\sin(\pi x/2)$  را، به ازای جمیع مقادیر  $x$  متعلق به بازه  $[-1, 1]$ ، برآورد کنید.

۵-۷۰۱ باقیمانده  $R_4(x, y)$  در مثال (۳۸۰۱) را پیدا و حداکثر مقدار آن را در ناحیه  $D$  که به وسیله

$$\left[ 0 \leq x \leq \frac{\pi}{2}, \quad 0 \leq y \leq \frac{\pi}{2} \right]$$

تعریف می شود، تعیین کنید.

۶-۷۰۱ ثابت کنید که جمله باقیمانده در (۳۵۰۱) را می توان به صورت زیر نیز نوشت

$$\frac{h^{n+1}}{(n+1)!} f^{(n+1)}(c) + o(h^{n+1})$$

[اگر  $f^{(n+1)}(x)$  در  $x=c$  پیوسته باشد].

۷-۷۰۱ صحت حکم تمرین (۶-۷۰۱) را از راه محاسبه

$$R_4(h) = e^h - \left( 1 + h + \frac{h^2}{2} \right), \quad \frac{h^3}{3!} f^{(3)}(0) = \frac{h^3}{3!}$$

به ازای  $f(x) = e^x$  و  $c=0$  و به ازای مقادیر مختلف  $h$ ، مثلاً  $h = 2^{-k}$ ،  $k = 1, 2, \dots$  و مقایسه  $R_4(h)$  با  $(h^3/3!) f^{(3)}(0)$  نشان دهید.

۸-۷۰۱ قضیه ۹۰۱ را با استفاده از قضیه های ۷۰۱ و ۸۰۱ ثابت کنید.

۹-۷۰۱ فرمول اولر یعنی

$$e^{i\theta} = \cos \theta + i \sin \theta$$

را ( $i = \sqrt{-1}$ ) به وسیله مقایسه سری نمایی  $e^x$ ، به ازای  $x = i\theta$ ، با حاصل جمع سری توانی  $\cos \theta$  و  $i \sin \theta$  برای سری توانی اثبات کنید.



## درونیابی به وسیلهٔ بسجمله‌ایها

بسجمله‌ایها تقریباً در همهٔ زمینه‌های آنالیز عددی به‌عنوان وسیلهٔ اصلی تقریب زدن به‌کار برده می‌شوند. بسجمله‌ایها در حل معادلات و تقریب زدن توابع، تقریب زدن انتگرالها و مشتقها، تقریب زدن جوابهای انتگرال و معادلات دیفرانسیل و غیره به‌کار برده می‌شوند. مقبولیت بسجمله‌ایها به‌علت ساخت سادهٔ آنهاست که امر ساختن تقریبهای مؤثر و استفاده از آنها را آسان می‌سازد.

بدین دلیل، نمایش و ارزیابی بسجمله‌ایها یک مبحث اساسی در آنالیز عددی است. ما این موضوع را در این فصل، ضمن گفتگو از درونیابی بسجمله‌ای، که ساده‌ترین و مطمئناً گسترده‌ترین تکنیکی است که برای به‌دست آوردن تقریبهای بسجمله‌ای به‌کار برده می‌شود، مورد بحث قرار خواهیم داد. برای به‌دست آوردن تقریبهای خوب به کمک بسجمله‌ایها و توابع تقریبی دیگر، روشهای پیشرفته‌تری در فصل ۶ داده شده‌اند. اما در همان‌جا نشان داده خواهد شد که حتی بهترین تقریب، نتایجی چندانی بهتر از یک طرح مناسب درونیابی بسجمله‌ای عاید ما نمی‌سازد.

پایهٔ عمل ما برای بسجمله‌ای درونیاب، تفاضلات منقسم است. این خود به‌ما امکان می‌دهد که با درونیابی بوسانی (یا هرمیتی<sup>۱</sup>) به صورت یک حالت حدی خاص درونیابی بسجمله‌ای در نقاط متمایز، عمل کنیم.

## ۱۰۲ صورتهای بسجمله‌ای

در این بخش، اشاره می‌کنیم که راه معمولی برای توصیف یک بسجمله‌ای ممکن است همیشه

بهترین راه برای محاسبات نباشد، و لذا ما راه دیگر، به‌ویژه صورت نیوتنی را پیشنهاد می‌کنیم. همچنین نشان می‌دهیم که یک بسجمله‌ای را که به‌صورت نیوتنی داده شده چگونه باید ارزیابی کنیم. سرانجام، در مقدمهٔ درونیابی بسجمله‌ای، چگونگی شمارش ریشه‌های یک بسجمله‌ای را نیز مورد بحث قرار می‌دهیم.

بنابر تعریف، یک بسجمله‌ای  $p(x)$  از درجهٔ  $n$  بزرگتر از  $n$ ، تابعی است به‌صورت

$$p(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n \quad (1.2)$$

با ضرایب معین  $a_0, a_1, \dots, a_n$ . و در مواردی که ضریب جمله پیشرو  $a_n$ ، غیر صفر باشد، بسجمله‌ای دقیقاً از درجهٔ  $n$  خواهد بود.

در مباحث ریاضی، راه استاندارد مشخص کردن یک بسجمله‌ای، داشتن صورت توانی<sup>۲</sup> (۱.۲) آن است. برای مشتکگیری و انتگرالگیری یک بسجمله‌ای این صورت بسیار مناسب است. اما، در زمینه‌های خاص مختلف صورتهای دیگر مناسبترند.

□ مثال ۱.۲: صورت توانی ممکن است به ازدست دادن ارقام با معنی منجر شود. اگر صورت توانی یک خط مستقیم  $p(x)$  را که برای آن داریم،  $p(6001) = -2/3$  و  $p(6000) = 1/3$  را در نظر بگیریم، آنگاه در حساب با ممیز شناور با پنج رقم اعشاری خواهیم داشت  $p(x) = 6000r3 - x$ . اگر مقدار این خط راست را در همان حساب ارزیابی کنیم، خواهیم دید که  $p(6000) = 0r3$  و  $p(6001) = -0r7$ ، که فقط اولین رقم مقادیر تابع داده شده را به دست می‌دهد و در نتیجه موجب از دست رفتن چهار رقم اعشاری می‌شود. □

یک راه چاره برای جلوگیری از این گونه از دست دادن ارقام با معنی، استفاده از صورت توانی منتقله<sup>۳</sup> به شرح زیر است:

$$p(x) = a_0 + a_1(x-c) + a_2(x-c)^2 + \dots + a_n(x-c)^n \quad (2.2)$$

اگر مرکز  $c$  را مساوی ۶۰۰۰ انتخاب کنیم، آنگاه در مثال فوق خواهیم داشت  $p(x) = 0r33333 - (x - 6000r0)$ ، و با محاسبهٔ مقدار تابع در حساب با ممیز شناور با پنج رقم اعشاری، خواهیم داشت

$$p(6000) = 0r33333 \quad \text{و} \quad p(6001) = -0r66667$$

یعنی در این صورت مقادیر تابع تا پنج رقم صحیح اند.

هرگاه یک بسجمله‌ای در یک بازهٔ  $[a, b]$  مورد نظر باشد، بهتر است از صورت توانی منتقله با مرکز  $c \in [a, b]$  استفاده کنیم. یک راه چارهٔ عاقلانه‌تر برای جلوگیری از به‌در رفتن ارقام با معنی (یا بدشرطی<sup>۴</sup>) بسط به بسجمله‌ایهای چیبیش<sup>۵</sup> یا بسجمله‌ایهای متعامد

- 
1. leading
  2. Power form
  3. Shifted power form
  4. illconditioning
  5. Chebyshev

دیگر است. (به بخش ۳.۶ نگاه کنید).

اگر  $p(x)$  به وسیلهٔ (۲.۲) داده شده باشد، ضرایب صورت توانی منتقله (۲.۲) مقادیر مشتق زیر را به دست می‌دهند، یعنی

$$a_i = \frac{p^{(i)}(c)}{i!} \quad i = 0, \dots, n$$

در حقیقت صورت توانی منتقله همان بسط تیلور  $p(x)$  پیرامون مرکز  $c$  است. يك تعمیم دیگر صورت توانی منتقله، صورت نیوتنی آن است:

$$\begin{aligned} p(x) = & a_0 + a_1(x-c_1) + a_2(x-c_1)(x-c_2) \\ & + a_3(x-c_1)(x-c_2)(x-c_3) + \dots \\ & + a_n(x-c_1)(x-c_2)\dots(x-c_n) \end{aligned} \quad (3.2)$$

این صورت نقش عمده‌ای در ساختن يك بسجمله‌ای درونیاب ایفا می‌کند. اگر مراکز  $c_1, \dots, c_n$  همگی مساوی  $c$  باشند، این صورت نیز به همان صورت توانی بدل می‌شود، و اگر مراکز  $c_1, \dots, c_n$  همگی مساوی صفر باشند، صورت توانی به دست می‌آید. لذا بحث زیر دربارهٔ ارزیابی صورت نیوتنی، مستقماً برای این صورت‌های ساده نیز به کار خواهد رفت.

محاسبهٔ جداگانهٔ هر يك از  $n+1$  جملهٔ  $(3.2)$ ، و بعد جمع کردن آنها با هم، بیهوده است. این کار متضمن  $\frac{n+1}{2}n$  عمل جمع و  $\frac{n+1}{2}n$  عمل ضرب است. در عوض ملاحظه می‌شود که عامل  $(x-c_1)$  در همهٔ جملات هست، بجز در جمله اول، یعنی

$$\begin{aligned} p(x) = & a_0 + (x-c_1) \{ a_1 + a_2(x-c_2) + a_3(x-c_2)(x-c_3) \\ & + \dots + a_n(x-c_2)(x-c_3)\dots(x-c_n) \} \end{aligned}$$

و به همین ترتیب، تمامی جملات دیگر داخل ابرو، به استثنای اولین جملهٔ آن، شامل عامل  $(x-c_2)$  می‌باشند، یعنی

$$\begin{aligned} p(x) = & a_0 + (x-c_1) \{ a_1 + (x-c_2) [ a_2 + a_3(x-c_3) \\ & + \dots + a_n(x-c_3)\dots(x-c_n) ] \} \end{aligned}$$

با ادامهٔ عملیات فوق،  $p(x)$  را به صورت تودرتو<sup>۳</sup> به دست خواهیم آورد:

$$\begin{aligned} p(x) = & a_0 + (x-c_1) \{ a_1 + (x-c_2) [ a_2 + (x-c_3) \{ a_3 + \dots \\ & + (x-c_{n-1}) (a_{n-1} + (x-c_n) a_n) \dots \} ] \} \end{aligned}$$

که محاسبهٔ آن به ازای هر مقدار خاص  $\alpha$  تنها شامل  $2n$  عمل جمع و  $n$  عمل ضرب است. برای مثال اگر  $p(x)$  به صورت

$$p(x) = 1 + 2(x-1) + 3(x-1)(x-2) + 4(x-1)(x-2)(x-3)$$

باشد و بخواهیم  $p(4)$  را محاسبه کنیم، محاسبات به صورت زیر انجام می گیرند:

$$\begin{aligned} p(4) &= 1 + (4-1) \{ 2 + (4-2) [ 3 + (4-3) \langle 4 \rangle ] \} \\ &= 1 + (4-1) \{ 2 + (4-2) [ 7 ] \} \\ &= 1 + (4-1) \{ 16 \} \\ &= 49 \end{aligned}$$

این شیوهٔ عمل در الگوریتم زیر، به شکل رسمی بیان شده است.

**الگوریتم ۱.۲:** ضرب تودرتو برای صورت نیوتنی  $n+1$  ضریب  $a_0, \dots, a_n$  برای صورت نیوتنی (۳.۲) از بسجمله ای  $p(x)$ ، همراه با مراکز  $c_1, \dots, c_n$  داده شده اند. عدد  $z$  نیز مفروض است

$$\begin{aligned} a'_n &:= a_n \\ &\left\{ \begin{array}{l} \text{به ازای } i = n-1, n-2, \dots, 0 \text{ می شود:} \\ a'_i := a_i + (z - c_{i+1}) a'_{i+1} \end{array} \right. \end{aligned}$$

پس،  $a'_0 = p(z)$ . بعلاوه، کمیتهای کمکی  $a'_1, \dots, a'_n$  مستقلاً مورد علاقه ما هستند. زیرا، داریم

$$\begin{aligned} p(x) &= a'_0 + a'_1(x-z) + a'_2(x-z)(x-c_1) \\ &\quad + a'_3(x-z)(x-c_1)(x-c_2) \\ &\quad + \dots + a'_n(x-z)(x-c_1)(x-c_2) \dots (x-c_{n-1}) \quad (4.2) \end{aligned}$$

یعنی  $a'_1, \dots, a'_n$ ، ضرایب نیوتنی برای  $p(x)$  با مراکز  $c_1, c_2, \dots, c_{n-1}$  هستند.

حال به اثبات (۴.۲) می پردازیم. بنا بر الگوریتم فوق داریم

$$\begin{aligned} a_n &= a'_n \\ a_i &= a'_i + a'_{i+1}(c_{i+1} - z) \quad i = n-1, n-2, \dots, 0 \end{aligned}$$

از قرارداد این روابط در (۳.۲) به دست می آوریم

$$\begin{aligned}
 p(x) &= a_0 + a_1(x-c_1) + a_2(x-c_1)(x-c_2) \\
 &\quad + \dots + a_n(x-c_1)\dots(x-c_n) \\
 &= a'_0 + a'_1(c_1 - z) \\
 &\quad + [a'_1 + a'_2(c_2 - z)](x - c_1) \\
 &\quad + [a'_2 + a'_3(c_3 - z)](x - c_1)(x - c_2) \\
 &\quad \vdots \\
 &\quad + [a'_{n-1} + a'_n(c_n - z)](x - c_1)\dots(x - c_{n-1}) \\
 &\quad + a'_n(x - c_1)\dots(x - c_{n-1})(x - c_n) \\
 &= a'_0 + a'_1(x - z) + a'_2(x - z)(x - c_1) \\
 &\quad + \dots + a'_n(x - z)(x - c_1)\dots(x - c_{n-1})
 \end{aligned}$$

که مؤید صحت رابطه (۴.۲) است.

علاوه بر تعیین مقدار بسجمله ای (۳.۲) در هر نقطه خاص  $z$ ، که از نظر اقتصادی بسیار مؤثر است، الگوریتم ضرب تودرتو در عبور از يك صورت نیوتنی به صورت دیگر، بسیار مفید است. برای مثال، فرض کنید می خواهیم بسجمله ای

$$p_2(x) = 195709 + 000006(x-1) + 000012(x-1)(x-4)$$

را بر حسب توانهای  $x$ ، یعنی به صورت نیوتنی که کلیه مراکز آن صفرند در آوریم. در این حال، اگر از الگوریتم ۱.۲ با  $z=0$  (و  $n=2$ ) استفاده کنیم، خواهیم داشت

$$a'_2 = a_2 = 000012$$

$$a'_1 = a_1 + (z - c_2)a'_2 = 000006 + (0 - 4)(000012) = 0000012$$

$$a'_0 = a_0 + (z - c_1)a'_1 = 195709 + (0 - 1)(0000012) = 1957078$$

بنابراین

$$p_2(x) = 1957078 + 000012(x-0) + 000012(x-0)(x-1)$$

با به کار گیری الگوریتم ۱.۲ برای این بسجمله ای، دوباره به ازای  $z=0$  داریم

$$a'_2 = a_2 = 000012$$

$$a'_1 = a_1 + (z - c_2)a'_2 = 0000012 + (0 - 1)(000012) = 000$$

$$a'_0 = a_0 + (z - c_1)a'_1 = 1957078 + (0 - 0)(000) = 1957078$$

بنا بر این

$$\begin{aligned} p_1(x) &= 1957078 + 0.0(x-0) + 0.000012(x-0)(x-0) \\ &= 1957078 + 0.000012x^2 \end{aligned}$$

در این مثال ساده، می‌توانیم صحت نتیجه را سریعاً از ضرب جمله‌های رابطه اولیه، تحقیق کنیم.

$$\begin{aligned} p_2(x) &= 195709 + 0.00006(x-1) + 0.000012(x^2 - 5x + 4) \\ &= [195709 - 0.00006 + (0.000012)(4)] \\ &\quad + [0.00006 + (0.000012)(-5)]x + 0.000012x^2 \\ &= 1957078 + 0.000012x^2 \end{aligned}$$

استفاده مکرر از الگوریتم ضرب تودر تو در محاسبه مشتقهای یک بسجمله‌ای که به صورت نیوتنی داده شده باشد (به تمرین  $1.2-2$  تا  $1.2-5$  نگاه کنید) نیز مفید است. الگوریتم فوق در اثبات مطلب مسلم پایه‌ای زیر نیز مؤثر است.

لم  $1.2$  اگر  $z_1, \dots, z_k$  ریشه‌های متمایز بسجمله‌ای  $p(x)$  باشند، آنگاه

$$p(x) = (x - z_1)(x - z_2) \dots (x - z_k)r(x)$$

که در آن  $r(x)$  يك بسجمله‌ای است.

برای اثبات این لم،  $p(x)$  را به صورت توانی  $(1.2)$ ، یعنی به صورت نیوتنی که همه مراکز آن صفرند می‌نویسیم و سپس یکبار الگوریتم  $(1.2)$  را به کار می‌بریم و داریم

$$p(x) = p(z) + (x - z)q(x)$$

[زیرا  $a'_0 = p(z)$  که  $a'_0 + a'_1x + \dots + a'_n x^{n-1}$  که  $q(x)$  يك بسجمله‌ای با درجه کمتر از  $n$  است. در واقع،  $p(x)$  را بر بسجمله‌ای خطی  $(x - z)$  تقسیم کرده‌ایم و  $q(x)$  بسجمله‌ای خارج قسمت و عدد  $p(z)$  باقیمانده آن شده‌است. اکنون به طور اخص  $z = z_1$  را انتخاب می‌کنیم و آنگاه طبق فرض مسئله  $p(z_1) = 0$ ، یعنی

$$p(x) = (x - z_1)q(x)$$

بدین ترتیب به ازای  $k=1$  اثبات تمام شده‌است. بعد به ازای  $k > 1$ ، نتیجه می‌شود که  $z_1, \dots, z_k$  الزاماً ریشه‌های  $q(x)$  هستند، زیرا بنا بر فرض  $p(x)$  در این نقاط برابر با صفر می‌شود، در حالی که  $x - z_1$  در این نقاط برابر با صفر نمی‌شود. بنا بر این برای اتمام برهان می‌توان از استقراء نسبت به  $k$ ، تعداد ریشه‌ها، استفاده کرد.

فرض: اگر  $p(x)$  و  $q(x)$  دو بسجمله‌ای از درجه نایز رگتر از  $k$  باشند که در  $k+1$  نقطه

متمايز  $z_0, \dots, z_k$  برهم منطبق باشند، آنگاه  $p(x) = q(x)$ .

زیرا، تفاضل آنها  $d(x) = p(x) - q(x)$  يك بسجمله‌ای است از درجهٔ نایز بزرگتر از  $k$ ، و بنا بر لم ۱۰۲ می‌توان آن را به صورت

$$d(x) = (x - z_0) \dots (x - z_k) r(x)$$

که در آن  $r(x)$  يك بسجمله‌ای است، نوشت. فرض کنید که به‌ازای ضرایبی مانند  $c_0, \dots, c_m$  با شرط  $c_m \neq 0$ ، داشته باشیم  $r(x) = c_0 + \dots + c_m x^m$ . در این صورت رابطهٔ

$$k \geq d = k + 1 + m$$

برقرار می‌گردد که بی‌معنی است. بنا بر این  $r(x) = 0$ ، و لذا  $p(x) = q(x)$ . این فرع به‌این سؤال که «چند بسجمله‌ای از درجهٔ نایز بزرگتر از  $k$  وجود دارد که در  $k+1$  نقطهٔ بخصوص مقادیر خاصی داشته باشند؟» جواب «حداکثر يك» را می‌دهد. این بررسی‌های مربوط به‌ریشه‌های بسجمله‌ای را می‌توان در طی بحثی دربارهٔ مفهوم تکرار اصلاح نمود. و این مطلب در بحث‌های بعدی‌ما دربارهٔ درون‌یابی بوسانی بسیار مهم خواهد بود. گوییم  $z$  را يك ریشهٔ (دقیقاً) با تکرار  $\alpha$  یا از مرتبهٔ  $\alpha$  تا بس  $f(x)$  گوییم به‌شرط آنکه داشته باشیم:

$$f(z) = f'(z) = \dots = f^{(\alpha-1)}(z) = 0 \neq f^{(\alpha)}(z)$$

□ مثال: برای مثال، بسجمله‌ای

$$(x-z)^{\alpha}$$

دارای يك ریشهٔ  $z$  از مرتبهٔ  $\alpha$  است. قابل قبول است که يك چنین ریشه‌ای،  $\alpha$  مرتبه به‌حساب آید، زیرا که می‌تواند به‌عنوان حد بسجمله‌ای

$$(x-z_1) \dots (x-z_j)$$

بیان شود که دارای  $\alpha$  ریشهٔ متمايز و یا ساده است که همگی در نقطهٔ  $z$  برهم منطبق و یا باهم متحد می‌شوند.

مثال دیگر آنکه به‌ازای  $0 < A < 1$ ، تابع  $\sin x - Ax$  در بازهٔ  $-\pi < x < \pi$  دارای سه ریشهٔ (ساده) است که وقتی  $A \rightarrow 1$ ، ریشه‌ها به‌سمت عدد صفر همگرا می‌شوند. همین‌طور، تابع (حدی)  $\sin x - x$  در صفر دارای يك ریشهٔ سه‌گانه است. □

با این برداشت از تکرار يك ریشه، می‌توان به‌لم ۱۰۲ صورت قوی‌تری بخشید.

لم ۲.۲ اگر  $z_1, \dots, z_k$  يك دنباله از ریشه‌های بسجمله‌ای  $p(x)$  با احتساب تکرر آنها باشد، آنگاه

$$p(x) = (x - z_1)(x - z_2) \dots (x - z_k)r(x)$$

که در آن  $r(x)$  يك بسجمله‌ای است.

برای اثبات لم ۲.۲ به تمرین ۱-۰۲ نگاه کنید. باید توجه داشت در حالتی که  $z$  يك ریشهٔ  $p(x)$  از مرتبهٔ  $z$  ممتد در دنبالهٔ  $z_1, \dots, z_k$  به تعداد  $z$  مرتبه ظاهر شود. از لم ۲.۲ و بحث قبلی نتیجه می‌گیریم

فرض: اگر  $p(x)$  و  $q(x)$  دو بسجمله‌ای از درجهٔ نابزرگتر از  $k$ ، بسا  $k+1$  نقطهٔ  $z_1, \dots, z_k$  منطبق برهم باشند، یعنی تفاضل آنها،  $r(x) = p(x) - q(x)$ ،  $k+1$  ریشهٔ  $z_1, \dots, z_k$  (با احتساب تکرر آنها) داشته باشد، آنگاه  $p(x) \equiv q(x)$ .

### تمرین

۱-۰۲ مقدار بسجمله‌ای درجهٔ سوم

$$p(x) = (x - 99\pi)(x - 100\pi)(x - 101\pi)$$

را در  $x = 314r15$  حساب کنید، سپس با استفاده از ضرب تودرتو برای به دست آوردن  $p(x)$  به شکل توانی مقدار بسجمله‌ای حاصل را به شکل توانی در  $x = 314r15$  به دست آورید و نتایج را با هم مقایسه کنید!

۲-۱.۲ گیریم  $p(x) = a_0 + (x - c_1)(a_1 + \dots + (x - c_n)(a_n) \dots)$  يك بسجمله‌ای به صورت نیوتنی بسا باشد. ثابت کنید که اگر  $c_1 = c_2 = \dots = c_{r+1}$ ، آنگاه به ازای  $z = 0, \dots, r$ ، رابطهٔ  $p^{(j)}(c_1) = j! a_j$  را برقرار است. [داهنمایی: در شرایط فوق،  $p(x)$  را می‌توانید به صورت

$$p(x) = \sum_{j=0}^r a_j (x - c_1)^j + (x - c_1)^{r+1} q(x)$$

که در آن  $q(x)$  يك بسجمله‌ای است بنویسید. حال از  $p(x)$  مشتق بگیرید.

۳-۱.۲ مشتق اول بسجمله‌ای

$$p(x) = 3 - (x - 1)(4 - (x + 1)\{5 - x[6 - (x + 2)]\})$$

را در  $x = 2$  پیدا کنید. [داهنمایی: برای به دست آوردن صورت نیوتنی بسا مراکز  $2, 1, -1$ ، الگوریتم ۱.۲ را دوبار به کار گیرید و سپس طبق تمرین ۲-۱.۲ عمل کنید.]



۱۰۲-۴ مشتق دوم بسجمله‌ای  $p(x)$  در تمرین ۱۰۲-۳ را نیز در  $x=2$  پیدا کنید.

۱۰۲-۵ بسط تیلر بسجمله‌ای  $p(x)$  در تمرین ۱۰۲-۳ را پیرامون  $c=3$  پیدا کنید. [داهنمایی: بسط تیلر برای يك بسجمله‌ای پیرامون نقطه  $c$ ، همان صورت نیوتنی آن بسجمله‌ای است با مراکز  $c, c, c, \dots, c$ ].

۱۰۲-۶ لم ۲۰۲ را ثابت کنید. [داهنمایی: به موجب الگوریتم ۱۰۲، داریم  $p(x) = (x - z_1)q(x)$ . اما برای اتمام برهان به استقراء نسبت به عدد  $k$ ، تعداد ریشه‌ها در دنباله داده شده، باید ثابت نماییم که  $z_1, z_2, \dots, z_k$  الزاماً دنباله‌ای از ریشه‌های  $q(x)$  (با احتساب تکرار) است. بدین منظور فرض کنید که عدد  $z$  دقیقاً  $j$  مرتبه در دنباله  $z_1, z_2, \dots, z_k$  ظاهر شود و موارد  $z = z_1$  و  $z \neq z_1$  متمایز باشند. و نیز از رابطه

$$p^{(j)}(x) = (x - z_1)q^{(j)}(x) + jq^{(j-1)}(x)$$

استفاده کنید.]

۱۰۲-۷ ثابت کنید که، به زبان فرع لم ۲۰۲، بسجمله‌ای تیلر  $\sum_{i=0}^{j-1} f^{(i)}(a)(x-a)^i/i!$  در نقطه  $x=a$ ،  $j$  مرتبه  $1$  بر تابع  $f(x)$  منطبق است (یعنی  $a$  يك ریشه مرتبه  $j$  از تفاضل آنهاست).

۱۰۲-۸ فرض کنید يك تابع  $F(X)$  که قرار است مقدار يك بسجمله‌ای خاص از درجه کمتر از  $r$  را در نقطه  $X$  به ما بدهد، داده شده است. بعد، فرض کنید که از راه تجسس در یافتید که این تابع واقعاً مقدار يك بسجمله‌ای از درجه کمتر از  $r$  را به شما می‌دهد (مثلاً می‌بینید که فقط حاصلجمع (یا تفاضل) جملاتی پیدا شده‌اند که از ضرب اعداد در توانهای مختلف  $X$  کمتر از  $r$  بار به دست آمده‌اند). مقادیر چند تابع را قبلاً باید امتحان کنید تا بتوانید مطمئن شوید که این تابع آنچه را که قرار بوده انجام دهد واقعاً انجام می‌دهد یا نه؟ (البته با فرض نبودن خطای گرد کردن.)

۱۰۲-۹ برای هر يك از سریهای توانی زیر از ضرب تودرتو استفاده کنید تا راه مؤثری برای ارزیابی آنها پیدا کنید. (البته مجبورید فرض کنید که این سریها فقط باید به ازای مقداری از  $N$  که از قبل تعیین شده، روی دامنه  $N \leq n$  جمع شوند.)

$$e^x = \sum_n x^n/n! \quad (\text{الف})$$

$$\ln x = 2 \sum_{\text{فرد } n} \frac{1}{n} \left( \frac{x-1}{x+1} \right)^n \quad (\text{ب})$$

$$\sin^{-1} x = \sum_{\text{فرد } n} \frac{1 \cdot 3 \cdots (n-2)}{2 \cdot 4 \cdots (n-1)} x^n/n \quad (\text{پ})$$

## ۲.۲ وجود و یکتایی<sup>۱</sup> بسجمله‌ای درونیاب

گیریم  $x_0, \dots, x_n, n+1$  نقطهٔ متمایز روی محور حقیقی  $x$ ها باشند و  $f(x)$  تابعی با مقدار حقیقی باشد که در بازهٔ  $I = [a, b]$ ، که این نقاط را دربردارد، تعریف شده است. می‌خواهیم بسجمله‌ای  $p(x)$  از درجهٔ نایبتر از  $n$  را طوری به دست آوریم که تابع  $f(x)$  را در نقاط  $x_0, \dots, x_n$  درونیابی کند، یعنی شرایط  $p(x_i) = f(x_i)$ ،  $i = 0, 1, \dots, n$  برقرار باشند.

به طوری که مشاهده خواهیم کرد راههای زیادی برای نوشتن چنین بسجمله‌ای وجود دارد، و بنا بر این لازم است در آغاز به خواننده یادآوری کنیم که به موجب فرع ۱.۲ حداکثر یک بسجمله‌ای از درجهٔ نایبتر از  $n$  وجود دارد که  $f(x)$  را در  $n+1$  نقطهٔ متمایز  $x_0, \dots, x_n$  درونیابی می‌کند.

سپس نشان می‌دهیم که حداقل یک بسجمله‌ای از درجهٔ نایبتر از  $n$  وجود دارد که  $f(x)$  را در  $n+1$  نقطهٔ متمایز  $x_0, \dots, x_n$  درونیابی می‌کند. برای این امر، باز از شکل دیگری از بسجمله‌ای به نام صورت لاگرانژی<sup>۲</sup> استفاده کنیم

$$p(x) = a_0 l_0(x) + a_1 l_1(x) + \dots + a_n l_n(x) \quad (5.2)$$

که در آن بسجمله‌ایهای

$$l_k(x) = \prod_{\substack{i=0 \\ i \neq k}}^n \frac{x - x_i}{x_k - x_i} \quad k = 0, \dots, n \quad (6.2)$$

بسجمله‌ای لاگرانژی برای نقاط  $x_0, \dots, x_n$  نامیده می‌شوند. تابع  $l_k(x)$  حاصلضرب  $n$  عامل خطی است، بنا بر این یک بسجمله‌ای است دقیقاً از درجهٔ  $n$ . لذا رابطهٔ (۵.۲) واقعاً معرف یک بسجمله‌ای است از درجهٔ نایبتر از  $n$ . وانگهی  $l_k(x)$ ، به ازای جمیع مقادیر  $i \neq k$  در  $x_i$  برابر صفر و در  $x_k$  برابر با یک است. یعنی

$$l_k(x_i) = \begin{cases} 1 & i = k \\ 0 & i \neq k \end{cases} \quad i = 0, \dots, n$$

این، نشان می‌دهد که

$$p(x_i) = \sum_{k=0}^n a_k l_k(x_i) = a_i \quad i = 0, \dots, n$$

یعنی ضرایب  $a_0, \dots, a_n$  در صورت لاگرانژی، مقادیر بسجمله‌ای  $p(x)$  در نقاط  $x_0, \dots, x_n$  هستند. در نتیجه برای یک تابع دلخواه  $f(x)$ ، بسجمله‌ای

$$p(x) = \sum_{k=0}^n f(x_k) l_k(x) \quad (۷.۲)$$

از درجهٔ نایبتر از  $n$  است که  $f(x)$  را در نقاط  $x_0, \dots, x_n$  درونیابی می‌کند و این مطلب قضیهٔ زیر را ثابت می‌کند.

**قضیه ۱۰۲** يك تابع حقیقی مقدار  $f(x)$  و  $n+1$  نقطهٔ متمایز  $x_0, \dots, x_n$  داده شده‌اند. دقیقاً يك بسجمله‌ای با درجهٔ نایبتر از  $n$  وجود دارد که  $f(x)$  را در نقاط  $x_0, \dots, x_n$  درونیابی می‌کند.

معادلهٔ (۷.۲) فرمول لاجرانژ برای بسجمله‌ای درونیاب نامیده می‌شود. به‌عنوان يك کاربرد ساده، مورد  $n=1$  را بررسی می‌کنیم. یعنی تابع  $f(x)$  و دو نقطهٔ متمایز  $x_0, x_1$  داده شده‌اند. آنگاه داریم

$$l_0(x) = \frac{x-x_1}{x_0-x_1} \quad l_1(x) = \frac{x-x_0}{x_1-x_0}$$

و

$$\begin{aligned} p(x) &= f(x_0)l_0(x) + f(x_1)l_1(x) = f(x_0) \frac{x-x_1}{x_0-x_1} + f(x_1) \frac{x-x_0}{x_1-x_0} \\ &= \frac{f(x_0)(x-x_1) - f(x_1)(x-x_0)}{x_0-x_1} \\ &= f(x_0) + \frac{f(x_1) - f(x_0)}{x_1-x_0} (x-x_0) \end{aligned}$$

این همان مورد آشنای درونیابی خطی<sup>۱</sup> است که به یکی از صورتهای هم ارزش نوشته شده‌است.

□ **مثال ۲.۲:** انتگرالی مربوط به انتگرال بیضوی کامل به شرح زیر تعریف می‌شود

$$K(k) = \int_0^{\pi/2} \frac{dx}{[1 - (\sin k)^2 \sin^2 x]^{1/2}} \quad (۸.۲)$$

از جدول مقادیر این انتگرالها به ازای مقادیر مختلف  $k$  که بر حسب درجه اندازه گیری شده‌اند، نتیجه می‌گیریم

$$K(1) = ۱.۵۷۰۹$$

$$K(۴) = ۱۷۵۷۲۷$$

$$K(۶) = ۱۷۵۷۵۱$$

با استفاده از بسجمله‌ای درونیاب درجهٔ دوم، مقدار  $K(۳۷۵)$  را محاسبه کنیم.  
داریم:

$$I_0(۳۷۵) = \frac{(۳۷۵-۴)(۳۷۵-۶)}{(۱-۴)(۱-۶)} = \frac{۱۷۲۵}{۱۵} = ۰۰۸۳۳۳$$

$$I_1(۳۷۵) = \frac{(۳۷۵-۱)(۳۷۵-۶)}{(۴-۱)(۴-۶)} = \frac{-۶۷۲۵}{-۶} = ۱۱۰۴۱۶۷$$

$$I_2(۳۷۵) = \frac{(۳۷۵-۱)(۳۷۵-۴)}{(۶-۱)(۶-۴)} = \frac{-۱۷۲۵}{۱۰} = -۰۰۱۷۲۵۰۰$$

لذا

$$K(۳۷۵) \approx (۱۷۵۷۰۹)(۰۰۸۳۳۳) + (۱۷۵۷۲۷)(۱۱۰۴۱۶۷) \\ + (۱۷۵۷۵۱)(-۰۰۱۷۲۵۰۰) = ۱۷۵۷۲۲۵$$

□

رقم آخر این تقریب توأم با خطاست.

صورت لاگرانژی (۷.۲) برای بسجمله‌ای درونیاب، اثبات وجود يك بسجمله‌ای درونیاب را آسان می‌سازد. اما محاسبهٔ مقدار آن در يك نقطه مستلزم حداقل  $۲(n+۱)$  عمل ضرب یا تقسیم و  $(۲n+۱)$  عمل جمع و تفریق می‌باشد، این بعد از آن است که مخرج بسجمله‌ای لاگرانژی یکبار برای همیشه محاسبه و بر مقادیر تابع متناظر تقسیم شده‌است. این تعداد عملیات باید با  $n$  عمل ضرب و  $n$  عمل جمع که برای محاسبهٔ مقدار يك بسجمله‌ای درجهٔ  $n$ ام به‌صورت توانی، به‌وسیلهٔ ضرب تودرتو لازم است، مقایسه شود (به الگوریتم ۱.۲ نگاه کنید).

يك ایراد جدیتر به‌روش لاگرانژی ایراد زیر است: در عمل اغلب معلوم نیست که چه تعداد نقاط درونیابی باید به‌کار گرفته شود. بنا بر این به‌کمک  $p_j(x)$  که معرف يك بسجمله‌ای از درجهٔ نایبتر از  $n$  است و  $f(x)$  را در نقاط  $x_0, \dots, x_n$  درونیابی می‌کنند، مقادیر  $p_0(x), p_1(x), \dots, p_{n-1}(x)$  را محاسبه می‌کنیم و تعداد نقاط درونیابی و در نتیجه درجهٔ بسجمله‌ای درونیاب را افزایش می‌دهیم، به این امید که  $p_k(x)$  تقریب موردقبولی از  $f(x)$  را به‌دست دهد. در چنین موردی به‌کار گرفتن صورت لاگرانژی بیفایده به‌نظر می‌رسد. زیرا که در محاسبهٔ  $p_k(x)$  از این واقعیت که  $p_{k-1}(x)$  در دست است هیچگونه استفاده‌ای

نمی‌شود. بدین دلیل و به دلایل دیگر صورت نیوتنی برای بسجمله‌ای درونیاب بسیار مناسبتر است.

در حقیقت، با استفاده از نقاط درونیابی  $x_0, \dots, x_{n-1}$  به عنوان مراکز، بسجمله‌ای درونیاب  $p_n(x)$  را به صورت نیوتنی، یعنی به شرح زیر می‌نویسیم

$$p_n(x) = A_0 + A_1(x-x_0) + A_2(x-x_0)(x-x_1) \\ + \dots + A_n(x-x_0)\dots(x-x_{n-1}) \quad (9.2)$$

به ازای هر عدد صحیح  $k$  بین  $0$  و  $n$ ، گیریم تابع  $q_k(x)$  مجموع  $k+1$  جمله اول این صورت باشد،

$$q_k(x) = A_0 + A_1(x-x_0) + A_2(x-x_0)(x-x_1) \\ + \dots + A_k(x-x_0)\dots(x-x_{k-1})$$

در این صورت هر کدام از جمله‌های باقیمانده در (۹.۲) دارای عامل  $(x-x_0)\dots(x-x_k)$  است و می‌توانیم رابطه (۹.۲) را به صورت

$$p_n(x) = q_k(x) + (x-x_0)\dots(x-x_k)r(x)$$

که در آن، صورت بسجمله‌ای  $r(x)$  چندان مهم نیست بنویسیم. نکته مهم این است که این جمله آخر یعنی  $(x-x_0)\dots(x-x_k)r(x)$  در نقاط  $x_0, \dots, x_k$  صفر می‌شود. از این رو خود  $q_k(x)$  می‌باید  $f(x)$  را در  $x_0, \dots, x_k$  درونیابی کند [زیرا که  $p_n(x)$  درونیابی می‌کند]. از آنجا که  $q_k(x)$  نیز یک بسجمله‌ای از درجه نایبتر از  $k$  است، نتیجه می‌شود که  $q_k(x) = p_k(x)$ ، یعنی  $q_k(x)$  می‌باید بسجمله یکتایی از درجه نایبتر از  $k$  باشد که  $f(x)$  را در نقاط  $x_0, \dots, x_k$  درونیابی می‌کند.

این نشان می‌دهد که صورت نیوتنی (۹.۲) برای بسجمله‌ای درونیاب  $p_n(x)$  را می‌توان مرحله به مرحله همان طوری که دنباله  $(p_0(x), p_1(x), p_2(x), \dots)$  ساخته می‌شود بنا کرد و  $p_k(x)$  را با افزودن جمله بعدی در صورت نیوتنی (۹.۲)، یعنی

$$p_k(x) = p_{k-1}(x) + A_k(x-x_0)\dots(x-x_{k-1})$$

از  $p_{k-1}(x)$  به دست آورد. همچنین رابطه بالا نشان می‌دهد که ضریب  $A_k$  در صورت نیوتنی (۹.۲)، برای بسجمله‌ای درونیاب ضریب جمله پیشرو، یعنی ضریب  $x^k$  در بسجمله‌ای  $p_k(x)$  است که  $p_k(x)$  از درجه نایبتر از  $k$  و در  $x_0, \dots, x_k$  بر  $f(x)$  منطبق است. این ضریب فقط به مقادیر  $f(x)$  در نقاط  $x_0, \dots, x_k$  بستگی دارد، این ضریب  $k$  امین تفاضل منقسم  $f(x)$  در نقاط  $x_0, \dots, x_k$  نامیده شده (که دلایل آن در بخش بعدی آمده است) و

با نماد

$$f[x_0, \dots, x_k]$$

شان داده شده است. با این تعریف، به صورت نیوتنی برای بسجمله‌ای درونیاب می‌رسیم

$$p_n(x) = f[x_0] + f[x_0, x_1](x-x_0) + f[x_0, x_1, x_2](x-x_0)(x-x_1) \\ + \dots + f[x_0, x_1, \dots, x_n](x-x_0)(x-x_1)\dots(x-x_{n-1})$$

فرمول بالا را می‌توان به صورت فشرده‌تر:

$$p_n(x) = \sum_{i=0}^n f[x_0, \dots, x_i] \prod_{j=0}^{i-1} (x-x_j) \quad (10.2)$$

نوشت. اگر از قرارداد

$$\prod_{m=r}^s a_m = \begin{cases} a_r a_{r+1} \dots a_s & r \leq s \\ 1 & r > s \end{cases} \quad \text{برای } r \leq s \\ \text{برای } r > s$$

استفاده کنیم، به ازای  $n=1$  رابطه (10.2) به صورت

$$p_1(x) = f[x_0] + f[x_0, x_1](x-x_0)$$

درمی‌آید و در مقایسه با فرمول  $p_1(x) = f(x_0) + [f(x_1) - f(x_0)] / [x_1 - x_0]$  که قبلاً به دست آمده است، دیده می‌شود که

$$f[x_0] = f(x_0)$$

$$f[x_0, x_1] = [f(x_1) - f(x_0)] / (x_1 - x_0) = \frac{f(x_0) - f(x_1)}{x_0 - x_1} \quad (11.2)$$

در هر حال نخستین تفاضل منقسم، نسبت تفاضلهاست.

### تمرین

۱-۲۰۲ با در نظر گرفتن رابطه  $w(x) = (x-x_0)(x-x_1)\dots(x-x_n)$ ، ثابت کنید که  $f[x_0, x_1, \dots, x_n] = \sum_{i=0}^n f(x_i) / w'(x_i)$ . (دانهمایی: ضریب پیشرو بسجمله‌ای (۷.۲) را پیدا کنید.)

۲-۲۰۲ حد فرمولی راکسه برای  $f[x_0, x_1, \dots, x_n]$  در تمرین ۱-۲۰۲ داده شده، در حالتی که  $x_p \rightarrow x_1$ ، و کلیه نقاط دیگر ثابت بمانند، محاسبه کنید.

۳-۲۲ در صورتی که  $f(x)$  یک بسجمله‌ای نایبتر از درجه  $n$  باشد، ثابت کنید آن بسجمله‌ای از درجه نایبتر از  $n$  که  $f(x)$  را در  $n+1$  نقطه متمایز درون‌یابی کند، همان خود  $f(x)$  است.

۴-۲۲ ثابت کنید که  $k$  امین تفاضل منقسم  $p[x_0, \dots, x_k]$  از بسجمله‌ای درجه نایبتر از  $k$  می  $p(x)$ ، به نقاط درون‌یابی  $x_0, x_1, \dots, x_k$  بستگی ندارد.

۵-۲۲ ثابت کنید که  $k$  امین تفاضل منقسم یک بسجمله‌ای با درجه کمتر از  $k$ ، برابر صفر است.

### ۳.۲ جدول تفاضل منقسم

تفاضلهای منقسم از مرتبه بالاتر را می توان از فرمول

$$f[x_0, \dots, x_k] = \frac{f[x_1, \dots, x_k] - f[x_0, \dots, x_{k-1}]}{x_k - x_0} \quad (12.2)$$

به دست آورده که صحت آن را می توان به شرح زیر بررسی کرد.

گیریم  $p_i(x)$  یک بسجمله‌ای از درجه نایبتر از  $i$  باشد که در نقاط  $x_0, \dots, x_i$  بر  $f(x)$  منطبق است. فرض می کنیم  $q_{k-1}(x)$  یک بسجمله‌ای از درجه نایبتر از  $k-1$  باشد که در نقاط  $x_1, \dots, x_k$  بر  $f(x)$  منطبق است. در این صورت

$$p(x) = \frac{x - x_0}{x_k - x_0} q_{k-1}(x) + \frac{x_k - x}{x_k - x_0} p_{k-1}(x) \quad (13.2)$$

یک بسجمله‌ای از درجه نایبتر از  $k$  است و به آسانی دیده می شود که به ازای  $p(x_i) = f(x_i)$ ،  $i = 0, \dots, k$  در نتیجه بدعلت یکتا بودن بسجمله‌ای درون‌یاب، باید داشته باشیم  $p(x) = p_k(x)$ . بنا بر این

ضریب پیشرو  $f[x_0, \dots, x_k] = p_k(x)$  به موجب تعریف

$$= \frac{\text{ضریب پیشرو } q_{k-1}(x)}{x_k - x_0} \quad \text{بنا بر (13.2)}$$

$$\frac{\text{ضریب پیشرو } p_{k-1}(x)}{x_k - x_0}$$

$$= \frac{f[x_1, \dots, x_k] - f[x_0, \dots, x_{k-1}]}{x_k - x_0} \quad \text{بنا بر تعریف}$$

که اثبات فرمول مهم (۱۲.۲) است.

□ مثال ۳.۲: با استفاده از فرمول نیوتن مثال ۲.۲ را حل کنید.  
در این مثال باید بسجمله ای  $p_2(x)$  از درجهٔ نایبتر از ۲ را طوری تعیین کنیم که روابط زیر برقرار باشند

$$p_2(1) = ۱۵۷۰۹ \quad p_2(۴) = ۱۵۷۲۷ \quad p_2(۶) = ۱۵۷۵۱$$

بموجب (۱۱.۲) می توانیم محاسبات زیر را انجام دهیم

$$K[۱, ۲] = \frac{۱۵۷۰۹ - ۱۵۷۲۷}{۱ - ۲} = ۰۰۰۰۰۶$$

$$K[۲, ۶] = \frac{۱۵۷۲۷ - ۱۵۷۵۱}{۲ - ۶} = ۰۰۰۰۱۲$$

بنا بر این، بموجب رابطه (۱۲.۲) داریم

$$K[۱, ۲, ۶] = \frac{۰۰۰۰۰۶ - ۰۰۰۰۱۲}{۱ - ۶} = ۰۰۰۰۰۱۲$$

و بنا بر (۱۰.۲) خواهیم داشت

$$p_2(x) = ۱۵۷۰۹ + ۰۰۰۰۰۶(x-1) + ۰۰۰۰۰۱۲(x-1)(x-۴)$$

با قراردادن  $x = ۳۵$  در این فرمول خواهیم داشت

$$\begin{aligned} p_2(۳۵) &= ۱۵۷۰۹ + (۰۰۰۰۰۶)(۳۵) + (۰۰۰۰۰۱۲)(۳۵)(-۰۵) \\ &= ۱۵۷۲۲۵ \end{aligned}$$

□ که با نتایج به دست آمده در مثال ۲.۲ یکی است.

معادله (۱۲.۲) نشان می دهد که  $k$  امین تفاضل منقسم همان تفاضل خارج قسمت  $(k-1)$  امین تفاضل منقسم است، که مؤید مناسب بودن نام «تفاضل منقسم» است. همچنین معادله (۱۲.۲) به ما امکان می دهد که کلیهٔ تفاضلهای منقسم را که بر ای فرمول نیوتن (۱۰.۲) لازم است با یک روش ساده به کمک جدولی موسوم به جدول تفاضل منقسم تولید کنیم.  
چنین جدولی در شکل ۱.۲، به ازای  $n = ۴$  داده شده است.



$x_1$	$f\{\} = f()$	$f\{i\}$	$f\{i, j\}$	$f\{i, j, k\}$	$f\{i, j, k, l\}$
$x_0$	$f[x_0]$				
$x_1$	$-f[x_1]$	$f[x_0, x_1]$			
$x_2$	$f[x_2]$	$f[x_1, x_2]$	$f[x_0, x_1, x_2]$		
$x_3$	$f[x_3]$	$f[x_2, x_3]$	$f[x_1, x_2, x_3]$	$f[x_0, x_1, x_2, x_3]$	
$x_4$	$-f[x_4]$	$f[x_3, x_4]$	$f[x_2, x_3, x_4]$	$f[x_1, x_2, x_3, x_4]$	$f[x_0, x_1, x_2, x_3, x_4]$

شکل ۱۰۲ جدول تفاضل منقسم.

درایه‌های جدول، مثلاً ستون به‌ستون، طبق الگوریتم زیر محاسبه شده‌اند.

**الگوریتم ۲۰۲:** جدول تفاضل منقسم. دو ستون اول جدول شامل  $x_0, x_1, \dots, x_n$  و مقادیر متناظر آنها،  $f[x_0], f[x_1], \dots, f[x_n]$  داده شده‌اند.

```

For k = 1, ..., n, do:
  For i = 0, ..., n - k, do:
     $f[x_i, \dots, x_{i+k}] := \frac{f[x_{i+1}, \dots, x_{i+k}] - f[x_i, \dots, x_{i+k-1}]}{x_{i+k} - x_i}$ 
  
```

اگر عملیات مربوط به این الگوریتم با دست انجام پذیرد، دستورالعمل‌های زیر ممکن است مفید واقع شوند. از درایه‌ای که باید محاسبه شود دو قطر از دو درایه مجاورش به طرف چپ رسم کنید. اگر این خطوط به ترتیب در  $f[x_i]$  و  $f[x_j]$  خاتمه پذیرند، تفاضل دو درایه مجاور را بر تفاضل متناظر آنها  $x_j - x_i$  تقسیم کنید تا درایه مطلوب را به دست آورید. این امر در شکل ۱۰۲ برای درایه  $f[x_1, \dots, x_4]$  نشان داده شده است.

پس از آنکه جدول تفاضل منقسم تکمیل گردید به ازای  $n, \dots, 0, z$  ضرایب  $f[x_0, \dots, x_i]$  را می‌توان برای فرمول نیوتن (۱۰۲) درصدر ستونهای مربوطه پیدا کرد.

به دلیل نیاز به حافظه و همچنین به علت آنکه در بسیاری از زبانهای فورترون تنها متغیرهای موجود در عبارت DO می‌توانند زیاد شوند، الگوریتم ۲۰۲ را به صورت اصلاح شده در برنامه‌های فورترون به کار می‌برند. اولاً برای ارزیابی صورت نیوتنی به سه موجب الگوریتم ۱۰۲ آسانتر است که از صورت

$$p_n(x) = \sum_{i=0}^n f[x_i, \dots, x_n] \prod_{j=i+1}^n (x - x_j)$$

یعنی، صورت نیوتنی با مراکز  $x_n, x_{n-1}, \dots, x_1$  استفاده کنیم، زیرا مقدار  $v = p_n(z)$

را می‌توان بعداً طبق الگوریتم ۱.۲ و

$$v := f[x_0, \dots, x_n]$$

$$\text{For } i = 1, \dots, n, \text{ do:}$$

$$v := f[x_i, \dots, x_n] + (z - x_i)v$$

محاسبه کرد. ثانیاً از آنجا که به ازای  $i = 0, \dots, n$  تنها اعداد  $f[x_i, \dots, x_n]$  مورد نظر هستند نیازی به ذخیره کردن تمامی جدول تفاضل منقسم (که مستلزم یک آرایهٔ دوبعدی است که در آن تقریباً از نیمی از درایه‌ها در هر حال، به علت ویژگی مثلثی جدول تفاضل منقسم، استفاده نخواهد شد) نیازی نیست. زیرا اگر از نماد اختصاری

$$d_{ij} = f[x_i, \dots, x_{i+j}]$$

استفاده کنیم، محاسبات الگوریتم ۲.۲ به صورت زیر درمی‌آید

$$\text{For } k = 1, \dots, n, \text{ do:}$$

$$\text{For } i = 0, \dots, n - k, \text{ do:}$$

$$d_{ik} := (d_{i+1, k-1} - d_{i, k-1}) / (x_{i+k} - x_i)$$

و مخصوصاً به محض اینکه  $d_{ik}$  محاسبه شد، عدد  $d_{i, k-1}$  دیگر مورد استفاده قرار نخواهد گرفت، لذا با اطمینان می‌توان  $d_{ik}$  را روی  $d_{i, k-1}$  ذخیره کرد.

**الگوریتم ۳.۲:** محاسبه ضرایب فرمول نیوتن.  $n+1$  نقطهٔ متمایز  $x_0, \dots, x_n$  و اعداد متناظر آنها  $f(x_0), \dots, f(x_n)$ ، که  $f(x_i)$ ،  $(i=0, \dots, n)$ ، در بردار  $d_i$  ذخیره شده‌است

$$\text{For } k = 1, \dots, n, \text{ do:}$$

$$\text{For } i = 0, \dots, n - k, \text{ do:}$$

$$d_i := (d_{i+1} - d_i) / (x_{i+k} - x_i)$$

مفروض‌اند. پس به ازای  $i = 0, \dots, n$  داریم  $d_i = f[x_i, \dots, x_n]$

□ **مثال ۴.۲:** گیریم  $f(x) = (1+x^2)^{-1}$ . به ازای  $n = 2, 4, \dots, 16$  مطلوب است محاسبهٔ بسجمله‌ای  $p_n(x)$  از درجهٔ نایبتر از  $n$  که  $f(x)$  را در  $n+1$  نقطهٔ

$$x_i = i \frac{10}{n} - 5, \quad i = 0, \dots, n$$

با فواصل مساوی درونیابی کند. سپس حداکثر خطای درونیابی

$$E_n = \max_{-5 \leq x \leq 5} |f(x) - p_n(x)| \quad n = 2, 4, \dots, 16$$

را در بازه  $[-۵, ۵]$  از راه محاسبه

$$E_n \approx \max_i |f(y_i) - p_n(y_i)|$$

که در آن

$$y_i = \frac{i}{10} - 5 \quad i = 0, \dots, 100$$

به دست آورید. برای حل این مسئله در برنامهٔ فورترن زیر، از الگوریتم ۱.۲ و ۳.۲ استفاده شده است.

### برنامهٔ فورترن برای مثال ۴.۲

```

C PROGRAM FOR EXAMPLE 2.4
  INTEGER I,J,K,N,NP1
  REAL D(17),ERRMAX,H,PNOFY,X(17),Y
C POLYNOMIAL INTERPOLATION AT EQUALLY SPACED POINTS TO THE FUNCTION
  F(Y) = 1./(1. + Y*Y)
C
  PRINT 600
600 FORMAT('1 N',5X,'MAXIMUM ERROR')
DO 40 N=2,16,2
  NP1 = N+1
  H = 10./FLOAT(N)
  DO 10 I=1,NP1
    X(I) = FLOAT(I-1)*H - 5.
    D(I) = F(X(I))
  10 CONTINUE
C CALCULATE DIVIDED DIFFERENCES BY ALGORITHM 2.3
DO 20 K=1,N
  DO 20 I=1,NP1-K
    D(I) = (D(I+1) - D(I))/(X(I+K) - X(I))
  20 CONTINUE
C ESTIMATE MAXIMUM INTERPOLATION ERROR ON (-5,5)
  ERRMAX = 0.
DO 30 J=1,101
  Y = FLOAT(J-1)/10. - 5.
C CALCULATE PN(Y) BY ALGORITHM 2.1
  PNOFY = D(1)
DO 29 K=2,NP1
  PNOFY = D(K) + (Y - X(K))*PNOFY
  29 CONTINUE
  ERRMAX = MAX(ABS(F(Y) - PNOFY), ERRMAX)
  30 CONTINUE
  PRINT 630, N,ERRMAX
630 FORMAT(15,E10.7)
40 CONTINUE
                                STOP
END

```

### برونداد (نتایج) کامپیوتری برای مثال ۴.۲

N	MAXIMUM ERROR
2	6.4615385E - 01
4	4.3813387E - 01
6	6.1666759E - 01
8	1.0451739E + 00
10	1.9156431E + 00
12	3.6052745E + 00
14	7.1920080E + 00
16	1.4051542E + 01

درونیابی به وسیلهٔ بسجمله‌ایها ۶۳

ملاحظه می‌کنید که اگر چه در روند درونیابی خود از اطلاعات موجود در باب  $f(x)$  بیش از پیش استفاده کردیم، ولی خطای درونیابی به سرعت با درجهٔ روبه‌افزایشی ازدیاد پیدا کرد. علت آن، استفادهٔ ما از نقاط درونیابی با فاصلهٔ یکنواخت است. به تمرین ۱۰۶-۱۲ و معادلهٔ (۲۰.۶) نگاه کنید.

□

## تمرین

۱-۳.۲ از يك جدول لگاریتم، مقادیر  $\log x$  در نقاط  $x$  طبق جدول زیر به دست آمده است

$x$	$\log x$
۱۰۰	۰۰۰
۱۰۵	۰۰۱۷۶۰۹
۲۰۰	۰۰۳۰۱۰۳
۳۰۰	۰۰۴۷۷۱۲
۳۰۵	۰۰۵۲۴۰۷
۴۰۰	۰۰۶۰۲۰۶

جدول تفاضل منقسم این نقاط را تشکیل دهید.

۲-۳.۲ با استفاده از جدول تفاضل منقسم در تمرین ۱-۳.۲ و به کار گرفتن صورت نیوتنی يك بسجمله‌ای درونیاب از درجهٔ سوم، مقادیر  $\log ۲۰۵$ ،  $\log ۱۰۲۵$ ،  $\log ۳۰۲۵$  را درونیابی کنید.

۳-۳.۲ خطای حاصل در نتیجهٔ محاسبهٔ  $\log ۲۰۵$  در تمرین ۲-۳.۲ را از راه محاسبهٔ جملهٔ بعدی در بسجمله‌ای درونیاب، برآورد کنید. همچنین از مقایسهٔ مقدار تقریبی  $\log ۲۰۵$  با حاصل جمع  $\log ۲$  و مقدار تقریبی  $\log ۱۰۲۵$ ، این خطا را تخمین بزنید.

۴-۳.۲ فرمول

$$f[x_0, \dots, x_j] = f[x_1, \dots, x_{j+1}] + (x_0 - x_{j+1})f[x_0, \dots, x_{j+1}]$$

را به دست آورید. سپس از آن به عنوان راهی برای محاسبهٔ  $p[z, x_0, \dots, x_{n-1}]$ ،  $p[z, x_0, \dots, x_{n-2}]$ ،  $\dots$ ،  $p[z, x_0]$  و  $p[z]$ ، یعنی راهی برای به دست آوردن قطر دیگری در جدول تفاضل منقسم برای  $p(x)$ ، در تعبیر الگوریتم ضرب تودرتوی ۱۰.۲ که در بسجمله‌ای (۱۰.۲) به کار برده شده استفاده کنید.

۳-۲-۵ طبق تمرین ۲-۲-۳، یک بسجمله‌ای از درجه نایبتر از  $k$  که یک تابع  $f(x)$  را در نقاط  $x_0, \dots, x_k$  درونیابی کند، خود  $f(x)$  است هر گاه بسجمله‌ای از درجه نایبتر از  $k$  باشد. این حقیقت را می‌توان برای بررسی دقت بسجمله‌ای درونیابی که حساب شده به‌کار گرفت. برنامه فورترن در مثال ۲-۲-۴ را برای انجام این بررسی به‌طریق زیر تغییر دهید. به‌ازای  $n = 4, 8, 12, \dots, 32$ ، بسجمله‌ای  $p_n(x)$  از درجه نایبتر از  $n$  را چنان پیدا کنید که تابع  $f_n(x) = \prod_{j=1}^n (x - z - 1/3)$  را در  $0, 1, 2, \dots, n$  درونیابی کند. سپس  $E_n = \max_{x_0 \leq x \leq x_n} |f_n(x) - p_n(x)|$  را به‌وسیله  $\max |f_n(y_i) - p_n(y_i)|$  برآورد کنید.  $y_i$  در اینجا، تعداد زیادی نقطه مناسب در  $[0, n]$  است.

۳-۲-۶ ثابت کنید که  $p'_k(x)$  مشتق سهمی درونیاب تابع  $f(x)$  در  $x_0 < x_1 < x_2$ ، خط مستقیمی است که مقدار آن در نقطه  $(x_{i-1} + x_i)/2$  به‌ازای  $i = 1, 2$  برابر  $f[x_{i-1}, x_i]$  است. در موردی که  $p_n(x)$  تابع  $f(x)$  را در  $x_0 < x_1 < \dots < x_n$  درونیابی کند، از تعمیم این مطلب برای بیان  $p'_n(x)$  به‌عنوان درونیاب داده‌های  $\{f(x_i), x_i, x_{i+1}\}$  به‌ازای  $i$  مناسب در  $x_{i+1} \leq x \leq x_i$  استفاده نمایید.

### ۴-۲-۴ درونیابی در یک عدد روه‌افزایش از نقاط درونیابی

اکنون با استفاده از درونیابی بسجمله‌ای در نقاط متمایز  $x_0, x_1, x_2, \dots$  به‌مسئله برآورد  $f(x)$  در نقطه  $x = \bar{x}$  می‌پردازیم. با استفاده از بسجمله‌ای  $p_k(x)$  از درجه نایبتر از  $k$  که  $f(x)$  را در نقاط  $x_0, \dots, x_k$  درونیابی می‌کند، متوالیاً مقادیر  $p_0(\bar{x}), p_1(\bar{x}), \dots, p_k(\bar{x})$  را تا آنجا محاسبه می‌کنیم که احتمالاً، تفاضل بین  $p_k(\bar{x})$  و  $p_{k+1}(\bar{x})$  به‌اندازه کافی کوچک شود. صورت نیوتنی برای بسجمله درونیاب

$$p_k(x) = \sum_{i=0}^k f[x_0, \dots, x_i] \psi_i(x)$$

با

$$\psi_i(x) = (x - x_0) \dots (x - x_{i-1})$$

بخصوص برای چنین محاسباتی طرح شده است. اگر  $p_k(\bar{x})$  و  $\psi_k(\bar{x})$  در دست باشند، آنگاه می‌توانیم  $p_{k+1}(\bar{x})$  را به‌وسیله زیر محاسبه کنیم.

$$p_{k+1}(\bar{x}) = p_k(\bar{x}) + f[x_0, \dots, x_{k+1}] \psi_k(\bar{x})(\bar{x} - x_k)$$

الگوریتم ۴-۲: درونیابی با استفاده از یک عدد روه‌افزایش از نقاط درونیاب. نقاط متمایز  $x_0, x_1, x_2, \dots$  و مقادیر  $f(x_0), f(x_1), f(x_2), \dots$  از تابع  $f(x)$  در این نقاط و همچنین یک نقطه  $\bar{x}$  داده شده‌اند.

$$f[x_0] := f(x_0), p_0(\bar{x}) := f(x_0), \psi_0(\bar{x}) := 1$$

For  $k = 0, 1, 2, \dots$ , until satisfied, do:

$$f[x_{k+1}] := f(x_{k+1})$$

For  $i = k, \dots, 0$ , do:

$$f[x_i, \dots, x_{k+1}] := \frac{f[x_{i+1}, \dots, x_{k+1}] - f[x_i, \dots, x_k]}{x_{k+1} - x_i}$$

$$\psi_{k+1}(\bar{x}) := \psi_k(\bar{x})(\bar{x} - x_k)$$

$$p_{k+1}(\bar{x}) := p_k(\bar{x}) + f[x_0, \dots, x_{k+1}]\psi_{k+1}(\bar{x})$$

این الگوریتم، هر بار درایه‌های یک قطر از جدول تفاضل منقسم برای  $f(x)$  را در  $x_0, x_1, x_2, \dots$  تولید می‌کند. هنگام محاسبهٔ  $p_{k+1}(\bar{x})$  قطر بالاروندهٔ ماربر  $f[x_{k+1}]$  با استفاده از عدد  $f[x_{k+1}] = f(x_{k+1})$  و درایه‌های قبلاً حساب شدهٔ  $f[x_k], \dots, f[x_{k-1}, x_k]$  در قطر پیشین، تا با عدد  $f[x_0, \dots, x_{k+1}]$  محاسبه شد. حتی اگر فقط جدیدترین قطر محاسبه شده (مثلاً، در یک برنامهٔ فورترن)، حفظ شود، این الگوریتم اتفاقاً ضرایب لازم برای صورت نیوتنی تابع  $p_{k+1}(x)$  با مراکز  $x_0, x_1, \dots, x_{k+1}$  را به دست می‌دهد:

$$p_{k+1}(x) = \sum_{i=0}^{k+1} f[x_i, \dots, x_{k+1}] \prod_{j=i+1}^{k+1} (x - x_j) \quad (14.2)$$

□ مثال ۵.۲: الگوریتم ۴.۲ را برای مسئلهٔ مثال ۲.۲ و ۳.۲ به کار می‌بریم و از  $x_0 = ۳$ ،  $x_1 = ۴$ ،  $x_2 = ۶$  و نیز از  $x_3 = ۰$  استفاده می‌کنیم. در این مثال داریم  $\bar{x} = ۳.۵$ . از آنجا خواهیم داشت  $K[x_0] = ۱۵۷۰۹$  و  $p_0(\bar{x}) = ۱$  و  $\psi_0(\bar{x}) = ۱$ . سپس با  $K[x_0, x_1] = ۰۰۰۰۶$ ،  $K[x_1] = ۱۵۷۲۷$  و  $\psi_1(\bar{x}) = (\bar{x} - x_0)\psi_0(\bar{x}) = ۲۵$  را به دست می‌آوریم

$$p_1(\bar{x}) = ۱۵۷۰۹ + ۰۰۰۰۱۵ = ۱۵۷۲۴$$

با افزودن نقطهٔ  $x_2 = ۶$  داریم  $K[x_2] = ۱۵۷۵۱$  در نتیجه  $K[x_0, x_1, x_2] = ۰۰۰۰۱۲$ ،  $K[x_1, x_2] = ۰۰۰۰۱۲$  و  $\psi_2(\bar{x}) = (-۰.۵)(۲.۵) = -۱.۲۵$  لذا

$$p_2(\bar{x}) = ۱۵۷۲۴ - ۰۰۰۰۱۵ = ۱۵۷۲۲.۵$$

که همان عددی است که قبلاً در مثال ۳.۲ محاسبه شده بود. به منظور بررسی خطا در این تقریب برای  $K(۳.۵)$ ، نقطهٔ  $x_3 = ۰$  را به نقاط موجود اضافه می‌کنیم. با  $K[x_3] = ۱۵۷۰۸$  و اعداد  $K[x_2, x_3] = ۰۰۰۰۷۱۷$  و  $K[x_1, x_2, x_3] = -۰۰۰۰۱۲۱$  و  $K[x_0, x_1, x_2, x_3] = -۰۰۰۰۰۰۱$  و به ازای  $\psi_3(\bar{x}) = (-۲.۵)(-۱.۲۵) = ۳.۱۲۵$  مقدار  $p_3(\bar{x})$  را به دست می‌آوریم

$$p_3(\bar{x}) = 1.57225 - 0.000005x$$

مقدار فوق نشان دهنده این است که  $1.57223$  یا  $1.57224$  احتمالاً مقدار  $K(3.5)$  در چارچوب دقت مقادیر داده شده برای مقدار  $K(x)$  است.

اگر این محاسبات با دست انجام گیرد به آسانی می‌تواند به صورت جدولی که در شکل ۲.۲ نشان داده شده تنظیم شود، و در ضمن چگونگی تولید تدریجی تفاضل منقسم بر اثر الگوریتم (۴.۲) را نیز نشان می‌دهد. □

در زیر فهرستی از يك تابع فورترن به نام تابع TABLE (تابع جدولی) داده شده است، که برای درونیایی در يك جدول مفروض از طول و عرضهای  $X(I), F(I)$  به ازای  $NTABLE$ ،  $I = 1, \dots, NTABLE$ ، با  $F(I) = f(X(I))$  در پیدا کردن يك تقریب خوب از  $f(x)$  در نقطه  $x = XBAR$  از الگوریتم (۴.۲) مورد استفاده قرار می‌گیرد. این برنامه  $p_0(XBAR)$  و  $p_1(XBAR)$  و ... را تا برقراری شرط

$$|(p_k(XBAR) - p_{k-1}(XBAR))| \leq TOL$$

(TOL حد خطای مطلوب) یا تا برقراری  $k+1 = \min(20, NTABLE)$  تولید می‌کند و آنگاه عدد  $p_k(XBAR)$  را باز می‌گرداند. دنباله  $x_0, x_1, x_2, \dots$  از نقاط درونیایی جدولی  $X(1), X(2), \dots, X(NTABLE)$  به شرح زیر انتخاب می‌شوند. اگر  $X(I) < XBAR \leq X(I+1)$  آنگاه  $x_0 = X(I+1)$  و  $x_1 = X(I)$  و  $x_2 = X(I+2)$  و  $x_3 = X(I-1)$  و ... انتخاب می‌شوند، جز نزدیک به ابتدا و انتهای جدول داده شده، که سرانجام فقط نقاط سمت راست یا سمت چپ  $XBAR$  به کار گرفته می‌شوند. برای حفاظت برنامه (و استفاده کننده!) در قبال انتخاب مقدار غیر معقول برای  $XTOL$ ، برنامه باید طوری اصلاح شود که تعیین کند که هنگام افزایش  $k$ ، چه موقع تفاضلهای متوالی  $|p_{k+1}(XBAR) - p_k(XBAR)|$  شروع به افزایش می‌کنند و آیا این افزایش صورت می‌گیرد یا نه. به ترمین ۱-۴.۲ نیز نگاه کنید.)

k	$p_k(\bar{x})$	$\psi_k(\bar{x})$	$x_k$	$K[1]$	$K[2]$	$K[3]$	$K[4]$
0	1.5709 + 15	1.	1	1.5709			
1	1.5724 - 15	2.5	4	1.5727	0.0006		
2	1.57225 - 3	-1.25	6	1.5751	0.0012	0.00012	-0.000001
3	1.572247	3.125	0	1.5708	0.000717	0.000121	

شکل ۲.۲

## زیور نامه فورترن برای درونیابی در یک تابع TABLE (تابع جدولی)

```

REAL FUNCTION TABLE (XBAR, X, F, NTABLE, TOL, IFLAG)
C RETURNS AN INTERPOLATED VALUE TABLE AT XBAR FOR THE FUNCTION
C TABULATED AS (X(I),F(I)), I=1,...,NTABLE.
      INTEGER IFLAG,NTABLE, J,NEXT,NEXTL,NEXTR
      REAL F(NTABLE),TOL,X(NTABLE),XBAR, A(20),ERROR,PSIK,XK(20)
C***** I N P U T *****
C XBAR POINT AT WHICH TO INTERPOLATE .
C X(I), F(I), I=1,...,NTABLE CONTAINS THE FUNCTION TABLE .
C A S S U M P T I O N ... X IS ASSUMED TO BE INCREASING.
C NTABLE NUMBER OF ENTRIES IN FUNCTION TABLE.
C TOL DESIRED ERROR BOUND .
C***** O U T P U T *****
C TABLE THE INTERPOLATED FUNCTION VALUE .
C IFLAG AN INTEGER,
C   = 1 , SUCCESSFUL EXECUTION ,
C   = 2 , UNABLE TO ACHIEVE DESIRED ERROR IN 20 STEPS,
C   = 3 , XBAR LIES OUTSIDE OF TABLE RANGE. CONSTANT EXTRAPOLATION IS
C   USED.
C***** M E T H O D *****
C A SEQUENCE OF POLYNOMIAL INTERPOLANTS OF INCREASING DEGREE IS FORMED
C USING TABLE ENTRIES ALWAYS AS CLOSE TO XBAR AS POSSIBLE. EACH IN-
C TERPOLATED VALUE IS OBTAINED FROM THE PRECEDING ONE BY ADDITION OF A
C CORRECTION TERM (AS IN THE NEWTON FORMULA). THE PROCESS TERMINATES
C WHEN THIS CORRECTION IS LESS THAN TOL OR, ELSE, AFTER 20 STEPS.
C
      LOCATE XBAR IN THE X-ARRAY.
      IF (XBAR .GE. X(1) .AND. XBAR .LE. X(NTABLE)) THEN
        DO 10 NEXT=2,NTABLE
          IF (XBAR .LE. X(NEXT)) GO TO 12
10      CONTINUE
      END IF
      IF (XBAR .LT. X(1)) THEN
        TABLE = F(1)
      ELSE
        TABLE = F(NTABLE)
      END IF
      PRINT 610,XBAR
610  FORMAT(E16.7,' NOT IN TABLE RANGE.')
```

RETURN

```

12  XK(1) = X(NEXT)
     NEXTL = NEXT-1
     NEXTR = NEXT+1
     A(1) = F(NEXT)
     TABLE = A(1)
     PSIK = 1.
C     USE ALGORITHM 2.4, WITH THE NEXT XK ALWAYS THE TABLE
C     ENTRY NEAREST XBAR OF THOSE NOT YET USED.
     KPLMAX = MIN(20,NTABLE)
     DO 20 KPL=2,KPLMAX
       IF (NEXTL .EQ. 0) THEN
         NEXT = NEXTR
         NEXTR = NEXTR+1
       ELSE IF (NEXTR .GT. NTABLE) THEN
         NEXT = NEXTL
         NEXTL = NEXTL-1
       ELSE IF (XBAR - X(NEXTL) .GT. X(NEXTR) - XBAR) THEN
         NEXT = NEXTR
         NEXTR = NEXTR+1
       ELSE
         NEXT = NEXTL
         NEXTL = NEXTL-1
       END IF
       XK(KPL) = X(NEXT)
       A(KPL) = F(NEXT)
       DO 13 J=KPL-1,1,-1
         A(J) = (A(J+1) - A(J))/(XK(KPL) - XK(J))
13      CONTINUE
C     FOR I=1,...,KPL, A(I) NOW CONTAINS THE DIV.DIFF. OF
C     F(X) OF ORDER K-I AT XK(1),...,XK(KPL).
     PSIK = PSIK*(XBAR - XK(KPL-1))
     ERROR = A(1)*PSIK
```



```

C          TEMPORARY PRINTOUT
          PRINT 613,KP1,XX(KP1),TABLE,ERROR
613      FORMAT(I10,3E17.7)
          TABLE = TABLE + ERROR
          IF (ABS(ERROR) .LE. TOL) THEN
              IFLAG = 1
              RETURN
          END IF
20      CONTINUE
          PRINT 620,KP1MAX
620      FORMAT(' NO CONVERGENCE IN ',I2,' STEPS. ')
          IFLAG = 2
          RETURN
          END
END

```

### تمرین

۱-۴۰۲ تابع TABLE (تابع جدولی) فورترن که در متن داده شده است به محض اینکه شرط  $|p_{k+1}(XBAR) - p_k(XBAR)| \leq TOL$  برقرار گردد، تمام می شود. با مثالهای زیر نشان دهید که این امر قرار گرفتن مقدار  $p_{k+1}(XBAR)$  را، که به وسیله TABLE برگردانده می شود، در محدوده TOL از عدد مطلوب  $f(XBAR)$  تضمین نمی کند.

(الف)  $X(I) = -10$ ,  $X(I+1) = 10$ ,  $XBAR = 0$ ,  $TOL = 10^{-5}$

به ازای  $I$ ;  $f(x) = x^2$

(ب)  $X(I) = -100$ ,  $X(I+1) = 0$ ,  $X(I+2) = 100$

$XBAR = -50$ ,  $TOL = 10^{-5}$

۲-۴۰۲ درونیابی خطی بارسته<sup>۱</sup> بر پایه ملاحظات زیر که به نوبت<sup>۲</sup> منسوب است استوار می باشد: بسجمله ای از درجه<sup>۳</sup> نایبتر از  $i - z$  را که  $f(x)$  را به ازای  $z \leq i$  در نقاط  $x_i, \dots, x_{i+1}, x_i$  درونیابی می کند با  $p_{i,j}(x)$  نشان می دهیم. پس داریم

$$P_{ij}(x) = \frac{x - x_i}{x_j - x_i} P_{i+1,j-1}(x) + \frac{x_j - x}{x_j - x_i} P_{i,j-1}(x)$$

صحت این همانی<sup>۳</sup> را تحقیق کنید. [دانهمایی: ما از این همانی در بخش ۳.۲ استفاده کرده ایم. به معادله (۱۳.۲) نگاه کنید.]

۳-۴۰۲ درونیابی خطی بارسته (ادامه) - همانی نوبل، که در تمرین ۲-۴۰۲ اثبات شد این امکان را به دست می دهد که درایه های جدول مثالی زیر

$$\begin{array}{ccc}
 f(x_0) = p_{00}(\bar{x}) & & \\
 & p_{01}(\bar{x}) & \\
 f(x_1) = p_{11}(\bar{x}) & & p_{02}(\bar{x}) \\
 & p_{12}(\bar{x}) & \\
 f(x_2) = p_{22}(\bar{x}) & & \vdots \quad ; \quad p_{0n}(\bar{x}) \\
 & \vdots & \\
 & & p_{n-2,n}(\bar{x}) \\
 \vdots & \vdots & \\
 & p_{n-1,n}(\bar{x}) & \\
 f(x_n) = p_{nn}(\bar{x}) & & 
 \end{array}$$

را ستون به ستون تولید کنیم، این عمل به وسیلهٔ عملیاتی نظیر درونیابی خطی انجام می‌گیرد و در نهایت عدد مورد نظر  $p_{0,n}(\bar{x})$  را، که با مقدار بسجمله‌ای درونیاب در  $\bar{x}$  برابر و در  $n+1$  نقطهٔ  $x_0, \dots, x_n$  بر  $f(x)$  منطبق است به دست می‌دهد. این الگوریتم، الگوریتم نویل نامیده می‌شود. تفاوت الگوریتم ایتکن<sup>۱</sup> با الگوریتم نویل در این است که در الگوریتم ایتکن یک جدول مثلثی تولید می‌شود که ستون  $j$ ام آن شامل اعداد

$$p_{0,1,\dots,j,j+1}(\bar{x})$$

$$p_{0,1,\dots,j,j+2}(\bar{x})$$

$$p_{0,1,\dots,j,n}(\bar{x})$$

می‌باشد که در آن  $p_{0,1,\dots,j,r}(x)$  (به ازای  $r > j$ ) یک بسجمله‌ای است از درجهٔ نایبتر از  $j+1$ ، منطبق بر  $f(x)$  در نقاط  $x_0, x_1, \dots, x_j$  و  $x_r$ . با شمارش تعداد عملیات نشان دهید که الگوریتم نویل پرهزینه‌تر از الگوریتم ۴.۲ است. (همچنین باید توجه کرد که الگوریتم ۴.۲، بدون هزینهٔ اضافی، صورت نیوتنی بسجمله‌ای درونیاب را تولید می‌کند که برای محاسبهٔ مقادیر بعدی تابع در نقاط دیگر مورد استفاده قرار می‌گیرد. در حالی که اطلاعات حاصل از الگوریتم نویل یا ایتکن برای محاسبهٔ مقادیر در نقاط دیگر کمکی نمی‌کنند.)

**۴-۴.۲ در درونیابی معکوس<sup>۲</sup> در یک جدول، عدد  $\bar{y}$  داده شده است و می‌خواهیم نقطهٔ  $\bar{x}$  را طوری پیدا کنیم که داشته باشیم  $f(\bar{x}) = \bar{y}$ . در اینجا  $f(x)$  یک تابع جدول‌بندی شده است. اگر  $f(x)$  (پیوسته و) اکیداً یکنوای صعودی یا نزولی<sup>۳</sup> باشد، همواره می‌توان این مسئله را بدین صورت حل کرد که جدول داده شدهٔ  $x_i, f(x_i)$ ،  $(i=0, 1, 2, \dots)$ ، را با اختیار کردن  $y_i = f(x_i)$ ،  $g(y_i) = x_i$ ،  $(i=0, 1, 2, \dots)$ ، یک جدول  $g(y_i)$ ،  $(i=0, 1, 2, \dots)$  مجهول  $g(\bar{y})$  در این جدول درونیابی کرد. از تابع جدولی فورتون برای پیدا کردن  $\bar{x}$  به طوری که  $\sin \bar{x} = 0.6$ ، استفاده کنید.**

## ۵.۲ خطای بسجمله‌ای درونیاب

گیریم  $f(x)$  تابعی حقیقی مقدار در بازهٔ  $I = [a, b]$  و  $x_0, \dots, x_n$   $n+1$  نقطهٔ متمایز از بازهٔ  $I$  باشند. اگر  $p_n(x)$  یک بسجمله‌ای باشد از درجهٔ نایبتر از  $n$  که  $f(x)$  را در نقاط  $x_0, \dots, x_n$  درونیابی می‌کند، آنگاه خطای درونیابی در  $p_n(x)$  عبارت است از:

1. Aitken's Algorithm
2. Inverse Interpolation
3. strictly monotone-increasing or-decreasing

$$e_n(x) = f(x) - p_n(x) \quad (15.2)$$

اکنون فرض کنید که  $\bar{x}$  نقطه‌ای غیر از  $x_0, \dots, x_n$  باشد. اگر  $p_{n+1}(x)$  بسجمله‌ای باشد از درجهٔ نایب‌تر از  $(n+1)$  که  $f(x)$  را در  $x_0, \dots, x_n$  و در  $\bar{x}$  درونیابی می‌کند، آنگاه  $p_{n+1}(\bar{x}) = f(\bar{x})$ ، درحالی که طبق رابطهٔ (۱۵.۲) داریم

$$p_{n+1}(x) = p_n(x) + f[x_0, \dots, x_n, \bar{x}] \prod_{j=0}^n (x - x_j)$$

در نتیجه

$$f(\bar{x}) = p_{n+1}(\bar{x}) = p_n(\bar{x}) + f[x_0, \dots, x_n, \bar{x}] \prod_{j=0}^n (\bar{x} - x_j)$$

بنابراین به‌ازای تمام نقاط  $\bar{x} \neq x_0, \dots, x_n$

$$e_n(\bar{x}) = f[x_0, \dots, x_n, \bar{x}] \prod_{j=0}^n (\bar{x} - x_j) \quad (16.2)$$

رابطهٔ بالا نشان می‌دهد که خطا «مشابه جملهٔ بعدی» در صورت نیوتنی است.

بدون داشتن عدد  $f(\bar{x})$  نمی‌توانیم سمت راست رابطهٔ (۱۶.۲) را محاسبه کنیم. اما به‌طوری که اکنون ثابت می‌کنیم، عدد  $f[x_0, \dots, x_n, \bar{x}]$  با  $(n+1)$  امین مشتق  $f(x)$  ارتباط نزدیک دارد و با به‌کار گرفتن این اطلاعات اکنون می‌توانیم،  $e_n(\bar{x})$  را برآورد کنیم.

**قضیهٔ ۲.۲** گیریم  $f(x)$  تابعی است حقیقی مقدار که در بازهٔ  $[a, b]$  تعریف شده و  $k$  مرتبه مشتقپذیر است. اگر  $x_0, \dots, x_k$   $k+1$  نقطهٔ متمایز در  $[a, b]$  باشند، آنگاه نقطه‌ای مانند  $\xi \in (a, b)$  وجود دارد به‌طوری که رابطهٔ زیر برقرار است

$$f[x_0, \dots, x_k] = \frac{f^{(k)}(\xi)}{k!} \quad (17.2)$$

به‌ازای  $k=1$ ، رابطهٔ فوق، همان قضیهٔ مقدار میانگین برای مشتق است (به‌بخش ۷.۱ نگاه کنید). برای اثبات قضیهٔ بالا در حالت کلی، باید توجه داشت که تابع خطای  $e_k(x) = f(x) - p_k(x)$  در  $I = [a, b]$  دارای  $k+1$  ریشهٔ متمایز  $x_0, \dots, x_k$  است. بنا بر این اگر  $f(x)$  و در نتیجه  $e_k(x)$ ،  $k$  مرتبه روی  $(a, b)$  مشتقپذیر باشد از قضیهٔ رول<sup>۱</sup> نتیجه می‌شود (به‌بخش ۷.۱ نگاه کنید) که  $e'_k(x)$  حداقل  $k$  ریشه در  $(a, b)$  دارد و بنا بر این حداقل  $e''(x)$  دارای  $k-1$  ریشه در  $(a, b)$  است، و اگر عمل را به‌همین ترتیب ادامه دهیم سرانجام به‌این نتیجه می‌رسیم که  $e_k^{(k)}(x)$  حداقل یک ریشه در  $(a, b)$  دارد. گیریم  $\xi$  چنین ریشه‌ای باشد. بنا بر این داریم:

$$0 = e_k^{(k)}(\xi) = f^{(k)}(\xi) - p_k^{(k)}(\xi)$$

از طرف دیگر می‌دانیم که به ازای هر  $x$

$$p_k^{(k)}(x) = f[x_0, \dots, x_k]k!$$

زیرا بنا بر تعریف،  $f[x_0, \dots, x_k]$  ضریب پیشرو  $p_k(x)$  است و بدین ترتیب رابطه (۱۷.۲) به دست خواهد آمد.

به فرض آنکه  $a = \min_i x_i$ ،  $b = \max_i x_i$ ، نتیجه می‌شود که نقطه مجهول  $\xi$  در رابطه (۱۷.۲) در بین  $x_i$ ها واقع است.

اگر قضیه ۲.۲ را برای (۱۶.۲) به کار ببریم، قضیه ۳.۲ به دست می‌آید.

**قضیه ۳.۲** گیریم  $f(x)$  تابعی است حقیقی مقدار که روی  $[a, b]$  تعریف شده و در  $(a, b)$ ،  $n+1$  مرتبه مشتقپذیر است. اگر  $p_n(x)$  یک بسجمله‌ای باشد از درجه نایبتر از  $n$  که  $f(x)$  را در  $n+1$  نقطه متمایز  $x_0, \dots, x_n$  و متعلق به بازه  $[a, b]$  درونیابی می‌کند، آنگاه به ازای جمیع مقادیر  $\bar{x} \in [a, b]$  نقطه‌ای مانند  $\bar{x} \in (a, b)$   $\xi = \xi(\bar{x})$  وجود دارد به طوری که

$$e_n(\bar{x}) = f(\bar{x}) - p_n(\bar{x}) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{j=0}^n (\bar{x} - x_j) \quad (18.2)$$

لازم است توجه کنیم که  $\xi = \xi(\bar{x})$  به نقطه  $\bar{x}$ ، یعنی به جایی که بر آورد خط لازم است، بستگی دارد. این بستگی حتی نیاز ندارد که پیوسته باشد. از آنجا که در فصل ۷ به انتگرالگیری و مشتقگیری از  $e_n(x)$  نسبت به  $x$  احتیاج داریم، معمولاً برای چنین منظوری فرمول (۱۶.۲) ترجیح داده می‌شود. زیرا همان طوری که در بخش ۷.۲ نشان داده خواهد شد  $f[x_0, \dots, x_n, x]$  تابعی است خوش رفتار از  $x$ .

فرمول خطای (۱۸.۲) از فایده عملی محدودی برخوردار است، زیرا در حالت کلی به ندرت از مقدار  $f^{(n+1)}(x)$  آگاهیم. تقریباً هیچوقت جای نقطه  $\xi$  را نمی‌دانیم. اما هر گاه یک کران  $^3$  برای  $|f^{(n+1)}(x)|$  در تمامی بازه  $(a, b)$  در دست باشد، آنگاه می‌توانیم از (۱۸.۲) برای پیدا کردن کران خطای بسجمله‌ای درونیابی در آن بازه استفاده کنیم.

□ مثال ۶.۲: کران خطایی در درونیابی خطی پیدا کنید.

بسجمله خطی که  $f(x)$  را در  $x_0$  و  $x_1$  درونیابی می‌کند بدشرح زیر تعریف شده است

$$p_1(x) = f(x_0) + f[x_0, x_1](x - x_0) = \frac{(x_1 - x)f(x_0) + (x - x_0)f(x_1)}{x_1 - x_0}$$

1. real valued function

2. well-behaved function

3. bound

پس معادله (۱۸.۲) فرمول خطای

$$f(\bar{x}) - p_1(\bar{x}) = (\bar{x} - x_0)(\bar{x} - x_1) \frac{f''(\xi)}{2!}$$

را که در آن  $\xi$  بستگی به  $\bar{x}$  دارد، به دست می‌دهد. اگر  $\bar{x}$  نقطه‌ای بین  $x_0$  و  $x_1$  باشد،  $\xi$  نیز بین  $x_0$  و  $x_1$  خواهد بود. بنا بر این اگر بدانیم که در  $[x_0, x_1]$ ، نامساوی  $|f''(x)| \leq M$  برقرار است، آنگاه خواهیم داشت

$$|f(\bar{x}) - p_1(\bar{x})| \leq |(\bar{x} - x_0)(\bar{x} - x_1)| \frac{M}{2}$$

ماکسیمم مقدار  $|(\bar{x} - x_0)(\bar{x} - x_1)|$  به ازای  $\bar{x} \in [x_0, x_1]$  در  $\bar{x} = (x_0 + x_1)/2$  حاصل می‌شود و مقدار آن برابر  $(x_1 - x_0)^2/4$  است. در نتیجه به ازای هر  $\bar{x} \in [x_0, x_1]$  داریم:

$$\square \quad |f(\bar{x}) - p_1(\bar{x})| \leq (x_1 - x_0)^2 \frac{M}{8}$$

$\square$  مثال ۷.۲: بازه  $h$  را در جدولی که مقادیر تابع  $f(x) = \sqrt{x}$  را بین ۱ و ۲ در فواصل مساوی به دست می‌دهد، طوری تعیین کنید که درونیایی با بسجمله‌ای درجه دوم در این جدول دقت عمل مطلوب را نتیجه دهد.

بنا بر فرض، این جدول با شرط  $x_i = 1 + ih$  برای  $i = 0, \dots, N$  و  $N = (2 - 1)/h$  شامل  $f(x_i)$  است. اگر  $\bar{x} \in [x_i, x_{i+1}]$  آنگاه  $f(\bar{x})$  را با  $p_2(\bar{x})$  تقریب می‌زنیم که  $p_2(x)$  یک بسجمله‌ای است از درجه دوم که  $f(x)$  را در  $x_{i-1}, x_i, x_{i+1}$  درونیایی می‌کنسد. پس به موجب ۱۸.۲، به ازای نقطه‌ای مانند  $\xi$  در  $(x_{i-1}, x_{i+1})$ ، خطا برابر است با

$$f(\bar{x}) - p_2(\bar{x}) = (\bar{x} - x_{i-1})(\bar{x} - x_i)(\bar{x} - x_{i+1}) \frac{f'''(\xi)}{3!}$$

از آنجا که مقدار  $\xi$  را نمی‌دانیم فقط می‌توانیم  $f'''(\xi)$  را برآورد کنیم

$$|f'''(\xi)| \leq \max_{1 \leq i \leq 2} |f'''(x)|$$

طبق محاسبه داریم:  $f'''(x) = (3/8)x^{-5/2}$ ، بنا بر این  $|f'''(\xi)| \leq 3/8$ . بعلاوه با استفاده از تغییر متغیر خطای  $y = x - x_i$  داریم:

$$\begin{aligned} & \max_{x \in [x_{i-1}, x_{i+1}]} |(x-x_{i-1})(x-x_i)(x-x_{i+1})| \\ &= \max_{y \in [-h, h]} |(y+h)y(y-h)| \\ &= \max_{y \in [-h, h]} |y(y^2-h^2)| \end{aligned}$$

از آنجا که مقدار تابع  $\psi(y) = y(y^2 - h^2)$  در  $y = h$  و  $y = -h$  برابر صفر می شود، ما کسیمم  $|\psi(y)|$  در بازه  $[-h, h]$  باید در یکی از نقاط فرین  $\psi(y)$  پیش آید. این نقاط فرین را با حل معادله  $\psi'(y) = 3y^2 - h^2 = 0$  می توان پیدا کرد. که از این معادله  $y = \pm h/\sqrt{3}$  به دست می آید. بنابراین

$$\max_{x \in [x_{i-1}, x_{i+1}]} |x-x_{i-1})(x-x_i)(x-x_{i+1})| = \frac{2h^3}{3\sqrt{3}}$$

اکنون اگر  $p_2(x)$  به عنوان يك بسجملهٔ درجهٔ دومی که  $f(x) = \sqrt{x}$  را در سه نقطهٔ جدولی نزدیکتر به  $\bar{x}$  درونیابی می کند انتخاب شود، مطمئن هستیم که به ازای هر  $\bar{x} \in [1, 2]$  خواهیم داشت

$$|f(\bar{x}) - p_2(\bar{x})| \leq \frac{(2h^3/[3\sqrt{3}])(3/8)}{6} = \frac{h^3}{24\sqrt{3}}$$

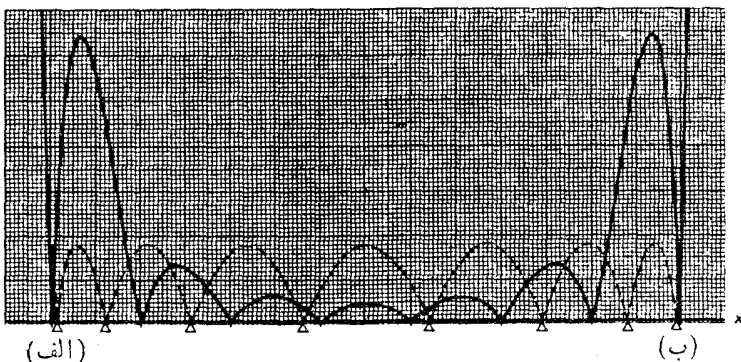
اگر بخواهیم با روش فوق، دقتی تا هفت رقم به دست آوریم می بایستی  $h$  را طوری انتخاب کنیم که

$$\frac{h^3}{24\sqrt{3}} < 5 \times 10^{-8}$$

□

و در نتیجه خواهیم داشت  $h \approx 0.0128$  یا  $h \approx 79$ .

البته تابع  $\psi_{n+1}(x) = \prod_{j=0}^n (x-x_j)$  که در (۱۸.۲) پیدا می شود، قویاً به تعیین جای نقاط درونیابی بستگی دارد. این امکان وجود دارد که نقاط مذکور به ازای مقدار مفروض  $n$  در بازهٔ مفروض  $a \leq x \leq b$  طوری انتخاب شوند که مقدار  $|\psi_{n+1}(x)|$  در حد امکان کوچک شود. این انتخاب نقاط، که نقاط چیبیشف<sup>۲</sup> نامیده می شوند، مشروحاً در بخش ۱۰.۶ مورد بحث قرار خواهد گرفت. در صورتی که نقاط درونیابی با فواصل مساوی با انتخاب معمولی صورت گیرند، هنگامی که از وسط بازه به طرف راست آن حرکت می کنیم ما کسیمم های موضعی  $|\psi_{n+1}(x)|$  افزایش می یابند، و این افزایش با افزایش  $n$  محسوستر می شود (به شکل ۳.۲ نگاه کنید). به موجب (۱۸.۲) مناسبتر است (حداقل وقتی که



شکل ۳.۲ تابع  $|\psi_{\pi} + 1(x)|$  به ازای  $n=8$  و (الف) نقاط درونیابی با فواصل مساوی (خط پر)؛ (ب) نقاط چبیشف برای همان بازه (خط چین)

داده‌هایی با فواصل یکنوا درونیابی می‌شوند که از بسجمله‌ای درونیاب فقط نزدیک به وسط نقاط داده شده استفاده شود. با نزدیک شدن به انتهای سمت چپ یا انتهای سمت راست نقاط داده شده، درونیاب کمتر مورد اطمینان است. البته بیرون از فاصله مورد نظر درونیابی این عدم اطمینان باز هم بیشتر می‌شود. این عمل را پرونیابی<sup>۱</sup> نامند و فقط باید این عمل با احتیاط زیاد انجام گیرد.

### تمرین

۱-۵.۲ جدولی از مقادیر  $\cos x$  لازم است، تا درونیابی خطی به ازای هر مقدار از  $x$  در  $[0, \pi]$ ، دقتی تا شش رقم به دست دهد. به فرض اینکه نقاط جدولی با فواصل مساوی باشند، مینیمم تعداد درایه‌های لازم در جدول چقدر است؟

۲-۵.۲ تابعی که با ضابطه

$$f(x) = \int_0^x \sin s^2 ds$$

تعریف شده است، برای مقادیر مساوی الفاصله  $x$  به ازای  $h=0.05$  جدول بندی شده است. در صورتی که به ازای هر نقطه  $\bar{x}$  از بازه  $[0, \pi/2]$  مقدار  $f(\bar{x})$  به وسیله درونیابی درجه سوم محاسبه شود، ما کسیم خطایی که با آن مواجه می‌شویم چقدر است؟

۳-۵.۲ ثابت کنید که اگر مقادیر  $f(x_0), \dots, f(x_n)$  تنها اطلاعات ما درباره تابع  $f(x)$  باشند، آنگاه درباره خطای  $e_n(\bar{x}) = f(\bar{x}) - p_n(\bar{x})$  در نقطه  $\bar{x} \neq x_0, \dots, x_n$

هیچ‌گونه اظهار نظری نمی‌توانیم بکنیم، یعنی خطا ممکن است «بسیار بزرگ» یا «بسیار کوچک» باشد. [دانهمایی: درونیابی تابع  $f(x) = K(x - x_0) \dots (x - x_n)$  را که در آن  $K$  ثابت مجهولی است در نقاط  $x_0, x_1, \dots, x_n$  در نظر بگیرید.] این موضوع چه تأثیری روی برنامه‌هایی نظیر **تابع جدولی**<sup>۱</sup> در بخش ۴.۲ یا الگوریتم ۴.۲ دارد؟

۴-۵.۲ رابطه (۱۸.۲) را برای به‌دست آوردن کران پایین در خطای درونیابی  $|f(\bar{x}) - p_n(\bar{x})|$  به‌کار برید وقتی که  $f(x) = \ln x$ ،  $n = 3$ ،  $x_0 = 1$ ،  $x_1 = 4/3$ ،  $x_2 = 2$ ،  $x_3 = 5/3$ ،  $\bar{x} = 3/2$

### ۶.۲ درونیابی یک تابع جدولی بر اساس نقاط با فواصل مساوی

در بیشتر محاسبات مهندسی و علمی، توابعی مانند  $\sin x$ ،  $e^x$ ،  $J_n(x)$ ،  $\operatorname{erf}(x)$  و غیره مورد استفاده قرار می‌گیرند که تعریف آنها به‌وسیلهٔ یک سری نامتناهی، یا به‌صورت جواب یک معادلهٔ دیفرانسیل، یا به‌وسیلهٔ فرایندهای مشابیهی که مستلزم حد می‌باشند، و بنا بر این در حالت کلی نمی‌توانند در تعداد مراحل متناهی به‌دست آیند، صورت می‌گیرد. مؤسسات کامپیوتری برنامه‌هایی برای به‌دست آوردن این گونه توابع در اختیار دارند که به‌کمک بسجمله‌ایها یا نسبت آنها تقریبهایی برای توابع فوق به‌دست می‌آورند. اما پیش از پیدایش کامپیوترهای با سرعت زیاد، تنها وسیلهٔ استفاده از چنین توابعی در محاسبات، تابع جدولی بوده است. جدولی از این گونه، شامل مقادیر تابع  $f(x_i)$  برای چند نقطهٔ  $x_0, \dots, x_n$  است و استفاده‌کننده با ایستی با درونیابی (که معنی تحت‌اللفظی آن «پر کردن رخنه‌ها و صاف کردن» و لذا «اثبات نادرستی» نیز هست) مقادیر مفروض، مقدار  $f(x)$  را در نقطه‌ای که در جدول موجود نیست به‌دست آورد. درونیابی بسجمله‌ای اصلاً برای تسهیل این روند گسترش پیدا کرده است. از آنجا که معمولاً در چنین جداولی  $f(x)$  به‌صورت یک دنبالهٔ روبه‌افزایشی از نقاط با فواصل مساوی داده شده‌اند، در محاسبهٔ بسجمله‌ای درونیاب ساده‌سازیهایی می‌تواند انجام گیرد که در این قسمت مورد بحث قرار خواهیم داد.

در سراسر این بخش، فرض بر این است که  $f(x)$  به‌ازای  $x = a(h)b$  جدولی‌بندی شده است، یعنی اعداد  $f(x_i)$  به‌ازای  $i = 0, \dots, N$  در اختیار ما هستند. در اینجا داریم

$$x_i = a + ih \quad i = 0, \dots, N \quad \text{با} \quad N = \frac{b-a}{h} \quad (19.2)$$

بهتر است که یک تبدیل متغیر خطی به‌صورت زیر را دخالت دهیم

$$x = x(s) = x_0 + sh \quad \text{به‌طوری‌که} \quad s = s(x) = \frac{x - x_0}{h} \quad (20.2)$$



و یا به اختصار

$$f(x) = f(x_0 + sh) = f_s \quad (21.2)$$

این امر موجب مقایسه کردن این حالات با حالتی می‌شود که در آن مقدار  $f(x)$  در اولین  $N+1$  عدد صحیح غیر منفی معلوم است، و در نتیجه موجب ساده تر شدن نمادگذاری می‌شود. باید توجه داشت که تبدیل متغیر خطی (20.2) بسجمله‌ای درجه  $n$  ام بر حسب  $x$  را به بسجمله‌ای درجه  $n$  ام بر حسب  $s$  بدل می‌کند.

در این حالت محاسبهٔ یک بسجمله‌ای از درجهٔ نایبتر از  $n$  که  $f(x)$  را در  $x_k, \dots, x_{k+n}$  درونیابی کند، به محاسبهٔ جدول تفاضل منقسم نیازی نیست. بلکه بهتر است یک جدول تفاضل را محاسبه کنیم. برای توضیح این مطلب، تفاضل پیشرو<sup>۲</sup> را وارد می‌کنیم

$$\Delta^i f_s = \begin{cases} f_s & i = 0 \\ \Delta(\Delta^{i-1} f_s) = \Delta^{i-1} f_{s+1} - \Delta^{i-1} f_s & i > 0 \end{cases} \quad (22.2)$$

تفاضل پیشرو و تفاضل منقسم به شرح زیر با هم در ارتباط اند.

لم 30.2 به ازای جميع مقادیر  $i \geq 0$

$$f[x_k, \dots, x_{k+i}] = \frac{1}{i! h^i} \Delta^i f_k \quad (23.2)$$

از آنجا که هر دو طرف رابطه (23.2) به استقرا نسبت به  $i$  تعریف شده‌اند، اثبات لم 30.2 باید با استقرا انجام پذیرد. به ازای  $i=0$ ، معادله (23.2) فقط مدعی صحت قرارداد

$$f[x_k] = f(x_k) = f_k = \Delta^0 f_k$$

و بنابراین درست است. با فرض اینکه (23.2) به ازای  $i \geq 0$  برقرار باشد، داریم

$$\begin{aligned} f[x_k, \dots, x_{k+n+1}] &= \frac{f[x_{k+1}, \dots, x_{k+n+1}] - f[x_k, \dots, x_{k+n}]}{x_{k+n+1} - x_k} \\ &= \frac{\Delta^n f_{k+1} / (n! h^n) - \Delta^n f_k / (n! h^n)}{(n+1)h} \\ &= \frac{\Delta^{n+1} f_k}{(n+1)! h^{n+1}} \end{aligned}$$

که نشان می‌دهد (۲۳.۲) به‌ازای  $i = n + 1$  نیز برقرار است.

بدین ترتیب بسجمله‌ای از درجهٔ نایبتر از  $n$  که  $f(x)$  را در نقاط  $x_{k+n}, \dots, x_k$  درونیایی می‌کند به‌صورت زیر درمی‌آید

$$p_n(x) = \sum_{i=0}^n \frac{1}{i!h^i} \Delta^i f_k \prod_{j=0}^{i-1} (x - x_{k+j}) \quad (24.2)$$

و بر حسب  $s$ ، داریم

$$x - x_{k+j} = x_0 + sh - [x_0 + (k+j)h] = (s - k - j)h$$

بنابراین

$$p_n(x) = p_n(x_0 + sh) = \sum_{i=0}^n \Delta^i f_k \prod_{j=0}^{i-1} \frac{s - k - j}{j + 1}$$

با یک تعریف نهایی با هم این عبارت ساده‌تر می‌شود. به‌ازای مقدار حقیقی  $y$  و عدد صحیح نامنفی  $i$ ، یک تابع دو جمله‌ای<sup>۱</sup> را به‌طریق زیر تعریف می‌کنیم

$$\binom{y}{i} = \begin{cases} 1 & i = 0 \\ \prod_{j=0}^{i-1} \frac{y - j}{j + 1} = \frac{(y)(y-1)\dots(y-i+1)}{1 \cdot 2 \dots i} & i > 0 \end{cases} \quad (25.2)$$

واژهٔ «دو جمله‌ای» قابل توجه است، زیرا وقتی  $y$  یک عدد صحیح باشد، (۲۵.۲) تنها

ضرب دو جمله‌ای  $\binom{y}{i}$  است. بدین ترتیب معادلهٔ (۲۴.۲) به‌شکل سادهٔ

$$\begin{aligned} p_n(x_0 + sh) &= \sum_{i=0}^n \Delta^i f_k \binom{s-k}{i} \\ &= f_k + (s-k)\Delta f_k + \frac{(s-k)(s-k-1)}{2} \Delta^2 f_k \\ &\quad + \dots + \frac{(s-k)\dots(s-k-n+1)}{n!} \Delta^n f_k \end{aligned} \quad (26.2)$$

نوشته می‌شود که فرمول نیوتنی تفاضل پیشرو برای بسجمله‌ای از درجهٔ نایبتر از  $n$  که  $f(x)$  را در  $x_k + ih$ ،  $i = 0, \dots, n$  درونیایی می‌کند نامیده می‌شود.

اگر در (۲۶.۲)،  $k$  را مساوی صفر قرار دهیم، که این یک امر معمولی است، فرمول نیوتنی تفاضل پیشرو به‌صورت:

$$p_n(x_0 + sh) = \sum_{i=0}^n \Delta^i f_0 \binom{s}{i} \quad (27.2)$$

درمی آید. اگر  $s$  يك عدد صحیح بین صفر و  $n$  باشد، آنگاه فرمول فوق به فرمول زیر بدل می شود

$$f_s = p_n(x_0 + sh) = \sum_{j=0}^n \binom{s}{j} \Delta^j f_0 = \sum_{j=0}^s \binom{s}{j} \Delta^j f_0 \quad (28.2)$$

شابهت عجیب این فرمول با قضیه دو جمله ای

$$(a+b)^s = \sum_{j=0}^s \binom{s}{j} a^j b^{s-j}$$

تصادفی نیست. زیرا با معرفی عملگر تعویض پیشرو

$$E f_i = f_{i+1} \quad \text{به ازای جميع مقادير } i,$$

می توانیم بنویسیم  $E + 1 = \Delta$ ، یعنی

$$(\Delta + 1) f_i = (f_{i+1} - f_i) + f_i = f_{i+1} = E f_i$$

بنابراین

$$f_s = E^s f_0 = (\Delta + 1)^s f_0 = \sum_{j=0}^s \binom{s}{j} \Delta^j 1^{s-j} f_0 = \sum_{j=0}^s \binom{s}{j} \Delta^j f_0$$

که همان رابطه (28.2) است.

ما، در اینجا، از دست زدن به حساب عملی وسیع برای تفاضلات، بر اساس فرمولهایی نظیر  $E + 1 = \Delta$ ، خودداری می کنیم، لیکن فرمولی را که بلاواسطه مورد استفاده قرار می گیرد، از آن نتیجه می گیریم. از آنجا که  $\Delta = E - 1$ ، از قضیه دو جمله ای نتیجه می گیریم

$$\Delta^s = (E - 1)^s = \sum_{j=0}^s \binom{s}{j} E^j (-1)^{s-j}$$

یا

$$\Delta^s f_i = \sum_{j=0}^s (-1)^{s-j} \binom{s}{j} f_{i+j} \quad (29.2)$$

ضرایب  $\Delta^s f_i$  در (29.2) به آسانی از جدول تفاضلهای (پیشرو) برای  $f(x)$  به دست

$x_{-4}$	$f_{-4}$			
$x_{-3}$	$f_{-3}$	$\Delta f_{-4}$		
$x_{-2}$	$f_{-2}$	$\Delta f_{-3}$	$\Delta^2 f_{-4}$	
$x_{-1}$	$f_{-1}$	$\Delta f_{-2}$	$\Delta^2 f_{-3}$	$\Delta^3 f_{-4}$
$x_0$	$f_0$	$\Delta f_{-1}$	$\Delta^2 f_{-2}$	$\Delta^3 f_{-3}$
$x_1$	$f_1$	$\Delta f_0$	$\Delta^2 f_{-1}$	$\Delta^3 f_{-2}$
$x_2$	$f_2$	$\Delta f_1$	$\Delta^2 f_0$	$\Delta^3 f_{-1}$
$x_3$	$f_3$	$\Delta f_2$	$\Delta^2 f_1$	$\Delta^3 f_0$
$x_4$	$f_4$	$\Delta f_3$	$\Delta^2 f_2$	$\Delta^3 f_1$

شکل ۴.۲ جدول تفاضل-پیشرو

می‌آید. این جدول در شکل ۴.۲ نشان داده شده است. به موجب معادله (۲۲.۲) هر درایه تنها تفاضل بین درایهٔ سمت چپ پایینی و درایهٔ سمت چپ بالایی است. تفاضلهایی که در (۲۷.۲) ظاهر می‌شوند، در طول قطری قرار دارند که در شکل ۴.۲ با علامت ① مشخص شده است.

جدولهای تفاضل برای بررسی همواری يك تابع مندرج در جدول، پیدا کردن خطاهای منفرد و تعیین درجهٔ بسجمله‌ای درونیاب مناسب برای جدول مورد نظر به کار می‌رود. این نکات را با مثال زیر روشن می‌کنیم.

□ **مثال ۸۰۲:** از کتابی در مختصات درون-سیاره‌ای، مختص  $x$  مریخ را در دستگانه مختصات خورشید- مرکزی در تاریخهای داده شده، به‌طور غلط (به‌منظور نشان دادن مطلبی) برداشته‌ایم. این مختصات در فواصل ده روزه داده شده‌اند و به‌توسط ستاره‌شناسان با وسایل مختلف به‌دست آمده‌اند. در شکل ۵.۲، يك جدول تفاضل (پیشرو) برای این داده‌ها درست کرده‌ایم.

در سه ستون اول تفاضلهای، علامتهای اعداد هر ستون یکی هستند، بنا بر این دو ستون اول تفاضلهای یکنوا هستند. از ستون سوم به بعد، تفاضلهای رفتار نوسانسی محسوس نشان

$t$	$x = f(t)$	$\Delta f$	$\Delta^2 f$	$\Delta^3 f$	$\Delta^4 f$	$\Delta^5 f$	$\Delta^6 f$
1,250.5	1.39140						
		-1444					
1,260.5	1.37696		-1469				
		-2913		55			
1,270.5	1.34783		-1414		-3		
		-4327		52		97	
1,280.5	1.30456		-1362		94		-302
		-5689		146		-205	
1,290.5	1.24767		-1216		-111		408
		-6905		35		203	
1,300.5	1.17862		-1181		92		-311
		-8086		127		-108	
1,310.5	1.09776		-1054		-16		128
		-9140		111		20	
1,320.5	1.00636		-943		4		
		-10083		115			
1,330.5	0.90553		-828				
		-10911					
1,340.5	0.79642						

شکل ۵.۲  $s$ -مختص مریخ، در دستگاه استوایی خورشیدمرکزی (تا حدی همراه با خطا)

می‌دهند. اگر بپذیریم که تابع جدول‌بندی شده هموارا، یعنی دارای تغییر آهسته و کند است، آنگاه این رفتار تفاضلهای درجات بالاتر باید نتیجه خطا باشد.

فرض کنید که خطا در  $i$  امین مقدار تابع برابر با  $\varepsilon_i$  (به ازای کلیه مقادیر  $i$ ها) باشد. در این صورت جدول شکل ۵.۲ شامل اعداد  $\Delta^s(f_i + \varepsilon_i)$  است که تفاوت این اعداد با اعداد دقیق  $\Delta^s f_i$ ، که بنا بر فرض تغییر آهسته داشتند، مقدار  $\Delta^s \varepsilon_i$  است. از رابطه (۲۹.۲) نتیجه می‌گیریم

$$\Delta^s \varepsilon_i = \sum_{j=0}^s \binom{s}{j} (-1)^{s-j} \varepsilon_{i+j} \quad (30.2)$$

و لذا با توجه به  $\varepsilon := \max_j |\varepsilon_j|$  داریم

$$|\Delta^s \varepsilon_i| \leq \sum_{j=0}^s \binom{s}{j} |\varepsilon_{i+j}| \leq \varepsilon \sum_{j=0}^s \binom{s}{j} = \varepsilon (1+1)^s = \varepsilon 2^s$$

اگر مقادیر مندرج در جدول، مقادیر دقیقاً گرد شده باشند، آنگاه  $\varepsilon \leq 0.0000005$ ، و خطاها در تفاضلهای چهارم نباید بزرگتر از ۸ واحد در آخرین رقم اعشاری باشند. اما اگر رفتار نوسانی را مربوط به خطا بدانیم، خطاها خیلی بیشتر از این است.

بررسی دقیقتر این تفاضلهای چهارم نشان می‌دهد که رفتار نوسانها منظم و اصولی است. اگر عدد ۱۰، یعنی میانگین ستون تفاضلهای چهارم را از هر يك از درایه‌های این ستون کم کنیم، خواهیم داشت

$$-۶ \quad -۲۶ \quad ۸۲ \quad -۱۲۱ \quad ۸۴ \quad -۱۳$$

که يك فرد با تجربه از اعداد فوق درمی‌یابد که اشتباهی تقریباً برابر ۲۰ واحد در آخرین رقم اعشاری، در درایهٔ جدول متناظر با عدد ۱۲۱-، یعنی در درایهٔ ۱۲۴۴۷۶۷ به ازای  $۲۹۰۵ = ۲$  رخ داده است. درحقیقت يك تغییر تنها به اندازهٔ ۲۰- واحد در آخرین رقم اعشاری این درایه، طبق رابطهٔ (۳۰.۲) ستون تفاضلهای چهارم را به صورت زیر تغییر می‌دهد

$$۰ \quad -۲۰ \quad ۸۰ \quad -۱۲۰ \quad ۸۰ \quad -۲۰$$

□ که این اساساً همان مقدار کل نوسانها در آن ستون است.

**خلاصه:** خطاهای منفرداً در يك جدول تابع به وسیلهٔ نوسانهای منظم و اصولی در تفاضلهای درجات بالاتر آشکار می‌شوند. از مقایسهٔ این نوسانها پیرامون مقدار میانگین (موضعی) با نوسانهای حاصل از يك خطای تنها، بنا بر (۳۰.۲)، بر آوردی از خطا را می‌توان به دست آورد و جدول را تصحیح نمود.

در مثال فوق، با اصلاح  $f(۱۲۹۰۵)$  به ۱۲۴۴۷۸۷، جدول تفاضل شکل ۶.۲ به دست می‌آید که در آن حتی تفاضلهای چهارم دارای يك علامت اند. تفاضلهای پنجم نوسان دارند، اما مقدار آنها کوچکتر از حداکثر خطای مربوط به گرد کردن مقادیر تابع یعنی  $۲۵/۲ = ۱۶$  واحد است. از اینجا نتیجه گیری می‌کنیم که تفاضلهای پنجم اصولاً اختلالات<sup>۲</sup> ناشی از گرد کردن مقادیر تابع هستند و درونیابی به وسیلهٔ يك بسجمله‌ای درجهٔ چهارم نتایج رضایتبخش (و قابل دفاع) می‌دهد.

به علت اهمیت قبلی جدولهای تابعی، مقدار نسبتاً زیادی از مطالب مربوط به درونیابی در جدولهای تابعی در طی سده‌های گذشته گسترش یافته است. عملگرهای تفاضلی غیر از عملگر تفاضل پیشرو  $\Delta$  (مانند تمویض پیشرو  $E$ ) برای نشانگذاری فزوده برای صورتهای مختلف بسجمله‌های درونیاب معرفی شده‌اند و تفاوت این شکلهای مختلف فقط در ترتیب نقاط درونیابی در آنهاست. این صورتهای اغلب نه به علت حقایق تاریخی بلکه بیشتر از لحاظ سنتی، با نامهایی مانند: نیوتن<sup>۳</sup>، گاوس<sup>۴</sup>، بسل<sup>۵</sup>، ستیرلینگ<sup>۶</sup>، گرگوری<sup>۷</sup>، اورت<sup>۸</sup> و غیره همراه‌اند. نحوهٔ استفادهٔ کامل از انواع صورتهای را می‌توان در کتاب هیلدبراند<sup>۹</sup> [۵] پیدا کرد.

- |                   |             |            |            |
|-------------------|-------------|------------|------------|
| 1. Isolated error | 2. noise    | 3. Newton  | 4. Gauss   |
| 5. Bessel         | 6. Stirling | 7. Gregory | 8. Everett |
| 9. Hildebrand     |             |            |            |

$t$	$x = f(t)$	$\Delta f$	$\Delta^2 f$	$\Delta^3 f$	$\Delta^4 f$	$\Delta^5 f$	$\Delta^6 f$
1,250.5	1.39140	-1444					
1,260.5	1.37696	-2913	-1469	55			
1,270.5	1.34783	-4327	-1414	72	17		
1,280.5	1.30456	-5669	-1342	86	14	-3	
1,290.5	1.24787	-6925	-1256	95	9	-5	8
1,300.5	1.17862	-8086	-1161	107	12	3	-11
1,310.5	1.09776	-9140	-1054	111	4	-8	8
1,320.5	1.00636	-10083	-943	115	4	0	
1,330.5	0.90553	-10911	-828				
1,340.5	0.79642						

شکل ۶.۲ مختص استوایی  $x$  سیاره مریخ در دستگاه خورشیدمرکزی.

در اینجا صورتهای مختلف را مورد بحث قرار نمی دهیم. احساس می کنیم که الگوریتم ۴.۲ و برنامه فرعی فورترن TABLE که در بخش ۴.۲ داده شده، برای چند موردی که احتمالاً دانشجویان از جداول استفاده می کنند کافی است.

### تمرین

۶-۶.۳ ثابت کنید که یک خطای منفرد در یک جدول تابع، میانگین چند ستون اول تفاضل را تغییر نمی دهد.

۶-۶.۴ مقادیر  $f(x)$  که ذیلا داده شده اند مقادیر یک بسجمله ای معین درجه چهارم هستند. یک جدول تفاضل برای آن تشکیل دهید و از این جدول  $f(5)$  را پیدا کنید (به تمرین ۶-۶.۲ نگاه کنید).

$x$	۰	۱	۲	۳	۴
$f(x)$	۱	۵	۳۱	۱۲۱	۳۴۱

۶-۶.۳ یک جدول تفاضل برای داده های زیر تشکیل دهید و درجه بسجمله ای درونیاب لازم برای ایجاد مقادیر درونیافته را دقیقاً تا تعداد ارقام اعشاری با معنی داده شده بر آورد کنید:

$x$	$f(x)$
۱۰۰	۱۰۵۷۰۹
۲۰۰	۱۰۵۷۱۳
۳۰۰	۱۰۵۷۱۹
۴۰۰	۱۰۵۷۲۷
۵۰۰	۱۰۵۷۳۸
۶۰۰	۱۰۵۷۵۱
۷۰۰	۱۰۵۷۶۷
۸۰۰	۱۰۵۷۸۵
۹۰۰	۱۰۵۸۰۵

۴-۶۰۲ با استفاده از جدول تفاضل در شکل ۶.۲ مقادیر زیر را به دست آورید و در هر دو مورد خطا را برآورد کنید

$$f(۱۲۵۲۰۵) \text{ (الف)} \quad f(۱۳۳۲۰۵) \text{ (ب)}$$

۵-۶۰۲ ثابت کنید که اگر  $p_n(x)$  یک بسجمله‌ای از درجهٔ  $n$ ام با ضرایب پیشرو  $a_n$  باشد و  $x_0$  یک نقطهٔ دلخواه، آنگاه

$$\Delta^n p_n(x_0) = a_n n! h^n$$

و

$$\Delta^{n+1} p_n(x_0) = 0$$

[دانهمایی: از تعریف (۲۲.۲) در مورد عملگر تفاضل پیشرو، یا لم ۱۰.۲ رابطه (۱۷.۲) استفاده کنید.]

۶-۶۰۲ فرض کنیم تساوی  $x_i = x_0 + ih$ ،  $i = 0, 1, 2, \dots$ ، درست باشد و فرض کنیم که اعداد  $\Delta^0 p_n(x_n)$ ،  $\Delta^1 p_n(x_{n-1})$ ،  $\dots$ ،  $\Delta^n p_n(x_0)$  برای یک بسجمله‌ای از درجهٔ  $n$  تا بیشتر از  $n$  معلوم باشد. نشان دهید که چگونه می‌توان با این اطلاعات مقادیر  $p_n(x_{n+1})$ ،  $p_n(x_{n+2})$ ،  $\dots$  را فقط با  $n$  عمل جمع برای هر مقدار به دست آورد. [دانهمایی: بنابر تمرین ۵-۶۰۲،  $\Delta^n p_n(x_i)$  به  $i$  بستگی ندارد در حالی که به موجب تعریف تفاضل پیشرو به ازای جمیع مقادیر  $i$  و  $j$  داریم  $\Delta^j p_n(x_i) = \Delta^j p_n(x_{i-1}) + \Delta^{j+1} p_n(x_{i-1})$ ] این روش برای رسم بسجمله‌ایها مفید است. چه رابطه‌ای بین این روش و الگوریتم ۱۰.۲ وجود دارد؟



۷-۶۰۴ چگونه صورت لاگرانژی بسجمله‌ای درونیاب را، وقتی نقاط متساوی الفاصله باشند، می‌توانید ساده کنید؟

۸-۶۰۴ فرمول تفاضل پسرو نیوتن

$$p_n(x_0 + sh) = \sum_{i=0}^n (-1)^i \Delta^i f_{-i} \binom{-s}{i}$$

رابر برای استفاده در نزدیکی انتهای راست جدول به دست آورید. در این فرمول از تفاضلهای واقع در امتداد قطری از شکل ۴.۲ که با (۲) مشخص شده، استفاده شده است.

### ۷.۲\* تفاضل منقسم به عنوان تابعی از شناسه‌ها و درونیابی بوسانی اش

تاکنون با تفاضلهای منقسم تنها در نقشی که به عنوان ضرایب صورت نیوتنی برای بسجمله‌ای درونیاب داشته‌اند، یعنی به عنوان ثابت‌هایی که به وسیله اعداد مفروض  $f(x_i)$ ،  $i = 0, \dots, n$ ، محاسبه می‌شوند، برخورد کرده‌ایم. اما پیدایش تابع  $g_n(x) = f[x_0, x_1, \dots, x_n, x]$  در عبارت خطای (۱۸.۲)، مربوط به درونیابی بسجمله‌ای (۱۸.۲)، مستلزم درک چگونگی رفتار تفاضل منقسم  $f[x_0, \dots, x_k]$  هنگام تغییر یک یا چندتا از نقاط  $x_0, \dots, x_k$  است. در ابتدا تعریف  $k$ امین تفاضل منقسم  $f[x_0, \dots, x_k]$  را برای کلیه انتخابهای  $x_0, \dots, x_k$  تعمیم می‌دهیم، یعنی از شرط دو به دو متمایز بودن نقاط  $x_0, \dots, x_k$  صرف نظر می‌کنیم. از آنجا که  $k$ امین تفاضل منقسم  $f[x_0, \dots, x_k]$  از  $f$  در نقاط  $x_0, \dots, x_k$  به عنوان ضریب پیشرو (یعنی ضریب  $x^k$ ) در بسجمله‌ای  $p_k(x)$ ، نایبتر از درجه  $k$ ، که با  $f(x)$  در  $k+1$  نقطه  $x_0, \dots, x_k$  منطبق است تعریف می‌شود. لذا باید منظور خود را از عبارت « $p_k(x)$  بر  $f(x)$  در نقاط  $x_0, \dots, x_k$  منطبق است»، وقتی برخی از این نقاط برهم منطبق‌اند، توضیح دهیم.

تعریف ما از این جمله چنین است. دو تابع  $f(x)$  و  $g(x)$  را در نقاط  $x_0, \dots, x_m$  برهم منطبق گوئیم هر گاه به ازای هر نقطه  $z$  که  $m$  بار در دنباله  $x_0, \dots, x_m$  آمده است، داشته باشیم

$$f^{(j)}(z) = g^{(j)}(z) \quad , j = 0, 1, \dots, m-1$$

در حقیقت  $f(x)$  و  $g(x)$  در نقاط  $x_0, \dots, x_m$  برهم منطبق‌اند هر گاه، با احتساب تکرار ریشه‌ها، تفاضل آنها ریشه‌های  $x_0, \dots, x_m$  داشته باشند (به بخش ۱.۲ نگاه کنید).

□ مثال:  $f(x)$  و  $g(x)$  در نقاط  $1, 2, 4, 5, 2, 4, 5, 1, 2$  در حالتی که

$$f(1) = g(1), \quad f(2) = g(2), \quad f'(2) = g'(2), \quad f''(2) = g''(2)$$

$$f(\varphi) = g(\varphi), \quad f'(\varphi) = g'(\varphi), \quad f(\delta) = g(\delta)$$

بنابراین تعریف، بسجمله‌ای تیلر

$$\hat{p}_n(x) = \sum_{j=0}^n f^{(j)}(c)(x-c)^j / j! \quad (31.2)$$

در نقطهٔ  $c$ ، مرتبه  $n+1$  بر  $f(x)$  منطبق است

$$(d/dx)^i (x-c)^j |_{x=c} = \begin{cases} j! & i=j \\ 0 & i \neq j \end{cases} \quad \text{چون}$$

بنابراین

$$\square \quad \hat{p}_n^{(i)}(c) = f^{(i)}(c) \quad i=0, \dots, n$$

زمانی از درونیابی بوسانی<sup>۱</sup> صحبت می‌کنیم که درجهٔ تماس بسجمله‌ای درونیاب با  $f(x)$  در یکی از نقاط درونیابی از مرتبهٔ اول<sup>۲</sup> بیشتر باشد.

صحبت از این بسجمله‌ای نایبتر از درجهٔ  $k$  که در  $k+1$  نقطه بر تابع  $f(x)$  منطبق باشد بجاست، زیرا که به موجب فرع ۲.۲ (در بخش ۱.۲) دو بسجمله‌ای نایبتر از درجهٔ  $k$  که در  $k+1$  نقطهٔ متمایز یا نامتمايز، ولی با احتساب تکرار) بر هم منطبق‌اند، بایستی یکی باشند. اگر یک چنین بسجمله‌ای درونیاب  $p_k(x)$  از درجهٔ نایبتر از  $k$  برای  $f(x)$  در نقاط  $x_0, \dots, x_k$  وجود داشته باشد، آنگاه بنا بر تعریف، ضریب پیشرو این بسجمله‌ای،  $k$ امین تفاضل منقسم  $f[x_0, \dots, x_k]$  خواهد شد. بنابراین بسجمله‌ای

$$p(x) = p_k(x) - f[x_0, \dots, x_k](x-x_0) \dots (x-x_{k-1})$$

نایبتر از درجهٔ  $k-1$  است. از آنجا که  $(x-x_0) \dots (x-x_{k-1})$  در  $x_0, \dots, x_{k-1}$  بر تابع صفر منطبق است، لذا نتیجه می‌شود که  $p(x)$  در نقاط  $x_0, \dots, x_k$  بر  $p_k(x)$  و در نتیجه بر  $f(x)$  منطبق است، یعنی  $p(x)$  می‌باید یک بسجمله‌ای از درجهٔ نایبتر از  $k-1$  باشد که در نقاط  $x_0, \dots, x_{k-1}$  بر  $f(x)$  منطبق است. بنابراین، استقرا روی  $n$  منجر به اثبات

$$p_n(x) = \sum_{j=0}^n (f[x_0, \dots, x_j] \prod_{i=0}^{j-1} (x-x_i)) \quad (32.2)$$

برای بسجمله‌ای نایبتر از درجهٔ  $n$  می‌شود که در  $x_0, \dots, x_n$  بر  $f(x)$  منطبق است. البته این فرمول همان فرمول (۱۰.۲) است، که تمامی موضوع این بخش را تشکیل می‌دهد.

در پایان، می‌بایستی مطمئن می‌شدیم که به‌ازای هر انتخابی از نقاط درونیایی  $x_0, \dots, x_k$  و تابع  $f(x)$ ، يك بسجمله‌ای نایبتر از درجه  $k$  وجود دارد که در این نقاط بر  $f(x)$  منطبق است ولی این امر را نمی‌توانیم تضمین کنیم، زیرا ممکن است  $f(x)$  به تعداد  $x_i$ هایی که بر هم منطبق می‌شوند، مشتق نداشته باشد. اما اگر  $f(x)$  به اندازه کافی دارای مشتق باشد، آنگاه به‌استقرار نسبت به  $k$  می‌توانیم وجود بسجمله‌ای درونیاب  $p_k(x)$  را ثابت کنیم و دوباره يك فرمول مفید [که اساساً همان فرمول ۱۲.۲ خواهد بود] برای تفاضل منقسم به‌دست آوریم.

قضیه ۴.۲ اگر  $f(x)$  دارای  $n$  مشتق پیوسته باشد و هیچ نقطه‌ای از دنباله  $x_0, \dots, x_n$  بیش از  $m+1$  مرتبه تکرار نشود، آنگاه دقیقاً يك بسجمله‌ای  $p_n(x)$  از درجه نایبتر از  $n$  وجود دارد که در نقاط  $x_0, \dots, x_n$  بر  $f(x)$  منطبق است. برای اثبات وجود چنین بسجمله‌ای، همچنین می‌توانیم فرض کنیم که دنباله نقاط درونیایی غیر نزولی است، یعنی

$$x_0 \leq \dots \leq x_n$$

به‌ازای  $n=0$  مسئله محرز است. فرض کنید که مطلب به‌ازای  $n=k-1$  صحیح باشد و صحت آن را به‌ازای  $n=k$  بررسی می‌کنیم، دو حالت وجود دارد.

حالت  $x_0 = x_k$ ، که در این حالت  $x_0 = \dots = x_k$ ، و بنا بر فرض بالا باید داشته باشیم  $m \geq k$ ، یعنی  $f(x)$  حداقل  $k$  مشتق پیوسته دارد. در این صورت، همان‌طور که قبلاً اشاره شد [به‌رابطه (۳۱.۲) نگاه کنید]، بسجمله‌ای تیلر  $\hat{p}_k(x)$  برای  $f(x)$  پیرامون مرکز  $c = x_0$  همان بسجمله‌ای مورد نظر است. باید توجه داشت که ضریب پیشرو آن عدد  $f^{(k)}(x_0)/k!$  است، بنابراین

$$f[x_0, \dots, x_k] = \frac{f^{(k)}(x_0)}{k!} \quad \text{اگر } x_0 = x_1 = \dots = x_k \quad (33.2)$$

حالت  $x_0 < x_k$ ، که در این حالت بنا بر فرض استقرار، می‌توانیم يك بسجمله‌ای  $p_{k-1}(x)$  از درجه نایبتر از  $k-1$  پیدا کنیم که در نقاط  $x_0, \dots, x_k$  بر  $f(x)$  منطبق باشد و يك بسجمله‌ای  $q_{k-1}(x)$  از درجه نایبتر از  $k-1$  که در  $x_0, \dots, x_{k-1}$  بر  $f(x)$  منطبق باشد. در این صورت بسجمله‌ای

$$p_k(x) = \frac{x-x_0}{x_k-x_0} q_{k-1}(x) + \frac{x_k-x}{x_k-x_0} p_{k-1}(x) \quad (34.2)$$

از درجه نایبتر از  $k$  خواهد بود، و اکنون ادعا می‌کنیم که این بسجمله‌ای همان بسجمله‌ای مطلوب است، یعنی  $p_k(x)$  در نقاط  $x_0, \dots, x_k$  بر  $f(x)$  منطبق است. داریم

$$p_k^{(j)}(x) = \frac{x-x_0}{x_k-x_0} q_{k-1}^{(j)}(x) + \frac{x_k-x}{x_k-x_0} p_{k-1}^{(j)}(x) + j \frac{q_{k-1}^{(j-1)}(x) - p_{k-1}^{(j-1)}(x)}{x_k-x_0}$$

(۳۵.۲)

فرض کنید  $z = x_i = \dots = x_{i+r}$ . اگر  $z = x_0$ ، آنگاه به‌ازای  $j = 0, \dots, r-1$  داریم

$$p_{k-1}^{(j)}(z) = f^{(j)}(z) \quad \text{و همچنین} \quad q_{k-1}^{(j)}(z) = f^{(j)}(z) = p_{k-1}^{(j)}(z) \quad (z < x_k \text{ زیرا که})$$

سپس با توجه به رابطه (۳۵.۲)، به‌ازای  $j = 0, \dots, r$  داریم

$$p_k^{(j)}(z) = 0 \times q_{k-1}^{(j)}(z) + f^{(j)}(z) + j [f^{(j-1)}(z) - f^{(j-1)}(z)] / (x_k - x_0) = f^{(j)}(z)$$

استدلال در حالت  $z = x_k$  نیز مشابه این است. بالاخره اگر  $x_k \neq x_0$ ، آنگاه به‌ازای  $j = 0, \dots, r$  داریم  $q_{k-1}^{(j)}(z) = f^{(j)}(z) = p_{k-1}^{(j)}(z)$  و نیز بنا بر (۳۵.۲) داریم

$$\begin{aligned} p_k^{(j)}(z) &= \frac{z - x_0}{x_k - x_0} f^{(j)}(z) + \frac{x_k - z}{x_k - x_0} f^{(j)}(z) + j \frac{f^{(j-1)}(z) - f^{(j-1)}(z)}{x_k - x_0} \\ &= f^{(j)}(z) \quad j = 0, \dots, r \end{aligned} \quad \text{به‌ازای}$$

که موضوع را برای  $n = k$  ثابت می‌کند.

از مقایسه ضرایب پیشرو در دو طرف رابطه (۳۴.۲) دوباره فرمول (۱۲.۲) را به‌دست خواهیم آورد، یعنی با شرط  $x_k \neq x_0$  داریم

$$f[x_0, \dots, x_k] = (f[x_1, \dots, x_k] - f[x_0, \dots, x_{k-1}]) / (x_k - x_0) \quad (۳۶.۲)$$

اگر  $x_k \neq x_0$

با تعمیم تعریف  $f[x_0, \dots, x_k]$  برای هر انتخاب دلخواه از  $x_0, \dots, x_k$  اکنون چگونگی بستگی  $f[x_0, \dots, x_k]$  را به نقاط  $x_0, \dots, x_k$  مورد بررسی قرار می‌دهیم. این بررسی آشکار می‌سازد که انگیزه تعمیم تعریف، ملاحظات پیوستگی بوده است. مسئله را با این ملاحظه آغاز می‌کنیم که  $f[x_0, \dots, x_k]$  نسبت به‌شناسه‌هایش یک تابع متقارن<sup>۱</sup> است، یعنی  $f[x_0, \dots, x_k]$  فقط به‌اعداد  $x_0, \dots, x_k$  بستگی دارد، نه به‌ترتیبی که این اعداد در فهرست شناسه‌ها ظاهر می‌شوند. این موضوع کاملاً واضح است، زیرا که کلیه بسجمله‌ایهای درونیاب  $p_k(x)$  مستقل از ترتیب نوشتن نقاط درونیابی هستند. این بدان معنی است که هر گاه مناسب باشد می‌توانیم بی‌آنکه خللی در استدلال وارد آید، فرض کنیم که ترتیب شناسه‌های  $x_0, \dots, x_k$  از  $f[x_0, \dots, x_k]$  صعودی است. در مرحله بعد نشان خواهیم داد که  $f[x_0, \dots, x_k]$  تابعی است پیوسته از شناسه‌هایش.

قضیه ۵.۲ فرض کنید که  $f(x)$ ،  $n$  بار به‌طور پیوسته روی  $[a, b]$  مشتقپذیر باشد و گیریم

## 1. Symmetric function

$y_n, \dots, y_0$  نقاط متمایز یا غیر متمایزی روی  $[a, b]$  باشند. در این صورت  
 (i) نقطه‌ای مانند  $\xi \in [\min_i y_i, \max_i y_i]$  وجود دارد به طوری که تساوی  
 $f[y_0, \dots, y_n] = f^{(n)}(\xi)/n!$  برقرار است.  
 (ii) اگر به ازای جميع مقادیر  $r, x_0^{(r)}, \dots, x_n^{(r)}$   $n+1$  نقطه روی  $[a, b]$  باشند

و به ازای  $i = 0, \dots, n$  رابطه  $\lim_{r \rightarrow \infty} x_i^{(r)} = y_i$  برقرار باشد. آنگاه داریم

$$\lim_{r \rightarrow \infty} f[x_0^{(r)}, \dots, x_n^{(r)}] = f[y_0, \dots, y_n]$$

اثبات به استقرا نسبت به  $n$  انجام می‌گیرد. به ازای  $n=0$  هر دو حکم به طور بدیهی صحیح اند. فرض کنید این احکام به ازای  $n=k-1$  نیز صحیح باشد و حالت  $n=k$  را بررسی کنید.

ابتدا به اثبات حالت (ii) در موردی که همه  $n+1$  نقطه  $y_0, \dots, y_n$  یکی نیستند می‌پردازیم. سپس با فرض  $y_0 \leq \dots \leq y_n$  داریم  $y_0 < y_n$  و بنابراین به ازای همه مقادیر بزرگ  $r$  رابطه  $x_0^{(r)} < x_n^{(r)}$  برقرار است و بنابراین طبق (۳۶.۲) داریم

$$\begin{aligned} \lim_{r \rightarrow \infty} f[x_0^{(r)}, \dots, x_n^{(r)}] &= \lim_{r \rightarrow \infty} \frac{f[x_1^{(r)}, \dots, x_n^{(r)}] - f[x_0^{(r)}, \dots, x_{n-1}^{(r)}]}{x_n^{(r)} - x_0^{(r)}} \\ &= \frac{\lim_{r \rightarrow \infty} f[x_1^{(r)}, \dots, x_n^{(r)}] - \lim_{r \rightarrow \infty} f[x_0^{(r)}, \dots, x_{n-1}^{(r)}]}{y_n - y_0} \\ &= \frac{f[y_1, \dots, y_n] - f[y_0, \dots, y_{n-1}]}{y_n - y_0} \end{aligned}$$

تساوی آخر از فرض استقرایی نتیجه می‌شود. اما رابطه آخر بنا بر (۳۶.۲) برابر  $f[y_0, \dots, y_n]$  است که در نتیجه برای این، حالت (ii) ثابت می‌شود.  
 سپس حالت (i) را ثابت می‌کنیم. اگر  $y_0 = y_1 = \dots = y_n$  آنگاه (i) فقط بیان دیگری از (۳۳.۲) است. در غیر این صورت فرض کنید که

$$y_0 \leq y_1 \leq \dots \leq y_n$$

و سپس  $y_0 < y_n$ . اما در این صورت خواهیم دید که برای کلیه  $r$ ها، می‌توان نقاط  $x_0^{(r)} < \dots < x_n^{(r)}$  را در  $[a, b]$  طوری پیدا کرد که به ازای  $i = 0, \dots, n$  تساوی  $\lim_{r \rightarrow \infty} x_i^{(r)} = y_i$  برقرار باشد. سپس، بنا بر قضیه ۲.۲ می‌توانیم  $\xi^{(r)} \in [x_0^{(r)}, x_n^{(r)}]$  را طوری پیدا کنیم که

$$f[x_0^{(r)}, \dots, x_n^{(r)}] = \frac{f^{(n)}(\xi^{(r)})}{n!} \quad r = 1, 2, 3, \dots$$

اما بنا بر حالت (ii) که در بالا برای آن ثابت شد، به ازای نقطه‌ای مانند

$$\xi \in [\lim x_0^{(r)}, \lim x_n^{(r)}] = [y_0, y_n]$$

داریم

$$\begin{aligned} f[y_0, \dots, y_n] &= \lim_{r \rightarrow \infty} f[x_0^{(r)}, \dots, x_n^{(r)}] = \frac{\lim_{r \rightarrow \infty} f^{(n)}(\xi^{(r)})}{n!} \\ &= \frac{f^{(n)}(\xi)}{n!} \end{aligned}$$

که به موجب پیوستگی  $f^{(n)}(x)$ ، (i) ثابت می‌شود.

و بالاخره اثبات (ii) در موردی که  $y_0 = y_1 = \dots = y_n$  اکنون از (i) برای اثبات وجود  $\xi^{(r)} \in [\min_i x_i^{(r)}, \max_i x_i^{(r)}]$  که به ازای کلیهٔ  $r$  ها در رابطهٔ  $f[x_0^{(r)}, \dots, x_n^{(r)}] = f^{(n)}(\xi^{(r)})/n!$  صادق باشد، استفاده می‌کنیم. در این صورت چون  $y_0 = \dots = y_n$  و به ازای کلیهٔ مقادیر  $i$ ، رابطهٔ  $\lim x_i^{(r)} = y_i$  برقرار است، داریم  $\lim \xi^{(r)} = y_0$  و بنا بر این طبق (۳۶.۲) و پیوستگی  $f^{(n)}(x)$  داریم

$$f[y_0, \dots, y_n] = \frac{f^{(n)}(y_0)}{n!} = \frac{\lim_{r \rightarrow \infty} f^{(n)}(\xi^{(r)})}{n!} = \lim_{r \rightarrow \infty} f[x_0^{(r)}, \dots, x_n^{(r)}]$$

که این رابطه هر دو حالت (i) و (ii) را به ازای  $n = k$  و برای کلیهٔ انتخابهای  $y_0, \dots, y_n$  در بازهٔ  $[a, b]$ ، ثابت می‌کند. این قسمت را با نتایج جالبی از قضیهٔ ۵.۲ خاتمه می‌دهیم. بلافاصله نتیجه‌گیری می‌شود که تابع

$$g_n(x) = f[x_0, \dots, x_n, x]$$

که در عبارت خطای مربوط به درونیابی بسجمله‌ای ظاهر می‌شود، برای کلیهٔ  $x$ ها تعریف شده و تابعی پیوسته از  $x$  است، اگر  $f(x)$  به اندازهٔ کافی هموار باشد. بنا بر این نتیجه می‌شود که

$$f(x) = \sum_{i=0}^n f[x_0, \dots, x_i] \prod_{j=0}^{i-1} (x - x_j) + f[x_0, \dots, x_n, x] \prod_{j=0}^n (x - x_j) \quad (37.2)$$

به ازای همهٔ  $x$ ها و نه فقط به ازای  $x \neq x_0, \dots, x_n$  [به ۱۶.۲ نگاه کنید] و نیز به ازای کلیهٔ نقاط متمایز یا نامتمايز  $x_0, \dots, x_n$ ، درحالی‌که  $f(x)$ ، به اندازهٔ کافی مشتق داشته باشد، برقرار است.

بعلاوه اگر  $f(x)$  به اندازهٔ کافی مشتقپذیر باشد، آنگاه  $g_n(x)$  مشتقپذیر خواهد بود. زیرا که طبق تعریف مشتق داریم

$$g'_n(x) = \lim_{h \rightarrow 0} g_n[x, x+h]$$

اگر این حد وجود داشته باشد. از طرف دیگر بنا بر قضیهٔ ۵.۲ داریم

$$g_n[x, x+h] = f[x_0, \dots, x_n, x, x+h] \xrightarrow{h \rightarrow 0} f[x_0, \dots, x_n, x, x]$$

بنابراین

$$\frac{d}{dx} f[x_0, \dots, x_n, x] = f[x_0, \dots, x_n, x, x] \quad (38.2)$$

در پایان، قضیهٔ بالا مبین تعریف درونیابی مماسی به صورت درونیابی همگرا است. زیرا نشان می‌دهد که بسجمله‌ای درونیاب در نقاط  $x_0, \dots, x_n$  همراه با  $y_i \rightarrow x_i$  برای کلیهٔ  $i$ ها به سمت بسجمله‌ای درونیاب در نقاط  $y_0, \dots, y_n$  همگرا می‌شود، بنا بر این  $k$  مرتبه درونیابی در یک نقطه، حد حالتی است که  $k$  نقطهٔ درونیابی متمایز یکی می‌شوند. دانشجویان با این پدیده در مورد  $m=1$ ، که همان درونیابی خطی است، آشنا هستند. در این حالت خط مستقیم  $p_1(x) = f(x_0) + f[x_0, x_1](x - x_0)$  قاطعی است برای (نمودار)  $f(x)$  که هر گاه دو نقطهٔ  $x_0$  و  $x_1$  به سمت نقطهٔ  $y$  نزدیک شوند، این قاطع به خط مماس  $\hat{p}_1(x) = f(y) + (x - y)f'(y)$  تبدیل می‌شود و مقدار  $\hat{p}_1(x)$  و ضریب زاویهٔ آن در نقطهٔ  $y = x$  با  $f(x)$  یکی می‌گردد.

□ مثال ۹.۲: با فرض  $f(x) = \ln x$ ، مطلوب است محاسبهٔ  $f(1.05)$  به توسط درونیابی تابع درجهٔ سوم، با استفاده از  $f(1) = 0.693147$ ،  $f(2) = 0.693147$ ،  $f'(1) = 1$ ،  $f'(2) = 0.5$ .

در این مورد چهار نقطهٔ درونیابی عبارت‌اند از:  $y_0 = y_1 = 1$ ،  $y_2 = y_3 = 2$ . محاسبات زیر را انجام می‌دهیم

$$f[y_0, y_1] = f'(y_0) = 1 \quad f[y_1, y_2] = 0.693147$$

$$f[y_2, y_3] = f'(y_2) = 0.5$$

$$f[y_0, y_1, y_2] = \frac{0.693147 - 1}{2 - 1} = -0.306853$$

$$f[y_1, y_2, y_3] = \frac{0.5 - 0.693147}{2 - 1} = -0.193147$$

$$f[y_0, y_1, y_2, y_3] = \frac{-0.193147 + 0.306853}{2-1} = 0.113706$$

جدول کامل تفاضل منقسم به صورت زیر نوشته می‌شود

$y_i$	$f[ ]$	$f[ , ]$	$f[ , , ]$	$f[ , , , ]$
۱	۰٫۰۰			
		۱٫۰۰		
۱	۰٫۰۰		-۰٫۳۰۶۸۵۳	
		۰٫۶۹۳۱۴۷		۰٫۱۱۳۷۰۶
۲	۰٫۶۹۳۱۴۷		-۰٫۱۹۳۱۴۷	
		۰٫۰۵		
۲	۰٫۶۹۳۱۴۷			

با استفاده از جدول فوق، تابع

$$p_3(x) = 0.00 + (1.00)(x-1) + (-0.306853)(x-1)^2 + (0.113706)(x-1)^2(x-2)$$

یک بسجمله‌ای درجهٔ سه است که از لحاظ مقدار و شیب با  $\ln x$  در دو نقطهٔ  $x=1$  و  $x=2$  یکی است. مشخصهٔ بوسانی تقریب  $\ln x$  به توسط  $p_3(x)$  از شکل ۷.۲ آشکار است. از به کار بستن الگوریتم ۱.۲ برای محاسبهٔ  $P_3(x)$  در ۱٫۵ خواهیم داشت

$$\ln 1.5 \approx p_3(1.5) = 0.409074$$

بر آورد زیر برای خطای  $e_3(x) = f(x) - p_3(x)$ ، با استفاده از (۳۷.۲) و حالت (i) قضیهٔ ۵.۲ به دست می‌آید

$$|e_3(1.5)| \leq \frac{1}{4!} \max_{1 \leq \xi \leq 2} |f^{(4)}(\xi)| (1.5-1)^2 (1.5-2)^2 = 0.015625$$

از آنجا که  $\ln 1.5 = 0.405465$ ، عملاً مقدار خطا فقط برابر با ۰٫۰۰۰۳۶۱ است. این مقدار بار دیگر نشان می‌دهد که نامعلوم بودن جای  $\xi$  خطا را بر اساس رابطهٔ (۱۸.۲) تقریباً محتاطانه - درست بگوئیم - بر آورد می‌کند.

□





### درونیایی به وسیلهٔ بسجمله ایها ۹۳

```

501 FORMAT (I3/2F10.3)
DO 30 J=1,NPOINT
  PNOFX = F(1)
  DO 29 I=2,NP1
    PNOFX = F(I) + (X - Y(I))*PNOFX
  29 CONTINUE
  PRINT 629,J,X,PNOFX
629 FORMAT (I10,2E20.7)
  X = X + DX
30 CONTINUE

```

STOP

END

اگر کلیهٔ نقاط درونیایی متمایز باشند، محاسبهٔ تفاضلهای منقسم، به الگوریتم ۳.۲ مربوط می‌شود. اگر بعضی از نقاط درونیایی برهم منطبق باشند، ورودی برنامه بایستی شامل مشتق تابع درونیاب باشد. به خصوص، فرض بر این است که ورودی از آرایه‌ای از نقاط درونیایی  $Y(I)$ ، به ازای  $I = 1, \dots, NP1 = n+1$  تشکیل شده است. برای ساده‌تر شدن برنامه نویسی فرض شده است که دنبالهٔ نقاط درونیایی در شرط زیر صدق می‌کند

$$Y(I) = Y(I+1) = \dots = Y(I+K) \text{ اگر } Y(I) = Y(I+K) \text{ آنگاه}$$

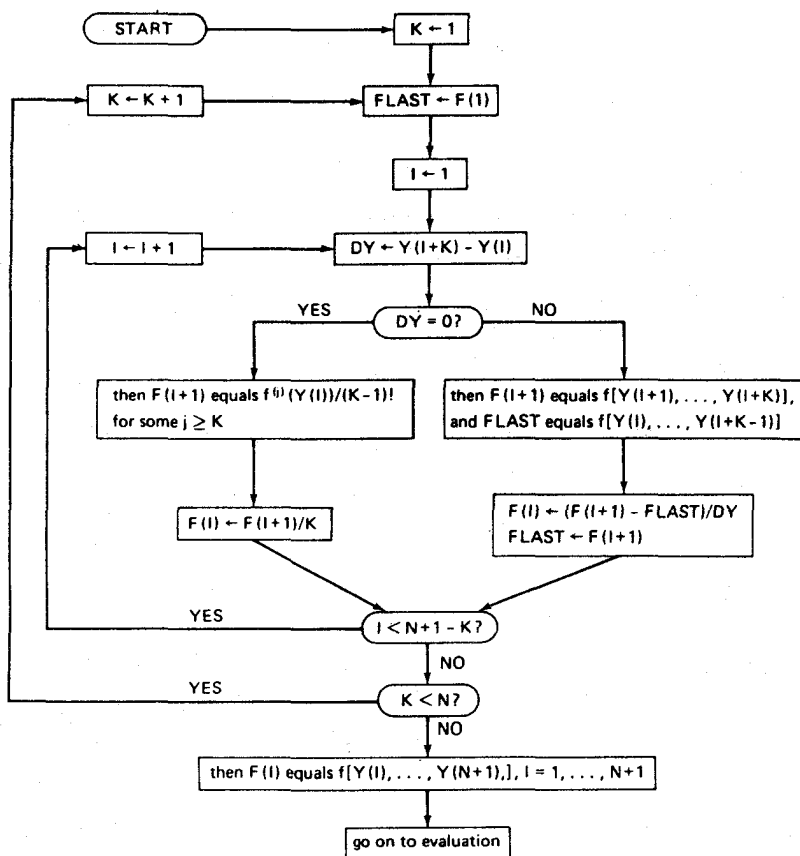
یعنی تمامی نقاط درونیایی تکراری با هم ظاهر می‌شوند. علاوه بر این شرط، فرض بر این است که به ازای هر  $I$

$$F(I) = f^{(j)}(Y(I)) \text{ آنگاه } Y(I) = Y(I-j) \neq Y(I-j-1) \text{ اگر}$$

بنا بر این با  $f(x) = 1/x$  و  $n=6$  و ورودیهای زیر صحیح خواهند بود، بدین معنی که یک بسجمله‌ای از درجهٔ نایبتر از ۶ تولید خواهد شد که  $f(x) = 1/x$  را در  $Y(I)$  داده شده به ازای  $I = 1, \dots, 7$ ، درونیایی می‌کند

I	۱	۲	۳	۴	۵	۶	۷
Y(I)	۲۰۰	۲۰۰	۲۰۰	۱۰۰	۴۰۰	۴۰۰	۵۰۰
F(I)	۰۰۵	-۰۰۲۵	۰۰۲۵	۱۰۰	۰۰۲۵	-۰۰۰۶۲۵	۰۰۲

دانشجویان را تشویق می‌کنیم که مثال مشابهی انتخاب و محاسبات را از راه برنامهٔ فورترن دنبال کنند. روند نمای<sup>۲</sup> زیر که معرف محاسبات تفاضلهای منقسم است ممکن است در مساعی فوق مفید واقع شود.



### تمرین

۲-۷۰۲ با استفاده از درونیاپی درجه دوم و مقادیر داده شده  $f(0) = 1$  و  $f'(0) = 1$  مقدار  $f(1) = 2.7183$ ، مقدار  $f(0.5)$  را برای تابع  $f(x) = e^x$  محاسبه کنید. این مقدار را با  $f(0.5) = 1.6487$  که به طور صحیح گرد شده است، مقایسه کنید.

۲-۷۰۲ مقادیر

$$f(0) = 0 \quad f'(0) = 1 \quad f(1) = 1.1752 \quad f'(1) = 1.8543$$

برای تابع  $f(x) = \sinh x$  داده شده اند. جدول تفاضل منقسم را تشکیل دهید و با استفاده از درونیاپی تابع درجه سوم، مقدار  $f(0.5)$  را محاسبه نمایید و آن را با  $\sinh 0.5 = 0.5211$  مقایسه کنید.

۳-۷۰۲ تابع  $f(x)$  یک ریشه دوگانه در  $x_1$  و یک ریشه سه گانه در  $x_2$  دارد. شکل بسجمله ای

نابیشتر از درجهٔ پنجمی را که  $f(x)$  را دوبار در  $z_1$  و سه بار در  $z_2$  و یک بار در نقطهٔ غیر مشخص  $z_3$  درونیایی کند، تعیین نماید.

۴-۷۰۲ ضرایب  $a_0, a_1, a_2, a_3$  از بسجمله‌ای درجهٔ سه

$$p_3(x) = a_0 + a_1(x-y) + a_2(x-y)^2 + a_3(x-y)^3$$

را چنان پیدا کنید که داشته باشیم

$$p_3(y) = f_y \quad p'_3(y) = f'_y \quad p_3(z) = f_z \quad p'_3(z) = f'_z$$

در اینجا،  $y, z, f_y, f'_y, f_z, f'_z$  اعداد مفروضی هستند ( $z \neq y$ ).

۵-۷۰۲ عبارت ساده‌ای برای  $p_3[(y+z)/2]$  بر حسب اعداد مفروض  $y, z, f_y, f'_y, f_z, f'_z$  به دست آورید، در اینجا  $p_3(x)$  بسجمله‌ای حاصل در تمرین ۴-۷۰۲ است.

۶-۷۰۲ گیریم  $f(x)$  و  $g(x)$  توابع هموار باشند. ثابت کنید که تابع  $f(x)$  در نقطهٔ  $x=c$   $k$  مرتبه بر  $g(x)$  منطبق است، اگر و فقط اگر، به ازای  $x$  نزدیک به  $c$  داشته باشیم

$$f(x) = g(x) + O(|x-c|^k)$$

۷-۷۰۲ گیریم  $g(x) = f[x_0, \dots, x_k, x]$ . ثابت کنید (با استفاده از استقرا) که

$$g[y_0, \dots, y_n] = f[x_0, \dots, x_k, y_0, \dots, y_n]$$

۸-۷۰۲ با استفاده از تمرین ۷-۷۰۲، ثابت کنید که اگر  $g(x) = f[x_0, \dots, x_k, x]$  آنگاه

$$g^{(n)}(x) = n! f[x_0, \dots, x_k, \overbrace{x, \dots, x}^{(n+1) \text{ مرتبه}}]$$

۹-۷۰۲ گیریم  $f(x) = g(x)h(x)$ . ثابت کنید که

$$f[x_0, \dots, x_k] = \sum_{i=0}^k g[x_0, \dots, x_i] h[x_i, \dots, x_k]$$

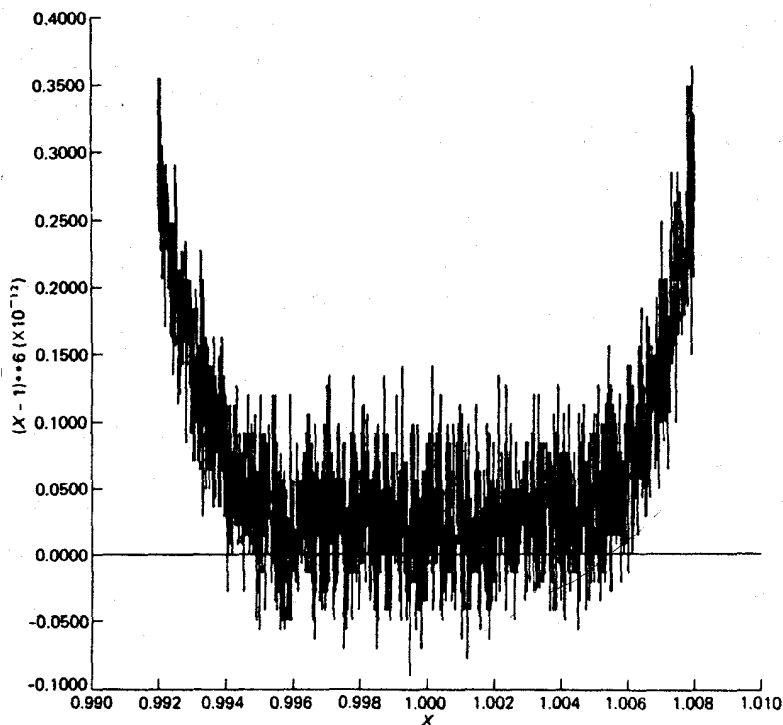
(از روش استقرا استفاده کنید، یا سمت راست را به صورت ضرب پیشرو بسجمله‌ای از درجهٔ نابیشتر از  $k$  که در نقاط  $x_0, \dots, x_k$ ،  $g(x)h(x)$  را درونیایی می‌کند، مشخص سازید.) این حالت  $x_0 = \dots = x_k$ ، کدام فرمول معروف دیفرانسیل و انتگرال را به دست می‌آورد؟

## حل معادلات غیر خطی

یکی از مسائلی که غالباً در کارهای عملی پیش می‌آید پیدا کردن ریشه‌های معادلاتی به صورت

$$f(x) = 0 \quad (1.3)$$

یعنی، صفرهای تابع  $f(x)$  است. تابع  $f(x)$  ممکن است صریحاً، مثلاً به صورت یک بسجمله‌ای بر حسب  $x$  یا به صورت یک تابع متعالی<sup>۱</sup> داده شده باشد. اما اغلب ممکن است  $f(x)$  فقط به صورت ضمنی داده شده باشد، بدین معنی که ممکن است قاعده‌ای برای ارزیابی  $f(x)$  به ازای هر شناسه در دست باشد، ولی صورت ساده آن معلوم نباشد. از این رو  $f(x)$  ممکن است معرف مقادیر جواب یک معادله دیفرانسیل در یک نقطه مشخص باشد، در حالی که  $x$  ممکن است بیانگر یک شرط اولیه آن معادله باشد. در موارد نادر ممکن است بتوان ریشه‌های دقیق معادله (۱.۳) را پیدا کرد، که یک مثال آن بسجمله‌ای‌هایی هستند که به صورت حاصلضرب چند عامل درمی‌آیند. لیکن در حالت کلی می‌توانیم امیدوار باشیم که، با اتکا به تکنیکهای محاسباتی، فقط جوابهای تقریبی را به دست آوریم. بسته به اینکه زمینه کار چه باشد، «جواب تقریبی» ممکن است بدین معنی باشد که یا  $x^*$  در معادله (۱.۳) به طور تقریبی صدق می‌کند یعنی،  $|f(x^*)|$  «بسیار کوچک» است، و یا نقطه  $x^*$  به جواب (۱.۳) «نزدیک» است. متأسفانه مفهوم یک «جواب تقریبی» تقریباً نامعلوم است. یک جواب تقریبی که از یک کامپیوتر به دست



شکل ۱۰۳

آمده است، به علت گرد کردن یا ناپایداری وی با به سبب حساب خاصی که به کار گرفته می شود، تقریباً همیشه دارای خطاست. در حقیقت ممکن است «جوابهای تقریبی» بسیاری وجود داشته باشند که آنها هم صحیح باشند ولو اینکه جواب مطلوب یگانه باشد. به منظور نشان دادن عدم اطمینانی که در پیدا کردن ریشه ها وجود دارد، در شکل ۱۰۳ نمودار تابع

$$P_6(x) = (1-x)^6 = 1 - 6x + 15x^2 - 20x^3 + 15x^4 - 6x^5 + x^6$$

را می آوریم. البته این تابع تنها دارای ریشه  $x=1$  است. یک برنامه فورترن برای ارزیابی منبسط این تابع نوشته شده است. این برنامه برای محاسبه مقدار  $P_6(x)$  در تعداد زیادی نقاط  $x_1 < x_2 < \dots < x_N$  نزدیک به  $x=1$  روی یک کامپیوتر CDC 6500 مورد استفاده قرار گرفته است. سپس یک رسم کامپیوتری برای رسم نمودار پاره خطهای تکه ای، که در شکل ۱۰۳ مشاهده می شود، به کار گرفته شده است. در این نمودار دیده می شود که  $P_6(x)$  آشکارا دارای صفرهای زیادی است زیرا که این نمودار چندین بار تغییر علامت

نشان می‌دهد. این صفرهای آشکار در محدودهٔ  $۰.۹۹۴$  تا  $۱.۰۰۶$  هستند. بنا بر این به کار گرفتن شکل منبسط برای برآورد صفر در  $x=1$ ، آشکارا به برآوردهای قابل قبولی منجر می‌شود که فقط تا دو رقم اعشاری صحیح‌اند، درحالی‌که کامپیوتر CDC ۶۵۰۰ در حساب با ۱۴ رقم ممیز شناور کار کرده‌است. دلیل این رفتار را باید در خطای گرد کردن، از دست رفتن و حذف ارقام با معنی در برنامهٔ فورترن هنگام محاسبهٔ  $P_6(x)$  جستجو کرد. این مثال برخی از مخاطرات موجود در ریشه‌یابی را نشان می‌دهد.

در بقیهٔ این فصل روشهای بارستی گوناگونی را برای پیدا کردن تقریبهای ریشه‌های سادهٔ معادلهٔ (۱.۳) مورد بررسی قرار می‌دهیم. به علت اهمیت کار بردهای معادلات بسجمله‌ای در مهندسی، در این فصل توجه خاصی به آنها خواهیم داشت.

### ۱.۳ بررسی اجمالی روشهای بارستی

در این بخش برخی از روشهای مقدماتی برای تعیین يك جواب معادله

$$f(x) = 0 \quad (1.3)$$

را مطرح می‌کنیم و مورد استفاده آنها را با به کار بردن این روشها در حل معادلهٔ بسجمله‌ای سادهٔ

$$x^2 - x - 1 = 0 \quad (2.3)$$

نشان می‌دهیم. در اینجا داریم  $f(x) = x^2 - x - 1$ . برای این مثال پیدا می‌کنیم که

$$f(1) = -1 < 0 < 5 = f(2) \quad (3.3)$$

بنا بر این، از آنجا که  $f(x)$  پیوسته‌است، با در نظر گرفتن قضیهٔ مقدار میانگین در توابع پیوسته (به بخش ۷.۱ نگاه کنید)،  $f(x)$  باید در بازهٔ  $[1, 2]$  در نقطه‌ای صفر شود. هر گاه  $f(x)$  در دو نقطه یا بیشتر در بازهٔ  $[1, 2]$  صفر شود، با توجه به قضیهٔ ژول (به بخش ۷.۱ نگاه کنید)، می‌بایستی  $f'(x)$  نیز در بازهٔ  $[1, 2]$  صفر شود. اما از آنجا که  $f'(x) = 3x^2 - 1$  روی  $[1, 2]$  مثبت است، بنا بر این  $f(x)$  در بازهٔ  $[1, 2]$  دقیقاً يك ریشه دارد. اگر این ریشه را  $\xi$  بنامیم، داریم:

$$0.8 \leq \xi \leq 1.8 \quad \text{با خطای مطلق}$$

برای شناخت بیشتر این ریشه،  $f(x)$  را در نقطهٔ میانی بازهٔ  $[1, 2]$  یعنی در نقطهٔ ۱.۵ ارزیابی می‌کنیم. داریم:

$$f(1.5) = 0.875 > 0 > -1 = f(1)$$

در نتیجه اکنون می‌دانیم که ریشهٔ  $\xi$  در بازهٔ کوچکتر  $[1, 1.5]$  یعنی

$$\xi = ۱۰۲۵ \quad \text{با خطای مطلق} \leq ۰.۰۲۵$$

قرار دارد. دوباره  $f(x)$  را در نقطه میانی ۱۰۲۵ آزمایش می‌کنیم. داریم

$$f(۱۰۲۵) = -۰.۰۲۹۶\dots < 0 < ۰.۰۸۷۵ = f(۱۰۵)$$

بنابراین می‌دانیم که  $\xi$  حتی در بازه کوچکتر  $[۱۰۲۵, ۱۰۵]$  قرار دارد یعنی

$$\xi \approx ۱۰۳۷۵ \quad \text{با خطای مطلق} \leq ۰.۰۱۲۵$$

این شیوه جایابی یک ریشه از معادله  $f(x) = 0$  در یک دنباله از بازه‌ها، که اندازه آنها سیر نزولی دارد، به روش تصنیف معروف است.

**الگوریتم ۱۰.۳: روش تصنیف.** فرض می‌کنیم تابع  $f(x)$  در بازه  $[a_0, b_0]$  پیوسته باشد به طوری که  $f(a_0)f(b_0) < 0$ .

For  $n = 0, 1, 2, \dots$ , until satisfied, do:

$$\text{Set } m = (a_n + b_n)/2$$

$$\text{If } f(a_n)f(m) \leq 0, \text{ set } a_{n+1} = a_n, b_{n+1} = m$$

$$\text{Otherwise, set } a_{n+1} = m, b_{n+1} = b_n$$

Then  $f(x)$  has a zero in the interval  $[a_{n+1}, b_{n+1}]$

از این به بعد غالب الگوریتمها را به صورت فشرده فوق بیان خواهیم کرد. برای دانشجویانی که با زبان الگوریتم آشنا نیستند، نشانگذاریهای فوق باید کاملاً طبیعی باشد. علاوه بر الگوریتم بالا، اصطلاح «until satisfied» به کار برده شده است تا تأکید شود که این شرح الگوریتم کامل نیست، و استفاده کننده از الگوریتم بایستی ملاک دقیق ختم آن را خود تعیین کند. بخشی از این ملاکها به مسئله خاصی که با الگوریتم فوق حل می‌شود بستگی دارد. بعضی از ملاکهای اختتامی ممکن در بخش بعد مورد بحث قرار خواهند گرفت.

در هر مرحله از الگوریتم تصنیف ۱۰.۳، طول بازه‌ای که می‌دانیم شامل یک صفر  $f(x)$  است، نصف می‌شود. بنابراین هر مرحله یک رقم دودویی صحیح جدید دیگری برای ریشه  $\xi$  از  $f(x) = 0$  تولید می‌کند. در مثال بالا با همان مقادیر اولیه  $a_0 = ۱$  و  $b_0 = ۲$  و پس از اجرای ۲۰ مرحله از این الگوریتم خواهیم داشت:

$$۱۰۳۲۴۷۱۷۵\dots = a_{20} \leq \xi \leq b_{20} = ۱۰۳۲۴۷۱۸۴\dots$$

$$f(a_{20}) = (-۱۰۸۵۷\dots) \cdot 10^{-6} < 0 < (۲۰۲۰۹\dots) \cdot 10^{-6} = f(b_{20})$$

آشکار است که با این الگوریتم و کوشش کافی همیشه می‌توان یک ریشه را تا دقت عمل مورد



نظری جایابی کرد. اما در مقایسه با روشهای دیگری که بحث خواهد شد، روش تصنیف دارای سرعت همگرایی نسبتاً اندکی است.

می توان امیدوار بود که با به کار گرفتن کاملتر اطلاعاتی که در هر مرحله از  $f(x)$  در دست است، سریعتر به ریشه رسید. در مثال (۲.۳) با اطلاعات زیر شروع کردیم:

$$f(1) = -1 < 0 < 5 = f(2)$$

از آنجا که  $|f(1)|$  نسبت به  $|f(2)|$  نزدیکتر به صفر است، ریشه  $\xi$  احتمالاً به ۱ نزدیکتر است تا ۲ (دست کم اگر  $f(x)$  «تقریباً» خطی باشد). بنابراین به جای اینکه نقطه میانی یا متوسط ۱ و ۲ را که نقطه ۱.۵ است آزمایش کنیم، اکنون  $f(x)$  را در میانگین وزنی یعنی

$$w = \frac{|f(2)| \times 1 + |f(1)| \times 2}{|f(2)| + |f(1)|} \quad (4.3)$$

آزمایش می کنیم. از آنجا که علامتهای  $f(1)$  و  $f(2)$  مختلف هستند، می توانیم رابطه (۴.۳) را به صورت ساده تر زیر بنویسیم:

$$w = \frac{f(2) \times 1 - f(1) \times 2}{f(2) - f(1)} \quad (5.3)$$

و برای مثال فوق جواب

$$w = \frac{5 \times 1 + 1 \times 2}{6} = 1.166666666\dots$$

و

$$f(w) = -0.578703\dots < 0 < 5 = f(2)$$

به دست می آید. از این رو  $\xi$  در بازه  $[1.166666666\dots, 2]$  قرار دارد. با تکرار روند فوق برای این بازه خواهیم داشت:

$$w = \frac{5 \times (1.166666666\dots) + (0.578703\dots) \times 2}{5.578703\dots} = 1.253112\dots$$

$$f(w) = -0.285363\dots < 0 < 5 = f(2)$$

در نتیجه  $f(x)$  در بازه  $[1.253112\dots, 2]$  یک ریشه دارد. این الگوریتم به روش تصحیح خطا معروف است.

الگوریتم ۲.۳: تصحیح خطا. تابع  $f(x)$  که در بازه  $[a_0, b_0]$  پیوسته است و  $f(a_0)f(b_0) < 0$  داده شده است.

For  $n = 0, 1, 2, \dots$ , until satisfied, do:

$$\text{Calculate } w = [f(b_n)a_n - f(a_n)b_n] / [f(b_n) - f(a_n)].$$

$$\text{If } f(a_n)f(w) \leq 0, \text{ set } a_{n+1} = a_n, b_{n+1} = w$$

$$\text{Otherwise, set } a_{n+1} = w, b_{n+1} = b_n$$

بعد از ۱۶ بار استفاده از این الگوریتم در مثال قبلی و با همان نقاط شروع  $a_0 = 1$  و  $b_0 = 2$  خواهیم داشت:

$$1.32247174\dots = a_{16} \leq \xi \leq b_{16} = 2$$

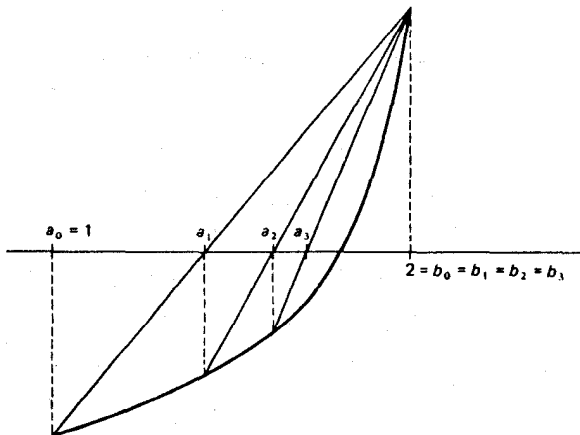
$$f(a_{16}) = (-1.95\dots)10^{-6} < 0 < 5 = f(b_{16})$$

بنابراین، گرچه روش تصحیح خطا تا حدی سریعتر از روش تصفیه نقطه‌ای راکه در آن نقطه مقدار  $|f(x)|$  «کوچک» است به دست می‌دهد، اما به طور قطع نمی‌تواند بازه «کوچکی» راکه بدانیم ریشه در آن وجود دارد پدید آورد.

نگاهی اجمالی به شکل ۲.۳ دلیل این امر را نشان می‌دهد. همان‌طور که به آسانی دیده می‌شود، میانگین وزنی

$$w = \frac{f(b_n)a_n - f(a_n)b_n}{f(b_n) - f(a_n)}$$

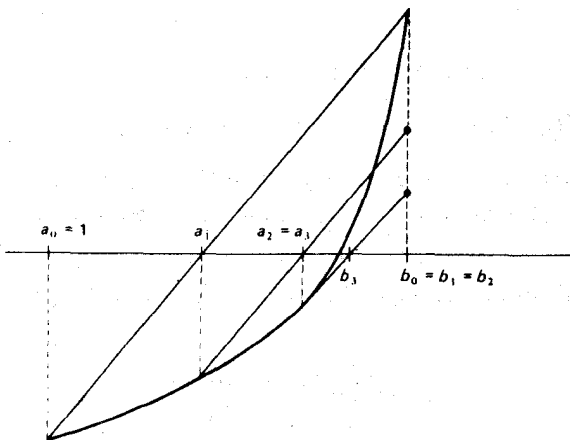
نقطه تقاطع خط مستقیم ماربر نقاط  $\{a_n, f(a_n)\}$  و  $\{b_n, f(b_n)\}$  است با محور  $x$ ها. این خط مستقیم قاطع  $f(x)$  است، و در مثال فوق  $f(x)$  در بازه مورد نظر  $[1, 2]$  صعودی و



شکل ۲.۳ موقعیت اشتباهی.

دارای تقرری به سمت بالاست. بنابراین قاطع همیشه در بالای نمودار  $f(x)$  قرار دارد. در نتیجه، در این مثال  $w$  همیشه در طرف چپ صفر قرار دارد. اگر  $f(x)$  صعودی و دارای تقرری به سمت پایین می‌بود،  $x$  همیشه در سمت راست ریشه قرار می‌گرفت.

الگوریتم تصحیح خطا را می‌توان به چندین طریق بهبود بخشید که دو طریقۀ از آنها را ذیلاً مورد بحث قرار می‌دهیم. طریقۀ اول «تصحیح خطای اصلاح شده» نامیده می‌شود که در آن خطهای مستقیمی که شیب آنها پیوسته کمتر می‌شود جایگزین قاطع می‌شوند تا جایی که  $w$  در طرف دیگر ریشه قرار می‌گیرد. این نکته در نمودار ۳.۳ نشان داده شده است.



شکل ۳.۳ تصحیح خطای اصلاح شده.

الگوریتم ۳.۳: تصحیح خطای اصلاح شده. فرض می‌کنیم  $f(x)$  روی  $[a_0, b_0]$  پیوسته باشد و  $f(a_0)f(b_0) < 0$ ، پس

Set  $F = f(a_0)$ ,  $G = f(b_0)$ ,  $w_0 = a_0$

For  $n = 0, 1, 2, \dots$ , until satisfied, do:

Calculate  $w_{n+1} = (Ga_n - Fb_n)/(G - F)$

If  $f(a_n)f(w_{n+1}) \leq 0$ , set  $a_{n+1} = a_n$ ,  $b_{n+1} = w_{n+1}$ ,  $G = f(w_{n+1})$

If also  $f(w_n)f(w_{n+1}) > 0$ , set  $F = F/2$

Otherwise, set  $a_{n+1} = w_{n+1}$ ,  $F = f(w_{n+1})$ ,  $b_{n+1} = b_n$

If also  $f(w_n)f(w_{n+1}) > 0$ , set  $G = G/2$

Then  $f(x)$  has a zero in the interval  $[a_{n+1}, b_{n+1}]$

اگر این الگوریتم برای مثال قبلی با  $a_0 = 1$  و  $b_0 = 2$  به کار گرفته شود، بعد از ۶ مرحله جواب زیر به دست می آید:

$$1.322471795 \dots = a_6 \leq \xi \leq b_6 = 1.322471796 \dots$$

$$f(a_6) = (-1.736 \dots) 10^{-8} < 0 < (1.730 \dots) 10^{-8} = f(b_6)$$

که بهبود قابل توجهی نسبت به روش تنصیف نشان می دهد.

یک طریقه دیگر (روش دوم) تصحیح خطا، که بسیار متداول است، روش «خط قاطع» نامیده می شود، که در عمل همه جا از خطوط قاطع استفاده، ولی از تخصیص حاشیه به ریشه، هرگز فراموش نمی کنند.

**الگوریتم ۴.۳:** روش خط قاطع. تابع  $f(x)$  و دو نقطه  $x_0$ ،  $x_{-1}$  داده شده اند

For  $n = 0, 1, 2, \dots$ , until satisfied, do:

$$\left[ \text{Calculate } x_{n+1} = [f(x_n)x_{n-1} - f(x_{n-1})x_n] / [f(x_n) - f(x_{n-1})] \right]$$

اگر دومین روش اصلاح شده را برای مثال قبلی با  $x_0 = 2$  و  $x_{-1} = 1$  به کار ببریم، پس از ۶ مرحله خواهیم داشت:

$$x_6 = 1.32247179 \dots, \quad f(x_6) = (3.458 \dots) 10^{-8}$$

ظاهراً روش خط قاطع سریعاً نقطه ای را که در آن  $|f(x)|$  «کوچک» است جایابی می کند اما در حالت کلی هیچگونه احساسی نسبت به فاصله این نقطه از ریشه  $f(x)$  به شخص نمی دهد. همچنین از آنجا که مختلف علامه بودن توابع  $f(x_n)$  و  $f(x_{n-1})$  الزامی نیست، لذا عبارت

$$x_{n+1} = \frac{f(x_n)x_{n-1} - f(x_{n-1})x_n}{f(x_n) - f(x_{n-1})} \quad (6.3)$$

مهمی نتایج خطای گرد کردن است. در موارد استثنایی ممکن است تساوی  $f(x_n) = f(x_{n-1})$  برقرار باشد، که در این صورت محاسبه  $x_{n+1}$  غیر ممکن می شود. از این رو - اگر چه این کار نیز چاره ساز نیست - برای محاسبه  $x_{n+1}$  بهتر است از عبارت معادل زیر استفاده شود

$$x_{n+1} = x_n - f(x_n) \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} \quad (7.3)$$

که طبق این فرمول  $x_{n+1}$  با افزودن «جمله تصحیحی» زیر از  $x_n$  به دست می آید:

$$\frac{-f(x_n)}{[f(x_n) - f(x_{n-1})]/(x_n - x_{n-1})} \quad (۸.۳)$$

دانشجویان به یاد دارند که نسبت  $[f(x_n) - f(x_{n-1})]/(x_n - x_{n-1})$  همان اولین تفاضل منقسم  $f(x)$  است، و طبق (۱۵.۲) این نسبت برابر است با شیب خط قاطع ماربر نقاط  $\{x_{n-1}, f(x_{n-1})\}$  و  $\{x_n, f(x_n)\}$  برای خم  $f(x)$ . بعلاوه بنا بر (۱۷.۲) ملاحظه می کنیم که اگر  $f(x)$  مشتقپذیر باشد، نسبت فوق برابر با شیب  $f(x)$  در نقطه ای بین  $x_{n-1}$  و  $x_n$  خواهد بود. بنا بر این بجاست که در نقطه ای «نزدیک» به  $x_{n-1}$  و  $x_n$  به جای نسبت فوق،  $f'(x)$  گذاشته شود البته در صورتی که  $f'(x)$  قابل محاسبه باشد.

اگر  $f(x)$  مشتقپذیر باشد، در این صورت اگر در رابطه (۷.۳) به جای شیب قاطع، شیب مماس بر منحنی در نقطه  $x_n$  راقرار دهیم. فرمول بارستی روش نیوتن به دست می آید:

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \quad (۹.۳)$$

**الگوریتم ۵.۳: روش نیوتن.** تابع پیوسته-مشتقپذیر  $f(x)$  و نقطه  $x_0$  داده شده اند.

For  $n = 0, 1, 2, \dots$ , until satisfied, do:

$$\left[ \text{Calculate } x_{n+1} = x_n - f(x_n)/f'(x_n) \right.$$

اگر این الگوریتم برای مثال قبلی با  $x_0 = 1$  به کار برده شود، بعد از چهار مرحله خواهیم داشت:

$$x_4 = 1.3247181\dots, \quad f(x_4) = (9.24\dots)10^{-7}$$

سرانجام به روش «بارستی نقطه ثابت» که مورد خاصی از روش نیوتن است اشاره می کنیم. اگر قرار دهیم

$$g(x) = x - \frac{f(x)}{f'(x)} \quad (۱۰.۳)$$

آنگاه فرمول بارستی روش نیوتن (۹.۳) شکل ساده زیر را پیدا خواهد کرد:

$$x_{n+1} = g(x_n) \quad (۱۱.۳)$$

اگر دنباله  $x_1, x_2, x_3, \dots$  که این گونه تولید می شود به سمت نقطه  $\xi$  همگرا شود و  $g(x)$  پیوسته باشد، آنگاه خواهیم داشت:

$$\xi = \lim_{n \rightarrow \infty} x_{n+1} = \lim_{n \rightarrow \infty} g(x_n) = g(\lim_{n \rightarrow \infty} x_n) = g(\xi) \quad (۱۲.۳)$$

و یا  $\xi = g(\xi)$ ، یعنی  $\xi$  نقطه ثابتی است از  $g(x)$ . واضح است که اگر  $\xi$  یک نقطه ثابت تابع

بارستی  $g(x)$  برای روش نیوتن باشد، آنگاه یک ریشه معادله  $f(x) = 0$  نیز خواهد بود. اما برای معادله داده شده  $f(x) = 0$  این امکان وجود دارد که توابع بارستی مختلف،  $g(x)$  با این ویژگی انتخاب شوند که هر نقطه ثابت  $g(x)$  یک صفر  $f(x)$  باشد. برای هر چنین انتخابی می‌بایستی دنباله  $x_1, x_2, x_3, \dots$  به وسیله فرمول

$$x_{n+1} = g(x_n) \quad n = 0, 1, 2, \dots$$

محاسبه شود و امیدوار بود که این دنباله همگرا باشد. اگر چنین باشد، حد این دنباله ریشه معادله  $f(x) = 0$  خواهد بود. در بخشهای ۳.۳ و ۴.۳ بارست نقطه ثابت را به طور مشروح مورد بحث قرار می‌دهیم.

□ مثال ۱.۳: تابع  $f(x) = x - 0.2 \sin x - 0.5$  دقیقاً یک ریشه بین  $x_0 = 0$  و  $x_1 = 1$  دارد، زیرا که رابطه  $0 < f(1) f(0)$  برقرار است در حالی که  $f'(x)$  بر بازه  $[0, 1]$  صفر نمی‌شود. با استفاده از الگوریتمهای ۱.۳، ۳.۳، ۴.۳ و ۵.۳ ریشه معادله فوق را تا ۶ رقم با معنی تعیین کنید.

محاسبات زیر با کامپیوتر IBM ۷۰۹۴، ۲۷ بیت، با ممیز شناور و دقت ساده انجام گرفته است.

n	Algorithm 3.1		Algorithm 3.3		Algorithm 3.4	Algorithm 3.5
	$x_n$	$\epsilon_n$	$x_n$	$\epsilon_n$	$x_n$	$x_n$
-1					1.	
0	0.75	$3 \cdot 10^{-1}$	0.75	$3 \cdot 10^{-1}$	0.5	0.5
1	0.625	$2 \cdot 10^{-1}$	0.80606124	$2 \cdot 10^{-1}$	0.61212248	0.61629718
2	0.5625	$6 \cdot 10^{-2}$	0.61534080	$3 \cdot 10^{-3}$	0.61549349	0.61546820
3	0.59375	$3 \cdot 10^{-2}$	0.61701328	$2 \cdot 10^{-3}$	0.61546816	0.61546816
4	0.609375	$2 \cdot 10^{-2}$	0.61701363	$2 \cdot 10^{-3}$		
5	0.6171875	$8 \cdot 10^{-3}$	0.61546816	0		
6	0.61328125	$4 \cdot 10^{-3}$				
...	.....	.....				
10	0.61547852	$4 \cdot 10^{-4}$				
...	.....	.....				
19	0.61546850	$5 \cdot 10^{-7}$				

در الگوریتم ۱.۳ و ۳.۳،  $x_n$  نقطه وسط بین کرانه‌های پایینی و بالایی  $a_n$  و  $b_n$ ، بعد از  $n$  بارست است، در حالی که  $\epsilon_n$  معرف کران مربوط به خطا در  $x_n$  است که به وسیله الگوریتم محاسبه شده است. همگرایی سریع و منظم الگوریتمهای ۴.۳ و ۵.۳ قابل توجه است. روش تصحیح خیلی به کندی اما به طور یکنوا همگرا می‌شود. در حالی که ظاهراً روش تصحیح خطای اصلاح شده با «پرش» همگرا می‌شود، اگرچه جواب صحیح را نسبتاً سریع به دست می‌دهد. □

## تمرین

۱-۱۰۳ مطلوب بازه‌ای است که ریشه حقیقی و مثبت معادله  $f(x) = x^2 - 2x - 2$  در آن واقع باشد. سپس با استفاده از الگوریتمهای ۱.۳ و ۲.۳ این ریشه را تا دو رقم با معنی صحیح محاسبه کنید. آیا می‌توانید برآورد کنید که در هر روش چند مرحله لازم است تا ۶ رقم با معنی صحیح نتیجه شود؟

۲-۱۰۳ برای مثالی که در متن داده شده است دومر حله از روش تصحیح خطای اصلاح شده، (الگوریتم ۳.۳) را انجام دهید.

۳-۱۰۳ بسجمله‌ای  $x^3 - 2x - 1$ ، يك ریشه بين ۱ و ۲ دارد. با استفاده از روش خط قاطع (الگوریتم ۴.۳) این ریشه را تا سه رقم با معنی صحیح محاسبه نمایید.

۴-۱۰۳ در الگوریتم ۱.۳، گیریم  $M$  معرف طول بازه اولیه  $[a_0, b_0]$  باشد. همچنین نقاط وسط متوالی تولید شده به توسط روش تنصیف را با  $\{x_0, x_1, x_2, \dots\}$  نشان می‌دهیم. نشان دهید که

$$|x_{i+1} - x_i| = \frac{M}{2^{i+1}}$$

همچنین نشان دهید که تعداد بارستهای لازم،  $I$ ، برای تضمین يك تقریب از يك ریشه با دقت  $\varepsilon$  با فرمول زیر تعیین می‌شود:

$$I > -2 - \frac{\log(\varepsilon/M)}{\log 2}$$

۵-۱۰۳ روش تنصیف راهنگامی می‌توان به‌کار برد که نامساوی  $f(a)f(b) < 0$  برقرار باشد. در صورتی که  $f(x)$  در بازه  $(a, b)$  بیش از يك ریشه داشته باشد، الگوریتم ۱.۳ معمولاً جای کدامیک را مشخص می‌کند؟

۶-۱۰۳ با  $a = 0$  و  $b = 1$  کلیه توابع زیر در بازه  $(a, b)$  تغییر علامت می‌دهند یعنی  $f(a)f(b) < 0$ . الگوریتم ۱.۳ (روش تنصیف) چه نقطه‌ای را مشخص می‌کند؟ آیا این نقطه يك ریشه  $f(x)$  هست؟

$$f(x) = (3x - 1)^{-1} \quad f(x) = \cos 10x$$

$$f(x) = \begin{cases} 1 & x \geq 0.3 \\ -1 & x < 0.3 \end{cases}$$

۷-۱۰۳ تابع  $f(x) = e^{2x} - e^x - 2$  دارای يك ریشه در بازه  $[0, 1]$  است. با استفاده از روش نیوتن (الگوریتم ۵.۳) این ریشه را تا چهار رقم با معنی صحیح پیدا کنید.

۱-۱۰۳ تابع  $f(x) = 4 \sin x - e^x$  فروی بازه  $[0, 0.5]$  يك ریشه دارد. با استفاده از الگوریتم خط قاطع (۴.۳) این ریشه را تا چهار رقم با معنی صحیح پیدا کنید.

۹-۱۰۳ با استفاده از الگوریتم قاطع کوچکترین ریشه مثبت بسجمله‌ای

$$p(x) = 2x^3 - 3x - 4$$

را تا سه رقم با معنی صحیح پیدا کنید.

## ۲.۳ برنامه‌های فورترن برای بعضی از روشهای بارستی

هر گاه الگوریتمهای مذکور در بخشهای پیشین برای محاسبات به کار گرفته شوند، به جای اصطلاح مبهم «until satisfied» باید ملاک دقیقی برای پایان کار گذاشته شود. در این بخش تعدادی از راههای ممکن برای تعیین بارست پایانی معقول مورد بحث قرار می‌گیرند و ترجمه‌های الگوریتمهای ۱.۳ و ۳.۳ به زبان فورترن داده می‌شوند.

### زیر برنامه فورترن برای الگوریتم تنصیف ۱.۳

```

SUBROUTINE BISECT ( F, A, B, XTOL, IFLAG )
C***** I N P U T *****
C F NAME OF FUNCTION WHOSE ZERO IS SOUGHT. NAME MUST APPEAR IN AN
C E X T E R N A L STATEMENT IN THE CALLING PROGRAM.
C A,B ENDPOINTS OF THE INTERVAL WHEREIN A ZERO IS SOUGHT.
C XTOL DESIRED LENGTH OF OUTPUT INTERVAL.
C***** O U T P U T *****
C A,B ENDPOINTS OF INTERVAL KNOWN TO CONTAIN A ZERO OF F .
C IFLAG AN INTEGER,
C   = -1, FAILURE SINCE F HAS SAME SIGN AT INPUT POINTS A AND B
C   = 0, TERMINATION SINCE ABS(A-B)/2 .LE. XTOL
C   = 1, TERMINATION SINCE ABS(A-B)/2 IS SO SMALL THAT ADDITION
C       TO (A+B)/2 MAKES NO DIFFERENCE .
C***** M E T H O D *****
C THE BISECTION ALGORITHM 3.1 IS USED, IN WHICH THE INTERVAL KNOWN TO
C CONTAIN A ZERO IS REPEATEDLY HALVED .
C
C   INTEGER IFLAG
C   REAL A,B,F,XTOL, ERROR,FA,FM,XM
C   FA = F(A)
C   IF (FA*F(B) .GT. 0.) THEN
C     IFLAG = -1
C     PRINT 601,A,B
601   FORMAT(' F(X) IS OF SAME SIGN AT THE TWO ENDPOINTS',2E15.7)
C     RETURN
C   END IF
C
C   ERROR = ABS(B-A)
C   DO WHILE ERROR .GT. XTOL
6   ERROR = ERROR/2.
C     IF (ERROR .LE. XTOL) RETURN
C     XM = (A+B)/2.
C     CHECK FOR UNREASONABLE ERROR REQUIREMENT
C     IF (XM + ERROR .EQ. XM) THEN
C       IFLAG = 1
C       RETURN
C     END IF
C     FM = F(XM)
C     CHOOSE NEW INTERVAL
C     IF (FA*FM .GT. 0.) THEN
C       A = XM
C       FA = FM
C     ELSE
C       B = XM
C     END IF
C
C   GO TO 6
C
END

```



در برنامه زیر از این زیر برنامه برای پیدا کردن ریشهٔ معادلهٔ (۲.۳) که در بخش قبل مورد بحث قرار گرفت، استفاده می‌شود.

```
C MAIN PROGRAM FOR TRYING OUT BISECTION ROUTINE
  INTEGER IFLAG
  REAL A,B,ERROR,XI
  EXTERNAL FF
  A = 1.
  B = 2.
  CALL BISECT ( FF, A, B, 1.E-6, IFLAG )
  IF (IFLAG .LT. 0) STOP
  XI = (A+B)/2.
  ERROR = ABS(A-B)/2.
  PRINT 600, XI, ERROR
600 FORMAT(' THE ZERO IS ',E15.7, ' PLUS/MINUS ',E15.7)
      STOP
  END
  REAL FUNCTION FF(X)
  REAL X
  FF = -1. - X*(1. - X*X)
  PRINT 600,X,FF
600 FORMAT(' X, F(X) = ',2E15.7)
      RETURN
  END
```

اکنون راجع به زیر برنامه نصف کردن، مذکور در بالا، توضیح بیشتری می‌دهیم. در اینجا زیر نمایه‌هایی را که در الگوریتم ۱.۳ مورد استفاده قرار گرفته‌اند حذف کرده‌ایم. مقادیر اولیه به وسیلهٔ برنامهٔ اصلی به زیر برنامه داده می‌شوند و در هر مرحله از برنامه، متغیرهای  $A$  و  $B$  مقادیر جاری کرانه‌های پایینی و بالایی ریشهٔ مورد نظر هستند. به خصوص نقطهٔ وسط  $X_M = (A+B)/2$  همیشه بهترین برآورد جاری برای ریشه است و اختلاف مطلق آن با ریشه، همیشه به وسیلهٔ فرمول زیر محدود می‌شود:

$$\text{خطا} = \frac{|A-B|}{2}$$

و به محض اینکه نامساوی زیر برقرار شود محاسبه خاتمه می‌یابد:

$$\text{خطا} \leq XTOL$$

در رابطهٔ بالا  $XTOL$  کرانهٔ خطای مطلق مفروض است. سپس برنامهٔ اصلی برای برآورد ریشه مرتباً از مقادیر جاری  $A$  و  $B$  استفاده می‌کند. برنامهٔ اصلی بایستی علاوه بر  $A$ ،  $B$  و  $XTOL$ ، نام تابع  $f(x)$  را هم که ریشهٔ آن مورد نظر است در زبان فورترن ادامه دهد. از آنجا که فرض مختلف علامه بودن  $f(A)$  و  $f(B)$  برای اجرای الگوریتم شرط عمده است، بنا بر این در ابتدا یک آزمایش برای برقراری این شرط انجام می‌شود. در صورتی که  $f(A)$  و  $f(B)$  مختلف علامه نباشند، اجرای زیر برنامه فوراً خاتمه می‌یابد. متغیر خروجی  $IFLAG$  برای نشان دادن (فهماندن) این پیشامد نامطلوب به برنامهٔ اصلی به کار می‌رود. زیر برنامه، به علت مذکور، هرگز مقادیر تابع را برای یک شناسه بیش از یک بار، محاسبه نمی‌کند، اما مقادیری از تابع را که ممکن است در گامهای بعدی مورد نیاز باشند، ضبط می‌کند. این خط مشی معقولی است زیرا ممکن است این زیر برنامه برای محاسبهٔ مقادیر

توابعی که بسیار وقت گیرند به کار رود. و بالاخره زیر برنامه در مقابل خطای مجاز نامعقول محافظت می شود: برای سادگی فرض کنید که کلیه محاسبات با حساب با ممیز شناور و چهار رقم اعشاری انجام شوند و کرانه های  $A$  و  $B$  به حدی به ریشه نزدیک شده باشند که

$$A = 13196 \quad \text{و} \quad B = 13197$$

در نتیجه

$$\text{خطا} = \frac{B-A}{2} = 0.05$$

در این صورت

$$XM = \frac{A+B}{2} = \frac{26393}{2} = 13196 \quad \text{یا} \quad 13197$$

دو جواب فوق به چگونگی گرد کردن تا چهار رقم بستگی دارد. در هر صورت داریم:

$$XM = A \quad \text{یا} \quad XM = B$$

به طوری که در آخر این مرحله نه  $A$  و نه  $B$  هیچ کدام تغییری نکرده اند. اگر اغماض  $XTOL$  کمتر از  $0.05$  باشد، این زیر برنامه هیچگاه پایان نمی یابد. زیرا  $|B-A|/2$  هرگز کوچکتر از  $0.05$  نخواهد شد. برای جلوگیری از این گونه حلقه های بی پایان که از خطای مجاز نامعقول ناشی می شوند (نامعقول، زیرا این امر ایجاب می کند که  $A$  و  $B$ ، بیش از آنچه که برای دقت با آن دقت ممکن است، به هم نزدیک شوند بی آنکه برهم منطبق شوند) این روال مقدار جاری خطا را به طریق زیر محاسبه می کند. در ابتدا داریم:

$$\text{خطا} = |B-A|$$

اما در آغاز هر مرحله خطا نصف می شود، زیرا این مقدار کاهش خطا در هر مرحله در روش تنصیف وجود دارد. به محض اینکه خطا به اندازه ای کوچک شود که جمع ممیز شناور آن با مقدار جاری  $XM$ ، مقدار  $XM$  را تغییر ندهد، در این صورت برنامه خاتمه می یابد.

در قسمت بعدی، روش تصحیح خطای اصلاح شده (الگوریتم ۳.۳) را مورد بررسی قرار می دهیم. برعکس روش تنصیف، روش تصحیح خطای اصلاح شده ایجاد یک بازه تا حد ممکن کوچک را با استفاده از حساب دقیق متناهی، تضمین نمی کند (به تمرین ۲.۳-۱ نگاه کنید). بنابراین برای این الگوریتم باید ملاک اختتام دیگری به کار گرفته شود.

## برنامه فورترن برای الگوریتم تصحیح خطای اصلاح شده ۳.۳

```

SUBROUTINE MRGFLS ( F, A, B, XTOL, FTOL, NTOL, W, IFLAG )
C***** I N P U T *****
C F NAME OF FUNCTION WHOSE ZERO IS SOUGHT. NAME MUST APPEAR IN AN
C E X T E R N A L STATEMENT IN THE CALLING PROGRAM .
C A,B ENDPOINTS OF INTERVAL WHEREIN ZERO IS SOUGHT.
C XTOL DESIRED LENGTH OF OUTPUT INTERVAL
C FTOL DESIRED SIZE OF P(W)
C NTOL NO MORE THAN NTOL ITERATION STEPS WILL BE CARRIED OUT.
C***** O U T P U T *****
C A,B ENDPOINTS OF INTERVAL CONTAINING THE ZERO .
C W BEST ESTIMATE OF THE ZERO .
C IFLAG AN INTEGER,
C   =-1, FAILURE, SINCE F HAS SAME SIGN AT INPUT POINTS A, B .
C   = 0, TERMINATION BECAUSE ABS(A-B) .LE. XTOL .
C   = 1, TERMINATION BECAUSE ABS(F(W)) .LE. FTOL .
C   = 2, TERMINATION BECAUSE NTOL ITERATION STEPS WERE CARRIED OUT .
C***** M E T H O D *****
C THE MODIFIED REGULA FALSI ALGORITHM 3.3 IS USED. THIS MEANS THAT,
C AT EACH STEP, LINEAR INTERPOLATION BETWEEN THE POINTS (A,FA) AND
C (B,FB) IS USED, WITH FA*FB .LT. 0 ,TO GET A NEW POINT (W,F(W))
C WHICH REPLACES ONE OF THESE IN SUCH A WAY THAT AGAIN FA*FB .LT. 0 .
C IN ADDITION, THE ORDINATE OF A POINT STAYING IN THE GAME FOR MORE
C THAN ONE STEP IS CUT IN HALF AT EACH SUBSEQUENT STEP.
C
      INTEGER IFLAG,NTOL, N
      REAL A,B,F,FTOL,W,XTOL, FA,PB,PW,SIGNFA,PRVSW
      FA = F(A)
      SIGNFA = SIGN(1., FA)
      PB = F(B)
      IF (SIGNFA*PB .GT. 0.) THEN
1001      FORMAT(' F(X) IS OF SAME SIGN AT THE TWO ENDPOINTS',2E15.7)
          IFLAG = -1
          RETURN
      END IF
C
      W = A
      FW = FA
      DO 20 N=1,NTOL
C
          CHECK IF INTERVAL IS SMALL ENOUGH.
          IF (ABS(A-B) .LE. XTOL) THEN
              IFLAG = 0
              RETURN
          END IF
C
          CHECK IF FUNCTION VALUE AT W IS SMALL ENOUGH .
          IF (ABS(FW) .LE. FTOL) THEN
              IFLAG = 1
              RETURN
          END IF
C
          GET NEW GUESS W BY LINEAR INTERPOLATION .
          W = (FA*B - PB*A)/(FA - FB)
          PRVSW = SIGN(1.,FW)
          FW = F(W)
C
          CHANGE TO NEW INTERVAL
          IF (SIGNFA*FW .GT. 0.) THEN
              A = W
              FA = FW
              IF (PW*PRVSW .GT. 0.) FB = FB/2.
          ELSE
              B = W
              PB = FW
              IF (PW*PRVSW .GT. 0.) FA = FA/2.
          END IF
20      CONTINUE
          PRINT 620,NTOL
620      FORMAT(' NO CONVERGENCE IN ',I5,' ITERATIONS')
          IFLAG = 2
          RETURN
      END

```

اولاً این راهواره هنگامی پایان می پذیرد که قدرمطلق آخرین مقدار محاسبه شده تابع، بزرگتر از اغماض مفروض FTOL نباشد. در نتیجه به نظر می رسد که يك «ریشه تقریبی» معادله  $f(x) = 0$  نقطه ای مثل  $x$  است که در آن  $|f(x)|$  «کوچک» باشد. همچنین، از آنجایی که این راهواره مرتباً عمل تقسیم بر مقدار تابع را انجام می دهد، چنین پایان داندنی برای برنامه ضروری به نظر می رسد تا درحالتهای نهایی، تقسیم بر صفر پیش نیاید. ثانیاً، هنگامی که تعداد مرحله های بارستی بیشتر از عدد داده شده NTOL باشد، راهواره ختم می شود. به تعبیری NTOL مقدار محاسبه ای را - که برنامه نویس مایل است برای حل مسئله خود به کار برد - مشخص می کند. به کار گرفتن این ملاک اختتام، برنامه نویس را در مقابل خطای غیر معقول و خطاهای برنامه نویسی و همچنین احتمال اینکه مسئله ای را که می خواهد حل کند به طور کامل نفهمیده باشد، محافظت می کند. بنا بر این بهتر است که چنین ملاک اختتام در هر روش بارستی مورد استفاده قرار گیرد.

زیر برنامه MRGFLS مانند راهواره روش تصنیف، عدد صحیح IFLAG را برگشت می دهد، که نشان دهنده علت اختتام بارست است و نیز جدیدترین مقدار کرانهای A و B را برای ریشه مورد نظر به دست می دهد. بالاخره این راهواره، مانند راهواره تصنیف هرگز مقدار تابع مفروض را برای يك شناسه بیشتر از يك بار محاسبه نمی کند.

الگوریتمهای ۴.۳ و ۵.۳ به ترتیب برای روش خط قاطع و روش نیوتن الزاماً ریشه را محصور نمی کنند، بلکه هر دو، دنباله  $x_0, x_1, x_2, \dots, x_p$  را تولید می کنند که انتظار می رود به سمت ریشه مطلوب  $f(x) = 0$  از تابع داده شده همگرا باشد. بنا بر این باید مقدمتاً به هر دو الگوریتم به عنوان روشهایی برای پیدا کردن نقاطی که در آنها قدرمطلق  $f(x)$  «کوچک» باشد نگاه کرد. بارست هنگامی خاتمه می یابد که قدرمطلق آخرین مقدار محاسبه شده تابع، کمتر از يك FTOL داده شده باشد.

همچنین بارست ممکن است زمانی خاتمه یابد که قدرمطلق تفاضل مقادیر بارستهای متوالی کمتر از يك عدد مفروض XTOL باشد. بنا بر این عادت بر این است که برای اختتام روش خط قاطع یا روش نیوتن از يك یا هر دو ملاک زیر استفاده می کنند:

$$|f(x_n)| \leq \text{FTOL} \quad \text{یا} \quad |x_n - x_{n-1}| \leq \text{XTOL} \quad (13.3)$$

اگر اندازه اعداد مورد نظر را از قبل ندانیم، معمولاً روش بهتر آن است که ملاک خطای نسبی مجاز به کار گرفته شود یعنی بارستها هنگامی خاتمه یابند که داشته باشیم:

$$\frac{|f(x_n)|}{\text{FSIZE}} \leq \text{FTOL} \quad \text{یا} \quad |x_n - x_{n-1}| \leq \text{XTOL} * |x_n| \quad (14.3)$$

در فرمول بالا FSIZE بر آوردی است از مقدار  $f(x)$  در همسایگی ریشه ای که می خواهیم از راه بارست به آن برسیم.

در بخش ۶.۱ درباره خطر این نتیجه گیری که دنباله مفروضی همگرا است، فقط به علت اینکه اختلاف دو جمله متوالی آن «بسیار کم» است، بحث کرده ایم. با این حال دقیقاً از

همین ملاک در هر دو راهوارهٔ فوق استفاده کرده‌ایم، زیرا در روش خط قاطع چنین ملاکی برای اجتناب از تقسیم بر صفر ضروری است. همچنین در هر دو روش، هنگامی که بارستها «به اندازهٔ کافی نزدیک» به ریشهٔ «معادله» باشند، تفاضل مابین آخرین دو بارست محاسبه شده، یک کران نسبتاً مناسبی برای خطای آخرین بارست است. به بیان ساده: اگر بارستهای متوالی اختلاف زیادی با هم نداشته باشند، آنگاه دلیل کمی برای ادامهٔ بارستها وجود دارد. زیر بر نامه‌های مربوط به روش نیوتن و روش خط قاطع در متن داده نشده‌اند و به‌عنوان تمرین به‌عهدهٔ دانشجویان واگذار شده‌اند.

□ مثال ۲.۳ (الف): ریشهٔ حقیقی و مثبت معادلهٔ زیر را پیدا کنید:

$$x^3 - x - 1 = 0$$

نتایج الگوریتمهای ۱.۳، ۳.۳، ۴.۳ و ۵.۳ در جدول زیر که شبیه جدول مثال ۱.۳ است، داده شده‌اند.

n	Bisection		Modified regula falsi		Secant	Newton
	$x_n$	$\epsilon_n$	$x_n$	$\epsilon_n$	$x_n$	$x_n$
0	1.5	$5 \cdot 10^{-1}$	1.5	$5 \cdot 10^{-1}$	1.0	
1	1.25	$3 \cdot 10^{-1}$	1.5833333	$4 \cdot 10^{-1}$	2.0	1.0
2	1.375	$1 \cdot 10^{-1}$	1.6616541	$3 \cdot 10^{-1}$	1.1666667	1.5
3	1.3125	$6 \cdot 10^{-2}$	1.3249256	$2 \cdot 10^{-1}$	1.2531120	1.3478261
4	1.34375	$3 \cdot 10^{-2}$	1.3256293	$2 \cdot 10^{-3}$	1.3372064	1.3252004
5	1.328125	$2 \cdot 10^{-2}$	1.3256305	$9 \cdot 10^{-4}$	1.3238501	1.3247182
6	1.3203125	$8 \cdot 10^{-3}$	1.3247180	$9 \cdot 10^{-4}$	1.3247079	1.3247180
...	.....	.....				
10	1.3247070	$5 \cdot 10^{-4}$				
...	.....	.....				
20	1.3247180	$5 \cdot 10^{-7}$				

□

□ مثال ۲.۴ (ب): مسئلهٔ به‌اصطلاح پیشقدری<sup>۱</sup> در طرح مدارهای الکتریکی مستلزم حل معادله‌ای به‌شکل  $f(v) = 50 I(e^{qv} - 1) + v - 20$  است، که در آن  $v$  ولتاژ،  $I$  شدت جریان و  $q$  پارامتری است که بار الکتریکی را با دمای مطلق مرتبط می‌سازد. در یک مسئلهٔ مهندسی نمونه‌ای بایستی این معادله به‌ازای مقادیر مختلف  $I$  و  $q$  حل شود تا چگونگی تغییر کوچکترین ریشهٔ مثبت  $f(v)$  با تغییر پارامترها معلوم شود.

با استفاده از روش نیوتن کوچکترین ریشهٔ مثبت  $f(v)$  را برای دو مجموعهٔ مختلف از مقادیر پارامتر  $(I, q) = (10^{-8}, 40)$  و  $(I, q) = (10^{-6}, 20)$  پیدا کنید. فرض کنید  $XTOL = 10^{-8}$  و  $FTOL = 10^{-7}$ .

نتایج با مقادیر اولیهٔ مذکور در زیر داده شده‌اند:

## 1. Biasing

$I = 10^{-8}, q = 40$		$I = 10^{-6}, q = 20$	
XN	FXN	XN	FXN
0.35000000	-1.0486984E + 00	0.55000000	1.5436571E + 00
0.39186072	1.6002561E + 00	0.52464183	3.2740616E - 01
0.37948785	3.3539839E - 01	0.51580645	2.6559725E - 02
0.37525497	2.6518162E - 02	0.51495562	2.1750479E - 04
0.37485947	2.0554242E - 04	0.51494853	1.4904231E - 08
0.37485636	1.2605184E - 08		

□ در این مثال انتخاب مقادیر اولیه کمتر موجب واگرایی خواهد شد.

### تمرین

۱-۲-۴ کوشش کنید که ریشه  $x = ۱۳۳۳۳$  از معادله  $x^3 = (۱۳۳۳۳ - x)$  را با استفاده از الگوریتم تصحیح خطای اصلاح شده  $۳.۳$ ، تا پنج رقم اعشار صحیح پیدا کنید. این محاسبه را با بازه  $[۱, ۲]$  شروع کنید. توضیح دهید که چرا در این مورد این روش نمی تواند یک بازه «کوچکی» را که در برگیرنده ریشه باشد به دست دهد؟

۲-۲-۳ در زیر برنامه تصنیف، حتی با وجودی که مقادیر تابع FA و FM اعداد با ممیز شناور به خوبی تعریف شده ای هستند، اما به علت استفاده از حاصل ضرب  $FA * FM$  ممکن است طی اجرای این زیر برنامه سرریز یا پی ریز رخ دهد. این نقص زیر برنامه را با به کار بستن تابع فورترن SIGN تصحیح کنید. آیا لازم است هر دفعه که A تغییر می کند، مقدار FA نیز به تناسب تغییر داده شود؟

۳-۲-۳ ثابت کنید تابع  $f(x) = e^x - 1 - x - x^2/2$  دقیقاً یک ریشه  $= 0$  دارد. (دانهایی: از باقیمانده بسط تیلر  $e^x$  پیرامون صفر استفاده کنید.) سپس تابع فورترن

$$F(X) = \text{EXP}(X) - 1 - X - X * X / 2$$

را به ازای مقادیر مختلف از شناسه X «نزدیک» به صفر ارزیابی کنید تا نشان دهید این تابع تغییر علامت زیادی می دهد، بنا بر این دارای ریشه های زیادی «نزدیک»  $X = 0$  می باشد. از این حقیقت چند نتیجه گیری می توانید بکنید؟ بخصوص در رابطه با روش تصنیف و به طور کلی در رابطه با مفهوم (نظری) «ریشه یک تابع»؟

۴-۲-۳ فرض کنید می خواهید نزدیکترین ریشه معادله  $\tan x = x$  به عدد  $50$  را با استفاده از روش خط قاطع و ممیز شناور با نه رقم اعشاری محاسبه کنید. آیا منطقی خواهد بود که برای ختم محاسبات، از ملاک  $|f(x_n)| \leq 10^{-8}$  استفاده کنیم؟

۵-۲-۳ «تجسس دودویی» مسئله جدول مراجعه ای عبارت است از تعیین یک عدد [به ازای

يك  $X$  مفروض به طوری که  $X$  مابین  $TABLE(I)$  و  $TABLE(I+1)$  قرار گیرد، که در اینجا  $TABLE$  آرایه‌ای است يك بعدی متضمن يك دنباله صعودی (یا نزولی)، يك زیر برنامه فور تون بنویسید که در آن برای انجام این تجسس به نحو مؤثری از تصیف استفاده شود. با فرض آنکه  $TABLE$  دارای  $n$  درایه باشد، در این راهواره چند بار  $X$  با یکی از درایه‌های  $TABLE$  مقایسه می‌شود؟

۶-۲۰۳ يك زیر برنامه برای روش خط قاطع که بر پایه فرم (۷.۳) باشد بنویسید. برای ختم عملیات از یکی از ملاکهای خطای نسبی (۱۴.۳) استفاده کنید. همچنین برای مقایسه خطای نسبی  $|x_n| * XTOL \leq |x_n - x_{n-1}|$  تفاضل  $x_n - x_{n-1}$  را دوباره محاسبه نکنید بلکه مقدار تصحیحی آن از بارست قبلی را به کار گیرید.

۷-۲۰۳ يك زیر برنامه برای روش نیوتن بنویسید. مطمئن شوید که يك خروجی برای موردی که  $f'(x_n)$  مساوی صفر است تهیه می‌شود. بعلاوه با به کار گرفتن ملاکهای (۱۳.۳) یا (۱۴.۳) باید ترتیبی داده شود که در صورت برقرار نبودن همگرایی، برنامه بعد از تعداد معینی بارست  $NTOL$  خاتمه یابد.

۸-۲۰۳ کوچکترین ریشه مثبت هر يك از معادلات زیر را با استفاده از الگوریتمهای ۱۰.۳، ۳.۳، ۴.۳ و ۵.۳ تا حداکثر دقت موجود روی کامپیوتر مورد استفاده‌تان پیدا کنید. نتایج خود، و تعداد بارستها و دقت عمل حاصله از الگوریتمها را با هم مقایسه کنید.

$$e^{-x} - \sin x = 0 \quad (\text{الف})$$

$$x - e^{-x^2} = 0 \quad (\text{ب})$$

$$x^3 - x - 2 = 0 \quad (\text{پ})$$

۹-۲۰۳ معادله داده شده در (مثال ۲.۳ ب) را با استفاده از روش نیوتن و به ازای مقادیری از پارامترها برابر با  $(I, q) = (10^{-7}, 30)$  حل نمایید. سعی کنید این معادله را برای مقادیر اولیه مختلف بین ۰ و ۴ حل کنید و به اثرات این مقادیر اولیه روی همگرایی و واگرایی توجه نمایید.

### ۳.۳ بارست نقطه ثابت

در بخش ۱۰.۳ به بارست نقطه ثابت به عنوان يك روش ممکن برای به دست آوردن ریشه معادله

$$f(x) = 0 \quad (15.3)$$

اشاره کردیم. در این روش از معادله (۱۵.۳) معادله‌ای به صورت

$$x = g(x) \quad (16.3)$$

به دست می آوریم به گونه ای که هر جواب (۱۶.۳)، یعنی هر «نقطه ثابت» از  $g(x)$ ، يك ریشه (۱۵.۳) باشد. برای انجام این عمل راههای مختلفی وجود دارد. برای مثال، اگر داشته باشیم

$$f(x) = x^2 - x - 2 \quad (17.3)$$

آنگاه نمونه هایی از انتخابهای ممکن برای تابع  $g(x)$  عبارت اند از

$$g(x) = x^2 - 2 \quad (\text{الف})$$

$$g(x) = \sqrt{2+x} \quad (\text{ب})$$

$$g(x) = 1 + \frac{2}{x} \quad (\text{پ})$$

(۱۸.۳)

به ازای يك ثابت غیر صفر  $m$

$$g(x) = x - \frac{x^2 - x - 2}{m} \quad (\text{ت})$$

هر يك از  $g(x)$  های فوق يك تابع بارست<sup>۱</sup> برای حل معادله (۱۵.۳) نامیده می شود.  $f(x)$  به توسط (۱۷.۳) داده شده است. پس از انتخاب تابع بارست  $g(x)$ ، برای حل معادله (۱۵.۳) الگوریتم زیر را اجرا می کنیم.

**الگوریتم ۶.۳:** بارست نقطه ثابت. يك تابع بارست  $g(x)$  و يك نقطه شروع  $x_0$  داده شده اند.

For  $n = 0, 1, 2, \dots$ , until satisfied, do:

Calculate  $x_{n+1} = g(x_n)$

برای اینکه این الگوریتم مفید واقع شود باید ثابت کنیم که:

(i) به ازای يك نقطه شروع مفروض  $x_0$ ، می توانیم  $x_1, x_2, x_3, \dots$  را متوالیاً محاسبه کنیم.

(ii) دنباله  $x_1, x_2, \dots$  به سمت نقطه ای مانند  $\xi$  همگراست.

(iii) حد  $\xi$ ، يك نقطه ثابت  $g(x)$  است، یعنی  $\xi = g(\xi)$ .

مثال مربوط به تابع حقیقی مقدار  $g(x) = -\sqrt{x}$  نشان می دهد که (i) يك شرط پیش پا افتاده ساده ای نیست، زیرا در این حالت  $g(x)$  فقط به ازای  $x \geq 0$  تعریف شده است. با شروع از يك نقطه دلخواه  $x_0 > 0$  خواهیم داشت  $x_1 = g(x_0) < 0$ . بنابراین  $x_2$  را



نمی‌توانیم محاسبه کنیم. برای پرداختن به (i) از اصل (۱.۳) استفاده می‌کنیم:

اصل ۱.۳ يك بازهٔ  $I = [a, b]$  وجود دارد چنانکه به ازای جميع نقاط  $x \in I$ ،  $g(x)$  معین است و  $g(x) \in I$ ، یعنی تابع  $g(x)$ ،  $I$  را بر روی خودش می‌نگارد.

به استقراء نسبت به  $n$ ، از این اصل نتیجه می‌شود که اگر  $x_0 \in I$ ، آنگاه به ازای جميع مقادیر  $n$  داریم  $x_n \in I$ ؛ بنا بر این  $g(x_n) = x_{n+1}$  معین و در  $I$  است.

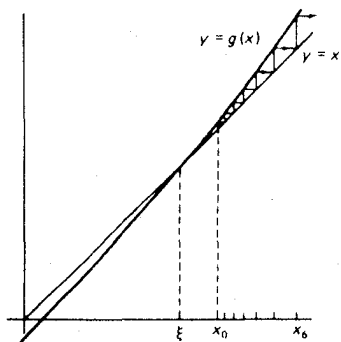
در بخش ۱.۳، قبلاً (iii) را مورد بحث قرار دادیم و ثابت کردیم که اگر  $g(x)$  پیوسته باشد، شرط (iii) صادق است. برای پرداختن به (iii) از اصل ۲.۳ استفاده می‌کنیم.

اصل ۲.۳ تابع بارست  $g(x)$  در  $I = [a, b]$  پیوسته است.

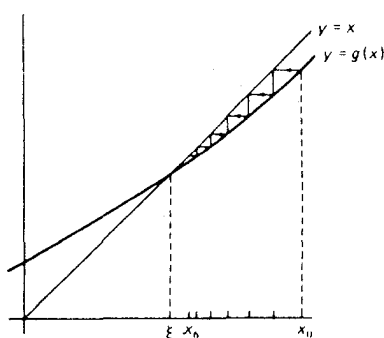
ملاحظه می‌کنیم که اصلهای ۱.۳ و ۲.۳ با هم، ایجاب می‌کنند که  $g(x)$  دارای يك نقطهٔ ثابت در  $I = [a, b]$  باشد. زیرا اگر  $g(a) = a$  و یا  $g(b) = b$ ، آنگاه وجود نقطهٔ ثابت روشن است. در غیر این صورت داریم  $g(a) \neq a$  و  $g(b) \neq b$ . اما به موجب اصل ۱.۳،  $g(a)$  و  $g(b)$  هر دو در  $I = [a, b]$  قرار می‌گیرند، بنابراین  $g(a) > a$  و  $g(b) < b$ . این امر ایجاب می‌کند که تابع  $h(x) = g(x) - x$  در سرطهای  $h(a) > 0$  و  $h(b) < 0$  صدق کند. از آنجا که طبق اصل ۲.۳، روی  $I$  پیوسته است، بنابراین طبق قضیهٔ مقدار میانگین برای توابع پیوسته (به بخش ۷.۱ نگاه کنید)،  $h(x)$  باید در نقطه‌ای واقع در  $I$  صفر شود. اما این امر مبین این است که  $g(x)$  يك نقطهٔ ثابتی در  $I$  دارد و این حکم ما را ثابت می‌کند.

برای بحث دربارهٔ (ii) که مربوط به همگرایی است، بهتر است که بارست را به طور نموداری انجام دهیم. این کار به نحو زیر صورت می‌گیرد. از آنجا که  $x_n = g(x_{n-1})$ ، نقطهٔ  $\{x_{n-1}, x_n\}$  بر نمودار  $g(x)$  قرار دارد. برای تعیین  $\{x_n, x_{n+1}\}$  از روی  $\{x_{n-1}, x_n\}$ ، از نقطهٔ اخیر خطی به موازات محور  $x$ ها رسم می‌کنیم. این خط، خط  $y = x$  را در نقطهٔ  $\{x_n, x_n\}$  قطع می‌کند. از این نقطه خطی به موازات محور  $y$ ها رسم می‌کنیم تا نمودار  $y = g(x)$  را در نقطهٔ  $\{x_n, g(x_n)\}$  قطع کند. اما از آنجا که  $g(x_n) = x_{n+1}$ ، این همان نقطهٔ مطلوب  $\{x_n, x_{n+1}\}$  است. در شکل ۴.۳، چند مرحلهٔ اول از بارست نقطهٔ ثابت را برای چهار مورد نمونه، نشان داده‌ایم. توجه کنید که  $\xi$  يك نقطهٔ ثابت  $g(x)$  است اگر، و تنها اگر،  $y = g(x)$  و  $y = x$  همدیگر را در  $(\xi, \xi)$  قطع کنند.

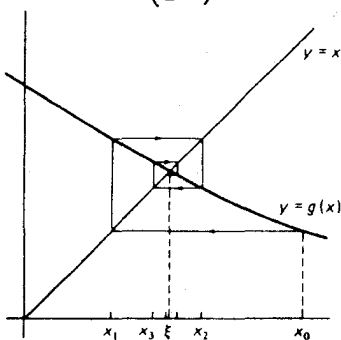
به طوری که شکل ۴.۳ نشان می‌دهد، بارست نقطهٔ ثابت ممکن است اصلاً همگرا نباشد مانند همگرا نبودن در شکلهای ۴.۳ (الف) و ۴.۳ (د). به نظر می‌رسد که همگرا بودن یا نبودن بارستها (در صورتی که  $g(x)$  نقطهٔ ثابتی داشته باشد) بستگی به شیب  $g(x)$  دارد. اگر قدر مطلق شیب  $g(x)$  در نزدیکی نقطهٔ ثابت  $\xi$  از  $g(x)$  خیلی بزرگ باشد، نمی‌توان به همگرایی امیدوار بود. بنابراین اصل ۳.۳ را اختیار می‌کنیم.



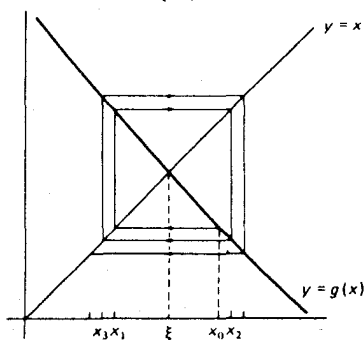
(الف)



(ب)



(ج)



(ت)

### شکل ۴.۳ بارست نقطه ثابت.

**اصل ۳.۳** تابع بارست در  $I = [a, b]$  مشتقپذیر است. بعلاوه يك عدد ثابت غیر منفی مانند  $K < 1$  وجود دارد به طوری که

$$|g'(x)| \leq K \quad x \in I$$

توجه شود که اصل ۳.۳، اصل ۲.۳ را ایجاب می کند. زیرا يك تابع مشتقپذیر حتماً پیوسته است.

**قضیه ۱.۳** گیریم  $g(x)$  تابع بارستی باشد که اصلهای ۱.۳ و ۳.۳ در آن صدق کنند. در این صورت  $g(x)$  دقیقاً يك نقطه ثابت  $\xi$  در  $I$  دارد و با شروع از هر نقطه  $x_0$  در  $I$ ، دنباله  $x_1, x_2, x_3, \dots$  که به توسط بارست نقطه ثابت الگوریتم ۶.۳ تولید می شود، در  $\xi$  همگراست.

برای اثبات این قضیه یادآور می شویم که قبلاً وجود يك نقطه ثابت برای  $g(x)$  را در بسازه  $I$  ثابت کرده ایم. اکنون فرض کنید  $x$  نقطه دلخواهی در  $I$  باشد. پس چنانکه قبلاً

اشاره شد، بارست نقطه ثابت دنباله  $x_1, x_2, \dots$  را تولید می کند که بنا بر اصل ۱.۳، همه این نقاط در  $I$  قرار دارند. اگر خطا در بارستن  $n$ ام را با  $e_n$  نشان دهیم

$$e_n = \xi - x_n \quad n = 0, 1, 2, \dots$$

در این صورت چون  $\xi = g(\xi)$  و  $x_n = g(x_{n-1})$ ، بنا بر قضیه مقدار میانگین برای مشتق (به بخش ۷.۱ نگاه کنید) مقداری مثل  $\eta_n$  بین  $\xi$  و  $x_{n-1}$  وجود دارد به طوری که

$$e_n = \xi - x_n = g(\xi) - g(x_{n-1}) = g'(\eta_n)e_{n-1} \quad (19.3)$$

بنابر این طبق اصل ۳.۳، داریم

$$|e_n| \leq K |e_{n-1}|$$

با استقراء نسبت به  $n$  نتیجه می شود که

$$|e_n| \leq K |e_{n-1}| \leq K^2 |e_{n-2}| \leq \dots \leq K^n |e_0|$$

از آنجا که  $0 < K < 1$ ، داریم  $\lim_{n \rightarrow \infty} K^n = 0$ ، بنا بر این صرف نظر از خطای اولیه  $e_0$

داریم

$$\lim_{n \rightarrow \infty} |e_n| = \lim_{n \rightarrow \infty} K^n |e_0| = 0$$

اما این تساوی بیانگر این است که  $x_1, x_2, \dots$  به سمت  $\xi$  همگراست. همچنین ثابت می کند که  $\xi$  تنها نقطه ثابت  $g(x)$  در بازه  $I$  است. زیرا اگر  $\zeta$  نیز يك نقطه ثابت از  $g(x)$  در  $I$  باشد، آنگاه با  $\xi = x_0$  خواهیم داشت  $\zeta = g(x_0) = x_1$ ، بنا بر این  $|e_1| \leq K |e_0|$ ، چون  $K < 1$ ، نتیجه می گیریم که  $|e_0| = 0$  یا  $\xi = \zeta$ ، که اثبات کامل می شود.

غالباً تعیین صحت و سقم اصل ۱.۳، کار بسیار مشکلی است. در چنین وضعیتی حکم ضعیفتر زیر ممکن است، حداقل در صورتی که بارست «به اندازه کافی نزدیک» به نقطه ثابت شروع شده باشد، موفقیت را تضمین کند.

فرض. اگر  $g(x)$  در يك بازه  $I$  متضمن نقطه ثابت  $\xi$  است به طور پیوسته مشتق پذیر باشد، و هرگاه  $|g'(\xi)| < 1$ ، آنگاه مقداری چون  $\varepsilon > 0$  وجود دارد که هر وقت شرط  $|\xi - x_0| \leq \varepsilon$  برقرار باشد، بارست نقطه ثابت با  $g(x)$  همگراست.

در حقیقت چون  $g'(x)$  نزدیک به  $\xi$  پیوسته است و  $|g'(\xi)| < 1$ ، به ازای هر  $K$  با شرط  $1 < K < |g'(\xi)|$ ، مقداری چون  $\varepsilon > 0$  وجود دارد. چنانکه به ازای هر  $x$  به شرط برقراری  $|\xi - x| \leq \varepsilon$  نامساوی  $|g'(x)| \leq K$  برقرار است. با تثبیت این  $K$  و  $\varepsilon$  مربوطه اش، به ازای  $I = [\xi - \varepsilon, \xi + \varepsilon]$  اصل ۳.۳ برقرار است. همچنین برای اصل ۱.۳، گیریم  $x$  نقطه دلخواهی از  $I$  باشد، لذا  $|\xi - x| \leq \varepsilon$ . پس مانند برهان قضیه ۱.۳ به ازای نقطه ای چون  $\eta$  بین  $x$  و  $\xi$  داریم

$$g(x) - \xi = g(x) - g(\xi) = g'(\eta)(x - \xi)$$

اما در این صورت داریم

$$|g(x) - \xi| \leq |g'(\eta)| |x - \xi| \leq K\varepsilon < \varepsilon$$

این رابطه نشان می‌دهد که اگر  $x$  در  $I$  باشد،  $g(x)$  نیز در  $I$  قرار دارد. این امر صحت اصل ۱.۳ را نشان می‌دهد و نتیجه از قضیه ۱.۳ حاصل می‌شود.

به علت این فرع، يك نقطه ثابت  $\xi$  برای  $g(x)$ ، که به ازای آن نامساوی  $g'(\xi) < 1$  برقرار باشد، غالباً يك نقطه جذب<sup>۱</sup> [برای بارست با  $g(x)$ ] نامیده می‌شود.

دو باره تابع درجه دوم  $f(x) = x^2 - x - 2$  (۱۷.۳) یعنی  $f(x) = x^2 - x - 2$  را در نظر می‌گیریم. ریشه‌های این تابع ۲ و ۱- هستند. فرض کنید می‌خواهیم ریشه  $\xi = 2$  را با استفاده از بارست نقطه ثابت محاسبه کنیم. اگر از تابع بارستی که به وسیله (۱۸.۳) الف داده شده بود، یعنی

$$g(x) = x^2 - 2$$

استفاده کنیم، در این صورت به ازای  $x > 1/2$ ، خواهیم داشت  $g'(x) > 1$ . از اینجا نتیجه می‌شود که اصل ۳.۳ برای هر بازه متضمن  $\xi = 2$  برقرار نیست. یعنی  $\xi = 2$  يك نقطه جذب نیست. در حقیقت برای این مثال می‌توان ثابت کرد که با شروع از يك نقطه دلخواه  $x_0$ ، دنباله  $x_1, x_2, \dots$  که به وسیله بارست نقطه ثابت تولید شده، تنها در صورتی به سمت  $\xi = 2$  همگرا می‌شود که به ازای مقداری چون  $n_0$  و به ازای جمیع مقادیر  $n \geq n_0$  تساوی  $x_n = 2$  برقرار باشد، یعنی اگر «تصادفاً»  $\xi = 2$  گرفته شده باشد (به ترمین ۱-۳.۳ نگاه کنید).

از سوی دیگر، اگر (۱۸.۳) ب، یعنی  $g(x) = \sqrt{2+x}$  را انتخاب کنیم، بنا بر این

$$g'(x) = \frac{1}{2\sqrt{2+x}}$$

اما  $x \geq 0$  ایجاب می‌کند که داشته باشیم  $g(x) \geq 0$  و  $1/\sqrt{8} < 1$  و  $g'(x) \leq 1/\sqrt{8}$  در حالی که مثلاً،  $x \leq 7$  ایجاب می‌کند که  $g(x) = \sqrt{2+x} \leq \sqrt{2+7} = 3$ ، بنا بر این یا  $I = [0, 7]$  هر دو اصل ۱.۳ و ۳.۳ صحیح‌اند و هر نقطه  $x_0 \in [0, 7]$  يك دنباله همگرا به دست می‌دهد. زیرا اگر  $x_0 = 0$  را انتخاب کنیم داریم

$$x_1 = \sqrt{2} = 1.41421$$

$$x_2 = \sqrt{3.41421} = 1.84775$$

$$x_3 = \sqrt{3.84775} = 1.96157$$

$$x_4 = \sqrt{3.96157} = 1.99036$$

$$x_5 = \sqrt{3.99036} = 1.99759$$

## 1. Point of attraction

که آشکارا دیده می‌شود که به سمت ریشهٔ  $\epsilon = 2$  همگراست. به عنوان مثال واقعگرایانه تر دیگر، معادلهٔ غیرجبری زیر را در نظر می‌گیریم

$$f(x) = x - 2 \sin x = 0 \quad (20.3)$$

طبیعیترین ترتیب در اینجا

$$x = 2 \sin x$$

است یعنی  $g(x) = 2 \sin x$ . مطالعهٔ خمهای  $y = x$  و  $y = g(x)$  نشان می‌دهد که یک ریشه بین  $\pi/3$  و  $2\pi/3$  وجود دارد. بعلاوه اگر  $\pi/3 \leq x \leq 2\pi/3$ ، آنگاه  $2 \leq g(x) \leq \sqrt{3}$ . بنا بر این اگر  $\pi/3 \leq a \leq \sqrt{3}$  و  $2 \leq b \leq 2\pi/3$ ، آنگاه اصل ۱.۳، صحیح است. بالاخره  $g'(x) = 2 \cos x$  اکیداً از ۱ تا ۱- تنزل می‌کند وقتی که  $x$  از  $\pi/3$  تا  $2\pi/3$  ترقی کند. در نتیجه هنگامی که  $\pi/3 < a \leq \sqrt{3}$  و  $2 \leq b < 2\pi/3$ ، اصل ۳.۳ صادق است و در نتیجه با رست نقطهٔ ثابت، با  $g(x) = 2 \sin x$  به طرف ریشهٔ یکنای (۲۰.۳) در  $[\pi/3, 2\pi/3]$  همگرا می‌شود، هر گاه  $x_0 \in (\pi/3, 2\pi/3)$ .

□ مثال ۳.۳: یک برنامهٔ کامپیوتری بنویسید که در آن برای پیدا کردن کوچکترین ریشهٔ مثبت تابع  $f(x) = e^{-x} - \sin x$  از با رست نقطهٔ ثابت استفاده شود. اولین گام، انتخاب یک تابع با رست و یک مقدار آغازینی است که به یک با رست همگرا منجر شود. تابع  $f(x) = 0$  را دوباره به شکل زیر می‌نویسیم

$$x = x + e^{-x} - \sin x =: g(x)$$

اما چون داریم  $f(0.5) = 0.127\dots$  و  $f(0.7) = -0.147\dots$ ، لذا، کوچکترین ریشهٔ مثبت در بازهٔ  $I = [0.5, 0.7]$  قرار دارد. برای اینکه تحقیق کنیم  $g(x)$  یک تابع با رست همگراست ملاحظه می‌کنیم که چون داریم

$$g'(x) = 1 - e^{-x} - \cos x$$

لذا  $\dots 0.48 = g'(0.5)$ ،  $\dots 0.26 = g'(0.7)$  و چون  $g'(x)$  یک تابع یکنوا روی  $I$  است، به ازای  $x \in I$  داریم  $|g'(x)| < 1$ . به همین ترتیب می‌توان ثابت کرد که به ازای جمیع مقادیر  $x \in I$  داریم  $0.7 < g(x) < 0.5$ . بنا بر این اگر  $x_0$  در  $I$  انتخاب شود با رست نقطهٔ ثابت همگرا خواهد بود.

برنامهٔ زیر با یک کامپیوتر CDC ۶۵۰ اجرا شده است. توجه کنید که اختتام نتیجه بخش این برنامه مستلزم برقرار بودن هر دو رابطهٔ زیر برای خطاست

$$|x_n - x_{n-1}| < XTOL |x_n|$$

حل معادلات غیر خطی ۱۲۱

$$|f(x_n)| < FTOL$$

همچنین اگر آزمونهای همگرایی پس از ۲۰ بارست صدق نکنند، برنامه مختومه تلقی می‌شود.

```

C PROGRAM FOR EXAMPLE 3.3
  INTEGER J
  REAL ERROR, FTOL, XNEW, XOLD, XTOL, Y
C THIS PROGRAM SOLVES THE EQUATION
C EXP(-X) = SIN(X)
C BY FIXED POINT ITERATION, USING THE ITERATION FUNCTION
  G(X) = EXP(-X) - SIN(X) + X
C
  DATA XTOL, FTOL / 1.E-8, 1.E-8 /
  PRINT 600
600 FORMAT(9X, 'XNEW', 12X, 'F(XNEW)', 10X, 'ERROR')
  XOLD = .6
  Y = G(XOLD) - XOLD
  PRINT 601, XOLD, Y
601 FORMAT(3X, 3E16.8)
  DO 10 J=1, 20
    XNEW = G(XOLD)
    Y = G(XNEW) - XNEW
    ERROR = ABS(XNEW - XOLD)/ABS(XNEW)
    PRINT 601, XNEW, Y, ERROR
    IF (ERROR .LT. XTOL .OR. ABS(Y) .LT. FTOL) STOP
    XOLD = XNEW
  10 CONTINUE
  PRINT 610
610 FORMAT(' FAILED TO CONVERGE IN 20 ITERATIONS')
  STOP
END

```

### برونداد برای مثال ۳.۳

XN	F(XN)	ERROR
6.00000000E - 01	-1.58308373E - 02	2.70997483E - 02
5.84169163E - 01	6.06240576E - 03	1.02712326E - 02
5.90231568E - 01	-2.35449276E - 03	4.00507667E - 03
5.87877076E - 01	9.09583240E - 04	1.54484349E - 03
5.88786659E - 01	-3.52118178E - 04	5.98398213E - 04
5.88434541E - 01	1.36203144E - 04	2.31413378E - 04
5.88570744E - 01	-5.27011849E - 05	8.95489706E - 05
5.88518043E - 01	2.03892661E - 05	3.46438992E - 05
5.88538432E - 01	-7.88865463E - 06	1.34039851E - 05
5.88530543E - 01	3.05208415E - 06	5.18591321E - 06
5.88533595E - 01	-1.18084550E - 06	2.00642389E - 06
5.88532415E - 01	4.56865632E - 07	7.76278869E - 07
5.88532871E - 01	-1.76760146E - 07	3.00340402E - 07
5.88532695E - 01	6.83880224E - 08	1.16200876E - 07
5.88532763E - 01	-2.64591478E - 08	4.49578182E - 08
5.88532737E - 01	1.02369739E - 08	1.73940600E - 08
5.88532747E - 01	-3.96065403E - 09	6.72970888E - 09
5.88532743E - 01	1.53236357E - 09	

تمرین

۱-۳۰۳ تحقیق کنید که بارست

$$x_{i+1} = x_i^2 - 2$$

به سمت جواب  $x = 2$  از معادله

$$x^2 - x - 2 = 0$$

همگراست، فقط اگر به ازای مقداری مانند  $n_0$  همهٔ بارسته‌های  $x_n$  با  $n \geq n_0$  مساوی ۲ باشند، یعنی فقط «به‌طور تصادفی» همگراست.

۳-۳.۳ برای هر یک از معادلات زیر یک تابع بارست (ویک بازهٔ  $I$ ) تعیین کنید به طوری که شرایط قضیهٔ ۱.۳ برقرار باشد (فرض کنید می‌خواهیم کوچکترین ریشهٔ مثبت را پیدا کنیم).

$$(الف) \quad x^2 - x - 1 = 0 \quad (ب) \quad x - \tan x = 0$$

$$(پ) \quad e^{-x} - \cos x = 0$$

۳-۳.۴ برنامه‌ای بر مبنای الگوریتم ۶.۳ بنویسید و از این برنامه برای محاسبهٔ کوچکترین ریشه‌های معادله‌های مذکور در تمرین ۳-۳.۳ استفاده کنید.

۴-۳.۳ مطلوب است تعیین بزرگترین بازهٔ  $I$  با ویژگی زیر: به ازای جميع مقادیر  $x_0 \in I$ ، بارست نقطهٔ ثابت با تابع بارست

$$g(x) = x(2 - ax)$$

در صورت شروع با  $x_0$ ، همگرا باشد. آیا اصلهای ۱.۳ و ۳.۳ که برای این  $I$  انتخاب شده‌اند صادق‌اند؟ حدود این بارست چه عددی می‌توانند باشند؟ آیا می‌توانید دلیل خوبی برای استفاده از این تابع بارست خاص بیاورید؟ (توجه کنید که بازهٔ  $I$  به ثابت  $a$  بستگی دارد).

۵-۳.۳ همان تمرین ۴-۳.۳ را برای تابع بارست  $g(x) = (x + a/x)/2$  حل کنید.

۶-۳.۳ تابع  $g(x) = \sqrt{1+x^2}$  از اصل ۱.۳ برای بازهٔ  $I = (-\infty, \infty)$  و در اصل ۳.۳ برای هر بازهٔ متناهی صادق است. ولی بارست نقطهٔ ثابت با این تابع بارست همگرا نیست. چرا؟

۷-۳.۳ معادلهٔ  $e^x - 4x^2 = 0$  یک ریشه بین  $x = 4$  و  $x = 5$  دارد. نشان دهید که این ریشه را با استفاده از بارست نقطهٔ ثابت، با تابع بارست «طبیعی»  $x = (1/2)e^{x/2}$  نمی‌توان پیدا کرد. آیا می‌توانید یک تابع بارستی پیدا کنید که جای صحیح این ریشه را تعیین کند؟

۸-۳.۳ معادلهٔ  $e^x - 4x^2 = 0$  نیز یک ریشه بین  $x = 0$  و  $x = 1$  دارد. نشان دهید که اگر  $x_0$  در بازهٔ  $[0, 1]$  انتخاب شود تابع بارست  $x = (1/2)e^{x/2}$  همگرا خواهد بود.

### ۴.۳ شتاب همگرایی برای بارست نقطه ثابت

در این بخش میزان همگرایی  $۱$  بارست نقطه ثابت بررسی، و نشان داده می‌شود که چگونه گاهی اوقات اطلاعات مربوط به میزان همگرایی می‌توان برای شتاب دادن به همگرایی استفاده کرد.

فرض کنیم که تابع بارست  $g(x)$  به‌طور پیوسته مشتق‌پذیر باشد و با شروع از هر نقطه  $x_0$  دنباله تولید شده  $x_1, x_2, \dots$  به‌وسیله بارست نقطه ثابت به سمت يك نقطه  $\xi$  همگرا باشد. در این صورت نقطه  $\xi$  يك نقطه ثابت  $g(x)$  خواهد بود و طبق (۱۹.۳) نقطه‌ای مانند  $\eta_n$  بین  $\xi$  و  $x_n, \dots, 2, 1, n$  وجود دارد به‌طوری که داریم

$$e_{n+1} = \xi - x_{n+1} = g'(\eta_n)e_n \quad (21.3)$$

از آنجا که  $\lim_{n \rightarrow \infty} x_n = \xi$ ، در نتیجه  $\lim_{n \rightarrow \infty} \eta_n = \xi$ . بنابراین، چون  $g'(x)$  بنا بر فرض پیوسته است، داریم

$$\lim_{n \rightarrow \infty} g'(\eta_n) = g'(\xi)$$

و در نتیجه

$$e_{n+1} = g'(\xi)e_n + \varepsilon_n e_n \quad (22.3)$$

که در اینجا  $\lim \varepsilon_n = 0$ . بنابراین اگر  $g'(\xi) \neq 0$ ، آنگاه به‌ازای  $n$  «به اندازه کافی بزرگ» داریم

$$e_{n+1} \approx g'(\xi)e_n \quad (23.3)$$

یعنی، خطای  $e_{n+1}$  در بارست  $(n+1)$ ام (کمابیش) به‌طورخطی به‌خطای  $e_n$  در بارست  $n$ ام بستگی دارد. از این‌رو می‌گوییم دنباله  $x_0, x_1, x_2, \dots$  به‌طورخطی به‌سمت  $\xi$  همگراست.

حال توجه کنید که می‌توان (۲۱.۳) را نسبت به  $\xi$  حل کرد. زیرا از

$$\xi - x_{n+1} = g'(\eta_n)(\xi - x_n) \quad (24.3)$$

نتیجه می‌شود

$$\begin{aligned} \xi(1 - g'(\eta_n)) &= x_{n+1} - g'(\eta_n)x_n \\ &= [1 - g'(\eta_n)]x_{n+1} + g'(\eta_n)(x_{n+1} - x_n) \end{aligned}$$

بنابراین



$$\xi = x_{n+1} + \frac{g'(\eta_n)(x_{n+1} - x_n)}{1 - g'(\eta_n)} = x_{n+1} + \frac{x_{n+1} - x_n}{g'(\eta_n)^{-1} - 1} \quad (25.3)$$

البته، عدد  $g'(\eta_n)$  بر ما معلوم نیست. اما بنا بر قضیه مقدار میانگین برای مشتق، نسبت

$$r_n := \frac{x_n - x_{n-1}}{x_{n+1} - x_n} = \frac{x_n - x_{n-1}}{g(x_n) - g(x_{n-1})} = g'(\zeta_n)^{-1} \quad (26.3)$$

به ازای مقداری مانند  $\zeta_n$  بین  $x_n$  و  $x_{n-1}$  بر ما معلوم است. بنا بر این به ازای  $n$  «به اندازه کافی بزرگ» داریم

$$r_n = \frac{1}{g'(\zeta_n)} \approx \frac{1}{g'(\xi)} \approx \frac{1}{g'(\eta_n)}$$

و در این صورت نقطه

$$r_n = \frac{x_n - x_{n-1}}{x_{n+1} - x_n} \quad \text{با شرط} \quad \hat{x}_n = x_{n+1} + \frac{x_{n+1} - x_n}{r_n - 1} \quad (27.3)$$

تقریبی به مراتب بهتر از  $x_n$  یا  $x_{n+1}$  برای  $\xi$  است.

این مطلب را از لحاظ نموداری نیز می توان دید. زیرا در واقع، ما معادله (27.3) را از حل (24.3) نسبت به  $\xi$ ، پس از قراردادن عدد  $g[x_{n-1}, x_n]$  به جای  $g'(\eta_n)$  و نشان دادن جواب حاصل با  $\hat{x}_n$ ، به دست آورده ایم. بنا بر این داریم

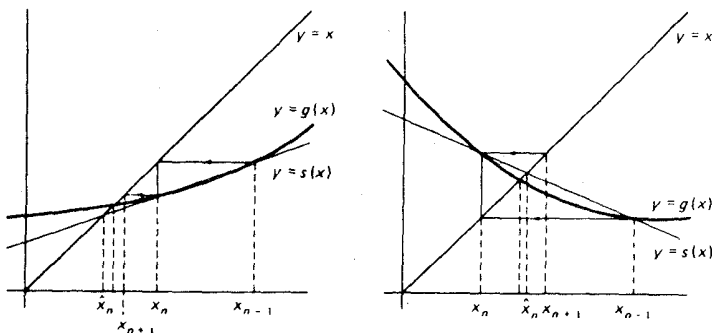
$$\hat{x}_n - x_{n+1} = g[x_{n-1}, x_n](\hat{x}_n - x_n)$$

از آنجا که  $x_{n+1} = g(x_n)$ ، این رابطه نشان می دهد که  $\hat{x}_n$  يك نقطه ثابت از خط مستقیم

$$s(x) = g(x_n) + g[x_{n-1}, x_n](x - x_n)$$

است. ما این رابطه را به عنوان درونیاب خطی برای  $g(x)$  در  $x_{n-1}$  و  $x_n$  می شماریم. اما اگر شیب  $g(x)$  بین نقاط  $x_{n-1}$  و  $\xi$  اندکی تغییر کند، یعنی اگر  $g(x)$  تقریباً يك خط مستقیم واصل بین دو نقطه  $x_{n-1}$  و  $\xi$  باشد، در این صورت و تر  $s(x)$  می باید تقریب بسیار خوبی برای  $g(x)$  در آن بازه باشد، از این رو نقطه ثابت  $\hat{x}_n$  از این خط قاطع می باید يك تقریب بسیار خوبی برای نقطه ثابت  $\xi$  از  $g(x)$  باشد؛ به شکل 5.3 نگاه کنید.

در عمل نمی توان ثابت کرد که يك  $x_n$  خاصی، «به اندازه کافی نزدیک» به  $\xi$  است تا بتوان  $\hat{x}_n$  را به عنوان تقریبی بهتر از  $x_n$  یا  $x_{n+1}$  برای  $\xi$  گرفت. اما می توان با بررسی نسبت های  $r_n$  و  $r_{n-1}$  این فرض را که  $x_n$  «به اندازه کافی نزدیک» است، آزمایش کرد. در صورتی که این نسبت ها تقریباً ثابت باشند، این فرض را که شیب  $g(x)$  در بازه مورد نظر خیلی کم تغییر می کند، می پذیریم. بنا بر این می توانیم عقیده داشته باشیم که و تر  $s(x)$  تقریب «به اندازه کافی خوبی» برای  $g(x)$  است تا بتواند  $\hat{x}_n$  را تقریبی خیلی بهتر از  $x_n$  برای  $\xi$



شکل ۵.۳ شتاب همگرایی برای بارست نقطه ثابت.

بسازد. به ویژه در این حالت قبول می‌کنیم که  $|\hat{x} - x_n|$  برآورد خوبی برای خطای  $|e_n|$  است.

□ مثال ۴.۳: معادله

$$1.25x - \tan x = 0.1 \quad (28.3)$$

دارای یک ریشه  $0.020592169510 \dots$  است. تابع بارست زیر را انتخاب

$$g(x) = \frac{0.1 + \tan x}{1.25}$$

و با شروع از  $x_0 = 0$  دنباله  $x_1, x_2, \dots$  را، با استفاده از بارست نقطه ثابت تولید می‌کنیم. بعضی از  $x_n$ ها در جدول زیر فهرست شده‌اند. به نظر می‌رسد که این دنباله آهسته اما به‌طور مطمئن به سمت  $\xi$  همگراست. دنباله نسبت‌های  $r_n$  را نیز محاسبه می‌کنیم. این دنباله نیز در همین جدول داده شده‌است.

$n$	$x_n$	$\hat{x}_n$	$r_n$	$r_n$
	0.0			0.0
1	0.0666667	0.2005954	1.4978	0.2005954
2	0.1111771	0.2024021	1.4879	0.2059125
3	0.1410916	0.2039180	1.4761	0.2059217
4	0.1613569	0.2048536	1.4659	
5	0.1751813	0.2053721	1.4579	
...	.....	.....	.....	.....
10	0.2009941	0.2059055	1.4408	
15	0.2051206	0.2059213	1.4379	
20	0.2057911	0.2059217	1.4374	
25	0.2059004		1.4373	
30	0.2059182		1.4372	

بخصوص با مشاهده

$$r_1 = 1.3978\dots \quad r_2 = 1.4879\dots \quad r_3 = 1.4761\dots$$

فکر می‌کنیم که این مقادیر «به اندازه کافی» ثابت هستند تا نتیجه بگیریم که به ازای جمیع مقادیر  $n \geq 1$ ،  $\hat{x}_n$  تقریبی بهتر از  $x_n$  برای  $\xi$  است. این امر با درج مقادیر  $\hat{x}_n$  در جدول نیز تأیید شده است.  $\square$

خواه یک  $\hat{x}_n$  خاص تقریبی بهتر از  $x_n$  برای  $\xi$  باشد یا نه، می‌توان ثابت کرد که دنبالهٔ  $\hat{x}_1, \hat{x}_2, \dots$  سریعتر از دنبالهٔ اصلی  $x_0, x_1, x_2, \dots$  به سمت  $\xi$  همگرا می‌شود، یعنی داریم

$$\hat{x}_n = \xi + o(e_n) \quad (29.3)$$

[برای تعریف  $o(\ )$  به بخش ۶.۱ نگاه کنید.]

این روند استخراج یک دنبالهٔ سریعتر همگرای  $\hat{x}_1, \hat{x}_2, \dots$  با استفاده از رابطه (۲۷.۳) از یک دنبالهٔ به طور خطی همگرای  $x_0, x_1, x_2, \dots$  را معمولاً روند  $\Delta^2$  ایتمکن نامند. با استفاده از علامات اختصار

$$\Delta x_k = x_{k+1} - x_k \quad \Delta^2 x_k = \Delta(\Delta x_k) = \Delta x_{k+1} - \Delta x_k$$

از بخش (۶.۲)، رابطه (۲۷.۳) را می‌توان به شکل زیر نوشت

$$\hat{x}_n = x_{n+1} - \frac{(\Delta x_n)^2}{\Delta^2 x_{n-1}} \quad (30.3)$$

و از این روند که نام «روند  $\Delta^2$ » بر آن نهاده شده است. این روند برای هر دنبالهٔ به طور خطی همگرا کاربرد دارد، خواه این دنباله بر اثر بارست نقطهٔ ثابت تولید شده یا نشده باشد.

الگوریتم ۷.۳: روند  $\Delta^2$  ایتمکن. دنبالهٔ  $x_0, x_1, x_2, \dots$  که به سمت  $\xi$  همگراست داده شده است. با استفاده از (۳۰.۳) دنبالهٔ  $\hat{x}_1, \hat{x}_2, \dots$  را محاسبه کنید. اگر دنبالهٔ  $x_0, x_1, x_2, \dots$  به طور خطی به سمت  $\xi$  همگرا شود، یعنی اگر به ازای مقداری مانند  $K \neq 0$  داشته باشیم

$$\xi - x_{n+1} = K(\xi - x_n) + o(\xi - x_n)$$

در این صورت

$$\hat{x}_n = \xi + o(\xi - x_n)$$

بعلاوه اگر از مقدار معینسی از  $k$  به بعد، دنباله نسبت‌های تفاضلی، یعنی  $\Delta x_{k-1}/\Delta x_k, \Delta x_k/\Delta x_{k+1}, \dots$  تقریباً ثابت باشد، آنگاه می‌توان فرض کرد که  $\hat{x}_k$  تقریبی بهتر از  $x_k$  برای  $\xi$  است. در این صورت بخصوص  $|\hat{x}_k - x_k|$  برآورد مطلوبی برای خطای  $|\xi - x_k|$  است.

اگر در مورد بارست نقطه ثابت، تشخیص داده شود که  $\hat{x}_k$  تقریبی بهتر از  $x_k$  برای  $\xi$  است، در این صورت قطعاً محاسبه دنباله  $x_{k+1}, x_{k+2}, \dots$  اتلاف وقت خواهد بود. در این حالت بجا به نظر می‌رسد که بارست نقطه ثابت تازه، با حدس اولیه  $\hat{x}_k$  شروع شود. این امر به الگوریتم زیر منجر می‌شود.

**الگوریتم ۸.۳:** بارست استغفنسن<sup>۱</sup>. تابع بارست  $g(x)$  و نقطه  $y_0$  داده شده‌اند.

For  $n = 0, 1, 2, \dots$ , until satisfied, do:

$$x_0 := y_n$$

$$\text{Calculate } x_1 = g(x_0), x_2 = g(x_1)$$

$$\text{Calculate } d = \Delta x_1, r = \Delta x_0/d$$

$$\text{Calculate } y_{n+1} = x_2 + d/(r-1)$$

یک مرحله این الگوریتم شامل دومرحله از بارست نقطه ثابت و به دنبال آن یکبار کاربرد (۲۷.۳) می‌باشد. نتیجه این سه بارست موجود برای محاسبه مقدار شروع برای مرحله بعدی به کار می‌رود.

در جدول فوق  $y_n$  حاصل از کاربرد این الگوریتم برای مثال (۴.۳) درج شده است. مشاهده می‌شود که کلیه رقمهای اعشاری  $y_n$  صحیح‌اند.

## تمرین

۱-۴.۳ فرض کنید خطای بارست نقطه ثابتی، به ازای مقداری از  $k$ ،  $|k| < 1$ ، در رابطه بازگشتی

$$e_{n+1} = ke_n$$

صدق کند. برای تعداد بارستها،  $N$ ، عبارتی چنان پیدا کنید که خطای اولیه  $e_0$  را به اندازه مضربی از  $10^{-m}$  ( $m > 0$ ) کاهش دهد.

۲-۴.۳ بارست نقطه ثابت برای معادله

$$f(x) = 0.5 - x + 0.2 \sin x = 0$$

به کار رفته و تقریبهای متوالی زیر که به ترتیب در جدول زیر داده شده‌اند، پدید آمده‌اند. با استفاده از الگوریتم ۷.۳ میتوان دنباله شتابدار  $\hat{x}_k$  و نسبتهای  $r_k$  را محاسبه کنید. با استفاده از نسبتهای  $r_k$  مقدار تقریبی  $g'(x)$  را نیز محاسبه کنید.

$k$	$x_k$
0	0.5000 0000
1	0.5958 8511
2	0.6122 4830
3	0.6149 4176
4	0.6153 8219
5	0.6154 5412
6	0.6154 6587
7	0.6154 6779
8	0.6154 6810
9	0.6154 6815

۴-۴.۳ برنامه‌ای برای روش شتابدار استنسن (الگوریتم ۸.۳) بنویسید. با استفاده از این برنامه کوچکترین ریشه مثبت تابع مذکور در تمرین ۴-۴.۳ را با استفاده از تابع بارست  $g(x) = 0.5 + 0.2 \sin x$  و  $x_0 = 0.5$  پیدا کنید.

۴-۴.۴ در بخش ۳.۳ نشان داده شد که بارست نقطه ثابت برای معادله

$$x_{i+1} = \sqrt{2 + x_i}$$

دنباله تقریبهای زیر را برای ریشه مثبت تابع  $f(x) = x^2 - x - 2$  تولید می‌کند:

$$x_0 = 0 \quad x_1 = 1.41421 \quad x_2 = 1.84776$$

$$x_3 = 1.96157 \quad x_4 = 1.99037 \quad x_5 = 1.99759$$

با استفاده از الگوریتم ایتکن ۷.۳ دنباله شتابدار مربوطه را محاسبه و بسه به‌بودمیزان همگرایی به طرف ریشه  $2 = x$  توجه کنید.

۵-۴.۳ تابع بارست  $g(x) = x - x^3$  را در نظر بگیرید. نقطه ثابت و یکنای  $g(x)$  را پیدا کنید. ثابت کنید که اگر  $x_0 \in (-1, 1)$ ، بارست نقطه ثابت بسا این تابع بارست به طرف نقطه ثابت یکنای  $g(x)$  همگرا می‌شود. (دانه‌مایی: از این واقعیت که اگر روابط  $c \leq x_n \leq x_{n+1} \leq c$  به ازای یک مقدار ثابت  $c$  برقرار باشد، آنگاه این دنباله همگرا می‌شود استفاده کنید). آیا نامساوی  $|e_n| \leq k |e_{n-1}|$  به ازای  $k < 1$  و جمیع مقادیر  $n$  صحیح است؟

### ۵.۳\* همگرایی نیوتن و روشهای خط قاطع

در بخش قبل، ثابت شد که برای بارست نقطه ثابت، خطای  $e_n$  در  $n$ مین بارستن  $x_n$  در رابطه

$$e_{n+1} \approx g'(\xi)e_n \quad (۳۱.۳)$$

صدق می‌کند، به شرط آنکه  $n$  به اندازه کافی بزرگ و  $g(x)$  به طور پیوسته مشتقپذیر باشد. ظاهراً هرچه  $|g'(\xi)|$  کوچکتر باشد وقتی  $n \rightarrow \infty$ ،  $e_n$  سریعتر به سمت صفر میل می‌کند. بنابراین حداکثر سرعت همگرایی با درست نقطه ثابت باید هنگامی باشد که  $g'(\xi) = 0$ . اگر  $g(x)$  دوبار مشتقپذیر باشد، به ازای مقداری از  $\xi_n$  بین  $\xi$  و  $x_n$ ، از فرمول تیلر خواهیم داشت

$$\begin{aligned} e_{n+1} &= \xi - x_{n+1} = g(\xi) - g(x_n) \\ &= -g'(\xi)(x_n - \xi) - \frac{1}{2}g''(\xi_n)(x_n - \xi)^2 \end{aligned}$$

یعنی

$$e_{n+1} = g'(\xi)e_n - \frac{1}{2}g''(\xi_n)e_n^2 \quad (۳۲.۳)$$

بنابراین اگر  $g'(\xi) = 0$  و  $g''(x)$  در  $\xi$  پیوسته باشد، آنگاه داریم

$$e_{n+1} \approx -\frac{1}{2}g''(\xi)e_n^2 \quad \text{برای } n \text{ به اندازه کافی بزرگ} \quad (۳۳.۳)$$

در این حالت،  $e_{n+1}$  (کمابیش) تابع درجه دومی از  $e_n$  است. بنابراین گوییم که دنباله  $e_n, x_n, \dots$  به طور مربعی به سمت  $\xi$  همگرا می‌شود. بدیهی است که یک چنین تابع بارستی بسیار مطلوب است. مقبولیت روش نیوتن را می‌توان به این حقیقت منتسب دانست که تابع بارست آن یعنی

$$g(x) = x - \frac{f(x)}{f'(x)} \quad (۳۴.۳)$$

معمولاً از این نوع (مربعی) است.

پیش از اینکه ثابت کنیم روش نیوتن (در صورت همگرا بودن) به طور مربعی همگرا می‌شود، مثال ساده زیر را در نظر می‌گیریم.

□ مثال: تعیین جذر مثبت یک عدد مثبت  $A$ ، با پیدا کردن جواب معادله

$$f(x) = x^2 - A = 0$$

هم‌ارز است. در اینجا  $f'(x) = 2x$ ، و با قراردادن آن در (۳۴.۳)، تابع بارست

$$g(x) = x - \frac{x^2 - A}{2x} = \frac{1}{2} \left( x + \frac{A}{x} \right) \quad (۳۵.۳)$$

به دست می آید. پیدا کردن جذر  $A$  به تابع بارست زیر منجر می شود

$$x_{n+1} = \frac{1}{4} \left( x_n + \frac{A}{x_n} \right) \quad (۳۶.۳)$$

به ویژه چنانچه  $A=۲$  و  $x_0=۲$  اختیار شود، نتیجهٔ بارست نقطهٔ ثابت با (۳۶.۳) به صورت زیر در می آید

$$x_0 = ۲.۰$$

$$x_1 = ۱.۵$$

$$x_2 = ۱.۴۱۶۶۶۶۶۶۶\dots$$

$$x_3 = ۱.۴۱۴۲۱۵۶۸\dots$$

$$x_4 = ۱.۴۱۴۲۱۳۵۶\dots$$

$$x_5 = ۱.۴۱۴۲۱۳۵۶\dots$$

$$r_1 = ۶$$

$$r_2 = ۳۴$$

$$r_3 = ۱,۱۵۲$$

$$r_4 = ۱,۳۳۱۷۱۴$$

دنبالهٔ بارستها آشکارا به سرعت همگرا می شود، و دنبالهٔ متناظر آن  $r_1, r_2, r_3, \dots$  که از نسبت های  $r_n = \Delta x_{n-1} / \Delta x_n$  به دست می آید، به سمت  $\infty$  همگرا می شود. از آنجا که همگرایی بارست نقطهٔ ثابت مستلزم برقراری تساوی  $\lim_{n \rightarrow \infty} r_n^{-1} = g'(\xi)$  است، مثال فوق این ادعا را روشن می کند و همگرایی سریع و مطلوب روش نیوتن را نشان می دهد.  $\square$

می توان همگرایی مربعی روش نیوتن را این گونه نشان داد که اگر  $f(\xi) = 0$  و  $f'(\xi) \neq 0$ ، آنگاه تابع بارست

$$g(x) = x - \frac{f(x)}{f'(x)}$$

در يك بازهٔ باز در همسایگی  $\xi$  به طور پیوسته مشتق پذیر است و  $g'(\xi) = 0$ . در نتیجه با در نظر گرفتن فرع قضیهٔ ۱.۳ مقداری مانند  $\varepsilon > 0$  وجود دارد، به طوری که به ازای هر انتخاب  $x_0$  که در رابطه  $|\xi - x_0| \leq \varepsilon$  صدق کند بارست نقطهٔ ثابت با  $g(x)$  به سمت  $\xi$  همگرا می شود. اما کاراتر به نظر می رسد که اثبات همگرایی مربعی نیوتن، مستقیماً صورت گیرد و همزمان با آن برهانی برای همگرایی روش خط قاطع فراهم گردد.

خط در روش نیوتن و در روش خط قاطع را می توان همزمان استخراج کرد. در هر دو روش تابع  $f(x)$  در دو نقطه مثل،  $\alpha$  و  $\beta$ ، به وسیلهٔ يك خط مستقیم

$$p(x) = f(\alpha) + f[\alpha, \beta](x - \alpha)$$

که صفر آن یعنی

$$\hat{\xi} = \alpha - \frac{f(\alpha)}{f[\alpha, \beta]}$$

به عنوان تقریب بعدی برای ریشه دقیق  $f(x)$  مورد استفاده قرار می‌گیرد، درونیابی می‌شود. در روش خط قاطع با انتخاب  $\alpha = x_n$  و  $\beta = x_{n-1}$  مقدار  $\hat{\xi} = x_{n+1}$  را به دست می‌آوریم، در حالی که در روش نیوتن می‌گیریم  $\alpha = \beta = x_n$  در هر دو حالت به موجب (۳۷.۲) داریم

$$f(x) = f(\alpha) + f[\alpha, \beta](x - \alpha) + f[\alpha, \beta, x](x - \alpha)(x - \beta)$$

که این معادله به ازای جمیع مقادیر  $x$  برقرار است. اکنون اگر  $\xi = x$  انتخاب گردد که همان ریشه مطلوب است، در این صورت

$$0 = f(\xi) = f(\alpha) + f[\alpha, \beta](\xi - \alpha) + f[\alpha, \beta, \xi](\xi - \alpha)(\xi - \beta)$$

و بنابراین

$$f[\alpha, \beta](\xi - \alpha) = -f(\alpha) - f[\alpha, \beta, \xi](\xi - \alpha)(\xi - \beta)$$

حال از حل معادله فوق نسبت به  $\xi$ ، خواهیم داشت:

$$\xi = \alpha - \frac{f(\alpha)}{f[\alpha, \beta]} - \frac{f[\alpha, \beta, \xi]}{f[\alpha, \beta]}(\xi - \alpha)(\xi - \beta)$$

یا

$$\xi = \hat{\xi} - \frac{f[\alpha, \beta, \xi]}{f[\alpha, \beta]}(\xi - \alpha)(\xi - \beta) \quad (37.3)$$

اما معادله (۳۷.۳) را می‌توان برای به دست آوردن معادله‌های خطی در روشهای نیوتن و خط قاطع به کار گرفت. برای روش نیوتن قرار می‌دهیم  $\alpha = \beta = x_n$ ؛ و با توجه به اینکه  $e_j = \xi - x_j$  و  $\hat{\xi} = x_{n+1}$  از معادله (۳۷.۳) چنین به دست می‌آوریم:

$$e_{n+1} = -\frac{f[x_n, x_n, \xi]}{f[x_n, x_n]} e_n^2 \quad (38.3)$$

با توجه به اینکه روابط  $f[x_n, x_n] = f'(x_n)$  و  $f[x_n, x_n, \xi] = (1/2)f''(\eta_n)$  به ازای مقداری چون  $\eta_n$  بین  $x_n$  و  $\xi$  برقرار هستند، می‌توانیم (۳۸.۳) را دوباره به صورت زیر بنویسیم

$$e_{n+1} = -\frac{1}{2} \frac{f''(\eta_n)}{f'(x_n)} e_n^2 \quad (\text{الف } 38.3)$$

این معادله نشان می‌دهد که روش نیوتن به طور مربعی همگرا می‌شود، زیرا که  $e_{n+1}$  تقریباً



با مربع  $e_n$  متناسب است.

برای به دست آوردن معادله‌ای برای خطا در روش خط قاطع، در معادلهٔ (۳۷.۳) قرار می‌دهیم  $\alpha = x_n$  و  $\beta = x_{n-1}$  و خواهیم داشت

$$e_{n+1} = -\frac{f[x_{n-1}, x_n, \xi]}{f[x_{n-1}, x_n]} e_n e_{n-1} \quad (۳۹.۳)$$

این معادله نشان می‌دهد که خطا در  $(n+1)$  امین بارستن تقریباً متناسب با حاصلضرب خطا در مراحل  $n$  ام و  $(n-1)$  ام است. همچنین از آنجا که  $f[x_{n-1}, x_n, \xi] = (1/2)f''(\xi_n)$  و  $f[x_{n-1}, x_n] = f'(\eta_n)$  و  $\xi_n$  و  $\eta_n$  نقاطی بین  $x_{n-1}$  و  $x_n$  و  $\xi$  هستند. پس برای  $n$  به اندازهٔ کافی بزرگ معادلهٔ (۳۹.۳) به صورت زیر در خواهد آمد

$$e_{n+1} \approx -\frac{1}{2} \frac{f''(\xi)}{f'(\xi)} e_n e_{n-1} \quad (\text{الف } ۳۹.۳)$$

برای دقیقتر ساختن مفهوم مرتبه همگرایی، تعریف زیر آورده می‌شود:

تعریف ۱۰۳: مرتبه همگرایی. گیریم  $x_0, x_1, x_2, \dots$  دنباله‌ای باشد که به سمت عدد  $\xi$  همگرا می‌شود، و قرار می‌دهیم  $e_n = \xi - x_n$ . اگر عددی مانند  $p$  و ثابتی چون  $C \neq 0$  طوری وجود داشته باشد که داشته باشیم

$$\lim_{n \rightarrow \infty} \frac{|e_{n+1}|}{|e_n|^p} = C$$

آنگاه  $p$  مرتبه همگرایی این دنباله و  $C$ ، ثابت مجانبی خطا خوانده می‌شود. در حالت کلی برای بارست نقطهٔ ثابت که بر پایهٔ  $x = g(x)$  استوار است، داریم

$$\lim_{n \rightarrow \infty} \frac{|e_{n+1}|}{|e_n|} = |g'(\xi)|$$

لذا مرتبه همگرایی مساوی یک و ثابت مجانبی خطا برابر  $|g'(\xi)|$  است. از رابطهٔ (۳۸.۳ الف) برای روش نیوتن، در صورتی که  $f'(\xi) \neq 0$  داریم

$$\lim_{n \rightarrow \infty} \frac{|e_{n+1}|}{|e_n|^2} = \frac{1}{2} \left| \frac{f''(\xi)}{f'(\xi)} \right|$$

لذا بنا بر تعریف، مرتبه همگرایی برابر ۲ و ثابت مجانبی خطا برابر با  $|f''(\xi)/(2f'(\xi))|$  است.

برای تعیین مرتبه همگرایی روش خط قاطع، ابتدا ملاحظه می‌کنیم که به موجب

(الف ۳۹.۳) داریم

$$\lim_{n \rightarrow \infty} c_n = c_\infty = \frac{1}{2} \left| \frac{f''(\xi)}{f'(\xi)} \right| \quad \text{به‌ازای} \quad |e_{n+1}| = c_n |e_n e_{n-1}| \quad (۴۰.۳)$$

عدد  $p$  را طوری پیدا می‌کنیم که رابطه

$$\lim_{n \rightarrow \infty} \frac{|e_{n+1}|}{|e_n|^p} = C$$

به‌ازای یک مقدار غیر صفر  $C$  برقرار باشد.

اما از (۴۰.۳) نتیجه می‌گیریم

$$\frac{|e_{n+1}|}{|e_n|^p} = c_n |e_n|^{1-p} |e_{n-1}| = c_n \left( \frac{|e_n|}{|e_{n-1}|} \right)^{\alpha} \quad (۴۱.۳)$$

به‌شرط آنکه داشته باشیم  $\alpha = 1 - p$  و نیز  $\alpha p = -1$ ؛ یعنی در صورتی که

$$p - p^2 = \alpha p = -1$$

معادله  $0 = p^2 - p - 1 = 0$  دارای ریشه ساده و مثبت  $p = (1 + \sqrt{5})/2 = 1.618 \dots$  است. با این انتخاب  $p$  و  $\alpha = -1/p$  ملاحظه می‌کنیم که معادله (۴۱.۳) تعریفی است برای «بارست شبه نقطه ثابت»

$$y_{n+1} = c_n y_n^{-1/p}$$

که در آن

$$\lim_{n \rightarrow \infty} c_n = c_\infty \quad \text{و} \quad y_{n+1} = \frac{|e_{n+1}|}{|e_n|^p}$$

از اینجا نتیجه می‌شود که  $y_n$  به سمت نقطه ثابت معادله

$$x = c_\infty x^{-1/p}$$

که جواب آن  $c_\infty^{1/p}$  می‌باشد، همگرا می‌شود زیرا داریم  $p = 1/p + 1$ . این تساوی نشان می‌دهد که برای روش خط قاطع با  $p = 1.618 \dots$  داریم

$$\frac{|e_{n+1}|}{|e_n|^p} \approx \left| \frac{f''(\xi)}{2f'(\xi)} \right|^{1/p} \quad \text{به‌ازای مقدار بزرگ } n \quad (۴۲.۳)$$

یعنی مرتبه همگرایی روش خط قاطع برابر با  $p = 1.618 \dots$  است و ثابت مجانبی خطا برابر است با  $|f''(\xi)/(2f'(\xi))|^{1/p}$ . این بدان معنی است که روش خط قاطع سریعتر از روش معمول بارست نقطه ثابت همگرا می‌شود اما سرعت همگرایی آن کمتر از روش نیوتن است.

□ مثال ۵.۳: با استفاده از داده‌های مثال ۲.۳، درستی فرمولهای خطای (۳۹.۳ الف) و (۴۲.۳) را برای روش خط قاطع تحقیق کنید.

در مثال (۲.۳ الف) بارستهای روش خط قاطع ریشهٔ مثبت  $1 = x - x^3$  داده شده و در جدول زیر مفادیر  $|e_n|$  و  $|e_{n+1}|/|e_n e_{n-1}|$  به ازای  $n = 2, 3, 4$  محاسبه شده‌اند. فرض کنید که مقدار  $\xi$  تا هشت رقم اعشاری  $\xi = 1.03247180$  صحیح باشد.

$n$	$ e_n $	$ e_{n+1} / e_n e_{n-1} $	$ e_{n+1} / e_n ^{1.618}$
1	0.1580513		1.41684
2	0.0716060	1.1034669	0.88969
3	0.0124884	0.9705400	1.04325
4	0.0008679	0.9318475	0.90778
5	0.0000101		

اگر ثابت  $(\xi f'(\xi))/(f''(\xi))$  را مستقیماً حساب کنیم، مقدار آن برابر با  $0.93188000$  خواهد شد که با نسبت  $|e_{n+1}|/|e_n e_{n-1}|$  به ازای  $n = 4$  بسیار توافق دارد. □

می‌توان مستقیماً نشان داد که اگر  $f(\xi) = 0$ ،  $f'(\xi) \neq 0$ ، و  $f''(x)$  دوبار به‌طور پیوسته مشتق‌پذیر باشد، آنگاه داریم  $g'(\xi) = 0$ ، که

$$g(x) = x - \frac{f(x)}{f'(x)}$$

تابع بارسٔ نیوتن است. سپس بنا بر فرغ قضیهٔ ۱.۳ نتیجه می‌شود که اگر  $x_0$  «به اندازهٔ کافی نزدیک» به  $\xi$  انتخاب گردد، بارسٔ نیوتن همگرا می‌شود. اصطلاح «به اندازهٔ کافی نزدیک» خیلی دقیق تعریف نشده است، و در حقیقت روش نیوتن غالباً واگرا می‌شود و اگر هم همگرا شود به سمت ریشهٔ دیگری غیر از آن ریشه‌ای که خواسته شده همگرا خواهد شد. آنچه که ما می‌خواهیم برقراری شرایطی است که همگرایی را برای هر انتخاب بارسٔ اولیه در يك بازهٔ داده شده تضمین کند. مجموعه‌ای از چنین شرایط در قضیهٔ زیر داده شده است.

قضیهٔ ۲.۳ گیریم  $f(x)$  در بازهٔ بسته و معین  $[a, b]$  دوبار به‌طور پیوسته مشتق‌پذیر باشد و شرایط زیر برقرار

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

$$f'(x) \neq 0, x \in [a, b] \quad (ii)$$

$$f''(x) \geq 0 \quad \text{یا} \quad \leq 0 \quad x \in [a, b] \quad \text{تمام} \quad (iii)$$

$$\text{در دو نقطهٔ } a \text{ و } b \quad (iv)$$

$$\frac{|f(a)|}{|f'(a)|} < b-a \quad \frac{|f(b)|}{|f'(b)|} < b-a$$

در شرایط فوق روش نیوتن به سمت جواب یگانه  $\xi$  از معادله  $f(x) = 0$  در بازه  $[a, b]$  برای هر انتخابی از  $x_0 \in [a, b]$  همگرا می‌شود.

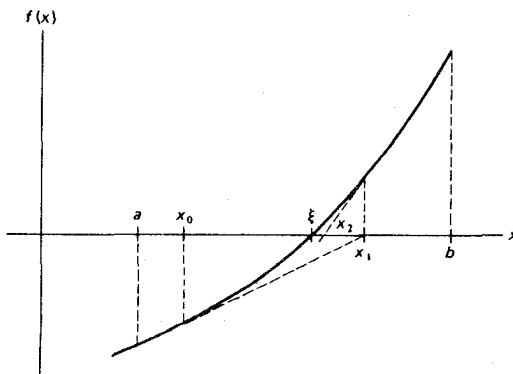
ذکر چند نکته درباره این شرایط مناسب به نظر می‌رسد.

شرایط (i) و (ii) تضمین می‌کنند که فقط یک جواب در  $[a, b]$  وجود دارد. شرط (iii) مبین آن است که نمودار  $f(x)$  تقریباً به سمت بالا یا تقریباً به سمت پایین دارد، و بعلاوه همراه با شرط (ii) ایجاب می‌کند که  $f'(x)$  روی  $[a, b]$  یکنوا باشد. علاوه بر اینها، شرط (iv) مبین این است که خط مماس بر منحنی در هر یک از دو نقطه  $a$  و  $b$  محور  $x$ ها را در بازه  $[a, b]$  قطع می‌کند. برهان این قضیه در اینجا داده نخواهد شد (به تمرین ۵.۳-۷ نگاه کنید). اما اشاره می‌کنیم که چرا این قضیه درست است. بی‌آنکه لطمه‌ای به کلیت وارد آید، فرض می‌کنیم که  $f(a) < 0$ ، پس دو حالت می‌توانیم تشخیص دهیم:

$$\text{حالت (الف)} \quad f''(x) \geq 0$$

$$\text{حالت (ب)} \quad f''(x) \leq 0$$

اگر به جای  $f$ ،  $-f$  — بگذاریم حالت (ب) به حالت (الف) تبدیل می‌شود. بنا بر این بررسی حالت (الف) کفایت می‌کند. در اینجا نمودار  $f(x)$ ، به گونه‌ای است که در شکل ۶.۳ داده شده است. از این نمودار آشکار است که به ازای  $\xi > x_0$ ، بارستن‌های حاصل به‌طور یکنوا به سمت  $\xi$  کاهش می‌یابند، در حالی که به ازای  $\xi < x_0$ ،  $x_1$  بین  $\xi$  و  $b$  قرار می‌گیرد و آنگاه بارستن‌های بعدی به‌طور یکنوا به سمت  $\xi$  همگرا می‌شوند.

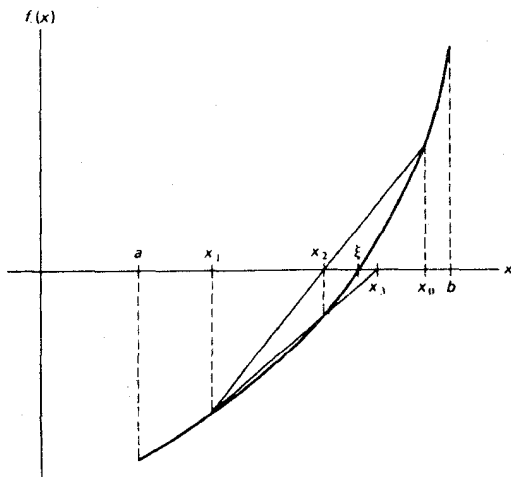


شکل ۶.۳ همگرایی نیوتن.

□ مثال ۶.۳: بازه‌ای پیدا کنید که کوچکترین ریشهٔ مثبت تابع  $f(x) = e^{-x} - \sin x$  را در برگیرد و در شرایط قضیهٔ ۲.۳ برای همگرایی روش نیوتن صدق نماید.

با توجه به عبارت  $f(x) = e^{-x} - \sin x$ ، داریم  $f'(x) = -e^{-x} - \cos x$  و  $f''(x) = e^{-x} + \sin x$ . بازهٔ خود را به صورت  $[a, b] = [0, 1]$  انتخاب می‌کنیم. پس چون  $f(0) = 1$ ،  $f(1) = -0.۴۷$ ، داریم  $f(a)f(b) < 0$  لذا شرط (i) برقرار است. از آنجا که برای کلیهٔ مقادیر  $x \in [0, 1]$ ، نامساوی  $f'(x) < 0$  برقرار است، شرط (ii) برقرار بوده و چون به ازای جمیع مقادیر  $x \in [0, 1]$ ، داریم  $f''(x) > 0$ ، شرط (iii) نیز برقرار می‌باشد. در نهایت چون  $f(0) = 1$ ،  $f'(0) = -2$ ، داریم  $|f(0)| / |f'(0)| = 1/2 < b - a = 1$ ، و چون  $f(1) = -0.۴۷$ ،  $f'(1) = -0.۹۰$ ، داریم  $|f(1)| / |f'(1)| = 0.۵۲ < 1$  که در شرط (iv) صدق می‌کند. بنا بر این باروست نیوتن برای هر انتخابی از  $x_0$  در  $[0, 1]$  همگرا می‌شود. □

همچنین شرایط قضیهٔ ۲.۳ برای اثبات همگرایی روش خط قاطع کافی است، ولو اینکه حالات ممکن همگرایی با حالات همگرایی نیوتن کاملاً متفاوت باشد. اگر دوباره فرض کنیم که در بازهٔ  $[a, b]$ ، همان طوری که در شکل (۶.۳ الف) نشان داده شده است،  $f'(x) > 0$  و  $f''(x) \geq 0$ ، آنگاه اساساً دو حالت همگرایی وجود خواهد داشت که بستگی به گزینش نقاط اولیهٔ  $x_0$  و  $x_1$  دارند. در حالت اول که ساده‌تر است، اگر  $x_0$  و  $x_1$  در بازهٔ  $[a, b]$  برگزیده شوند، مانند روش نیوتن همگرایی از طرف راست یکنوا خواهد بود. دانشجویان می‌توانند این امر را برای روش نیوتن با رسم منحنیهای نمونه‌ای که در شرایط قضیهٔ ۲.۳ صدق می‌کنند بررسی نمایند.



شکل ۶.۳ (الف) همگرایی خط قاطع.

از سوی دیگر اگر یکی از نقاط، مثلاً  $x_0$  را، در بازه  $[a, b]$  و نقطه دیگر  $x_1$  را در بازه  $[a, \xi]$  انتخاب کنیم، آنگاه بارست بعدی  $x_2$  نیز در بازه  $[a, \xi]$  خواهد بود و بارستن  $x_3$  در طرف راست  $\xi$  قرار خواهد گرفت. در این حالت دوباره دو نقطه بارست بیابیم  $x_4$  و  $x_5$  را خواهیم داشت که ریشه  $\xi$  را در برمی گیرند، و تمامی دنباله تکرار خواهد شد. بنا بر این همگرایی به اصطلاح، حالت والسی<sup>۱</sup> دارد، یعنی در بازه ای روی می دهد که یک بارست در یک طرف آن و دو بارست در طرف دیگر آن است، برای نشان دادن این نوع همگرایی به شکل (۶.۳ الف) نگاه کنید.

□ مثال ۷.۰۳: حالت همگرایی روش خط قاطع را هنگامی که برای حل تابع  $f(x) = e^x - 3$  به کار برده می شود، بررسی کنید.

بدیهی است که به ازای جميع مقادیر  $x$ ، نامساویهای  $f'(x) > 0$  و  $f''(x) > 0$  برقرارند. بعلاوه، برای مثال، در بازه  $[0, 5]$  شرایط نقاط انتهایی قضیه ۲.۳ برقرارند. بنا بر این  $f(x)$  یک ریشه در آن بازه دارد که عبارت است از  $x_0 = \ln 3 = 1.098612289$  و اگر  $\xi < 0$  و  $x_0 > \xi$  و  $x_1 = 5$  انتخاب شوند، آنگاه انتظار همگرایی می رود. پس بارستن های زیر به دست خواهند آمد، که حالت همگرایی به اصطلاح والسی را نشان می دهند:

$$x_0 = 0$$

$$x_1 = 5$$

$$x_2 = 0.06783654900$$

$$x_3 = 0.1324350609$$

$$x_4 = 1.13459843$$

$$x_5 = 0.7584862650$$

$$x_6 = 0.9868082782$$

$$x_7 = 1.119119918$$

$$x_8 = 1.097448653$$

$$x_9 = 1.098600396$$

$$x_{10} = 1.098612296$$

اگر بگیریم  $\xi > x_0 = 3$  و  $\xi > x_1 = 2$ ، آنگاه بارستهای زیر را خواهیم داشت

$$x_0 = 3$$

$$x_1 = 2$$

$$x_2 = 1.9654309240$$

$$x_3 = 1.9297433158$$

$$x_4 = 1.9147393259$$

$$x_5 = 1.9103265406$$

$$x_6 = 1.9098724772$$

$$x_7 = 1.9098612550$$

$$x_8 = 1.9098612289$$

□

که معرف حالت یکنوا بودن همگرایی هستند.

از دیدگاه محاسباتی، دقت حاصل با روش نیوتن بستگی به دقت عمل در محاسبهٔ  $f(x)/f'(x)$  دارد. مثلاً ممکن است که در همسایگی ریشه، مقدار  $f'(x)$  صفر نباشد اما بسیار کوچک باشد. در این صورت اگر  $f(x)$  با خطا همراه باشد انتظار می‌رود که مقدار این خطا در محاسبهٔ  $f(x)/f'(x)$  بزرگ شود. در چنین مواردی به دست آوردن دقت مناسب دشوار خواهد بود.

در روش نیوتن دو اشکال مهم وجود دارد. اولاً برای اینکه بارستها به سمت  $\xi$  که ریشهٔ  $f(x)$  است همگرا باشند، باید از ریشه‌ای که «به اندازهٔ کافی» به  $\xi$  نزدیک است شروع کنیم. (به تمرینهای ۵-۳ و ۶-۳ و همچنین ۴-۳ و ۵-۳ نگاه کنید.) از آنجا که معمولاً مقدار  $\xi$  معلوم نیست انجام این کار در عمل مشکل است، مگر آنکه با استفاده از روشهای دیگر یک تخمین بهتر  $\xi$  برای  $\xi$  برآورد شود. مثلاً اگر با استفاده از روش تنصیف بازه یا سایر روشهای بارستی، تقریب  $\xi$  برای  $\xi$  که تا دو یا سه رقم اعشاری صحیح باشد محاسبه شود، آنگاه می‌توان روش نیوتن را با تقریب بهتر  $\xi = x_0$  شروع کرد و پس از انجام دو یا سه بارست سریعاً تقریب صحیحتری از  $\xi$  به دست آورد. بدین ترتیب، روش نیوتن غالباً برای بهبود بخشی یک تقریب مناسب از ریشه که خود از روش دیگری به دست آمده، به کار گرفته می‌شود.

اشکال دوم روش نیوتن لزوم محاسبهٔ  $f'(x)$  است. در برخی موارد ممکن است  $f'(x)$  صریحاً در دست نباشد و یا حتی اگر مقدار  $f'(x)$  قابل ارزیابی باشد این کار ممکن است

مستلزم محاسبات پرکار و قابل ملاحظه‌ای باشد. در حالت اخیر ممکن است مقدار  $f'(x_n)$  در هر  $k$  مرحله یکبار از الگوریتم نیوتن محاسبه شود و در هر مرحله آخرین مقدار محاسبه شده به کار برده شود. اما در برخورد با هر دو حالت بهتر است که از روش خط قاطع به جای روش نیوتن استفاده شود.

در روش خط قاطع از مقادیر مختلف  $f(x)$  استفاده می‌شود و در هر بارست فقط محاسبه یک مقدار تابع لازم است، در صورتی که در روش نیوتن در هر مرحله به محاسبه مقدار دو تابع\* نیاز است. از طرف دیگر اگر روش خط قاطع همگرا باشد سرعت همگرایی آن به اندازه روش نیوتن نیست گرچه که معمولاً خیلی سریعتر از همگرایی خطی است. سرعت بیشتر همگرایی روش نیوتن نسبت به روش خط قاطع در مثال ۲.۳ نشان داده شده است.

در این فصل، برای پیدا کردن ریشه‌های توابع شش الگوریتم را در نظر گرفته‌ایم. در مقایسه بین الگوریتمها، برای استفاده کامپیوتری آنها، ملاکهای مختلفی را می‌باید در نظر گرفت، که مهمترین آنها اطمینان از همگرایی، میزان همگرایی و کارایی محاسباتی است. همواره نمی‌توان یک روش را برتر از روش دیگر دانست. برای مثال، درحالی که روش تنصیف از نظر سرعت همگرایی، آهسته‌است ولی اگر این روش به‌طور صحیح به کار گرفته شود، روش مطمئنی خواهد بود، درحالی که روش نیوتن غالباً واگراست، مگر آنکه تقریب اولیه با دقت انتخاب شود. اصطلاح «کارایی محاسباتی» که در بالا به کار برده شد، سعی دارد که مقدار کار لازم برای به دست آوردن دقت عمل مفروضی را در نظر بگیرد. اگر چه به‌طور کلی روش نیوتن نسبت به روش خط قاطع، سریعتر همگرا می‌شود، ولی معمولاً به اندازه این روش کارایی ندارد. زیرا روش نیوتن در هر بارست، نیاز به محاسبه هر دو تابع  $f(x)$  و  $f'(x)$  دارد. در مواردی که  $f'(x)$  در دسترس و به آسانی قابل محاسبه باشد، روش نیوتن نسبت به روش خط قاطع کارا تر است، اما برای کارهای معمولی، روش تنصیف معمولاً کارا تر و کار برد آن ارجح است.

الگوریتمهای ۱.۳ تا ۳.۳، این مزیت را دارند که ریشه را محصور کرده و بنا بر این تضمین کننده کرانهای خطا روی ریشه هستند. از میان این الگوریتمها، الگوریتم ۲.۳ (تصحیح خطا) هرگز نباید به کار رود. زیرا این الگوریتم فاصله‌ای را که ریشه را در بر می‌گیرد کوچکتر نمی‌کند. به‌طور کلی از بین این سه الگوریتم، روش تصحیح خطای اصلاح شده (الگوریتم ۳.۳) ارجح است.

بارست نقطه ثابت زمانی کار است که مانند روش نیوتن به‌طور مربعی همگرا شود. در حالت کلی بارست نقطه ثابت تنها به‌طور خطی همگرا می‌شود، از این رو رقابت واقعی بین این روش و روش خط قاطع یا روش تصحیح خطای اصلاح شده، وجود ندارد. حتی بسا برونمایی مکرر، مانند بارست استفنسن، الگوریتم ۸.۳، همگرایی فقط در بهترین حالت

\* دو تابع مورد نظر  $f(x)$  و  $f'(x)$  می‌باشند. م.



مربعی خواهد بود. از آنجا که هر مرحله از بارست استغفنسن معادل با دو ارزیابی تابع  $g(x)$  بارست است؛ بنابراین بارست استغفنسن با روش نیوتن قابل قیاس است. اما چون قسمت برونیابی در یک مرحله از بارست استغفنسن، با یک مرحله از روش خط قاطع که برای تابع  $f(x) = x - g(x)$  به کار گرفته شده مشا به است، به نظر می رسد که کار اثر است که روش استغفنسن را کنار بگذاریم و روش خط قاطع را برای تابع  $f(x) = x - g(x)$  به کار گیریم.

اساساً، هدف اصلی از بحث درباره بارست نقطه ثابت این بوده است که مدل ساده ای برای یک شیوه بارستی که بتوان به آسانی تجزیه و تحلیل کرد، به دست آورد. بینش حاصل در مبحث مربوط به چند معادله چند مجهولی که موضوع فصل ۵ است، بسیار مفید خواهد بود.

### تمرین

۱-۵۰۳ بنا بر تعریف بارست نقطه ثابت با تابع بارست  $g(x)$ ، باید خطا در بارستن  $m$  در معادله زیر صدق کند

$$e_n = \xi - x_n = g(\xi) - g(x_{n-1})$$

در متن نشان داده شد که اگر  $g'(\xi) = 0$  و  $g''(x) = 0$  در  $x = \xi$  پیوسته باشد، بارست  $x = g(x)$  به طور مربعی همگرا خواهد بود. شرایطی را که تحت آنها، می توان انتظار داشت که یک بارست همگرایی مکعبی داشته باشد پیدا کنید.

۲-۵۰۳ نشان دهید که در روش نیوتن، اگر  $f(\xi) = 0$  و  $f'(\xi) \neq 0$ ، و اگر  $f(x)$  دو بار به طور پیوسته مشتق پذیر باشد، آنگاه  $g'(\xi) = 0$ . همچنین نشان دهید که

$$g''(\xi) = \frac{f''(\xi)}{f'(\xi)}$$

۳-۵۰۳ برای هر یک از توابع زیر، بازه ای را که کوچکترین ریشه مثبت را در برگیرد تعیین کنید و نشان دهید که شرایط قضیه ۲۰۳ برقرار است

$$e^{-x} - x = 0 \quad (\text{الف})$$

$$x^3 - x - 1 = 0 \quad (\text{ب})$$

$$e^{-2x} - \cos x = 0 \quad (\text{پ})$$

۴-۵۰۳ هر یک از معادله های تمرین ۳-۵۰۳ را با دو روش خط قاطع و روش نیوتن حل و جوابها را مقایسه کنید.

۵-۵۰۳ اگر  $x = \xi$  یک ریشه  $f(x)$  از مرتبه ۲ باشد، آنگاه  $f'(\xi) = 0$  و  $f''(\xi) \neq 0$ .

نشان دهید که در این حالت، روش نیوتن به طور مربعی همگرا نخواهد بود [یعنی نشان دهید که  $g'(\xi) = 1/2 \neq 0$ ]. همچنین نشان دهید که اگر  $f'(\xi) = 0$ ،  $f''(\xi) \neq 0$  و  $f'''(x) \neq 0$  در همسایگی  $\xi$  پیوسته باشد، بارست

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

به طور مربعی همگرا می‌شود. {داهنمایی: برای محاسبه  $g'(\xi)$ ، از رابطه زیر و قانون لوپیتال، استفاده کنید.}

$$\lim_{x \rightarrow \xi} \frac{f(x)f''(x)}{[f'(x)]^2} = \lim_{x \rightarrow \xi} \frac{f(x)}{[f'(x)]^2} \lim_{x \rightarrow \xi} f''(x)$$

۶-۵.۳ با استفاده از روش نیوتن، نزدیکترین ریشه معادله زیر به عدد ۱۰۰ را محاسبه کنید:

$$x = \tan x$$

(توجه: چنانچه  $x_0$  با دقت کافی انتخاب نشود، روش نیوتن يك دنباله واگر تولید خواهد کرد.)

۷-۵.۳ جزئیات اثبات قضیه ۲.۳ را بیان کنید.

۸-۵.۳ ثابت کنید که در شرایط قضیه ۲.۳، به ازای جميع مقادیر  $x_0$  و  $x_1$  در بازه  $[a, b]$ ، روش خط قاطع همگراست. همچنین نشان دهید که بسته به محل دو بارست پیاپی حالت همگرایی یا یکنوا یا الوسی است. [داهنمایی: از معادله خط (۳۹.۳) استفاده کنید و مانند اثبات همگرایی در روش نیوتن کار را ادامه دهید.]

۹-۵.۳ نشان دهید که اگر  $x = \xi$  يك ریشه مکرر مرتبه  $m$  از تابع  $f(x)$  باشد، آنگاه بارست

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

در شرایط پیوستگی مناسب، به طور مربعی همگرا خواهد بود.

### ۶.۳ معادلات بسجمله‌ای: ریشه‌های حقیقی

گرچه معادلات بسجمله‌ایها را می‌توان با هر يك از روشهای بارستی که قبلا بحث شده حل

نمود، اما از آنجا که کراراً در کار بردهای فیزیکی پدید می‌آیند نیاز به بحثی خاص دارند. بخصوص ما چند الگوریتم کارآمد را برای پیدا کردن ریشه‌های حقیقی و همثافت بسجمله‌ایها ارائه خواهیم داد. در این بخش به دست آوردن اطلاعات (معمولاً تقریبی) درباره‌ی جای ریشه‌های بسجمله‌ایها را مورد بحث قرار می‌دهیم و سپس روش نیوتن را برای پیدا کردن ریشه مثبت بسجمله‌ایها ذکر خواهیم کرد.

یک بسجمله‌ای از درجه (دقیق)  $n$  معمولاً به شکل زیر نوشته می‌شود:

$$p(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n \quad a_n \neq 0 \quad (۴۳.۳)$$

قبل از شروع به بحث درباره‌ی روشهای ریشه‌یابی، چند تذکر در باره‌ی ریشه‌های بسجمله‌ای ممکن است مناسب باشد. به ازای  $n=2$ ،  $p(x)$  یک بسجمله‌ای از درجه دوم است و البته ریشه‌های آن را می‌توان صریحاً با استفاده از دستور درجه دوم، همان طوری که در فصل یک انجام شد، پیدا کرد. دستورهای مشابه ولی پیچیده‌تر دقیقی برای بسجمله‌ایها از درجه ۳ و ۴ وجود دارند. اما به ازای  $n \geq 5$  به طور کلی دستورهای صریحی برای ریشه‌ها وجود ندارند. بنابراین مجبوریم که روشهای یارستی را برای محاسبه‌ی ریشه بسجمله‌ایها مورد بررسی قرار دهیم. همه‌ی روشهایی را که در این فصل بررسی شدند می‌توان برای پیدا کردن ریشه‌های حقیقی بسجمله‌ایها به کار گرفت و برخی از آنها را می‌توان برای محاسبه‌ی ریشه‌های همثافت تعمیم داد. غالباً ما به پیدا کردن همه‌ی ریشه‌های بسجمله‌ای علاقه‌مندیم. عده‌ای از قضایای درس جبر برای رده‌بندی و جایابی انواع ریشه‌های بسجمله‌ایها مفیدند.

در ابتدا قضیه‌ی اساسی جبر (قضیه‌ی ۱۰.۱ را ببینید) را داریم که بیان می‌دارد هر بسجمله‌ای از درجه  $n$  با شرط  $a_n \neq 0$ ، دقیقاً  $n$  ریشه حقیقی یا همثافت دارد، به شرطی که ریشه‌های مکرر از مرتبه  $m$ ،  $m$  مرتبه شمرده شوند. اگر کلیه ضرایب  $a_p$  از بسجمله‌ای  $p(x)$ ، حقیقی باشند و اگر  $z = a + ib$  یک ریشه باشد، آنگاه عدد  $\bar{z} = a - ib$  نیز یک ریشه خواهد بود. یک روش مناسب برای تعیین تعداد ریشه‌های حقیقی یک بسجمله‌ای با ضرایب حقیقی، قاعده‌ی علامتهای دکارت است. این قانون بیان می‌کند که تعداد ریشه‌های مثبت یک بسجمله‌ای،  $n_p$ ، نایزتر از تعداد تغییر علامتهای ضرایب  $p(x)$ ،  $v$ ، است. بعلاوه تفاضل  $n_p - v$  یک عدد زوج است. برای تعیین تعداد تغییر علامتها، کافی است که تعداد تغییر علامتهای ضرایب غیر صفر  $p(x)$  را بشماریم. بنابراین اگر داشته باشیم  $p(x) = x^4 + 2x^2 - x - 1$ ، چون تعداد تغییر علامتها یکی است، بنابراین قانون دکارت،  $p(x)$  حداکثر دارای یک ریشه مثبت است، اما چون  $n_p - v$  می‌باید یک عدد زوج غیر منفی باشد، بنا بر این معادله فوق دقیقاً دارای یک ریشه مثبت است. به طور مشابه، تعداد ریشه‌های منفی حقیقی  $p(x)$  برابر با تعداد تغییر علامتها در ضرایب بسجمله‌ای

$$p(-x) = -x^4 - 2x^2 - x - 1$$

می باشد و چون تغییر علامتی در  $p(-x)$  وجود ندارد، بنا بر این ریشه منفی حقیقی نخواهیم داشت

□ مثال ۸.۳: تا حدی که می توانید در بساطه ریشه های حقیقی بسجمله ای زیر اطلاعاتی به دست دهید

$$p(x) = x^4 - x^3 - x^2 + x - 1$$

چون ضرایب  $p(x)$  دارای سه تغییر علامت می باشند بنا بر این یا سه ریشه مثبت حقیقی یا یک ریشه مثبت حقیقی وجود دارد، اکنون  $p(-x) = x^4 + x^3 - x^2 - x - 1$  را در نظر می گیریم و چون تنها یک تغییر علامت داریم بنا بر این می باید یک ریشه منفی حقیقی وجود داشته باشد. بنا بر این، یا سه ریشه مثبت حقیقی و یک ریشه منفی حقیقی و یا یک ریشه مثبت حقیقی، یک ریشه منفی حقیقی و دو ریشه مزدوج همگفت، داریم. □

اکنون به ذکر چندین قضیه، که کرانه های ریشه های بسجمله ایها را معین می کنند، می پردازیم. یکی از این قضایا این است که اگر  $p(x)$  یک بسجمله ای مانند (۴۳.۳) با ضرایب  $a_k$  باشد، آنگاه  $p(x)$  در داخل دایره ای که با  $\min\{\rho_1, \rho_n\}$  تعریف می شود حداقل یک ریشه دارد که:

$$\rho_1 = n \frac{|a_0|}{|a_1|} \quad (44.3)$$

و

$$\rho_n = \sqrt[n]{\frac{|a_0|}{|a_n|}}$$

□ مثال: اگر یک بسجمله ای به شکل

$$p(x) = x^5 - 377x^4 + 774x^3 - 1078x^2 + 1078x - 678 \quad (45.3)$$

باشد، آنگاه  $a_0 = -678$ ،  $a_1 = 1078$ ،  $a_5 = 1$  از رابطه (۴۴.۳) به دست می آوریم

$$\rho_1 = 5 \times \frac{678}{1078} = 314 \dots$$

$$\rho_n = \sqrt[5]{\frac{678}{1}} = 1746 \dots$$

بنابراین، حداقل یک ریشه حقیقی با همگفت در داخل دایره  $|x| \leq 1746 \dots$  وجود دارد.

در بخش بعدی، این بسجمله‌ای (۴۵.۳) را با جزئیات بیشتر مورد بررسی قرار و نشان می‌دهیم که ریشه‌های آن تحقیقاً برابر  $i \pm 1$ ،  $i \pm \sqrt{2}$  و  $1 \pm \sqrt{2}$  هستند.  $\square$

دومین قضیه مفیدی که به کوشی منسوب است، به ما امکان می‌دهد که کرانه‌های ریشه‌های  $p(x)$  را به شرح زیر مشخص سازیم. اگر  $p(x)$  يك بسجمله‌ای مانند (۴۳.۳) باشد، دو بسجمله‌ای جدید به صورت زیر تعریف می‌کنیم:

$$P(x) = |a_n| x^n - |a_{n-1}| x^{n-1} - \dots - |a_0| \quad (۴۶.۳)$$

$$Q(x) = |a_n| x^n + |a_{n-1}| x^{n-1} + \dots + |a_1| x - |a_0| \quad (\text{الف } ۴۶.۳)$$

بنابر قاعده علامتهای دکارت بسجمله‌ای (۴۶.۳) دقیقاً يك ریشه مثبت حقیقی  $R$  و بسجمله‌ای (الف ۴۶.۳)، دقیقاً يك ریشه مثبت حقیقی  $r$  دارد. در این حال قضیه کوشی بیان می‌کند که همه ریشه‌های  $p(x)$  در ناحیه طوقی  $r \leq |x| \leq R$  قرار دارند.

$\square$  مثال: مجدداً بسجمله‌ای (۴۵.۳) را در نظر می‌گیریم. پس داریم

$$P(x) = x^5 - 37x^4 - 74x^3 - 108x^2 - 108x - 68$$

$$Q(x) = x^5 + 37x^4 + 74x^3 + 108x^2 + 108x - 68$$

که ریشه مثبت آنها به ترتیب برابر با  $r = 0.63\dots$  و  $R = 5.6\dots$  است. از این رو کلیه ریشه‌های  $p(x)$  می‌باید در رابطه زیر صدق کنند

$$\square \quad 0.63\dots < |x| \leq 5.6\dots$$

آخرین قضیه از این نوع بیان می‌دارد که اگر  $p(x)$  يك بسجمله‌ای به شکل (۴۳.۳) باشد و اگر داشته باشیم

$$r := 1 + \max_{0 \leq k \leq n-1} \left| \frac{a_k}{a_n} \right|$$

آنگاه هر ریشه‌ای از  $p(x)$  در ناحیه مستدیری قرار دارد که با رابطه  $|x| \leq r$  تعریف شده است.

$\square$  مثال: اگر بسجمله‌ای (۲.۳) را به صورت

$$p(x) = x^3 - x - 1$$

در نظر بگیریم، آنگاه  $r = 1 + 1/1 = 2$ ، لذا همه ریشه‌های  $p(x)$  درون قرصی به شعاع ۲ و مرکز مبدأ مختصات قرار دارند. در بخش ۱.۳ دیدیم که يك ریشه حقیقی  $1.324\dots$  بود. دو ریشه دیگر مختلط‌اند ولی باز درون دایره  $|x| \leq 2$  قرار دارند.  $\square$

اکنون به بررسی روشهای بارستی برای یافتن ریشه‌های حقیقی بسجمله‌ایها می‌پردازیم. در هر روش بارستی، غالباً مقدار بسجمله‌ای می‌باید ارزیابی شود، و این کار باید تا حد ممکن به قدر کافی انجام پذیرد. همان گونه که در فصل ۲ نشان داده شد، مؤثرترین روش برای ارزیابی یک بسجمله‌ای روش ضرب تو در تو است که در الگوریتم ۱۰۲ شرح داده شد. در الگوریتم ۱۰۲، فرض شده بود که بسجمله‌ای به شکل نیوتنی (۳.۲) با مراکز  $c_1, \dots, c_n$  داده شده است. اگر این مراکز، همگی برابر صفر باشند، شکل نیوتنی (۳.۲) به شکل استاندارد توانی (۴۳.۳) بدل می‌شود. اکنون اگر یک نقطه معین  $z$  را به ما بدهند، الگوریتم ۱۰۲، برای تعیین  $p(z)$  به شکل خاص زیر در خواهد آمد

$$a'_n := a_n$$

$$a'_{n-1} := a_{n-1} + za'_n$$

...

$$a'_1 := a_1 + za'_2$$

$$p(z) = a'_0 := a_0 + za'_1 \quad (47.3)$$

کمیت‌های کمکی  $a'_1, a'_2, \dots, a'_n$  جداگانه مورد علاقه ما هستند زیرا که با دوباره صفر گذاردن تمام  $c_k$ ها در (۴.۲) داریم

$$\begin{aligned} p(x) &= a'_0 + (x-z) \{a'_1 + a'_2x + a'_3x^2 + \dots + a'_nx^{n-1}\} \\ &= a'_0 + (x-z)q(x) \end{aligned} \quad (48.3)$$

از این رو،  $a'_1, \dots, a'_n$  ضرایب بسجمله‌ای خارج قسمت  $q(x)$  هستند که از تقسیم  $p(x)$  بر بسجمله‌ای خطی  $(x-z)$  به دست می‌آیند و  $a'_0$  باقیمانده این تقسیم است. به ویژه اگر در (۴۸.۳) قرار دهیم  $x=z$ ، از نو به دست می‌آوریم  $p(z) = a'_0$ .

□ مثال ۹.۳: تبدیل یک عدد صحیح دودویی به یک عدد صحیح دهدهی. در بخش ۱۰.۱، برای تبدیل عدد صحیح دودویی به عدد صحیح دهدهی الگوریتم ۱۰.۱ ارائه گردیده بنا بر قرارداد، عدد صحیح دودویی

$$\alpha = (a_n a_{n-1} a_{n-2} \dots a_0)_2$$

که در آن  $a_i$  برابر با صفر یا یک است، نمایانگر عدد زیر است

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

بنابراین معادل دهدهی این عدد را می‌توان با ارزیابی بسجمله‌ای

$$p(x) = a_0 + a_1 x + \dots + a_n x^n$$

به ازای  $x = 2$  و با استفاده از الگوریتم ضرب تودرتو (۱۰۲) به دست آورد. این امر نشان می‌دهد که الگوریتم ۱۰۱، حالت خاصی از الگوریتم ۱۰۲ است. به عنوان مثال، عدد صحیح دودویی  $\alpha = (110011)_2$  به شرح زیر به هم ارز دهدهی خود تبدیل می‌شود

$$a'_5 = a_5 = 1$$

$$a'_4 = a_4 + 2a'_5 = 3$$

$$a'_3 = a_3 + 2a'_4 = 6$$

$$a'_2 = a_2 + 2a'_3 = 12$$

$$a'_1 = a_1 + 2a'_2 = 25$$

□

$$\alpha = a'_0 = a_0 + 2a'_1 = 51$$

هدف مستقیم ما تطابق روش نیوتن برای یافتن ریشه‌های حقیقی بسجمله ایهاست. بدین منظور باید بتوانیم که نه فقط مقدار  $p(x)$  بلکه مقدار  $p'(x)$  را نیز محاسبه کنیم. برای محاسبه  $p'(x)$  در نقطه  $x = z$ ، از (۴۸.۳) نسبت به  $x$  مشتق می‌گیریم و داریم

$$p'(x) = q(x) + q'(x)(x - z)$$

بنابراین با قراردادن  $x = z$  داریم

$$p'(z) = q(z)$$

از آنجا که  $q(x)$  خود نیز یک بسجمله ای است که ضرایب آن معلوم اند، می‌توانیم الگوریتم ۱۰۲ را یکبار دیگر برای محاسبه  $q(z)$  و نتیجتاً برای  $p'(z)$  به کار گیریم. این عمل الگوریتم زیر را به دست می‌دهد.

**الگوریتم ۹۰۳:** روش نیوتن برای پیدا کردن ریشه‌های حقیقی بسجمله ایها.  $n + 1$  ضریب  $a_0, \dots, a_n$  از بسجمله ای  $p(x)$  در (۴۳.۳) و یک نقطه شروع  $x_0$  داده شده‌اند.

For  $m = 0, 1, \dots$ , until satisfied, do:

$$z := x_m, a'_n := a_n, a''_n := a'_n$$

For  $k = n - 1, \dots, 1$ , do:

$$a'_k := a_k + za'_{k+1}$$

$$a''_k := a'_k + za''_{k+1}$$

$$a'_0 := a_0 + za'_1$$

$$x_{m+1} := x_m - a'_0/a''_1$$

□ مثال ۱۰.۳: کلیه ریشه‌های معادلهٔ بسجمله‌ای  $0 = x^3 + x - 3 = p(x)$  را پیدا کنید. این معادله یک ریشهٔ حقیقی و دو ریشهٔ همثافت دارد. از آنجا که  $p(1) = -1$  و  $p(2) = 7$ ، ریشهٔ حقیقی باید بین  $x = 1$  و  $x = 2$  باشد.  $x_0$  را مساوی ۱٫۱ انتخاب می‌کنیم و الگوریتم ۹.۳ را به کار می‌بندیم، با انجام محاسبات به‌توسط ماشین حساب دستی و احتساب پنج رقم اعشاری نتایج زیر حاصل می‌شود.

$x_0 = 1.1$				$x_1 = 1.1 - (-0.569)/4.63 = 1.22289$	
$k$	$a_k$	$a'_k$	$a''_k$	$a'_k$	$a''_k$
3	1	1	1	1	1
2	0	1.1	2.2	1.22289	2.44578
1	1	2.21	4.63	2.49546	5.48638
0	-3	-0.569		0.05167	
$x_2 = 1.22289 - a'_0/a''_1 = 1.21347$				$x_3 = 1.21347 - a'_0/a''_1 = 1.21341$	
$k$	$a_k$	$a'_k$	$a''_k$	$a'_k$	$a''_k$
3	1	1	1	1	1
2	0	1.21347	2.42694	1.21341	2.42682
1	1	2.42751	5.41753	2.47236	5.41709
0	-3	0.000317		-0.00001	

توجه کنید که  $a'_0$  به سمت ریشه نزدیک و  $a'_k$ ‌ها همگرا می‌شوند. با توجه به دقت عملی که در نظر گرفته شده است، بهسازی جواب و یا بهسازی  $a'_k$  ممکن نیست. بنابراین  $x_3 = 1.21341$  را که تا پنج رقم بسا معنی صحیح است به‌عنوان ریشهٔ مطلوب قبول می‌کنیم. برای محاسبهٔ بقیهٔ ریشه‌های همثافت، از دستور درجهٔ دوم برای حل بسجمله‌ای زیر استفاده می‌کنیم:

$$x^2 + a'_2 x + a'_1 = x^2 + 1.21341 x + 2.47236 = 0$$

که نتایج زیر حاصل می‌شود

$$x = \frac{-a'_2 \pm (a'^2_2 - 4a'_1)^{1/2}}{2}$$

$$= \frac{-1.21341 \pm 2.490122i}{2} = -0.606705 \pm 1.245061i$$

به منظور مقایسه، دربخش ۷.۳ ریشه‌های این بسجمله‌ای دوباره با روشهای خاص مربوط به محاسبهٔ ریشه‌های همثافت محاسبه خواهد شد. □

□ مثال ۱۱.۳: ریشهٔ مثبت حقیقی بسجمله‌ای زیر را محاسبه کنید



$$x^5 - 37x^4 + 74x^3 - 108x^2 + 108x - 64 = 0$$

به آسانی می‌توان تحقیق کرد که ریشه مورد نظر بین ۱ و ۲ قرار دارد. بنابراین  $x_0 = 1.5$  را انتخاب می‌کنیم. برنامه فورترن و نتایج حاصله از کامپیوتر در زیر داده شده‌اند. ریشه دقیق برابر با ۱.۷ و نتایج حاصله از کامپیوتر تا هشت رقم اعشاری صحیح است.

### برنامه فورترن برای مثال ۱۱.۳

```

C NEWTON'S METHOD FOR FINDING A REAL ZERO OF A CERTAIN POLYNOMIAL.
C THE COEFFICIENTS ARE SUPPLIED IN A DATA STATEMENT. A FIRST GUESS
C X FOR THE ZERO IS READ IN .
  PARAMETER N=6
  INTEGER J,K
  REAL A(N),B,C,DELTA,X
  DATA A /-6.8, 10.8, -10.8, 7.4, -3.7, 1./
  1 READ 500, X
  500 FORMAT(E16.8)
  PRINT 601
  601 FORMAT('NEWTONS METHOD FOR FINDING A REAL ZERO OF A POLYNOMIAL'
  * '//4X,'I',10X,'X',14X,'AP(0)',12X,'APP(1)'/)
  DO 10 J=1,20
    B = A(N)
    C = B
    DO 5 K=N,3,-1
      B = A(K-1) + X*B
      C = B + X*C
  5 CONTINUE
    B = A(1) + X*B
    PRINT 605,J,X,B,C
  605 FORMAT(15,3(1PE17.7))
    DELTA = B/C
    IF (ABS(DELTA) .LT. 1.E-7 .OR. ABS(B) .LT. 1.E-7) STOP
    X = X - DELTA
  10 CONTINUE
  PRINT 610
  610 FORMAT(' FAILED TO CONVERGE IN 20 ITERATIONS' )
                                GO TO 1
  END

```

### نتایج به دست آمده از کامپیوتر برای مثال ۱۱.۳

I	X	AP	APP
1	1.500000E 00	-1.0625001E - 00	3.7124998E 00
2	1.7861953E 00	7.2393334E - 01	9.6004875E 00
3	1.7107894E 00	8.0013633E - 02	7.5470622E 00
4	1.7001875E 00	1.3663173E - 03	7.2905675E 00
5	1.7000000E 00	4.7683716E - 07	7.2861013E 00
6	1.7000000E 00	-1.1920929E - 07	7.2860994E 00
7	1.7000000E 00	-5.9604645E - 08	7.2860998E 00

□

گرچه در مثال فوق برای به دست آوردن جوابهای دقیق به مشکل مهمی برخوردیم، اما دانشجویان را بر حدزمنی داریم از اینکه تصور کنند که همواره محاسبه ریشه‌های بسجمله‌ایها بدون اشکال صورت می‌گیرد. بعضی از دشواریهایی که ممکن است با آنها مواجه شوید در

زیر شرح می‌دهیم.

۱. در روش نیوتن دقت ریشه بستگی به دقت محاسبه جمله تصحیحی  $p(x_i)/p'(x_i)$  دارد. برای مثال اگر خطای ناشی از گرد کردن یا عوامل دیگر برابر  $\varepsilon$  باشد، آنگاه ریشه محاسبه شده می‌تواند حداکثر دقتی برابر با ریشه واقعی به علاوه  $\varepsilon/p'(x_i)$  داشته باشد. شکل ۱.۳، مقدار خطاهای غیرمنتظره را به وضوح نشان می‌دهد. خطاهای جانبی نیز به‌طور قابل ملاحظه‌ای افزایش می‌یابند اگر  $p(x)$  ریشه مکرری در نقطه  $\xi = x$  داشته باشد، زیرا در این صورت وقتی  $\xi \rightarrow x_i$ ،  $p'(x)$  به صفر نزدیک، و هر گونه خطای ناشی از گرد کردن در محاسبه  $p(x_i)$  بزرگ می‌شود.

برای روش کردن رفتار روش نیوتن پیرامون یک ریشه مکرر، بسجمله‌ای

$$p(x) = x^3 - 3x^2 + 4$$

را که یک ریشه مکرر در نقطه  $x = 2$  دارد، بررسی می‌کنیم. با انتخاب  $x_0 = 1.5$  و به‌کارگیری کامپیوتر IBM ۷۰۹۴، نتایج به‌دست آمده در جدول ۱.۳ داده شده‌اند.

جدول ۱.۳

$i$	$x_i$	$p(x_i)$	$p'(x_i)$	$p(x_i)/p'(x_i)$
0	1.5	0.625 E + 0	-0.224999999E + 1	-0.27777778E + 0
1	1.7777777	0.13717422E + 0	-0.11851852E + 1	-0.11574074E + 0
2	1.8935185	0.32807648E - 1	-0.60487403E + 0	-0.54238810E - 1
3	1.9477573	0.80453157E - 2	-0.30526827E - 0	-0.26354902E - 1
4	1.9741122	0.19932091E - 2	-0.15331630E - 0	-0.13000633E - 1
5	1.9871128	0.49611926E - 3	-0.76824840E - 1	-0.64577974E - 2
6	1.9935706	0.12376904E - 3	-0.38452353E - 1	-0.32187638E - 2
7	1.9967893	0.30934811E - 4	-0.19232938E - 1	-0.16084287E - 2
8	1.9983977	0.77188015E - 5	-0.96056932E - 2	-0.80356526E - 3
9	1.9992013	0.19371510E - 5	-0.47900614E - 2	-0.40441045E - 3
10	1.9996057	0.47683716E - 6	-0.23651228E - 2	-0.20161200E - 3
11	1.9998073	0.11920929E - 6	-0.11558611E - 2	-0.10313461E - 3
12	1.9999104	0.59604645E - 7	-0.53713301E - 3	-0.11096813E - 3
13	2.0000214	0.29802322E - 7	+0.12850899E - 3	+0.23190846E - 3
14	1.9997895	0.14901161E - 6	-0.12628894E - 2	-0.11799259E - 3
15	1.9999074	0.59604645E - 7	-0.55501277E - 3	-0.10739328E - 3

اعداد بعد از حرف E معرف توانهای ۱۰ هستند، و ارقامی که زیر آنها خط کشیده شده، به‌علت از دست رفتن ارقام با معنی طی محاسبه  $p(x_i)$  و  $p'(x_i)$ ، غیر صحیح می‌باشند. از این جدول می‌توان به نکات زیر پی برد. (در رابطه با این مسئله، تمرین ۵.۳-۵ را ببینید). الف. بارستها علی‌رغم این حقیقت که  $p'(2) = 0$ ، همگرا هستند.

ب. میزان همگرایی خطی است و مانند حالت عادی روش نیوتن مرعی نیست. بررسی عامل تصحیحی  $p'(x_i)/p(x_i)$ ، نشان می‌دهد که تا بارست دوازدهم، خطا در هر بارست با ضرب  $1/2$  کاهش می‌یابد.

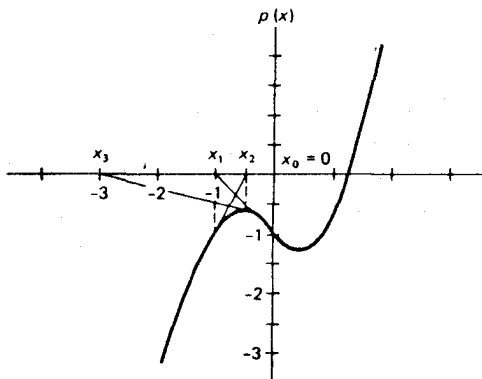
پ. بعد از ۱۳ بارست انتظار بهبود بیشتری در جواب را نمی‌توان داشت، علت آن این است که هیچ رقم صحیحی در  $p(x_i)$  باقی نمانده و در عین حال مقدار  $p'(x_i)$  از مرتبه  $3-10$  شده است. بنابراین خارج قسمت  $p(x_i)/p'(x_i)$  نتیجه غیر صحیحی در پنجمین رقم اعشاری تولید می‌کند و بهبود جواب را غیر ممکن می‌سازد.

۲. در برخی موارد انتخاب نامناسب تقریب اولیه موجب همگرایی به سمت ریشه‌ای غیر از ریشه مورد نظر می‌شود.

۳. برای برخی از بسجمله‌ایها انتخاب نامناسب  $x_0$  ممکن است به تولید يك دنباله واگرا منجر شود. به عنوان نمونه در مثال ۲.۳، اگر  $x_0 = 0$  انتخاب شود تقریبهای متوالی  $x_1 = -1/2$ ،  $x_2 = -3$ ،  $x_3 = -204$  و  $x_4 = -140$  به دست می‌آید، که قطعاً به سمت ریشه‌ای که قبلاً محاسبه شد همگرا نیست. بررسی نمودار بسجمله‌ای  $p(x) = x^3 - x - 1$  (شکل ۷.۳ را نگاه کنید) به توضیح این عملکرد کم خواهد کرد. بارستهای پیاپی ممکن است به‌طور نامتناهی پیرامون نقطه  $x = (-1/3)\sqrt{3}$ ، که در آن مقدار  $p(x)$  ماکسیمم است، نوسان کنند.

۴. برخی از بسجمله‌ایها به‌ویژه آنهایی که از درجات بالاتر هستند، بسیار ناپایدارند، بدین معنی که تغییرات جزئی و کوچک در ضرایب آنها منجر به تغییرات بزرگی در ریشه‌های آنها می‌شوند (به مثال ۱۲.۳، که در زیر آمده، نگاه کنید).

۵. وقتی ریشه‌ای از  $p(x)$  را پیدا کردیم، الگوریتم ضرب تودرتوی (۴۷.۳)، ضرایب بسجمله‌ای  $q(x)$ ، یعنی  $a'_1, \dots, a'_n$  را چنان معین می‌کند که بقیه ریشه‌های  $p(x)$  ریشه‌های  $q(x)$  باشند. لذا برای یافتن این ریشه‌ها بررسی بسجمله‌ای تقلیل یافته یا پایین



شکل ۷.۳

آمده  $g(x)$  آسانتر از بررسی  $p(x)$  خواهد بود. اما در مورد ریشه‌های اخیر، انتظار می‌رود که دقت آنها کاهش یابد، زیرا ضرایب بسجمله‌ای تقلیل یافته متضمن خطاهایی هستند که از همگرایی غیرکامل و گرد کردن این ریشه‌ها ناشی شده‌اند. برای سه‌حد اقل رساندن میزان کاهش دقت عمل، ریشه‌ها می‌باید به ترتیب صعودی محاسبه شوند (مثال ۱۲.۳ را ببینید). همچنین از بارستن بسجمله‌ای اصلی، می‌توان دقت ریشه به دست آمده از بسجمله‌ای تقلیل یافته را بهتر نمود.

□ مثال ۱۲.۳: به منظور تشریح برخی مخاطرات موجود در محاسبه ریشه بسجمله‌ایها، در زیر دو بسجمله‌ای را بررسی می‌کنیم

$$x^7 - 28x^6 + 322x^5 - 1960x^4 + 6769x^3 - 13132x^2 + 13068x - 5040 \quad (49.3)$$

و

$$x^5 - 1585x^4 + 7775x^3 - 155x^2 + 124x - 32 \quad (50.3)$$

برای یافتن کلیه ریشه‌های این بسجمله‌ایها، روش نیوتن روی یک کامپیوتر CDC ۶۵۰۰ اجرا و در هر مرحله بسجمله‌ایهای تقلیل یافته به کار گرفته شده، تقریبهای اولیه همان جواب دقیق با تقریباً ۱۰ درصد خطا در نظر گرفته شده و ملاک خاتمه  $|x_i - x_{i-1}| < 10^{-7}$  انتخاب شده است.

ریشه‌های بسجمله‌ای اولی (۴۹.۳)، به ترتیب صعودی ۱، ۲، ۳، ۴، ۵، ۶، و ۷ می‌باشند. ستون  $A$  در جدول زیر متضمن تقریبهای پیدا شده است که به ترتیب با حدسهای اولیه ۰ره، ۱ره، ۲ره، ۳ره، ۴ره، ۵ره و ۶ره شروع شده‌اند. تعداد بارستهای ضروری در سمت راست هر ریشه چاپ شده است.

ستون  $B$  شامل ریشه‌هایی است که پس از تبدیل ضریب  $x^2$  در (۴۹.۳) به  $131334x -$  به دست آمده‌اند، یعنی، بعد از آنکه یک تغییر یک واحدی در رقم پنجم اعشاری یکی از ضرایب داده شده است. فقط پنج ریشه محاسبه شده و برخی از آنها در رقم دوم اعشاری یا ریشه‌های متناظر خود در ستون  $A$  متفاوت‌اند. برای اینکه تأیید شود که این تغییرات فقط ناشی از خطای گرد کردن نیست و برای اینکه سرنوشت دو ریشه از دست رفته روشن شود، روش مولر<sup>۱</sup> (این روش در بخش بعد توضیح داده خواهد شد) به کار گرفته شده که این روش هر هفت ریشه چاپ شده در ستون  $C$  را تولید کرده است. تمام ارقام چاپ شده این اعداد، صحیح‌اند. توجه کنید که ریشه‌های ۵ و ۶ به یک جفت ریشه‌های مزدوج همثافت تبدیل شده‌اند. بنابراین تغییری معادل  $1/100$  در یکی از ضرایب به یک تغییر ۱۰ درصد در برخی از ریشه‌ها منجر شده است. هنگامی که ضرایب بسجمله‌ای به وسیله آزمایش به دست

آمده باشند، به آسانی خطا در چنین حدودی در ضرایب مشاهده می‌شود. بنا بر این باید ریشه‌های بسجمله‌ایهای درجهٔ بالا که این گونه پیدا شده‌اند، بسیار محتاطانه نگاه کرد، به ویژه هنگامی که در دقت ضرایب آنها تردید باشد.

ریشه‌های بسجمله‌ای دوم،  $(۵۰.۳)$ ، به ترتیب  $۵، ۴، ۳، ۲، ۱$  و  $۸$  هستند. با شروع عملیات از حدسهای اولیهٔ  $۵۰۴۵، ۵۰۹، ۱۰۸، ۳۰۶$  و  $۷۲$  ریشه‌های بسجمله‌ای را، به ترتیب صعودی، همان گونه که در ستون  $D$  می‌بینیم، محاسبه کردیم. سرانجام در ستون  $E$ ، نتایج محاسبهٔ این ریشه‌ها را به ترتیب نزولی ذکر کرده‌ایم، یعنی با حدسهای اولیهٔ  $۷۲$  شروع کرده و ریشهٔ  $۸$  را محاسبه کرده‌ایم، سپس با به کار گرفتن بسجمله‌ای تبدیل یافته و تقریب اولیهٔ  $۳$  ریشهٔ  $۴$  را به دست آوردیم، و هکذا. گرچه اولین ریشهٔ پیدا شده تا رقم اعشاری صحیح است، اما ریشه‌های بعدی فقط تا شش رقم اعشاری صحیح اند. بعلاوه در اینجا تعداد بارستهای مورد نیاز زیادتر است. این امر نشان می‌دهد که بهتر است در ابتدا ریشه‌هایی با قدر مطلق کوچکتر، محاسبه شوند.

نتایج کامپیوتری برای مثال ۱۲.۳

A	B	C	D	E
1.00000000 5	1.00139755 5	1.0013976	0.50000000 4	8.00000000 7
2.00000000 5	1.96892082 4	1.9689208	1.00000008 4	3.99999862 6
3.00000000 5	3.31832477 7	3.3183233	2.00000007 4	2.00000552 6
4.00000000 5	3.50505891 7	3.5050604	4.00000005 4	0.99999079 5
5.00000000 4		5.5731849	7.99999999 2	0.50000485 2
6.00000014 4		5.5731849		
6.99999993 2	7.05992816 40	7.0599281		

□

مهلی<sup>۱</sup> راهی برای استفاده از بسجمله‌ای تقلیل یافته پیشنهاد نموده است که در آن به مشکلات مذکور در بالا برخورد نمی‌شود. گیریم  $\xi_1, \xi_2, \dots, \xi_k$  ریشهٔ قبلا به دست آمده از یک بسجمله‌ای باشند. برای یافتن ریشهٔ بعدی، یک بارست نیوتن روی بسجمله‌ای تقلیل یافتهٔ  $\tilde{p}(x) = p(x) / (x - \xi_1)(x - \xi_2) \dots (x - \xi_k)$  انجام می‌دهند، اما  $\tilde{p}(x)$  را با تقسیم ترکیبی مکرر به دست نمی‌آورند. بلکه  $\tilde{p}(x)$  را به شکل فوق رها می‌کنند و در این حالت، تابع بارست به گونهٔ زیر می‌شود

$$\begin{aligned}
 x_{i+1} &= x_i - \frac{\tilde{p}(x_i)}{\tilde{p}'(x_i)} \\
 &= x_i - \frac{p(x_i)}{p'(x_i) - \sum_{r=1}^k \frac{p(x_i)}{x_i - \xi_r}}
 \end{aligned}$$

حل معادلات غیرخطی ۱۵۳

به نظر می‌رسد که این تکنیک در محاسبه ریشه‌های دقیق و متوالی کاملاً کارا باشد. تمرین ۶.۳-۷ را ببینید.

### تمرین

۶.۳-۱ با استفاده از الگوریتم ۹.۳ و به کمک یک ماشین حساب دستی، ریشه حقیقی معادله زیر را تا هفت رقم با معنی صحیح محاسبه کنید

$$x^3 + 2x - 1 = 0$$

بقیه ریشه‌ها را از بسجمله‌ای تقلیل یافته با استفاده از فرمول مربعی تعیین کنید. میزان دقت این ریشه‌ها چقدر است؟

۶.۳-۲ با استفاده از الگوریتم ۹.۳، ریشه‌های مثبت و حقیقی معادلات بسجمله‌ای زیر را به دست آورید

$$x^5 - 3x^3 + x^2 - 1 = 0 \quad (\text{الف})$$

$$x^3 + 3x - 1 = 0 \quad (\text{ب})$$

۶.۳-۳ بسجمله‌ای  $x^5 - 4x^4 - 63x^3 - 938x^2 - 380x^3 - 28x^4 + x^5$ ، دارای چهار ریشه حقیقی است. با استفاده از الگوریتم ۹.۳، این ریشه‌ها را محاسبه کنید.

۶.۳-۴ بسجمله‌ای

$$p(x) = x^8 - 170x^6 + 7392x^4 - 39712x^2 + 51200$$

دارای ریشه‌های  $\pm\sqrt{3}$ ،  $\pm 2$ ،  $\pm 8$  و  $\pm 10$  است. با انتخاب حدسهای اولیه‌ای که در حدود ۱۰ درصد با ریشه‌های دقیق فاصله داشته باشد، ریشه‌های بسجمله‌ای فوق را به ترتیب مقادیر صعودی با کامپیوتر پیدا کنید. سپس ضریب  $x^2$  را به  $39710 -$  تبدیل و مسئله را دوباره حل کنید. تغییر حاصل در جوابها را ملاحظه نمایید.

۶.۳-۵ با استفاده از قاعده‌علایم دکارت و قضایای مربوط به کرانه‌های متضمن ریشه بسجمله‌ایها و تا آنجا که امکان دارد، اطلاعاتی درباره محل و نوع ریشه‌های بسجمله‌ای زیر به دست آورید

$$p(x) = x^4 - x^3 + x^2 - x + 1$$

۶.۳-۶ بسجمله‌ای

$$p(x) = x^4 - 5885x^2 + 623504$$

دارای یک ریشه  $x_1 = 12$  است. ریشه مثبت حقیقی دیگری نزدیک  $x = 2$  وجود دارد.

با استفاده از تکنیک مهلی و با شروع عملیات از  $x_0 = ۲$  این ریشه را به دست آورید.

۷-۶-۳ بر مبنای روش مهلی برنامه‌ای برای یافتن ریشه‌های حقیقی و متوالی یک بسجمله‌ای  $p(x)$  بنویسید.

۸-۶-۳ با استفاده از روش مهلی، ریشه‌های بسجمله‌ای داده شده در مثال ۱۲-۳ را به دست آورید و با نتایج مربوط به آن مسئله مقایسه کنید.

### ۷-۳\* ریشه‌های هم‌تافت و روش مولر

روشهایی که تاکنون مورد بحث قرار گرفتند به ما امکان می‌دهند که یک ریشه تکین<sup>۱</sup> از یک تابع را، وقتی که یک تقریب اولیه از آن ریشه معلوم باشد، محاسبه کنیم. زمانی که کلیه ریشه‌های یک تابع لازم باشند و یا یک تقریب اولیه خوبی در دست نباشد، این روشها چندان رضایتبخش نیستند. برای توابع بسجمله‌ای روشهایی وجود دارند که تمامی ریشه‌ها را به طور همزمان به دست می‌دهند، و پس از آن می‌توان روشهای بارستی این فصل را برای محاسبه دقیقتر این ریشه‌ها به کار برد. از جمله این روشها می‌توان به الگوریتم تفاضلهای منقسم [۲] و روش گرافه<sup>۲</sup> [۵] اشاره کرد.

یک روش ریشه‌دار و تازه‌ای به وسیله مولر بیان و با موفقیت قابل ملاحظه‌ای روی کامپیوترها به کار برده شده است. این روش را می‌توان برای محاسبه هر تعداد ریشه حقیقی یا هم‌تافت خواسته شده از هر تابع دلخواهی به کار گرفت. این روش، روش بارستی است و پیرامون ریشه تقریباً به طور مریبی همگراست و احتیاج به محاسبه مشتق تابع ندارد و ریشه‌های حقیقی و هم‌تافت را، حتی در حالتی که این ریشه‌ها ساده نباشند، به دست می‌دهد. وانگهی، روشی است کلی، بدین معنی که احتیاج به تقریب اولیه ندارد. در این بخش با انصراف از بحث در باره همگرایی، مختصراً چگونگی به دست آوردن آن شرح داده می‌شود و سپس کاربرد آن برای محاسبه ریشه‌های حقیقی و هم‌تافت مورد بحث قرار می‌گیرد. به ویژه، روی محاسبه ریشه‌های هم‌تافت یک بسجمله‌ای با ضرایب حقیقی تکیه می‌شود، زیرا این مسئله با تمامی شاخه‌های مهندسی مربوط است.

روش مولر توسیع روش خط قاطع است. یادآوری می‌کنیم که در روش خط قاطع از  $x_{i-1}$  و  $x_i$  به مثابه تقریبهایی برای یک ریشه  $f(x) = 0$ ، تقریب بعدی  $x_{i+1}$  محاسبه می‌شود که یک ریشه بسجمله‌ای خطی  $p(x)$  مار بر نقاط  $\{x_i, f(x_i)\}$  و  $\{x_{i-1}, f(x_{i-1})\}$  است. در روش مولر تقریب بعدی  $x_{i+1}$ ، به عنوان ریشه یک سهمی مار بر سه نقطه  $\{x_i, f(x_i)\}$ ،  $\{x_{i-1}, f(x_{i-1})\}$  و  $\{x_{i-2}, f(x_{i-2})\}$  پیدا می‌شود. چنانچه در فصل ۲ نشان داده شد، تابع

$$p(x) = f(x_i) + f[x_i, x_{i-1}](x - x_i) + f[x_i, x_{i-1}, x_{i-2}](x - x_i)(x - x_{i-1})$$

سهمی منحصر به فردی است که با تابع  $f(x)$  در سه نقطه  $x_i$ ،  $x_{i-1}$  و  $x_{i-2}$  منطبق می باشد. و چون داریم

$$(x - x_i)(x - x_{i-1}) = (x - x_i)^2 + (x - x_i)(x_i - x_{i-1})$$

می توانیم  $p(x)$  را به صورت زیر هم بنویسیم

$$p(x) = f(x_i) + (x - x_i)c_i + f[x_i, x_{i-1}, x_{i-2}](x - x_i)^2 \quad (51.3)$$

در رابطه بالا  $c_i$  عبارت است از

$$c_i = f[x_i, x_{i-1}] + f[x_i, x_{i-1}, x_{i-2}](x_i - x_{i-1})$$

بنابراین هر ریشه  $\alpha$  از سهمی  $p(x)$  بر اساس یک شکل خاص از فرمول مربعی استاندارد (۲۰.۱) را ببینید، در رابطه زیر صدق می کند

$$\alpha - x_i = \frac{-2f(x_i)}{c_i \pm \{c_i^2 - 2f(x_i)f[x_i, x_{i-1}, x_{i-2}]\}^{1/2}} \quad (52.3)$$

اگر در رابطه (۵۲.۳) علامت به گونه ای انتخاب شود که مخرج از نظر قدر مطلق بزرگترین مقدار ممکنه را داشته باشد و اگر طرف راست رابطه (۵۲.۳) را با  $h_{i+1}$  نمایش دهیم، تقریب بعدی برای یک ریشه  $f(x)$  به صورت زیر درمی آید

$$x_{i+1} = x_i + h_{i+1}$$

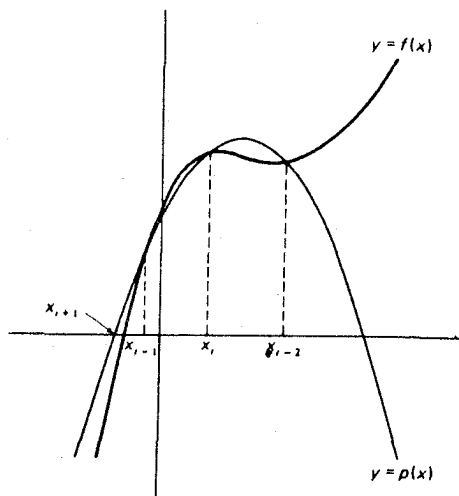
در این صورت با استفاده از  $x_i$ ،  $x_{i-1}$  و  $x_{i+1}$  به عنوان سه تقریب پایه، روند فوق تکرار می شود. اگر ریشه های به دست آمده از (۵۲.۳) حقیقی باشند، وضعیت به گونه ای خواهد شد که در شکل ۸.۳ نمایش داده شده است. ولی باید توجه داشت که حتی اگر ریشه مورد نظر حقیقی باشد امکان مواجه شدن با تقریبهای هم تافت نیز می رود. زیرا ممکن است که جوابهای به دست آمده از (۵۲.۳)، هم تافت باشند. در چنین حالتی مؤلفه هم تافت معمولاً از نظر مقدار آنقدر کوچک است که می تواند نادیده گرفته شود. در حقیقت در زیر برنامه ای که در زیر داده شده است، هنگام جستجو برای یک ریشه حقیقی، هر مؤلفه هم تافتی که به دست آید می تواند حذف شود.

در الگوریتم ۱۰.۳ دنباله مراحل خواسته شده در روش مولر به صورت رسمی بیان شده است.

### الگوریتم ۱۰.۳: روش مولر

۱. گیریم  $x_0$ ،  $x_1$ ،  $x_2$  سه تقریب برای یک ریشه  $f(x)$  باشد.  $f(x_0)$ ،  $f(x_1)$  و  $f(x_2)$  را محاسبه کنید.
۲. مقادیر زیر را محاسبه کنید





شکل ۸-۳

$$h_2 := x_2 - x_1, \quad h_1 := x_1 - x_0$$

$$f[x_2, x_1] = \frac{f(x_2) - f(x_1)}{h_2}$$

$$f[x_1, x_0] = \frac{f(x_1) - f(x_0)}{h_1}$$

۳. قرار دهید  $i = 2$ .

۴. مقادیر زیر را محاسبه کنید

$$f[x_i, x_{i-1}, x_{i-2}] = \frac{f[x_i, x_{i-1}] - f[x_{i-1}, x_{i-2}]}{(h_i + h_{i-1})}$$

$$c_i := f[x_i, x_{i-1}] + h_i f[x_i, x_{i-1}, x_{i-2}]$$

۵. مقدار زیر را محاسبه کنید

$$h_{i+1} := \frac{-2f(x_i)}{(c_i \pm \sqrt{c_i^2 - 4f(x_i)f[x_i, x_{i-1}, x_{i-2}]})}$$

علامت راطوری انتخاب کنید که مخرج رابطهٔ بالا از لحاظ قدرمطلق بیشترین مقدار را داشته باشد.

۶. قراردید  $x_{i+1} = x_i + h_{i+1}$

۷. مقادیر زیر را محاسبه کنید

$$f(x_{i+1}) \quad \text{و} \quad f[x_{i+1}, x_i] = \frac{f(x_{i+1}) - f(x_i)}{h_{i+1}}$$

۸. قراردید  $i = i + 1$ ، مراحل ۴ تا ۷ را تکرار کنید تا آنکه یکی از ملاکهای

تعیین شده زیر به ازای مقادیر مفروض  $\epsilon_1$  و  $\epsilon_4$  برقرار گردد:

$$|x_i - x_{i-1}| < \epsilon_1 \quad (\text{الف})$$

$$|f(x_i)| < \epsilon_4 \quad (\text{ب})$$

یا آنکه تعداد بارستها از یک حداکثر، بیشتر شود.

یک زیر برنامه کاملی بر اساس این الگوریتم در زیر داده شده است. پارامترهای لازم برای زیر برنامه در کارتهای توضیحی<sup>۱</sup> ارائه شده اند. ZEROS(I)، یک آرایه یک بعدی متضمن برآوردهای اولیه ریشههای مورد نظر است. زیر برنامه به طور خودکار دو تقریب اضافی برای ZEROS(I) را به صورت  $ZEROS(I) + \delta$  و  $ZEROS(I) - \delta$  حساب می کند و آنگاه بر طبق الگوریتم مولر پیش می رود.

```

SUBROUTINE MULLER ( FN, FNREAL, ZEROS, N, NPREV, MAXIT, EP1, EP2 )
C DETERMINES UP TO N ZEROS OF THE FUNCTION SPECIFIED BY FN, USING
C QUADRATIC INTERPOLATION, I.E., MUELLER'S METHOD.
EXTERNAL FN
LOGICAL FNREAL
INTEGER MAXIT, N, NPREV, KOUNT
REAL EP1, EP2, EPS1, EPS2
COMPLEX ZEROS(N), C, DEN, DIVDF1, DIVDF2, DVDF1P, FZR, FZRDFL
* FZRPRV, H, ZERO, SQR
C***** I N P U T *****
C FN NAME OF A SUBROUTINE, OF THE FORM FN(Z, FZ) WHICH, FOR GIVEN
C Z, RETURNS F(Z). MUST APPEAR IN AN E X T E R N A L STATE-
C MENT IN THE CALLING PROGRAM.
C FNREAL A LOGICAL VARIABLE. IF .TRUE., ALL APPROXIMATIONS ARE TAKEN
C TO BE REAL, ALLOWING THIS ROUTINE TO BE USED EVEN IF F(Z) IS
C ONLY DEFINED FOR REAL Z.
C ZEROS(1), ..., ZEROS(NPREV) CONTAINS PREVIOUSLY FOUND ZEROS (IF
C NPREV.GT. 0).
C ZEROS(NPREV+1), ..., ZEROS(N) CONTAINS FIRST GUESS FOR THE ZEROS TO BE
C FOUND. (IF YOU KNOW NOTHING, 0 IS AS GOOD A GUESS AS ANY.)
C MAXIT MAXIMUM NUMBER OF FUNCTION EVALUATIONS ALLOWED PER ZERO.
C EP1 ITERATION IS STOPPED IF ABS(H) .LT. EP1*ABS(ZR), WITH
C H = LATEST CHANGE IN ZERO ESTIMATE ZERO.
C EP2 ALTHOUGH THE EP1 CRITERION IS NOT MET, ITERATION IS STOPPED IF
C ABS(F(ZERO)) .LT. EP2.
C N TOTAL NUMBER OF ZEROS TO BE FOUND.
C NPREV NUMBER OF ZEROS FOUND PREVIOUSLY.
C***** O U T P U T *****
C ZEROS(NPREV+1), ..., ZEROS(N) APPROXIMATIONS TO ZEROS.
C
C
C INITIALIZATION
EPS1 = MAX(EP1, 1.E-12)
EPS2 = MAX(EP2, 1.E-20)

```

```

C
DO 100 I=NPREV+1,N
  KOUNT = 0
C      COMPUTE FIRST THREE ESTIMATES FOR ZERO AS
C      ZEROS(I)+5., ZEROS(I)-.5, ZEROS(I)
1  ZERO = ZEROS(I)
  H = .5
  CALL DFLATE(FN, ZERO+.5, I, KOUNT, FZR, DVDF1P, ZEROS, *1)
  CALL DFLATE(FN, ZERO-.5, I, KOUNT, FZR, FZRPRV, ZEROS, *1)
  HPREV = -1.
  DVDF1P = (FZRPRV - DVDF1P)/HPREV
  CALL DFLATE(FN, ZERO, I, KOUNT, FZR, FZRDFL, ZEROS, *1)
C  DO WHILE KOUNT.LE.MAXIT OR H IS RELATIVELY BIG
C      OR FZR = F(ZERO) IS NOT SMALL
C      OR FZRDFL = FDEFLATED(ZERO) IS NOT SMALL OR NOT MUCH
C      BIGGER THAN ITS PREVIOUS VALUE FZRPRV .
40  DIVDF1 = (FZRDFL - FZRPRV)/H
  DIVDF2 = (DIVDF1 - DVDF1P)/(H + HPREV)
  HPREV = H
  DVDF1P = DIVDF1
  C = DIVDF1 + H*DIVDF2
  SQR = C*C - 4.*FZRDFL*DIVDF2
  IF (FNREAL .AND. REAL(SQR) .LT. 0.) SQR = 0.
  SQR = SQRT(SQR)
  IF (REAL(C)*REAL(SQR)+AIMAG(C)*AIMAG(SQR) .LT. 0.) THEN
    DEN = C - SQR
  ELSE
    DEN = C + SQR
  END IF
  IF (ABS(DEN) .LE. 0.) DEN = 1.
  H = -2.*FZRDFL/DEN
  FZRPRV = FZRDFL
  ZERO = ZERO + H
  IF (KOUNT .GT. MAXIT) GO TO 99
C
C 70  CALL DFLATE(FN, ZERO, I, KOUNT, FZR, FZRDFL, ZEROS, *1)
C      CHECK FOR CONVERGENCE
C      IF (ABS(H) .LT. EPS1*ABS(ZERO)) GO TO 99
C      IF (MAX(ABS(FZR),ABS(FZRDFL)) .LT. EPS2) GO TO 99
C      CHECK FOR DIVERGENCE
C      IF (ABS(FZRDFL) .GE. 10.*ABS(FZRPRV)) THEN
C          H = H/2.
C          ZERO = ZERO - H
C              GO TO 70
C      ELSE
C          GO TO 40
C      END IF
99  ZEROS(I) = ZERO
100 CONTINUE
      RETURN
END
SUBROUTINE DFLATE ( FN, ZERO, I, KOUNT, FZERO, FZRDFL, ZEROS, * )
C TO BE CALLED IN M U L L E R
  INTEGER I,KOUNT, J
  COMPLEX FZERO,FZRDFL,ZERO,ZEROS(I), DEN
  KOUNT = KOUNT + 1
  CALL FN(ZERO, FZERO)
  FZRDFL = FZERO
  IF (I .LT. 2) RETURN
  DO 10 J=2,I
    DEN = ZERO - ZEROS(J-1)
    IF (ABS(DEN) .EQ. 0.) THEN
      ZEROS(I) = ZERO*1.001
      RETURN 1
    ELSE
      FZRDFL = FZRDFL/DEN
    END IF
  10 CONTINUE
      RETURN
END
END

```

روش مولر مانند الگوریتمهای دیگری که در این فصل مورد بحث قرار گرفته اند، در هر مرحله یک ریشه را محاسبه می کند. برای محاسبه بیش از یک ریشه، در این روش از شیوه

معروف پایین آوری<sup>۱</sup> استفاده می‌شود. برای مثال، اگر یک ریشه<sup>۱</sup>  $\xi_1$  قبلاً محاسبه شده باشد، زیر برنامه<sup>۱</sup> ریشه<sup>۱</sup> بعدی را با استفاده از تابع زیر محاسبه می‌کند

$$f_1(x) = \frac{f(x)}{x - \xi_1} \quad (۵۳.۳)$$

ما قبلاً هنگام حل معادله‌های بسجمله‌ای با روش نیوتن، با این تکنیک که در آن تابع پایین آمده یا تقلیل یافته<sup>۱</sup>  $f_1(x)$ ، نتیجه<sup>۱</sup> جانبی الگوریتم بود برخورد داشتیم. در روش مولر، اگر  $r$  ریشه<sup>۱</sup>  $\xi_1, \dots, \xi_r$  قبلاً محاسبه شده باشند، آنگاه ریشه<sup>۱</sup> بعدی با به کار گرفتن تابع تقلیل یافته<sup>۱</sup> زیر محاسبه می‌شود

$$f_r(x) = \frac{f(x)}{(x - \xi_1)(x - \xi_2) \dots (x - \xi_r)} \quad (۵۴.۳)$$

اگر هیچ تخمینی داده نشده باشد، زیر برنامه<sup>۱</sup> همواره ریشه را به ترتیب صعودی تجسس می‌کند، زیرا این عمل معمولاً رشد خطای ناشی از گرد کردن را حداقل می‌سازد. همچنین این زیر برنامه<sup>۱</sup> کلیه<sup>۱</sup> ریشه‌هایی را که از به کار گرفتن تابع تقلیل یافته<sup>۱</sup> به دست می‌آیند در تابع اصلی  $f(x)$  قرار می‌دهد و بدین وسیله دقت آنها را کنترل می‌کند. در عمل، استفاده از تابع تقلیل یافته<sup>۱</sup> برای تعیین ریشه<sup>۱</sup> تاحدی از دقت عمل می‌کاهد. ریشه‌های تقریبی حاصل از روش پایین آوری می‌توانند اصلاح شوند، هر گاه این ریشه‌ها به عنوان حدسهای اولیه<sup>۱</sup> روش نیوتن که به تابع اصلی اعمال می‌گردد، به حساب آیند. در کاربرد زیر برنامه<sup>۱</sup> مولر، استفاده کننده می‌تواند تعداد ریشه‌های مورد نظر را مشخص نماید. برای مثال برخی توابع ممکن است دارای بینهایت ریشه باشند که تنها چند ریشه<sup>۱</sup> اول آنها مورد نظر باشند.

□ مثال ۱۳.۳: تابع بسل<sup>۲</sup>  $J_0(x)$ ، به توسط سری نامتناهی زیر داده شده است

$$J_0(x) = 1 - \frac{x^2}{2^2 \times 1 \times 1} + \frac{x^4}{2^4 \times 2! \times 2!} - \frac{x^6}{2^6 \times 3! \times 3!} + \dots$$

می‌دانیم که  $J_0(x)$  دارای بینهایت ریشه<sup>۱</sup> حقیقی است. با استفاده از الگوریتم ۱۰.۳، سه تا از اولین ریشه‌های مثبت را پیدا کنید. نتایج کامپیوتری که در زیر ارائه می‌شود با استفاده از زیر برنامه<sup>۱</sup> استاندارد موجود در خود کامپیوتر، که بر اساس سری فسوق برای  $J_0(x)$  تهیه شده، روی یک کامپیوتر IBM ۷۰۹۴ به دست آمده است. در این برنامه<sup>۱</sup> مقادیر  $J_0(x)$  با حداکثر دقت عمل محاسبه شده‌اند.

با رسته‌ها همگی با تقریبهای اولیه<sup>۱</sup>  $x_0 = -1$ ،  $x_1 = 1$ ،  $x_2 = 0$  شروع شده، تا برقراری یکی از ملاکهای زیر ادامه یافته‌اند:

$$\frac{|x_{i+1} - x_i|}{|x_{i+1}|} < 10^{-6} \quad (\text{الف})$$

$$|J_0(x_i)| < 10^{-20} \quad (\text{ب})$$

مقادیر همگرا شده تا شش رقم با معنی صحیح اند. نکتهٔ دیگر اینکه مقادیر ریشه‌ها به ترتیب صعودی محاسبه شده‌اند.

نتایج کامپیوتری برای مثال ۱۳.۳

Zero 1	Zero 2	Zero 3
-0.099999999E 01	-0.099999999E 01	-0.099999999E 01
0.099999999E 01	0.099999999E 01	0.099999999E 01
0.	0.	0.
0.20637107E 01	0.36557332E 01	0.47983123E 01
0.23167706E 01	0.44416171E 01	0.59396663E 01
0.23970029E 01	0.50863190E 01	0.70758440E 01
0.24047983E 01	0.55024961E 01	0.88981197E 01
0.24048255E 01	0.55202182E 01	0.92976399E 01
0.24048255E 01	0.55200780E 01	0.86854592E 01
	0.55200780E 01	0.86529856E 01
		0.86537299E 01
		0.86537278E 01

کلیهٔ مثالهای زیر روی يك کامپیوتر CDC ۶۵۰۰ با استفاده از الگوریتم ۱۰.۳ اجرا شده‌اند. ملاکهای خطا برای این مثالها  $\epsilon_1 = \epsilon_2 = 10^{-8}$  می‌باشند و همگی از همان مقادیر اولیه (۰، ۰، ۰) و سپس پایین آوری استفاده کرده‌اند. گرچه نتایج تا هشت رقم با معنی چاپ شده‌اند ولی به‌یاد داشته باشید که در کامپیوتر CDC ۶۵۰۰ طول کلمه در محاسبات با ممیز شناور برابر با ۱۴ رقم اعشاری است. برون‌داد عبارت است از اجزای حقیقی و انگاری تقریبات همگرا به ریشه‌ها، و اجزای حقیقی و انگاری مقادیر تابع به ازای آن ریشه‌ها.

□ مثال ۱۴.۳: کلیهٔ ریشه‌های بسجمله‌ای زیر را پیدا کنید

$$p(x) = x^3 + x - 3$$

ROOT		F(ROOT)	
جزء حقیقی	جزء انگاری	جزء حقیقی	جزء انگاری
1.2134117E + 00	0.	-4.2632564E - 14	0.
-6.0670583E - 01	1.4506122E + 00	2.8421709E - 14	4.2632564E - 14
-6.0670583E - 01	-1.4506122E + 00	2.8421709E - 14	-2.6290081E - 13

نتایج فوق را با نتایج مثال ۱۰.۳، که با استفاده از ماشین حساب به دست آمده است،

## حل معادلات غیرخطی ۱۶۱

مقایسه کنید. توجه کنید که چون  $p(x)$  دارای ضرایب حقیقی است، ریشه‌های همتافت به صورت جفت‌های مزدوج ظاهر می‌شوند. همچنین توجه کنید که در اینجا هیچ نیازی به برآورد ریشه‌های همتافت نیست. در حالی که می‌توان روش نیوتن را برای محاسبه ریشه‌های همتافت به کار گرفت، ولی این روش احتیاج به یک برآورد مناسبی دارد که به دست آوردن آن معمولاً بسیار مشکل است. توجه شود که خطای موجود در  $F(\text{ROOT})$  همان گونه که ملاک خطا آن را تجویز کرده است، به طور قابل ملاحظه‌ای کوچکتر از  $10^{-8}$  است. درحقیقت در آخرین بارست می‌بایستی خطا از مقداری مانند  $10^{-7}$  به  $10^{-14}$  کاهش یابد و این امر نشان می‌دهد که روش مورد بحث تقریباً به‌طور مربعی همگر است.  $\square$

$\square$  مثال ۱۵.۳: ریشه‌های بسجمله‌ای زیر را پیدا کنید

$$f(x) = x^5 - 37x^4 + 774x^3 - 1088x^2 + 1088x - 688$$

این همان مثال ۱۱.۳ است که قبلاً به وسیله روش نیوتن حل شده بود. ریشه‌های دقیق این بسجمله‌ای،  $\pm i$ ،  $1 \pm \sqrt{3}i$ ، و  $107$  هستند. با وجود اینکه ریشه انگاری محض،  $\pm \sqrt{2}i$ ، دارای قسمت حقیقی کوچکی است، نتایج زیر تا هشت رقم با معنی صحیح اند.

ROOT		F(ROOT)	
جزء حقیقی	جزء انگاری	جزء حقیقی	جزء انگاری
1.0000000E + 00	-1.0000000E + 00	-2.8421709E - 14	2.8421709E - 14
1.0000000E + 00	1.0000000E + 00	1.4210855E - 13	-1.5631940E - 13
-9.0964472E - 12	1.4142136E + 00	-8.5330498E - 10	-1.2773698E - 09
8.3306265E - 11	-1.4142136E + 00	-8.5339025E - 10	-1.2785158E - 09
1.7000000E + 00	1.3036419E - 10	6.3431571E - 10	9.4984654E - 10

$\square$

$\square$  مثال ۱۶.۳: ریشه‌های بسجمله‌ای زیر را محاسبه کنید

$$f(x) = x^7 - 28x^6 + 322x^5 - 1960x^4 + 6769x^3 - 13132x^2 + 13068x - 5040$$

این مثال با روش نیوتن در مثال ۱۲.۳ حل شد و در به دست آوردن جواب‌های دقیق آن با اشکالاتی مواجه شدیم. ریشه‌های این بسجمله‌ای برابر با  $1, 2, 3, 4, 5, 6, 7$  هستند. نتایج زیر بسیار دقیق‌اند، گرچه طول نسبتاً بزرگ کلمه در کامپیوتر CDC ۶۵۰۰ در این دقت عمل نقش اساسی داشته است. همچنین باید توجه کرد که گرچه در حالت کلی روش مولر ریشه‌ها را به ترتیب صعودی تجسس می‌کند، اما در این مثال به چنین امری توفیق نیافته است.

ROOT		F(ROOT)	
جزء حقیقی	جزء انگاری	جزء حقیقی	جزء انگاری
2.0000000E + 00	-2.6080092E - 16	-5.8207661E - 11	3.1296110E - 14
3.0000000E + 00	7.4093893E - 11	-7.8580342E - 10	3.5565069E - 09
1.0000000E + 00	-1.7030067E - 16	0.	-1.2261648E - 13
6.0000000E + 00	1.6284031E - 15	-1.0710210E - 08	-1.9540837E - 13
5.0000000E + 00	-7.2393906E - 13	2.4156179E - 09	-3.4749075E - 11
4.0000000E + 00	-2.3682266E - 10	2.0954758E - 09	8.5256156E - 09
7.0000000E + 00	-8.1000834E - 20	4.0745363E - 10	-5.8320601E - 17

□ مثال ۱۷.۳: ریشه‌های بسجمله‌ای زیر را پیدا کنید

$$f(x) = x^8 - 170x^6 + 7392x^4 - 39712x^2 + 51200$$

ریشه‌های این بسجمله‌ای برابر با  $\pm 10$ ,  $\pm 2$ ,  $\pm 8$ ,  $\pm \sqrt{3}$  هستند. این برنامه در حالت همتافت اجرا شده و ریشه‌هایی با دقت هشت رقم بسا معنی صحیح به دست داده بود. این مثال نشان می‌دهد که الگوریتم فوق نمی‌تواند برای حل بسجمله‌ای از درجه نسبتاً بالا به کار برده شود و نتایج خوبی به دست دهد (تمرین ۳-۶ را نگاه کنید).

ROOT		F(ROOT)	
جزء حقیقی	جزء انگاری	جزء حقیقی	جزء انگاری
-1.4142136E + 00	0.	-2.3282064E - 10	0.
1.4142136E + 00	0.	-2.3283064E - 10	0.
2.0000000E + 00	0.	0.	0.
-2.0000000E + 00	0.	0.	0.
8.0000000E + 00	0.	2.9336661E - 08	0.
-8.0000000E + 00	0.	2.9336661E - 08	0.
1.0000000E + 01	0.	2.3352914E - 07	0.
-1.0000000E + 01	0.	1.8742867E - 07	0.

## تمرین

۱-۷.۳ ریشه‌های حقیقی و همتافت بسجمله‌ایهای زیر را با استفاده از روش مولر پیدا کنید

(الف)  $x^7 - 1$

(ب)  $x^4 - 7x^3 + 18x^2 - 20x + 8$

(پ)  $x^6 + 2x^5 + x^4 + 3x^3 + 5x^2 + x + 1$

۲-۷.۳ معادله  $x - \tan x = 0$  دارای بی‌نهایت ریشه حقیقی است. با استفاده از روش مولر اولین سه ریشه مثبت آن را پیدا کنید.

۳-۷.۳ انتگرال فونل  $C(x)$ ، با سری زیر تعریف می‌شود

حل معادلات غیرخطی ۱۶۳

$$C(x) = \sum_{n=0}^{\infty} \frac{(-1)^n \left(\frac{\pi}{2}\right)^{2n}}{(2n)!(2n+1)!} x^{2n+1}$$

با استفاده از روش مولر، اولین سه ریشه مثبت حقیقی این تابع را به دست آورید. عملیات را به ازای  $n=3$  و کنار گذاشتن سایر جملات در سری فوق، شروع کنید و سپس  $n$  را افزایش دهید تا آنکه ریشه‌های صحیحی که شما را قانع سازد به دست آید.

۴-۷.۳ تابع بسل از مرتبه ۱ با سری زیر تعریف می‌شود

$$J_1(z) = \frac{z}{2} \sum_{k=0}^{\infty} \frac{(-z^2/4)^k}{k!(k+1)!}$$

اولین چهار ریشه این تابع را مانند تمرین ۳-۷.۳ را به دست آورید.



## ماتریسها و دستگاههای معادلات خطی

بسیاری از مسائل آنالیز عددی را می‌توان به مسئله حل دستگاههای معادلات خطی تبدیل کرد. از آن جمله می‌توان مسئله حل معادلات دیفرانسیل معمولی یا معادلات بسا مشتقات جزئی به وسیله روشهای تفاضل متناهی، حل دستگاههای معادلات، مسائل ویژه مقصداری فیزیک ریاضی، برآوردن داده‌ها با روش کوچکترین توانهای دوم و تقریب با بسجمله‌ایها را نام برد. کاربرد نمادگذاری ماتریسی نه تنها مناسب، بلکه در نشر روابط بنیادی بسیار تواناست. در بخش ۱۰۴ برخی از ویژگیهای ساده ماتریسها را که در قسمتهای بعدی به کار برده می‌شوند معرفی خواهیم کرد. برخی از قضایا و ویژگیها بدون برهان بیان خواهند شد.

### ۱۰۴ ویژگیهای ماتریسها

یک دستگاه  $m$  معادله خطی  $n$  مجهولی دارای شکل کلی زیر است

$$a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1$$

$$a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2$$

$$\dots$$

$$a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n = b_m$$

(۱۰۴)

ضرایب  $a_{ij}$  ( $i=1, \dots, m, j=1, \dots, n$ ) و مقادیر سمت راست معادله،  $b_i$  ( $i=1, \dots, m$ )، اعدادی مفروض‌اند. مسئله ما عبارت است از پیدا کردن، اعداد  $x_j$  ( $j=1, \dots, n$ ) به طوری که  $m$  معادله (۱.۴) هم‌زمان برقرار باشند. بحث در باب این مسئله و درک آن زمانی بسیار ساده خواهد شد که مفاهیم جبری ماتریسها و بردارها را مورد استفاده قرار دهیم.

### تعریف ماتریس و بردار

یک ماتریس یک آرایه مستطیلی از اعداد (معمولاً حقیقی) است که در سطرها و ستونها، مرتب شده باشند. ضرایب (۱.۴) ماتریسی تشکیل می‌دهند که آن را  $A$  می‌نامیم. مرسوم است که ماتریس  $A$  به صورت زیر نمایش داده شود

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdot & \cdot & \cdot & \cdot \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \quad (2.4)$$

گاهی اوقات ماتریس فوق به صورت خلاصه‌تر زیر نوشته می‌شود

$$A = (a_{ij}) \quad (3.4)$$

ماتریس  $A$  در (۲.۴) دارای  $m$  سطر و  $n$  ستون و یا از مرتبه  $m \times n$  است. درایه  $a_{ij}$  درایه  $(i, j)$ ، در محل تقاطع سطر  $i$ ام و ستون  $j$ ام  $A$  قرار دارد. اگر  $A$  یک ماتریس  $n \times n$  باشد می‌گویند که  $A$  ماتریس مربع و از مرتبه  $n$  است. اگر ماتریسی فقط یک ستون داشته باشد آن را یک بردار ستونی و اگر فقط یک سطر داشته باشد، آن را یک بردار سطری می‌نامند. بردارهای ستونی را با حروف کوچک سیاه نشان می‌دهیم و بدین ترتیب آنها را از بردارهای سطری متمایز می‌سازیم و مختصراً آنها را بردار می‌گوییم. بنابراین ثابتهای سمت راست یعنی  $b_i$  ها ( $i=1, \dots, m$ ) و مجهولات  $x_j$  ( $j=1, \dots, n$ ) هر دو بردارهای زیر را تشکیل می‌دهند

$$\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_m \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \quad (4.4)$$

گوییم که  $\mathbf{b}$  یک  $m$ -بردار و  $\mathbf{x}$  یک  $n$ -بردار است.

## تساوی

اگر  $A = (a_{ij})$  و  $B = (b_{ij})$  دو ماتریس باشند، گوییم  $A$  مساوی  $B$  است یا  $A = B$ ، به شرط آنکه  $A$  و  $B$  هم مرتبه باشند و به ازای کلیه مقادیر  $i$  و  $j$ ، رابطه  $a_{ij} = b_{ij}$  برقرار باشد.

## ضرب ماتریسها

در اصطلاحاتی که تاکنون وارد شده اند، (۱.۴) گویای آن است که اگر ماتریس  $A$  به طریق خاصی با ماتریس یک ستونی، یا بردار  $\mathbf{x}$ ، ترکیب شود حاصل می‌بایستی مساوی ماتریس یک ستونی، یا بردار  $\mathbf{b}$ ، شود. روند ترکیب ماتریسها که در اینجا دخالت کرده‌اند ضرب ماتریس نام دارد و به طور کلی به صورت زیر تعریف می‌شود: گیریم  $A = (a_{ij})$  یک ماتریس  $m \times n$  و  $B = (b_{ij})$  یک ماتریس  $n \times p$  باشد، آنگاه ماتریس  $C = (c_{ij})$  حاصل ضرب (ماتریسهای)  $A$  و  $B$  (با همین ترتیب) و یا  $C = AB$  است، در صورتی که  $C$  از مرتبه  $m \times p$  باشد و داشته باشیم

$$c_{ij} = \sum_{k=1}^n a_{ik} b_{kj} \quad j = 1, \dots, p \quad \text{و} \quad i = 1, \dots, m \quad (5.4)$$

به طور لفظی درایه  $(i, j)$  ام حاصل ضرب  $A$  و  $B$  یا  $C = AB$ ، بدین صورت محاسبه می‌شود که  $n$  درایه سطر  $i$  ام ماتریس  $A$  و  $n$  درایه ستون  $j$  ام ماتریس  $B$  نظیر به نظیر در هم ضرب و این  $n$  حاصلضربها با هم جمع می‌شوند.

□ مثال:

$$B = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \quad \text{و} \quad A = \begin{bmatrix} 3 & 0 & 2 \\ 1 & 2 & 0 \\ 0 & 1 & 1 \end{bmatrix} \quad \text{اگر}$$

$$AB = \begin{bmatrix} 8 & 3 \\ 2 & 3 \\ 1 & 1 \end{bmatrix} \quad \text{آنگاه}$$

برای مثال درایه  $(2, 1)$  از  $AB$ ، همان طوری که به وسیله پیکانها نشان داده شده‌اند، از ترکیب سطر دوم  $A$  و ستون اول  $B$  به دست آمده است:

□

$$1 \times 2 + 2 \times 0 + 0 \times 1 = 2$$

با این تعریف از حاصلضرب ماتریسها و تعاریف (۲.۴) و (۴.۴)، می‌توان دستگاه (۱.۴) را به صورت ساده‌تر نوشت

$$AX = b \quad (۶.۴)$$

فعلاً، چنین به نظر می‌رسد که این سادگی بیان به‌بهای چندین تعریف به دست آمده، که یکی از آنها کاملاً پیچیده است، و لسی بسیاری از مزایای قرارداد ماتریسی طی این فصل آشکار می‌شوند.

ضرب ماتریسها به هیچ وجه مشابه با ضرب اعداد نیست. برای مثال، تعیین حاصلضرب ماتریس  $A$  در ماتریس  $B$  تنها وقتی امکان پذیر است که تعداد ستونهای  $A$  مساوی با تعداد سطرهای  $B$  باشد. بنابراین حتی وقتی که حاصلضرب  $AB$  تعریف شده باشد، حاصلضرب  $BA$  با  $A$  الزاماً تعریف شده نیست. وانگهی، حتی اگر هر دو حاصلضرب  $AB$  و  $BA$  تعریف شده باشند این دو الزاماً مساوی نخواهند بود.

□ مثال:

$$\text{اگر } A = \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix} \text{ و } B = \begin{bmatrix} 2 & 1 \\ 0 & 1 \end{bmatrix}, \text{ آنگاه}$$

$$\square \quad AB = \begin{bmatrix} 4 & 3 \\ 2 & 4 \end{bmatrix} \neq \begin{bmatrix} 5 & 5 \\ 1 & 3 \end{bmatrix} = BA$$

از طرف دیگر ضرب ماتریس ویژگی شرکتپذیری<sup>۱</sup> دارد: اگر  $A$ ،  $B$  و  $C$  به ترتیب ماتریسهایی از مرتبه  $m \times n$  و  $n \times p$  و  $p \times q$  باشند، آنگاه داریم

$$(AB)C = A(BC) \quad (۷.۴)$$

این مطلب را می‌توان به صورت زیر توجیه کرد: چون  $A$  از مرتبه  $m \times n$  و  $B$  از مرتبه  $n \times p$  است، پس  $AB$  تعریف شده و از مرتبه  $m \times p$  است و بنابراین  $(AB)C$  تعریف شده و از مرتبه  $m \times q$  است. به طریق مشابه می‌توان نشان داد که  $A(BC)$  نیز تعریف شده و از مرتبه  $m \times q$  است، بنابراین حداقل یک شرط برای تساوی برقرار است. بعلاوه:

$$(AB)C \text{ در } m(i, j) = \sum_{k=1}^p [AB \text{ در } m(k, i)] \times c_{kj}$$

$$\begin{aligned}
 &= \sum_{k=1}^p \left[ \sum_{r=1}^n a_{ir} \times b_{rk} \right] \times c_{kj} \\
 &= \sum_{r=1}^n a_{ir} \times \left[ \sum_{k=1}^p b_{rk} \times c_{kj} \right] \\
 &= \sum_{r=1}^n a_{ir} \times [BC \text{ در } (r, j)] \\
 &= A(BC) \text{ در } (i, j) \text{ درایه}
 \end{aligned}$$

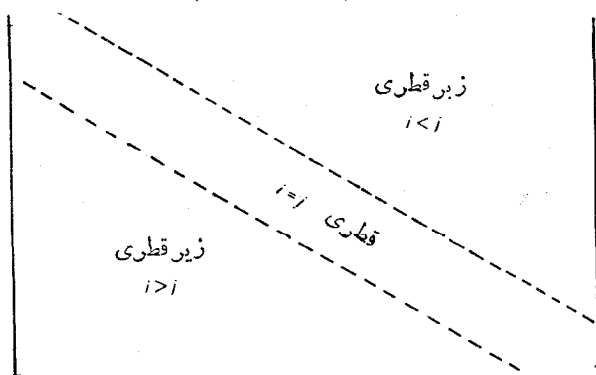
که ثابت می‌کند  $(AB)C = A(BC)$ . در حالت خاصی که  $C$  یک بردار (از مرتبهٔ مناسب باشد) کراً از تساوی بالا استفاده خواهیم کرد یعنی

$$(AB)x = A(Bx)$$

### ماتریسهای قطری و مثلثی

اگر  $A = (a_{ij})$  یک ماتریس مربعی از مرتبهٔ  $n$  باشد، درایه‌های  $a_{11}, a_{22}, \dots, a_{nn}$  را درایه‌های قطری و درایه‌های دیگر را درایه‌های غیرقطری گویند. تمامی درایه‌های  $a_{ij}$  به ازای  $i < j$  درایه‌های زیر قطری و تمام درایه‌های  $a_{ij}$  به ازای  $i > j$  درایه‌های زیر قطری نامیده می‌شوند (به شکل ۱۰۴ نگاه کنید).

اگر کلیهٔ درایه‌های غیرقطری ماتریس مربعی  $A$  صفر باشند،  $A$  ماتریس قطری نامیده می‌شود. اگر کلیهٔ درایه‌های زیر قطری ماتریس مربعی  $A$  صفر باشند، آنگاه  $A$  یک ماتریس بالامثلثی (یا راست مثلثی) نامیده می‌شود، در حالی که اگر تمام درایه‌های زیر قطری ماتریس مربعی  $A$  صفر باشند  $A$  را ماتریس پایین مثلثی (یا چپ مثلثی) گویند. واضح است که یک ماتریس قطری است اگر، و فقط اگر، هم بالامثلثی و هم پایین مثلثی باشد.



شکل ۱۰۴

□ مثال:

در نمونه‌های زیر ماتریسهای  $A$  و  $C$  قطری و  $A$  و  $B$ ،  $C$  و ماتریسهای بالامثلی و  $A$  و  $C$  و  $D$  ماتریسهای پایین‌مثلی هستند و  $E$  هیچکدام از این ویژگیها را ندارد

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad B = \begin{bmatrix} 0 & 1 & 2 \\ 0 & 3 & 1 \\ 0 & 0 & 4 \end{bmatrix} \quad C = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\square \quad D = \begin{bmatrix} 10 & 0 & 0 \\ 8 & 7 & 0 \\ 1 & 1 & 1 \end{bmatrix} \quad E = \begin{bmatrix} 1 & 0 & 0 \\ 3 & 2 & 6 \\ 0 & 0 & 4 \end{bmatrix}$$

### ماتریس واحد و وارون ماتریس

یک ماتریس قطری از مرتبه  $n$  را که تمام درایه‌های قطری آن یک باشند ماتریس واحد از مرتبه  $n$  گویند و با حرف مخصوص  $I$  یا  $I_n$  (در صورتی که مرتبه مهم باشد) نشان می‌دهند. این نام برای ماتریس واحد بدین منظور انتخاب شده است که

$$I_n A = A \quad \text{به‌ازای هر ماتریس } A \text{ از مرتبه } n \times p \text{ داریم:}$$

$$B I_n = B \quad \text{برای هر ماتریس } B \text{ از مرتبه } m \times n \text{ داریم:}$$

مشاهده می‌شود که ماتریس  $I$  دقیقاً همان نقش عدد یک را در ضرب معمولی ایفا می‌کند. تقسیم ماتریسها، در حالت کلی تعریف نشده است. لیکن برای ماتریسهای مربعی در ارتباط با مفهوم تقسیم، مفهومی موسوم به وارون ماتریس تعریف شده است. گویند ماتریس مربعی  $A$  از مرتبه  $n$  وارون‌پذیر است به شرط آنکه ماتریسی مانند  $B$  از مرتبه  $n$  وجود داشته باشد به طوری که

$$AB = I = BA \quad (۸.۴)$$

به‌عنوان مثال، ماتریس  $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$  وارون‌پذیر است، زیرا داریم

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

از طرف دیگر ماتریس  $A = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}$  وارون‌پذیر نیست، برای اینکه اگر  $B$  ماتریسی

باشد به طوری که  $BA=I$ ، آنگاه نتیجه می شود که

$$\begin{bmatrix} b_{11} + 2b_{12} & 2b_{11} + 4b_{12} \\ b_{21} + 2b_{22} & 2b_{21} + 4b_{22} \end{bmatrix} = BA = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

بنابراین می بایستی هم داشته باشیم  $b_{11} + 2b_{12} = 1$  و همزمان با آن

$$2(b_{11} + 2b_{12}) = 2b_{11} + 4b_{12} = 0$$

و این، غیرممکن است.

توجه داشته باشید که رابطه (۸.۴) می تواند حداکثر برای يك ماتریس  $B$  برقرار باشد. زیرا اگر داشته باشیم

$$AB=I \quad \text{و} \quad CA=I$$

که در آن  $B$  و  $C$  ماتریسهای مربعی و از مرتبه مشابه  $A$  باشند، آنگاه داریم

$$C = CI = C(AB) = (CA)B = IB = B$$

و این بدان معنی است که  $B$  و  $C$  باید مساوی باشند. بنا بر این اگر  $A$  وارونپذیر باشد، دقیقاً يك ماتریس  $B$  وجود دارد که رابطه (۸.۴) را برقرار می کند، این ماتریس، وارون  $A$  نامیده شده و با  $A^{-1}$  نشان داده می شود.

بلافاصله از (۸.۴) نتیجه می شود که اگر  $A$  وارونپذیر باشد،  $A^{-1}$  نیز چنین است و وارون آن  $A$  می باشد، یعنی

$$(A^{-1})^{-1} = A \quad (9.4)$$

بعلاوه اگر ماتریسهای  $A$  و  $B$  هر دو مربعی و وارونپذیر و از يك مرتبه باشند، حاصلضرب آنها نیز وارونپذیر است و داریم

$$(AB)^{-1} = B^{-1}A^{-1} \quad (10.4)$$

به تغییر در ترتیب ماتریسها توجه کنید! اثبات (۱۰.۴) به شرح کتبخیری ضرب ماتریسها منوط است

$$(AB)(B^{-1}A^{-1}) = A(BB^{-1})A^{-1} = AA^{-1} = I$$

$$(B^{-1}A^{-1})(AB) = B^{-1}(A^{-1}A)B = B^{-1}B = I$$

□ مثال: وارون ماتریس  $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$  ماتریس  $A^{-1} = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix}$  و وارون

ماتریس  $B = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$  ماتریس  $B^{-1} = \begin{bmatrix} 1 & 0 \\ -1 & 1 \end{bmatrix}$  است. بعلاوه  $AB = \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}$

و بنابراین با استفاده از (۱۰.۴) داریم  $B^{-1}A^{-1} = \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix} = (AB)^{-1}$ . از طرف دیگر

$$A^{-1}B^{-1} = \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix}$$

و

$$(AB)(A^{-1}B^{-1}) = \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix} = \begin{bmatrix} 3 & -1 \\ 1 & 0 \end{bmatrix} \neq I$$

□

یعنی  $A^{-1}B^{-1}$  نمی تواند وارون  $AB$  باشد.

### جمع ماتریسها و ضرب ماتریس در یک اسکالر

ممکن است یک ماتریس را در یک اسکالر  $\lambda$  (= عدد) ضرب و دو ماتریس از مرتبه‌های مساوی را به‌طور مناسبی باهم جمع کرد. اول اگر  $A = (a_{ij})$  و  $B = (b_{ij})$  دو ماتریس باشند و  $d$  یک عدد، گوئیم  $B$  حاصلضرب  $d$  در  $A$  یا  $B = dA$  در صورتی که  $A$  دارای مرتبه‌های مساوی باشند و به‌ازای جمیع مقادیر  $i$  و  $j$  داشته باشیم

$$b_{ij} = da_{ij}$$

بعلاوه اگر  $A = (a_{ij})$  و  $B = (b_{ij})$  ماتریسهایی از مرتبه‌های مساوی باشند و  $C = (c_{ij})$  یک ماتریس باشد، گوئیم  $C$  حاصلجمع  $A$  و  $B$  است، یا  $C = A + B$ ، در صورتی که مرتبه مساوی مرتبه‌های  $A$  و  $B$  باشد و به‌ازای جمیع مقادیر  $i$  و  $j$  داشته باشیم:

$$c_{ij} = a_{ij} + b_{ij}$$

بنابراین ضرب یک ماتریس در یک عدد و جمع ماتریسها، درایه به‌درایه انجام می‌گیرد. قواعد امر در رابطه با عملیات فوق و همچنین ضرب ماتریسها به‌آسانی قابل تحقیق هستند. فرض کنید  $A$ ،  $B$  و  $C$  ماتریسهایی باشند که تمام حاصلجمعها و حاصلضربهای مذکور در زیر برای آنها تعریف شده باشند و  $a$  و  $b$  نیز دو عدد باشند، آنگاه



$$(i) \quad A+B=B+A$$

$$(ii) \quad (A+B)+C=A+(B+C)$$

$$(iii) \quad a(A+B)=aA+aB$$

$$(iv) \quad (a+b)A=aA+bA$$

$$(v) \quad (A+B)C=AC+BC \quad (11.4)$$

$$(vi) \quad A(B+C)=AB+AC$$

$$(vii) \quad a(AB)=(aA)B=A(aB)$$

(viii) اگر  $a \neq 0$  و  $A$  وارونپذیر باشد، آنگاه  $aA$  وارونپذیر است و داریم

$$(aA)^{-1} = (1/a)A^{-1}$$

برای توضیح بیشتر برهان (vi) در زیر داده می‌شود. اگر ماتریس  $A$  از مرتبهٔ  $m \times n$  و  $B$  و  $C$  از مرتبهٔ  $n \times p$  باشند، هر دو طرف (vi)، ماتریسهایی تعریف شده و از مرتبهٔ  $m \times p$  خواهند بود. بعلاوه

$$\begin{aligned} A(B+C) \text{ در } (i, j) \text{ درایه} &= \sum_{k=1}^n a_{ik} \times [B+C \text{ در } (k, j)] \\ &= \sum_{k=1}^n a_{ik} \times [b_{kj} + c_{kj}] \\ &= \sum_{k=1}^n a_{ik} b_{kj} + \sum_{k=1}^n a_{ik} c_{kj} \\ &= [AB \text{ در } (i, j)] + [AC \text{ در } (i, j)] \\ &= (AB+AC) \text{ در } (i, j) \text{ درایه} \end{aligned}$$

و بالاخره اگر کلیهٔ درایه‌های ماتریس  $A$  از مرتبهٔ  $m \times n$  صفر باشند، آنگاه این ماتریس را ماتریس صفر از مرتبهٔ  $m \times n$  گویند و با حرف مخصوص  $O$  نشان می‌دهند. ماتریس صفر دارای ویژگی آشکار زیر است که به ازای کلیهٔ ماتریسهای هم‌مرتبهٔ خودش مانند  $B$ ، داریم

$$B+O=B$$

### ترکیبهای خطی

تعریف جمع ماتریسها و حاصلضرب اعداد در ماتریسها، بخصوص، امکان جمع  $n$ -بردارها

با هم و ضرب  $n$ -برداریهسا در اعداد را به وجود می آورد. اگر  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}$  و  $k$  برداری و  $b_1, b_2, \dots, b_k$  عدد باشند، آنگاه حاصلجمع وزن ۱ دار

$$b_1 \mathbf{x}^{(1)} + b_2 \mathbf{x}^{(2)} + \dots + b_k \mathbf{x}^{(k)}$$

نیز يك  $n$ -برداري است که به ترکیب خطی بردارهای  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}$  با وزنهای  $b_1, b_2, \dots, b_k$  موسوم است.

اکنون بار دیگر دستگاه معادلات (۱۰۴) را بررسی می کنیم. به ازای  $n, \dots, 2, 1, j$  گیرسیم  $\mathbf{a}_j$  معرف ستون  $j$ ام ماتریس ضرایب  $A$  از مرتبه  $m \times n$  باشد، یعنی  $\mathbf{a}_j$  يك  $m$ -برداري باشد که درایه  $j$ ام آن عدد  $a_{ij}$  به ازای  $i = 1, \dots, m$  باشد، در این صورت می توان  $m$ -برداري  $A\mathbf{x}$  را به صورت

$$A\mathbf{x} = x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \dots + x_n \mathbf{a}_n$$

یعنی به عنوان يك ترکیب خطی از  $n$  ستون  $A$  با درایه های وزن دار  $\mathbf{x}$  نوشت. بنا بر این حل مسئله (۱۰۴) با بیان زیر هم ارز است: وزنهای  $x_1, \dots, x_n$  را طوری پیدا کنید که ترکیب خطی  $n$  ستون  $A$  با این وزنهای مساوی  $m$ -برداري  $\mathbf{b}$  سمت راست شود. در رابطه با نشانگرهای سازگار فوق، ستون  $j$ ام ماتریس واحد  $I$  با نماد خاص

$$\mathbf{i}_j$$

نشان داده می شود. واضح است که تمامی درایه های  $\mathbf{i}_j$  صفرند به استثنای عنصر  $j$ ام که يك است. چنین مرسوم است که  $\mathbf{i}_j$  را  $j$ امین بردار واحد می نامند (همان گونه که با ماتریس واحد عمل کردیم، طول و یا مرتبه  $\mathbf{i}_j$  را صریحاً نشان نمی دهیم زیرا که مرتبه آن از متن مشخص می شود). با این نشانگرها برای هر  $n$ -برداري  $\mathbf{b} = (b_i)$  داریم

$$b_1 \mathbf{i}_1 + b_2 \mathbf{i}_2 + \dots + b_n \mathbf{i}_n = \mathbf{b}$$

بعلاوه، ستون  $j$ ام ماتریس  $A$ ، یعنی  $\mathbf{a}_j$  را می توان به وسیله ضرب  $A$  در  $\mathbf{i}_j$  نیز به دست آورد، یعنی

$$\mathbf{a}_j = A\mathbf{i}_j$$

بنا بر این اگر  $C = AB$ ، آنگاه

$$C\mathbf{i}_j = (AB)\mathbf{i}_j = A(B\mathbf{i}_j)$$

لذا ستون  $j$ ام حاصلضرب  $AB$  از ضرب اولین عامل  $A$  در ستون  $j$ ام عامل دوم  $B$  به دست می آید.

## 1. weighted sum

## وجود و یکتایی جوابهای دستگاه معادلات (۱.۴)

در بخشهای بعد، مختصراً با دستگاه معادلات خطی، که ماتریس ضرایب مربعی دارند، سروکار خواهیم داشت. این مطلب بدین صورت توجیه می‌شود که نشان می‌دهیم دستگاه معادلات (۱.۴) نمی‌تواند به‌ازای هر عامل سمت راست دقیقاً يك جواب داشته باشد مگر اینکه ماتریس ضرایب مربعی باشد.

لم ۱.۴ اگر  $x = x_1$  يك جواب دستگاه معادله خطی  $AX = b$  باشد، آنگاه هر جواب  $x = x_2$  از این دستگاه به‌شکل زیر است

$$x_2 = x_1 + y$$

در رابطه بالا  $x = y$  يك جواب از دستگاه همگن  $AX = 0$  است. درحقیقت اگر  $x_1$  و  $x_2$  هر دو، ریشه دستگاه  $AX = b$  باشند آنگاه داریم:

$$A(x_2 - x_1) = Ax_2 - Ax_1 = b - b = 0.$$

یعنی تفاضل آنها  $y = x_2 - x_1$  ریشه دستگاه همگن  $AX = 0$  است.

□ مثال: دستگاه خطی

$$x_1 + 2x_2 = 3$$

$$2x_1 + 4x_2 = 6$$

دارای جواب (ریشه)  $x_1 = x_2 = 1$  است. دستگاه همگن مربوطه به‌صورت

$$x_1 + 2x_2 = 0$$

$$2x_1 + 4x_2 = 0$$

است که دارای جوابهای  $x_2 = a$  و  $x_1 = -2a$  است که در آن  $a$  عددی است دلخواه. بنابراین هر جواب دستگاه اولیه دارای شکل  $x_1 = 1 - 2a$  و  $x_2 = 1 + a$  است که در آن  $a$  عددی است دلخواه. □

قضایای زیر از لم فوق نتیجه می‌شوند.

قضیه ۱.۴ دستگاه خطی  $AX = b$  حداکثر دارای يك جواب است (یعنی اگر جواب داشته باشد این جواب یکتاست) اگر، و فقط اگر، دستگاه همگن مربوطه  $AX = 0$  فقط دارای جواب «بدیهی»  $x = 0$  باشد.

بعد ثابت خواهیم کرد که نمی‌توان امیدوار بود که دستگاه خطی فوق يك جواب

یکتا داشته باشد، مگر اینکه تعداد معادلات دستگاه خطی حداقل برابر با تعداد مجهولات آن باشد.

**قضیه ۲.۴** هر دستگاه خطی همگن با تعداد معادلاتی کمتر از تعداد مجهولات، دارای جوابهای نامبدیهی (یعنی غیرصفر) است.

باید ثابت کنیم که اگر  $A$  یک ماتریس  $m \times n$  با  $m < n$  باشد آنگاه می توان  $\mathbf{y} \neq \mathbf{0}$  را طوری پیدا کرد که  $A\mathbf{y} = \mathbf{0}$ . این عمل با استقرا نسبت به  $n$  انجام می گیرد. اول حالت  $n=2$  را بررسی می کنیم. در این حالت می توانیم فقط یک معادله داشته باشیم

$$a_{11}x_1 + a_{12}x_2 = 0$$

و این معادله در صورت  $a_{12} = 0$ ، دارای جواب نامبدیهی  $x_1 = 0$ ،  $x_2 = 1$  و در غیر این صورت دارای جوابهای نامبدیهی  $x_1 = a_{12}$ ،  $x_2 = -a_{11}$  می باشد. بدین صورت قضیه را برای  $n=2$  اثبات کردیم. حال گیریم  $n > 2$ ، و فرض کنید ثابت شده است که هر دستگاه همگنی که تعداد معادلاتش کمتر از تعداد مجهولات و تعداد مجهولاتش کمتر از  $n$  باشد، دارای جوابهای نامبدیهی است و بعلاوه، گیریم  $A\mathbf{x} = \mathbf{0}$  یک دستگاه خطی همگن با  $n$  معادله و  $n$  مجهول،  $m < n$  باشد. باید ثابت کنیم که این دستگاه دارای جوابهای نامبدیهی است. اگر ستون  $n$ ام  $A$  صفر باشد، یعنی اگر  $\mathbf{a}_n = \mathbf{0}$ ، این دستگاه قطعاً جوابهای نامبدیهی دارد. زیرا در این حالت  $n$ -بردار غیر صفر  $\mathbf{x} = \mathbf{i}_n$  یک جواب دستگاه است. در غیر این صورت یکی از عناصر  $\mathbf{a}_n$  باید غیر صفر باشد، مثلاً فرض کنیم

$$a_{in} \neq 0$$

در این حالت، ماتریس  $B$  از مرتبه  $(n-1) \times m$  را، که ستون  $j$ ام آن به صورت

$$\mathbf{b}_j = \mathbf{a}_j - \frac{a_{ij}}{a_{in}} \mathbf{a}_n \quad j = 1, \dots, n-1$$

است مورد بررسی قرار می دهیم. اگر بتوانیم نشان دهیم که دستگاه همگن  $B\mathbf{x} = \mathbf{0}$  دارای جوابهای نامبدیهی است، اثبات انجام گرفته است. زیرا اگر بتوان اعداد  $x_1, x_2, \dots, x_{n-1}$  را که همگی صفر نیستند، به گونه ای یافت که رابطه

$$x_1 \mathbf{b}_1 + x_2 \mathbf{b}_2 + \dots + x_{n-1} \mathbf{b}_{n-1} = \mathbf{0}$$

برقرار باشد در این صورت از تعریف بردارهای  $\mathbf{b}_j$  نتیجه می شود که

$$x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \dots + x_{n-1} \mathbf{a}_{n-1} + \left( - \sum_{j=1}^{n-1} x_j \frac{a_{ij}}{a_{in}} \right) \mathbf{a}_n = \mathbf{0}$$

رابطه بالا یک جواب نامبدیهی برای معادله  $A\mathbf{x} = \mathbf{0}$  به دست می دهد. بنابراین تنها باید

نشان داد که  $B\mathbf{x} = \mathbf{0}$  دارای جوابهای نابدیهی است. برای اثبات این موضوع توجه کنید که به ازای هر  $j$ ، درایهٔ  $i$ ام بردار  $\mathbf{b}_j$  برابر است با

$$a_{ij} - \frac{a_{ij}}{a_{in}} a_{in} = 0$$

به طوری که معادلهٔ  $i$ ام دستگاه  $B\mathbf{x} = \mathbf{0}$  به صورت زیر درمی آید

$$0 \times x_1 + 0 \times x_2 + \dots + 0 \times x_{n-1} = 0$$

و بنابراین به ازای هر انتخابی از  $x_1, \dots, x_{n-1}$  برقرار است. از اینجا نتیجه می شود که  $\mathbf{x} = \mathbf{y}$ ، ریشهٔ  $B\mathbf{x} = \mathbf{0}$  است اگر، و فقط اگر،  $\mathbf{x} = \mathbf{y}$  ریشهٔ دستگاه همگن  $\hat{B}\mathbf{x} = \mathbf{0}$  باشد، که از  $B\mathbf{x} = \mathbf{0}$  فقط با صرف نظر کردن از معادلهٔ  $i$ ام به دست می آید. اما اکنون يك دستگاه خطی همگنی است با  $m-1$  معادله و  $n-1$  مجهول یعنی با معادلاتی کمتر از مجهولات و مجهولاتی کمتر از  $n$ . بنابراین با فرض استقرای  $\hat{B}\mathbf{x} = \mathbf{0}$  دارای جوابهای نابدیهی است که این امر خاتمهٔ اثبات است.

□ مثال: دستگاه خطی همگن  $A\mathbf{x} = \mathbf{0}$  را که به صورت زیر داده شده، بررسی کنید

$$x_1 + 2x_2 - x_3 = 0$$

$$x_1 - x_2 + x_3 = 0$$

برای این دستگاه  $m=2$  و  $n=3$ ، به دنبال برهان قضیهٔ ۲.۴ به طریق زیر يك جواب نابدیهی می سازیم: از آنجا که  $a_{23} \neq 0$ ،  $i=2$  را انتخاب می کنیم و داریم

$$\mathbf{b}_1 = \mathbf{a}_1 - \frac{a_{21}}{a_{23}} \mathbf{a}_2 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - (1/1) \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ 0 \end{bmatrix}$$

$$\mathbf{b}_2 = \mathbf{a}_2 - \frac{a_{22}}{a_{23}} \mathbf{a}_1 = \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} - (-1/1) \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

بنابراین دستگاه همگن کوچکتر  $B\mathbf{x} = \mathbf{0}$  عبارت است از

$$2x_1 + x_2 = 0$$

$$0 \times x_1 + 0 \times x_2 = 0$$

می توانیم از معادلهٔ دوم صرف نظر کنیم، در این صورت دستگاه همگن  $\hat{B}\mathbf{x} = \mathbf{0}$  فقط شامل يك معادلهٔ زیر خواهد شد

ماتریسها و دستگاههای معادلات خطی ۱۷۷

$$2x_1 + x_2 = 0$$

یک جواب نابدیهی برای این معادله،  $x_1 = 1$  و  $x_2 = -2$  است. بنا بر این با احتساب  
۳- برداری

$$x_3 = -[1(1/1) + (-2)(-1/1)] = -3$$

□ یک  $\mathbf{x} = (x_j)$  جواب نابدیهی از دستگاه اولیه است.

اکنون ثابت می‌کنیم که نمی‌توان انتظار داشت که به ازای هر مقدار دلخواه بردار سمت راست  $\mathbf{b}$ ، این دستگاه یک جواب داشته باشد، مگر اینکه تعداد معادلات بیشتر از تعداد مجهولات نباشد.

لم ۲.۴ اگر  $A$  یک ماتریس  $m \times n$  و دستگاه خطی  $A\mathbf{x} = \mathbf{b}$  دارای یک جواب به ازای هر  $m$ - برداری  $\mathbf{b}$  باشد، در این صورت یک ماتریس  $n \times m$  مانند  $C$  وجود دارد به طوری که

$$AC = I_m$$

چنین ماتریسی را می‌توان به صورت زیر تشکیل داد: طبق فرض قضیه، بدون توجه به چه بود  $\mathbf{b}$  می‌توانیم یک جواب برای  $A\mathbf{x} = \mathbf{b}$  پیدا کنیم. بنا بر این اگر  $\mathbf{b}$  را مساوی ستون  $j$ ام  $I$  انتخاب کنیم، می‌توانیم یک  $n$ - برداری  $\mathbf{c}_j$  طوری پیدا کنیم که

$$A\mathbf{c}_j = \mathbf{i}_j \quad j = 1, \dots, m$$

در این صورت با ماتریس  $C$  از مرتبه  $n \times m$  که ستون  $j$ ام آن عبارت است از  $\mathbf{c}_j$ ،  
 $j = 1, \dots, m$

$$(AC)\mathbf{i}_j = A(C\mathbf{i}_j) = A\mathbf{c}_j = \mathbf{i}_j = I\mathbf{i}_j \quad j = 1, \dots, m$$

رابطه بالا نشان می‌دهد که به ازای  $j = 1, \dots, m$ ، ستون  $j$ ام حاصلضرب  $AC$ ، با ستون  $j$ ام ماتریس  $I$  برابر است و این به معنی  $AC = I$  است.

لم ۳.۴ اگر  $B$  و  $C$  ماتریسهایی چنان باشند که تساوی

$$BC = I$$

برقرار باشد، آنگاه دستگاه همگن  $C\mathbf{x} = \mathbf{0}$  فقط دارای جواب بدیهی  $\mathbf{x} = \mathbf{0}$  است. زیرا، اگر  $C\mathbf{x} = \mathbf{0}$ ، آنگاه داریم:

$$\mathbf{x} = I\mathbf{x} = (BC)\mathbf{x} = B(C\mathbf{x}) = B\mathbf{0} = \mathbf{0}$$

قضیه ۳.۴ اگر  $A$  یک ماتریس  $m \times n$  و دستگاه خطی  $A\mathbf{x} = \mathbf{b}$  دارای یک جواب به ازای هر  $m$ - برداری دلخواه  $\mathbf{b}$  باشد، آنگاه  $m \leq n$ .

برای اثبات، از لم ۲.۴ نتیجه می‌گیریم که یک ماتریس  $C$  از مرتبه  $n \times m$  وجود دارد به طوری که

$$AC = I$$

اما از لم ۳.۴ نتیجه می‌شود که دستگاه همگن  $Cx = 0$  فقط یک جواب بدیهی  $x = 0$  دارد. بنابراین طبق قضیه ۲.۴،  $C$  باید حداقل به تعداد سطرها، ستون داشته باشد، یعنی  $n \geq m$  که این امر پایان اثبات است.

اما می‌دانیم که، نمی‌توان انتظار داشت که به ازای هر بردار ممکن سمت راست دقیقاً یک جواب برای دستگاه (۱.۴) به دست آورد، مگر اینکه دستگاه دقیقاً به تعداد معادلات مجهول داشته باشد؛ یعنی، مگر اینکه ماتریس ضرایب، مربعی باشد. بنابراین از این به بعد فقط دستگاه‌های خطی با ماتریس ضرایب مربعی بررسی خواهند شد. برای چنین ماتریس‌های مربعی یک قضیه نهایی را ثابت خواهیم کرد.

**قضیه ۴.۴** گیریم  $A$  یک ماتریس  $n \times n$  باشد. احکام زیر هم‌ارزند:

- (i) دستگاه همگن  $Ax = 0$  فقط یک جواب بدیهی  $x = 0$  دارد.  
 (ii) به ازای هر بردار سمت راست  $b$ ، دستگاه  $Ax = b$  یک جواب دارد.  
 (iii)  $A$  وارون‌پذیر است.

اول ثابت می‌کنیم که صحت (i) منجر به صحت (ii) می‌شود. گیریم  $b$  یک  $n$ -بردار داده شده باشد. باید ثابت کنیم که معادله  $Ax = b$  یک جواب دارد. برای این امر گیریم  $D$  یک ماتریس از مرتبه  $(n+1) \times m$  باشد که  $n$  ستون اول آن با ستون‌های  $A$  مساوی و ستون  $(n+1)$ ام آن بردار  $b$  باشد. از آنجا که تعداد ستون‌های  $D$  بیش از تعداد سطرهاى آن است، با استفاده از قضیه ۲.۴ می‌توان یک  $(n+1)$ -بردار غیر صفر  $y$  طوری پیدا کرد که تساوی  $Dy = 0$  برقرار باشد، یعنی

$$y_1 a_1 + y_2 a_2 + \dots + y_n a_n + y_{n+1} b = 0 \quad (12.4)$$

واضح است که  $y_{n+1}$  نمی‌تواند صفر باشد. زیرا اگر  $y_{n+1}$  صفر باشد، آنگاه به ازای  $y \neq 0$  حداقل یکی از اعداد  $y_1, \dots, y_n$  باید غیر صفر باشد، و درعین حال رابطه

$$y_1 a_1 + \dots + y_n a_n = 0$$

برقرار باشد. اما این رابطه بیان می‌کند که معادله  $Ax = 0$  به ازای  $i = 1, \dots, n$  جواب ناپذیری  $x_i = y_i$  را قبول می‌کند، که مغایر با (i) است. بنابراین چون  $y_{n+1} \neq 0$  می‌توانیم (۱۲.۴) را نسبت به  $b$  حل کنیم، که خواهیم داشت:

$$-\frac{y_1}{y_{n+1}} a_1 - \dots - \frac{y_n}{y_{n+1}} a_n = b$$

که این رابطه مبین آن است که معادله  $AX = \mathbf{b}$  يك جواب  $x_i = -(y_i/y_{n+1})$  به ازای  $i = 1, \dots, n$  دارد که این امر (ii) را ثابت می کند.

در مرحله بعد ثابت می کنیم که صحت (ii) منجر به صحت (iii) می شود. بنا فرض درست بودن (ii)، از لم ۲.۴ نتیجه می گیریم که يك ماتریس  $C$  از مرتبه  $n \times n$  وجود دارد به طوری که رابطه

$$AC = I$$

برقرار است. بنابراین با توجه به لم ۳.۴، معادله  $CX = \mathbf{0}$  تنها يك جواب بدیهی  $\mathbf{x} = \mathbf{0}$  دارد. این نکته مبین این است که ماتریس  $C$  از مرتبه  $n \times n$ ، در (i) صدق می کند و بنا بر برهان پیش،  $C$  در (ii) نیز صدق می کند. بنابراین، با استفاده از لم ۲.۴ يك ماتریس  $D$  از مرتبه  $n \times n$  وجود دارد به طوری که

$$CD = I$$

بدین ترتیب اثبات به انجام رسیده است. قبلا نشان داده شده بود که اگر رابطه زیر برای ماتریسهای مربعی  $A$  و  $C$  و  $D$  برقرار باشد

$$AC = I = CD$$

آنگاه  $C$  وارون پذیر است و

$$A = D = C^{-1}$$

بدین ترتیب  $A$  وارون يك ماتریس وارون پذیر است و بنا بر این خود نیز وارون پذیر است. سرانجام لم ۳.۴، نشان می دهد که از (iii) صحت (i) نتیجه می شود.

□ مثال: در مثال قبلی نشان داده شد که ماتریس  $2 \times 2$

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}$$

وارون پذیر نیست، و در مثال دیگری نشان دادیم که دستگاه همگن  $AX = \mathbf{0}$  دارای جواب نابديهی است. بنا بر این به موجب قضیه ۲.۴ دستگاه خطی  $AX = \mathbf{b}$  بایستی به ازای يك برداری قابل حل نباشد. در واقع با  $\mathbf{b} = \mathbf{i}_1$  دستگاه زیر به دست می آید

$$x_1 + 2x_2 = 1$$

$$2x_1 + 4x_2 = 0$$

این دستگاه جوابی ندارد، زیرا معادله دوم مستلزم  $0 = 2(x_1 + 2x_2)$  است، در حالی که معادله اول مستلزم  $2(x_1 + 2x_2) = 1$  است. □

اکنون به عنوان يك نمونه ساده از کاربرد قضیه ۳.۴، ثابت می کنیم که با داشتن ماتریس مربعی  $A$  و صحت رابطه  $AB = I$ ، درستی  $BA = I$  و  $B = A^{-1}$  نتیجه می شود. در واقع



اگر  $A$  يك ماتریس از مرتبه  $n \times n$  باشد، آنگاه رابطه  $AB = I$  بیان می‌کند که  $B$  از مرتبه  $n \times n$  است و به‌ازای هر  $n$ -برداري  $\mathbf{b}$ ، داریم  $A(B\mathbf{b}) = \mathbf{b}$ . اما این تساوی بیانگر این است که بدون توجه به  $\mathbf{b}$ ، می‌توان  $A\mathbf{x} = \mathbf{b}$  را نسبت به  $\mathbf{x}$  حل نمود و از این روطی قضیه ۴.۴،  $A$  وارونپذیر است و  $\mathbf{x} = B\mathbf{b}$  يك جواب آن. از این رو به‌ازای جمیع مقادیر  $\mathbf{b}$  خواهیم داشت  $B\mathbf{b} = A^{-1}\mathbf{b}$  یا  $B = A^{-1}$  و سرانجام  $BA = I$ .

### استقلال خطی و پایه‌ها

گیریم  $\mathbf{a}_1, \dots, \mathbf{a}_n$  تا  $n$ -برداري باشند و  $A$  يك ماتریس  $m \times n$  که ستون  $j$ ام آن  $\mathbf{a}_j$ ،  $j = 1, \dots, n$ ، است. این  $m$ -برداريها را مستقل خطی<sup>۱</sup> گوئیم اگر از رابطه  $x_1\mathbf{a}_1 + \dots + x_n\mathbf{a}_n = \mathbf{0}$  روابط  $x_1 = \dots = x_n = 0$  نتیجه شوند. در غیر این صورت بردارها را وابسته خطی گویند. واضح است که این  $n$  تا  $m$ -برداري، مستقل خطی هستند اگر، و فقط اگر، دستگاه همگن  $A\mathbf{x} = \mathbf{0}$  تنها دارای جواب بدیهی  $\mathbf{x} = \mathbf{0}$  باشد. بنابراین می‌توان از قضیه ۴.۴ نتیجه گرفت که هر مجموعه‌ای که بیش از  $m$  تا  $m$ -برداري داشته باشد باید وابسته خطی باشد.

گیریم  $\mathbf{a}_1, \dots, \mathbf{a}_n$  مستقل خطی باشند. اگر بتوان هر  $m$ -برداري  $\mathbf{b}$  را به‌صورت يك ترکیب خطی از این  $n$  تا  $m$ -برداري نوشت، آنگاه  $\mathbf{a}_1, \dots, \mathbf{a}_n$  را يك پایه (برای همه  $m$ -برداريها) می‌نامند. واضح است که  $\mathbf{a}_1, \dots, \mathbf{a}_n$  يك پایه است اگر، و فقط اگر، دستگاه خطی  $A\mathbf{x} = \mathbf{b}$  برای هر  $m$ -برداري  $\mathbf{b}$  فقط دارای يك جواب باشد، یعنی اگر، و فقط اگر، هر  $m$ -برداري را فقط به‌يك طریق بتوان به‌صورت يك ترکیب خطی از  $m$ -برداريهای  $\mathbf{a}_1, \dots, \mathbf{a}_n$  نوشت. به‌ویژه آنکه يك پایه (برای همه  $m$ -برداريها) دقیقاً شامل  $m$  تا  $m$ -برداري می‌شود (یعنی  $n = m$ ) و ماتریس مربوطه وارونپذیر است.

□ مثال: بردارهای

$$\mathbf{a}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{a}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

به‌طورخطی مستقل اند؛ اما يك پایه تشکیل نمی‌دهند، زیرا که فقط دو تا ۳-برداري هستند. علاوه بر ۲-برداري را می‌توان به‌صورت يك ترکیب خطی از سه تا ۲-برداري نوشت:

$$\mathbf{a}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{a}_2 = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}, \quad \mathbf{a}_3 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

اما این سه تا ۲-برداري تشكيل يك پایه نمی‌دهند، زیرا می‌باید به‌طورخطی مستقل باشند.

### 1. linearly independent

بالاخره سه تا ۳- برداری

$$\mathbf{a}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \mathbf{a}_2 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad \mathbf{a}_3 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

یک پایه تشکیل می‌دهند زیرا، ماتریس مربوطه وارونپذیر است. برای مشاهده این امر کافی است که بنا بر قضیه ۴.۴ ثابت کنیم که دستگاه

$$x_1 + x_2 + x_3 = 0$$

$$x_2 + x_3 = 0$$

$$x_3 = 0$$

فقط دارای جواب بدیهی  $x_1 = x_2 = x_3 = 0$  است، ولی این امر از خود معادلات آشکار است.  $\square$

### ماتریس ترانژاده

بالاخره، در ماتریسها عملی انجام پذیر است که هم ارز آن در حساب معمولی وجود ندارد و آن تشکیل ماتریس ترانژاده<sup>۱</sup> است. اگر  $A = (a_{ij})$  و  $B = (b_{ij})$  دو ماتریس باشند،  $B$  را ترانژاده  $A$  گویند و با  $B = A^T$  نشان می‌دهند، به شرط آنکه تعداد سطرهای  $B$  برابر تعداد ستونهای  $A$  و تعداد ستونهای آن برابر تعداد سطرهای  $A$  باشد و تساوی  $b_{ij} = a_{ji}$  به ازای جمیع مقادیر  $i$  و  $j$  درست باشد. به طور لفظی برای به دست آوردن  $A^T$ ، ترانژاده  $A$ ، «قرینه آن را نسبت به قطر» آن به دست می‌آوریم. اگر رابطه

$$A^T = A$$

برقرار باشد،  $A$  را متقارن<sup>۲</sup> نامند  
ماتریسهای

$$A = \begin{bmatrix} 1 & 3 & 2 \\ 3 & 0 & 4 \\ 2 & 4 & 5 \end{bmatrix} \quad B = \begin{bmatrix} 3 & -1 \\ 2 & 6 \\ 0 & 8 \end{bmatrix} \quad C = [3 \quad 4 \quad 7]$$

دارای ترانژاده‌های زیرند.

$$A^T = \begin{bmatrix} 1 & 3 & 2 \\ 3 & 0 & 4 \\ 2 & 4 & 5 \end{bmatrix} \quad B^T = \begin{bmatrix} 3 & 2 & 0 \\ -1 & 6 & 8 \end{bmatrix} \quad C^T = \begin{bmatrix} 3 \\ 4 \\ 7 \end{bmatrix}$$

به ویژه ترانزاده  $b^T$  ی یک بردار ستونی  $b$ ، یک بردار سطری است. به آسانی می توان در رابطه با ترانزش، صحت قواعد زیر را اثبات کرد:

۱. اگر  $A$  و  $B$  دو ماتریس چنان باشند که  $AB$  تعریف شده باشد، آنگاه  $B^T A^T$  نیز تعریف شده است و داریم:

$$(AB)^T = B^T A^T$$

۲. به ازای هر ماتریس  $A$  داریم  $(A^T)^T = A$

۳. اگر ماتریس  $A$  وارون پذیر باشد،  $A^T$  نیز وارون پذیر است و داریم

$$(A^T)^{-1} = (A^{-1})^T$$

برای اثبات قاعده (۱)، گوئیم اگر  $A$  یک ماتریس  $m \times n$  و  $B$  یک ماتریس  $n \times p$  باشد به طوری که  $AB$  یک ماتریس  $m \times p$  و  $(AB)^T$  یک ماتریس  $p \times m$  باشند، آنگاه  $A^T$  از مرتبه  $n \times m$  و  $B^T$  از مرتبه  $p \times n$  است، و بنا بر این حاصلضرب  $B^T A^T$  ماتریسی است به نحوی تعریف شده و از مرتبه  $p \times m$  است. سرانجام

$$\text{درایه } (j, i) \text{ ام در } AB = \text{درایه } (i, j) \text{ ام در } (AB)^T$$

$$= \sum_{k=1}^n [A \text{ در } (j, k) \text{ ام در } A] \times [B \text{ در } (k, i) \text{ ام در } B]$$

$$= \sum_{k=1}^n [B^T \text{ در } (i, k) \text{ ام در } B^T] \times [A^T \text{ در } (k, j) \text{ ام در } A^T]$$

$$= B^T A^T \text{ در } (i, j) \text{ ام در } B^T A^T$$

برای اثبات قاعده (۳) از قاعده (۱) نتیجه می گیریم که

$$A^T (A^{-1})^T = (A^{-1} A)^T = I^T = I$$

$$(A^{-1})^T A^T = (A A^{-1})^T = I^T = I$$

که قاعده (۳) را ثابت می کند.

اگر  $a$  و  $b$  دو  $n$ -برداری باشند، آنگاه  $b^T a$  یک ماتریس  $1 \times 1$  یا یک عدد است، که آن را حاصلضرب عددی  $a$  و  $b$  گویند در حالتی که  $a$  و  $b$  دو بردار حقیقی باشند.

دربارهٔ ماتریسهای که دارای درایه‌های همتافت (که در بحث مربوط به ویژه مقادیرها حائز اهمیت است) هستند، مفهوم ترانژادهٔ مزدوج یا هرمتیتی  $A^H$ ، پیش می‌آید. برای بیان این مطلب اشاره می‌کنیم که مزدوج يك عدد همتافت  $z$ ، یعنی  $\bar{z}$ ، از تغییر علامت بخش انگاری آن، به دست می‌آید. اگر  $z \neq 0$ ،  $\bar{z}$  آنگاه  $\bar{\bar{z}}$  عدد یگانه‌ای است مانند  $\alpha$  به طوری که  $\alpha z = |z|^2$ . تعیین هرمتیتی  $A$  یعنی  $A^H$ ، درست‌شده به تعیین ترانژادهٔ  $A^T$  است، جز آنکه به جای همهٔ درایه‌های  $A^T$ ، مزدوجهای همتافت آنها قرار داده می‌شوند. بنابراین در حالت

$$b_{ij} = \bar{a}_{ji} \quad j, i \text{ به ازای جمیع مقادیر}$$

رابطهٔ  $A^H = (b_{ij})$  برقرار است.

از این رو در حالتی که  $A$  يك ماتریس حقیقی\* باشد، خواهیم داشت  $A^H = A^T$ . توجه کنید که برای  $n$ -برداریهایی  $\mathbf{a}$  و  $\mathbf{b}$  که دارای درایه‌های همتافت هستند، حاصلضرب عددی معمولی بر  $\mathbf{a}$  برابر با عدد  $\mathbf{a}^H \mathbf{a}$  است نه  $\mathbf{b}^T \mathbf{a}$ ، زیرا عدد  $\mathbf{a}^H \mathbf{a}$  است که مربع طول بردار  $\mathbf{a}$  را به دست می‌دهد.

### جایگشتها و ماتریسهای جایگشتی

يك جایگشت<sup>۲</sup> از درجهٔ  $n$ ، بازآرایی است از اولین  $n$  عدد صحیح، یعنی يك توالی از  $n$  عدد صحیح که در آن هر عدد صحیح بین ۱ تا  $n$  حداقل یکبار، حداکثر نیز یکبار و بنا بر این دقیقاً یکبار در توالی مزبور ظاهر شود. يك جایگشت از درجهٔ  $n$  را می‌توان به شکل‌های مختلفی نوشت. برای اهداف ماکافی است (به تعبیری کاملاً دقیق) که يك جایگشت، به عنوان يك  $n$ -برداری  $\mathbf{p} = (p_i)$  با  $p_i \in \{1, 2, \dots, n\}$  به ازای جمیع مقادیر  $i$  و  $p_i \neq p_j$  به ازای  $j \neq i$ ، نگریسته شود. در اینجا تعداد  $n!$  جایگشت از درجهٔ  $n$  وجود دارد. بسته به اینکه تعداد برگردانها<sup>۳</sup>  $\mathbf{p}$  زوج یا فرد باشد، جایگشت  $\mathbf{p}$  را زوج یا فرد نامند. در اینجا منظور از برگردانهای يك جایگشت  $\mathbf{p} = (p_i)$ ، تعداد دفعاتی است که يك عدد صحیح بزرگتر پیش از يك عدد کوچکتر آمده باشد. برای مثال در جایگشت  $\mathbf{p}$  بسا

$$\mathbf{p}^T = [7, 2, 6, 3, 4, 1, 5]$$

۷ پیش از ۲، ۳، ۴، ۱، ۵	آمده که	۶	برگردانی به دست می‌دهد
۲ پیش از ۱	آمده که	۱	برگردانی به دست می‌دهد
۶ پیش از ۳، ۴، ۱، ۵	آمده که	۴	برگردانی به دست می‌دهد
۳ پیش از ۱	آمده که	۱	برگردانی به دست می‌دهد
۴ پیش از ۱	آمده که	۱	برگردانی به دست می‌دهد

بنابراین  $\mathbf{p}$  جمعاً ۱۳ برگردانی دارد

### 1. Hermitian

\* توضیح مترجم: منظور آن است که درایه‌های ماتریس  $A$  حقیقی باشند.

### 2. permutation

### 3. inversions

باید توجه داشت که هر تعویض دودرایه با هم در یک جایگشت، تعداد برگردانها را به تعداد فردی تغییر می‌دهد.

یک ماتریس جایگشتی از مرتبه  $n$ ، یک ماتریس  $P$  از مرتبه  $n \times n$  است که ستونها (سطرها)ی آن یک باز آرایش یا جایگشتی از ستونها (سطرها)ی ماتریس واحد از مرتبه  $n$  باشد. دقیقاً بگوییم، ماتریس  $P$  از مرتبه  $n \times n$  یک ماتریس جایگشتی است اگر به ازای یک جایگشت  $\mathbf{p} = (p_i)$  از درجه  $n$  داشته باشیم

$$P\mathbf{i}_j = \mathbf{i}_{p_j} \quad j = 1, \dots, n \quad (13.4)$$

قضیه ۵.۴ اگر  $P$  یک ماتریس جایگشتی باشد که در (۱۳.۴) صدق کند، آنگاه

$P^T(j)$  یک ماتریس جایگشتی است که به ازای  $n, \dots, 1, j$  در تساوی:

$$P^T\mathbf{i}_{p_j} = \mathbf{i}_j$$

صدق می‌کند. از این رو  $P^T P = I$  و بنابراین وارونپذیر است و داریم  $P^{-1} = P^T$ .

(ii) اگر  $A$  یک ماتریس  $m \times n$  باشد، آنگاه  $AP$  یک ماتریس  $m \times n$  است که ستون  $j$ ام

آن، به ازای  $n, \dots, 1, j$  مساوی است با ستون  $p_j$ ام ماتریس  $A$ .

(iii) اگر  $A$  یک ماتریس  $n \times m$  باشد، آنگاه  $P^T A$  ماتریسی است  $n \times m$  که سطر  $i$ ام

آن، به ازای  $n, \dots, 1, i$  مساوی است با سطر  $p_i$ ام ماتریس  $A$ .

□ مثال: ماتریس

$$P = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

یک ماتریس جایگشتی متناظر با جایگشت  $\mathbf{p}^T = [2 \ 3 \ 1]$  است، زیرا داریم

$$P\mathbf{i}_2 = \mathbf{i}_1, P\mathbf{i}_3 = \mathbf{i}_2, P\mathbf{i}_1 = \mathbf{i}_3$$

بعلاوه

$$P^T = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

از این رو تساویهای  $P^T\mathbf{i}_2 = \mathbf{i}_1$  و  $P^T\mathbf{i}_3 = \mathbf{i}_2$  برقرارند که قسمت (i) قضیه

(۵.۴) را نشان می‌دهد. بعلاوه می‌توان محاسبه کرد که مثلاً

$$AP = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 3 & 1 \\ 0 & 6 & 4 \\ 8 & 9 & 7 \end{bmatrix}$$

بنابراین ستون دوم ماتریس  $AP$  مساوی ستون سوم ماتریس  $A$ ، یعنی  $P_2 = P_3$  است که

قسمت (ii) قضیه (۵.۴) را نشان می‌دهد.  $\square$

### حل عددی دستگاه معادلات خطی

در اینجا فقط دستگاههای خطی

$$AX = b$$

را بررسی خواهیم کرد که به ازای هر مقدار سمت راست،  $b$ ، يك و فقط يك جواب دارند. بنا بر قضایای ۲.۴ و ۳.۴ باید خود را به دستگاههایی که تعداد مجهولات آنها دقیقاً با تعداد معادلات برابرند، یعنی به دستگاههایی که ماتریس ضرایبشان ماتریس مربعی است، محدود سازیم. برای این گونه دستگاهها، حکم قضیه ۴.۴ این است که شرط اینکه دستگاه به ازای جمیع مقادیر سمت راست،  $b$ ، فقط يك جواب داشته باشد،  $A$  می‌باید وارونپذیر باشد. بنا بر این فرض می‌کنیم که همه دستگاههای خطی مورد بحث ماتریس ضرایبی وارونپذیر دارند.

آزمونی که غالباً برای وارونپذیری به کار برده می‌شود مبتنی بر مفهوم دترمینان است. قضیه مربوطه چنین بیان می‌کند که ماتریس  $A$  وارونپذیر است اگر، و فقط اگر،  $\det(A) \neq 0$ . اگر  $\det(A) \neq 0$ ، حتی ممکن است جواب دستگاه  $AX = b$  بر حسب دترمینانها از راه به اصطلاح قاعده کرامرا، داده شود. ولی دترمینانها در حل دستگاههای خطی ارزش عملی ندارند، زیرا در حالت کلی محاسبه يك دترمینان، خود به اندازه حل يك دستگاه خطی دشواری دارد. به همین دلیل در اینجا برای حل دستگاههای خطی از دترمینان استفاده نمی‌شود و برای تعریف دترمینان نیز کوششی به عمل نمی‌آید. بسا این حال در بخش ۷.۴ برای ارزیابی دترمینانها (بر اساس روش مستقیم حل دستگاههای خطی) برای استفاده در موارد دیگر، روشی ارائه خواهد شد.

روشهای عددی برای حل دستگاههای خطی راممکن است به دو نوع، مستقیم و بارستی تقسیم کرد. روشهای مستقیم آنهایی هستند که در صورت نبودن خطای گرد کردن و خطاهای دیگر، جواب دقیق را با تعداد متناهی از عملیات مقدماتی حسابی به دست می‌دهند. در عمل، چون يك کامپیوتر با طول کلمه متناهی سروکار دارد، روشهای مستقیم معمولاً جوابهای دقیق نمی‌دهند. زیرا خطاهای ناشی از گرد کردن، ناپایداری، و از دست رفتن ارقام با معنی ممکن

است به نتایج خیلی نامطلوب و یا حتی بی معنی منجر شود. سروکار قسمت اعظم آنالیز عددی با چرا و چگونه پدید آمدن این خطاها و بررسی روشهایی است که این خطاها را بر روی هم حداقل می سازند. روش اساسی که برای حل مستقیم به کار گرفته می شود روش حذف گاوس است، ولی حتی در این دسته (روش مستقیم) نیز می توان روشهای متفاوتی انتخاب کرد که کارایی محاسباتی و دقت آنها متفاوت است. برخی از این روشها در بخشهای بعدی بررسی خواهند شد.

روشهای بارستی روشهایی هستند که با يك تقریب اولیه شروع می شوند و با به کار بستن يك الگوریتم انتخابی مناسب، متوالیاً به تقریبهایی بهتری دست می یابند. حتی اگر این روند همگرا باشد، می توان تنها امیدوار بود که از راه روشهای بارستی به يك حل تقریبی رسید. روشهای بارستی از لحاظ الگوریتمهای انتخابی و میزان همگرایی، متفاوت اند. برخی از روشهای بارستی ممکن است عملاً واگرا باشند و برخی دیگر ممکن است آنقدر آهسته همگرا شوند که کارایی محاسباتی نداشته باشند. امتیاز مهم روشهای بارستی در سادگی آنها و یکنواختی عملیاتی است که باید انجام شوند، و در نتیجه موجب می شوند برای کار برد کامپیوتری بسیار مناسب باشند و در قبال افزایش خطاهای گرد کردن عدم حساسیت نسبی داشته باشد.

ماتریسهای مربوط به دستگاههای خطی نیز به دو رده چگال<sup>۲</sup> و تنک<sup>۳</sup> تقسیم می شوند. در ماتریسهای چگال تعداد بسیار کمی از درایهها صفرند و مرتبه این گونه ماتریسها نسبتاً کم است. شاید از مرتبه ۱۰۰ یا کمتر. معمولاً کاراتر این است که مسائل مشتعل بر این نوع ماتریسها با روش مستقیم حل شوند. ماتریسهای تنک معمولاً تعداد خیلی کمی درایههای غیر صفر دارند. این گونه ماتریسها معمولاً در حل معادلات دیفرانسیل با روش تفاضلات متناهی به دست می آیند. مرتبه این گونه ماتریسها ممکن است بسیار بزرگ باشد و معمولاً در حل بارستی بسیار مناسب اند و این امتیازی است در قبال ماهیت تنک بودن ماتریسهای مورد نظر. روشهای بارستی برای دستگاههای خطی در فصل ۵ مورد بحث قرار خواهند گرفت.

## تمرین

۱-۱۰۴ گیریم

$$A = \begin{bmatrix} 2 & 3 & 1 \\ 1 & -1 & 1 \\ 0 & 2 & 2 \end{bmatrix} \quad B = \begin{bmatrix} 1 & 2 & 1 \\ -1 & 2 & -1 \\ 2 & 0 & 2 \end{bmatrix} \quad C = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 2 \end{bmatrix}$$

(الف) حاصل ضرب  $AB$  و  $BA$  را محاسبه کنید و نشان دهید که  $AB \neq BA$

(ب)  $(A+B)+C$  و  $A+(B+C)$  را پیدا کنید.

(پ) نشان دهید که  $A(BC) = (AB)C$ .

(ت) تحقیق کنید که  $(AB)^T = B^T A^T$ .

۲-۱۰۴ نشان دهید که ماتریس  $A$  که در زیر داده شده، وارونپذیر نیست (به قضیه ۴۰۴ نگاه کنید):

$$A = \begin{bmatrix} 1 & 2 & 1 \\ 2 & -1 & -1 \\ 6 & 2 & 0 \end{bmatrix}$$

۳-۱۰۴ برای ماتریس  $A$  که در زیر داده شده است، ماتریس جایگشتی  $P$  را طوری پیدا کنید که

(الف) ضرب  $A$  از سمت چپ در  $P$ ، ستونهای اول و چهارم  $A$  را با هم عوض

کند.

(ب) ضرب  $A$  از سمت راست در  $P$ ، سطرهای اول و سوم  $A$  را با هم عوض کند.

$$A = \begin{bmatrix} 4 & 1 & 2 & 1 \\ 3 & 2 & 1 & 1 \\ 1 & 2 & 0 & 1 \\ 1 & 1 & 0 & 1 \end{bmatrix}$$

۴-۱۰۴ در ماتریس  $A$ ی تمرین ۳-۱۰۴ يك دنباله از ماتریسهای جایگشتی پیدا کنید که  $A$  را به شکل زیر درآورد

$$A' = \begin{bmatrix} 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 2 \\ 1 & 1 & 3 & 2 \\ 1 & 2 & 2 & 1 \end{bmatrix}$$

۵-۱۰۴ دستگاه زیر را به شکل ماتریسی بنویسید و ماتریس  $A$  و بردار  $\mathbf{b}$  را مشخص کنید

$$2x_1 + 3x_2 + 4x_3 + x_4 = 1$$

$$x_1 + 2x_2 + x_4 = 0$$



$$2x_1 + 3x_2 + x_3 - x_4 = 2$$

$$x_1 - 2x_2 - x_3 - x_4 = 3$$

۶-۱۰۴ با اثبات عبارت زیر خود را متقاعد سازید که مفهوم وارونپذیری تنها برای ماتریسهای مربعی بامعنی است: گیریم  $A$  يك ماتریس  $m \times n$  باشد. اگر  $B$  و  $C$  ماتریسهای  $n \times m$  باشند به طوری که  $AB = I_m$  و  $CA = I_n$ ، آنگاه  $B = C = A^{-1}$ . به ویژه اینکه در این صورت  $m = n$ . [دانهجایی: اول ثابت کنید  $B = C$ ، سپس نشان دهید که  $m = n = \text{trace}(AB) = \text{trace}(BA)$  که  $\text{trace}$  (اثر) يك ماتریس به صورت حاصلجمع درایه‌های قطری آن تعریف شده است].

۷-۱۰۴ با استفاده از قضیهٔ ۴.۴ ثابت کنید که ماتریس جایگشتی وارونپذیر است.

۸-۱۰۴ با استفاده از قضیهٔ ۴.۴ ثابت کنید که اگر  $A$  و  $B$  ماتریسهای مربعی باشند به طوری که حاصلضرب آنها وارونپذیر باشد، آنگاه  $A$  و  $B$  باید هر دو وارونپذیر باشند.

۹-۱۰۴ آیا بردارهای زیر تشکیل يك پایه می‌دهند؟

$$\begin{bmatrix} 0 \\ -1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ -1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \\ -1 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

۱۰-۱۰۴ ثابت کنید که سه بردار

$$\begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ -1 \end{bmatrix}$$

تشکیل يك مجموعهٔ مستقل خطی می‌دهند. آیا يك پایه تشکیل می‌دهند؟

۱۱-۱۰۴ برای هر يك از سه عمل ماتریسی، یعنی جمع دو ماتریس، ضرب دو ماتریس و ضرب يك اسكالر در يك ماتریس، يك زیربرنامهٔ فورترن بنویسید که با ورودی مناسب محاسبه را انجام دهد و ماتریس حاصل را برگرداند.

۱۲-۱۰۴ اگر  $p(x) = c_0 + c_1x + c_2x^2 + \dots + c_kx^k$  يك بسجمله‌ای باشد و  $A$  يك ماتریس داده شده  $n \times n$ ، آنگاه ماتریس  $p(A)$  چنین تعریف می‌شود

$$p(A) = c_0A^0 + c_1A^1 + c_2A^2 + \dots + c_kA^k$$



به ویژه اینکه معادلهٔ آخر فقط شامل  $x_n$  است. بنابراین از آنجا که  $a_{nn} \neq 0$  باید داشته باشیم

$$x_n = \frac{b_n}{a_{nn}}$$

و چون حالا  $x_n$  در دست است معادلهٔ ماقبل آخر یعنی

$$a_{n-1, n-1} x_{n-1} + a_{n-1, n} x_n = b_{n-1}$$

تنها یک مجهول دارد که  $x_{n-1}$  است. چون  $a_{n-1, n-1} \neq 0$  خواهیم داشت

$$x_{n-1} = \frac{b_{n-1} - a_{n-1, n} x_n}{a_{n-1, n-1}}$$

با در دست داشتن  $x_n$  و  $x_{n-1}$  اکنون از معادلهٔ سوم از آخر، یعنی

$$a_{n-2, n-2} x_{n-2} + a_{n-2, n-1} x_{n-1} + a_{n-2, n} x_n = b_{n-2}$$

استفاده می‌کنیم که تنها یک مجهول واقعی  $x_{n-2}$  دارد. باز چون  $a_{n-2, n-2} \neq 0$  معادلهٔ فوق را می‌توان نسبت به  $x_{n-2}$  حل کرد

$$x_{n-2} = \frac{b_{n-2} - a_{n-2, n-1} x_{n-1} - a_{n-2, n} x_n}{a_{n-2, n-2}}$$

در حالت کلی که مقادیر  $x_{k+2}, x_{k+1}, x_k, \dots, x_n$  از قبل محاسبه شده‌اند و با توجه به  $a_{kk} \neq 0$  معادلهٔ  $k$ ام را می‌توان به تنهایی نسبت به  $x_k$  حل کرد، که در این صورت خواهیم داشت

$$x_k = \frac{b_k - \sum_{j=k+1}^n a_{kj} x_j}{a_{kk}}$$

این روند تعیین جواب دستگاه (۱۵.۴) را روند پسجایگذاری<sup>۱</sup> می‌نامند.

**الگوریتم ۱۰.۴:** پسجایگذاری ماتریس بالا مثلثی  $A$  از مرتبهٔ  $n \times n$  که تمامی درایه‌های قطری آن غیر صفرند و  $n$ -بردار  $\mathbf{b}$  داده شده‌اند. درایه‌های  $x_n, x_{n-1}, \dots, x_1$  متعلق به جواب  $\mathbf{x}$  از دستگاه  $\mathbf{Ax} = \mathbf{b}$  را می‌توان (به همین ترتیب) به وسیلهٔ

$$\left[ \begin{array}{l} \text{For } k = n, n-1, \dots, 1, \text{ do} \\ \quad x_k := \frac{b_k - \sum_{j=k+1}^n a_{kj} x_j}{a_{kk}} \end{array} \right.$$

## 1. back-substitution

به دست آورد.

در اینجا دو نکته قابل توجه است: وقتی  $k = n$ ، در این صورت حاصلجمع  $\sum_{j=k+1}^n$  به صورت  $\sum_{j=n+1}^n$  در می آید و به صورت حاصلجمع روی هیچ عضوی تعبیر می شود و طبق قرارداد مقدار آن صفر است. همچنین متذکر می شویم که با توجه به توضیح پسجایگذاری، تقریباً نتیجه زیر واضح است.

**قضیه ۶.۴** ماتریس بالامثلثی  $A$  وارونپذیر است اگر و فقط اگر همه درایه های قطری آن مخالف صفر باشند.

در واقع روش پسجایگذاری نشان می دهد که وقتی تمام درایه های قطری  $A$  غیر صفر باشند، دستگاه خطی  $AX = b$  به ازای بردار مفروض  $b$ ، حداکثر یک جواب دارد؛ از این رو بنا بر قضیه ۴.۴،  $A$  می باید وارونپذیر باشد. از سوی دیگر به ازای هر  $j = 1, \dots, n$ ، مقادیری مانند  $x_1, \dots, x_n$  که همه برابر صفر نیستند، وجود دارند به طوری که بنا بر قضیه ۲.۴ داریم

$$a_{11}x_1 + \dots + a_{1j}x_j = 0$$

$$\dots \dots \dots$$

$$a_{j-1,1}x_1 + \dots + a_{j-1,j}x_j = 0$$

اما اگر  $a_{jj} = 0$ ، بردار  $y = [x_1 \dots x_j \ 0 \dots 0]^T$  بردار صفر نیست، ولی در دستگاه  $Ay = 0$  صدق می کند که با توجه به قضیه ۴.۴، نتیجه می شود که  $A$  وارونپذیر نیست. بنابراین ما محق بوده ایم که بردار  $x$  را که با الگوریتم ۱.۴ محاسبه شده، جواب دستگاه (۱۵.۴) بنامیم.

□ مثال ۱۰.۴: دستگاه خطی زیر را در نظر می گیریم

$$2x_1 + 3x_2 - x_3 = 5$$

$$-2x_2 - x_3 = -7 \quad (16.4)$$

$$-5x_3 = -15$$

از معادله آخر نتیجه می گیریم  $x_3 = b_3/a_{33} = 15/5 = 3$ ، با این جواب، از معادله دوم چنین به دست می آوریم

$$x_2 = (b_2 - a_{23}x_3)/a_{22} = (-7 + 3)/(-2) = 2$$

لذا با توجه به معادله اول خواهیم داشت

$$\square \quad x_1 = (b_1 - a_{12}x_2 - a_{13}x_3)/a_{11} = (5 - 3 \times 2 + 3)/2 = 1$$

اما اگر ماتریس ضرایب دستگاه  $A\mathbf{x} = \mathbf{b}$  بالامتثلی نباشد، نخست روش حذف گاوس را بر این دستگاه اعمال می‌کنیم. احتمالاً دانشجویان در جبر مقدماتی با این روش آشنایی پیدا کرده‌اند. هدف از این روش تبدیل دستگاه مفروض به دستگاه هم‌ارزی است که ماتریس ضرایب آن هم‌ارز یک ماتریس بالامتثلی باشد. در این صورت دستگاه اخیر را می‌توان از راه پس‌جا یکنگداری حل کرد.

دو دستگاه خطی  $A\mathbf{x} = \mathbf{b}$  و  $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$  را هم‌ارز نامیم هر گاه هر جواب یک دستگاه یک جواب دیگری نیز باشد.

**قضیه ۷.۴** گیریم دستگاه خطی  $A\mathbf{x} = \mathbf{b}$  داده شده باشد، و فرض می‌کنیم که این دستگاه را تحت سلسله عملیاتی از نوع زیر قرار دهیم.

- (i) ضرب یک معادله در یک عدد ثابت مخالف صفر
- (ii) جمع مضربی از یک معادله با معادلهٔ دیگر
- (iii) تعویض جای دو معادله با هم

اگر این سلسله عملیات دستگاه جدید  $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$  را پدید آورد، آنگاه دستگاه‌های  $A\mathbf{x} = \mathbf{b}$  و  $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$  هم‌ارزند. به‌ویژه در این صورت،  $A$  وارون‌پذیر است اگر و فقط اگر،  $\tilde{A}$  وارون‌پذیر باشد.

برای دیدن یک برهان از این قضیه، به مسئلهٔ ۲.۴-۱۱ مراجعه کنید.

عمل حذف، بر پایهٔ قضیهٔ فوق و ملاحظات زیر انجام می‌گیرد: اگر  $A\mathbf{x} = \mathbf{b}$  یک دستگاه خطی باشد و اگر به‌ازای مقداری از  $k$  و  $j$  داشته باشیم  $a_{kj} \neq 0$ ، آنگاه می‌توان مجهول  $x_j$  را از یک معادلهٔ  $i \neq k$  با افزودن  $-(a_{ij}/a_{kj})$  برابر معادلهٔ  $k$ ام به معادلهٔ  $i$ ام، حذف کنیم. دستگاه حاصل  $\tilde{A}\mathbf{x} = \tilde{\mathbf{b}}$  با دستگاه اصلی هم‌ارز است.

روش حذف گاوس، در ساده‌ترین شکل خود، از یک دستگاه خطی مفروض  $A\mathbf{x} = \mathbf{b}$  از مرتبهٔ  $n$ ، به‌ازای  $k = 0, \dots, n-1$ ، یک سلسله از دستگاه‌های هم‌ارز  $A^{(k)}\mathbf{x} = \mathbf{b}^{(k)}$  را پدید می‌آورد. در اینجا  $A^{(0)}\mathbf{x} = \mathbf{b}^{(0)}$  درست همان دستگاه اصلی است. دستگاه  $(k-1)$ ام به‌شکل زیر است:

$$\begin{aligned} a_{11}^{(k-1)}x_1 + a_{12}^{(k-1)}x_2 + \dots + a_{1, k-1}^{(k-1)}x_{k-1} + a_{1k}^{(k-1)}x_k + \dots + a_{1n}^{(k-1)}x_n &= b_1^{(k-1)} \\ a_{22}^{(k-1)}x_2 + \dots + a_{2, k-1}^{(k-1)}x_{k-1} + a_{2k}^{(k-1)}x_k + \dots + a_{2n}^{(k-1)}x_n &= b_2^{(k-1)} \\ \dots & \dots \\ a_{k-1, k-1}^{(k-1)}x_{k-1} + a_{k-1, k}^{(k-1)}x_k + \dots + a_{k-1, n}^{(k-1)}x_n &= b_{k-1}^{(k-1)} \\ a_{kk}^{(k-1)}x_k + \dots + a_{kn}^{(k-1)}x_n &= b_k^{(k-1)} \\ \dots & \dots \\ a_{nn}^{(k-1)}x_n &= b_n^{(k-1)} \end{aligned}$$

به بیان لفظی،  $k$  معادلهٔ اول هم‌اکنون به‌شکل بالامتثلی درآمده‌اند، درحالی‌که  $n-k$  معادلهٔ

آخر تنها شامل مجهولات  $x_k, \dots, x_n$  هستند. از این رو دستگاه  $k$ ام  $A^{(k)}\mathbf{x} = \mathbf{b}^{(k)}$  در طی مرحله  $k$ ام روش حذف گاوس، به طریق زیر به دست آمده است:  $k$  معادله اول بدون تغییر مسانده اند، بعلاوه اگر ضریب  $x_k$ ، یعنی  $a_{kk}^{(k-1)}$  در معادله  $k$ ام صفر نباشد، آنگاه  $m_{ik} = a_{ik}^{(k-1)} / a_{kk}^{(k-1)}$  برابر معادله  $k$ ام از معادله  $i$ ام کم شده و از این طریق مجهول  $x_k$  از معادله  $i$ ام،  $i = k+1, \dots, n$  حذف شده است. روشن است که دستگاه حاصل  $A^{(k)}\mathbf{x} = \mathbf{b}^{(k)}$  هم ارز دستگاه  $A^{(k-1)}\mathbf{x} = \mathbf{b}^{(k-1)}$  است، لذا، با توجه به استقرا، هم ارز دستگاه اصلی است. بعلاوه اولین  $k+1$  معادله دستگاه  $k$ ام به شکل بالامثلثی است.

بعد از  $n-1$  مرحله از این عملیات، دستگاه  $A^{(n-1)}\mathbf{x} = \mathbf{b}^{(n-1)}$  به دست آمده است که ماتریس ضرایب آن به صورت بالامثلثی است و بنا بر این می توان این دستگاه را سریعاً از راه پسجایگذاری حل کرد.

□ مثال ۲۰۴: دستگاه خطی زیر را در نظر می گیریم

$$2x_1 + 3x_2 - x_3 = 5 \quad (\text{الف})$$

$$4x_1 + 4x_2 - 3x_3 = 3 \quad (\text{ب}) \quad (17.4)$$

$$-2x_1 + 3x_2 - x_3 = 1 \quad (\text{پ})$$

برای حذف  $x_1$  از معادلات (ب) و (پ)،  $2 = -4/2 = -2$ ، برابر معادله (الف) را با معادله (ب) جمع می کنیم تا معادله جدید زیر به دست آید

$$0x_1 - 2x_2 - x_3 = -7$$

همچنین  $1 = (-2)/2 = -1$  برابر معادله (الف) را با معادله (پ) جمع می کنیم، در نتیجه

$$0x_1 + 6x_2 - 2x_3 = 6$$

که این معادله دستگاه جدید  $A^{(1)}\mathbf{x} = \mathbf{b}^{(1)}$  را به صورت زیر به دست می دهد

$$2x_1 + 3x_2 - x_3 = 5 \quad (\text{الف})$$

$$-2x_2 - x_3 = -7 \quad (\text{ب}) \quad (18.4)$$

$$6x_2 - 2x_3 = 6 \quad (\text{پ})$$

این عملیات، مرحله اول روش حذف گاوس را برای این دستگاه کامل می کند. در مرحله دوم (و آخرین مرحله برای این مثال)،  $x_2$  را از معادله (پ) با جمع کردن  $3 = 6/(-2)$  برابر معادله (ب) با معادله (پ) حذف می کنیم، داریم

$$0x_2 - 5x_3 = -15$$

بنا بر این دستگاه جدید و نهایی

$$2x_1 + 3x_2 - x_3 = 5$$

$$-2x_2 - x_3 = -7$$

(۱۹.۴)

$$-5x_3 = -15$$

به دست می آید که بنا بر قضیهٔ ۷.۴ با دستگاه اولیهٔ (۱۷.۴) هم ارز است، اما ماتریس ضرایب آن به صورت ماتریس بالامثلثی است، از این رو می توان آن را مانند مثال ۱۰.۴ سریعاً با استفاده از پسجایگذاری حل کرد. □

در تشریح سادهٔ روش حذف گاوس که در بالا ذکر شد، ما از معادلهٔ  $k$ ام برای حذف  $x_k$  از معادلات  $1, \dots, m, \dots, m+1$ ، طی مرحلهٔ  $k$ ام عملیات، استفاده کردیم. البته این کار تنها در صورتی امکانپذیر است که در شروع مرحلهٔ  $k$ ام، ضریب  $x_k$  یعنی  $a_{kk}^{(k-1)}$  در معادلهٔ  $k$ ام صفر نباشد. متأسفانه به آسانی می توان دستگاهی خطی طرح کرد که در آن این شرط برقرار نباشد. برای مثال، اگر دستگاه  $AX = b$  با معادلات

$$x_2 + x_3 = 1 \quad (\text{الف})$$

$$x_1 + x_3 = 1 \quad (\text{ب}) \quad (20.4)$$

$$x_1 + x_2 = 1 \quad (\text{ب})$$

داده شده باشد، آنگاه غیرممکن است بتوان برای حذف  $x_1$  از معادله‌های دیگر، از معادلهٔ (الف) استفاده کرد. به منظور غلبه بر این مشکل به طوری که باز یک دستگاه بالامثلثی هم ارز با دستگاه مفروض نتیجه شود، باید در هر مرحله در انتخاب معادلهٔ لولایی<sup>۱</sup>، یعنی معادله‌ای که برای حذف یک مجهول از یکی از معادلات اختیار می شود، آزادی بیشتری داشت.

مثلاً در دستگاه (۲۰.۴) می توانستیم از معادلهٔ (ب) در مرحلهٔ اول حذف به عنوان معادلهٔ لولایی استفاده کنیم. برای حفظ شکل قبلی عملیات، نخست معادلهٔ (ب) را با عوض کردن با معادلهٔ (الف) در جای اول قرار می دهیم. اکنون در این ترتیب جدید دستگاه ضریب  $x_1$  در معادلهٔ (الف) غیر صفر بوده و مشابه بسا عملیات قبلی می توان پیش رفت تا دستگاه جدید  $A^{(1)}x = b^{(1)}$  به فرم زیر به دست آید:

$$x_1 + x_3 = 1 \quad (\text{الف})$$

$$x_2 + x_3 = 1 \quad (\text{ب})$$

$$x_2 - x_3 = 0 \quad (\text{ب})$$

بسا این دستگاه، مرحلهٔ دوم (و آخرین مرحله) از روش حذف گاوس، بدون هیچ مشکلی پیش می رود و دو دستگاه نهایی بالامثلثی به گونهٔ زیر حاصل می شود

$$x_1 + x_3 = 1$$

$$x_2 + x_3 = 1$$

$$-2x_3 = -1$$

که با پسجایگذاری جوابهای  $1/2$   $x_3 = x_2 = x_1 = 1/2$  به دست می آیند.

این آزادی عمل بیشتر در انتخاب معادلات لولایی نه فقط به دلیل امکان صفر بودن ضرایب ضروری است، بلکه تجربه نشان داده است که این امر برای مقابله با نتایج خطای گرد کردن نیز اساسی است (بهبخش ۳.۴ نگاه کنید). کار اضافی که باید انجام گیرد خیلی ناچیز است؛ در شروع مرحله  $k$ ام، بررسی در مورد ضریب غیر صفر  $x_k$  در معادله های  $k$ ام،  $(k+1)$ ام،  $\dots$ ،  $n$ ام انجام می گیرد، و اگر چنین ضریبی در حداقل یک معادله به ازای  $k > n$  به دست آمد، جای معادله های  $k$ ام با هم عوض می شود.

در حالتی که ماتریس  $A$  وارونپذیر باشد، لزوماً چنین ضریب غیر صفری می باید وجود داشته باشد. در غیر این صورت دستگاه خطی به دست آمده شامل  $(n-k+1)$  تا معادله

$$0 \times x_k + a_{i, k+1}^{(k-1)} x_{k+1} + \dots + a_{i, n}^{(k-1)} x_n = b_i^{(k-1)} \quad i = k, \dots, n \quad (21.4)$$

می شود که در واقع دارای  $n-k$  مجهول به صورت  $x_{k+1}, \dots, x_n$  است. لذا بنا بر قضیه ۳.۴، دستگاه معادلات (۲۱.۴) به ازای بعضی مقادیر سمت راست قابل حل نیستند؛ از این رو کل دستگاه به ازای بعضی از مقادیر سمت راست قابل حل نخواهد بود و بنا بر این به موجب قضیه ۴.۴ ماتریس ضرایب دستگاه معادلات فعلی وارونپذیر نیست. اما چون دستگاه معادلات فعلی با دستگاه معادلات اولیه هم ارز یعنی  $AX = b$  است، نتیجه می شود که ماتریس  $A$  وارونپذیر نیست و این همان اثبات ادعای ماست.

وقتی که این اعمال به کمک کامپیوتر انجام می شود،  $n$  معادله اصلی و تغییرات مختلفی که در آنها انجام می گیرند، باید به طریق مناسب و منظمی ثبت شوند. بدین منظور معمولاً از یک آرایه عملی  $1$  یا یک ماتریس  $(n+1) \times n$  که آن را  $W$  می نامند و در ابتدا شامل ضرایب و مقادیر سمت راست  $n$  معادله  $AX = b$  است استفاده می کنند. هرگاه مجهولی از یک معادله حذف شود، ضرایب تغییر یافته و سمت راست مربوط به این معادله محاسبه و در آرایه عملی  $W$  به جای ضرایب و مقادیر سمت راست قبلی ذخیره می شوند. به دلایلی که در زیر روشن خواهد شد، مضرب  $m_{ik} = a_{ik}^{(k-1)} / a_{kk}^{(k-1)}$  (که در حذف  $x_k$  از معادله  $i$ ام به کار می رود) به جای عدد  $a_{ik}^{(k)}$  در محل  $w_{ik}$  ذخیره می شود، زیرا به هر حال عدد  $a_{ik}^{(k)}$  صفر می شود (یا آنکه فرض آن است که صفر هست). تعویض سطرها نیز با استفاده از یک تعداد صحیح آرایه  $p$ ، ثبت می گردد.

## 1. working array



الگوریتم ۲۰۴: روش حذف گاوس. ماتریس  $W$  از مرتبه  $n(n+1)$  که  $n$  اولین ستون آن متضمن ماتریس مرتبه  $n \times n$   $A$  و ستون آخرش،  $n$ -برداری  $\mathbf{b}$  است، داده شده است.

$n$ -برداری  $\mathbf{p}$  را با  $p_i = i$ ، به ازای  $i = 1, \dots, n$ ، پر کنید

For  $k = 1, \dots, n-1$ , do:

کوچکترین مقدار  $i \geq k$  را چنان پیدا کنید که  $w_{ik} \neq 0$   
 اگر چنین مقدار  $i$  به دست نیامد، علامت بدهید که  $A$   
 وارونپذیر نیست و عملیات را متوقف سازید.  
 در غیر این صورت، محتوای  $p_i$  و  $p_k$  را با یکدیگر و  
 سطرها  $k$ ام و  $i$ ام را با هم در  $W$  عوض نموده، عملیات را  
 ادامه دهید

For  $i = k+1, \dots, n$ , do:

$$m := w_{ik} := w_{ik} / w_{kk}$$

For  $j = k+1, \dots, n+1$ , do:

$$w_{ij} := w_{ij} - mw_{kj}$$

اگر  $w_{nn} = 0$ ، علامت دهید که  $A$  وارونپذیر نیست و عملیات را متوقف سازید.

در غیر این صورت، دستگاه اولیه  $\mathbf{Ax} = \mathbf{b}$  هم ارز است با دستگاه  $\mathbf{Ux} = \mathbf{y}$  که در آن  $\mathbf{U}$  و  $\mathbf{y}$  بر حسب درایه‌های نهایی  $W$  به گونه زیر داده شده‌اند

$$u_{ij} = \begin{cases} w_{ij} & i \leq j \\ 0 & i > j \end{cases}, \quad y_i = w_{i, n+1}, \quad i = 1, \dots, n \quad (22.4)$$

به ویژه  $U$  یک ماتریس بالامثلثی است که کلیه درایه‌های قطری آن غیر صفرند. پس حالا الگوریتم ۱۰۴، را می‌توان برای محاسبه جواب  $\mathbf{x}$  به کار گرفت.

با توجه به جنبه‌های خاص، مانند تقارن یا تنگی ماتریس ضرایب  $A$ ، غالباً ممکن است عملیات ضروری محاسباتی برای حل  $\mathbf{Ax} = \mathbf{b}$  را کاهش داد. به عنوان مثال، اکنون حل دستگاه‌های سه قطری را مختصراً مورد بحث قرار می‌دهیم.

ماتریس مرتبه  $n \times n$   $A = (a_{ij})$  را سه قطری گوئیم اگر وقتی که  $|i-j| > 1$ ، آنگاه تساوی  $a_{ij} = 0$  برقرار باشد.

## 1. tridiagonal

به بیان لفظی،  $A$  زمانی سه قطری است، که تنها درایه‌های قطری آن یعنی  $a_{ii}$ ،  
 ( $i = 1, \dots, n$ ) یا درایه‌های زیر قطری آن یعنی  $a_{i, i-1}$ ، به ازای ( $i = 2, \dots, n$ ) و یا  
 درایه‌های زیر قطری آن یعنی  $a_{i, i+1}$ ، ( $i = 1, \dots, n-1$ ) غیر صفر باشند. بنابراین  
 ماتریسهای زیر سه قطری هستند

$$\begin{bmatrix} 3 & 1 & 0 & 0 \\ 1 & 3 & 1 & 0 \\ 0 & 2 & 2 & 1 \\ 0 & 0 & 1 & 6 \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \quad \begin{bmatrix} 4 & 7 & 0 & 0 \\ 0 & 8 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

فرض کنید که ماتریس ضرایب  $A$  از دستگاه خطی  $AX = b$  سه قطری باشد، بعلاوه  
 فرض می‌کنیم که به ازای هر  $k$  بتوان معادله  $k$ ام را به عنوان معادله لولایی در مرحله  $k$ ام  
 به کار گرفت. سپس طی مرحله  $k$ ام از الگوریتم ۲.۴، فقط متغیر  $x_k$  می‌باید از معادله  
 ( $k+1$ )ام، حذف شود. بعلاوه طی بسجایگذاری، تنها لازم است  
 که  $x_{k+1}$  در معادله  $k$ ام قرار داده شود تا بتوان  $x_k$  را،  $k = 1, \dots, n-1$ ، پیدا کرد.  
 سرانجام، احتیاجی به ذخیره کردن درایه‌های صفر ماتریس  $A$  نیست، بلکه فقط لازم است  
 سه بردار متضمن درایه‌های زیر قطری، قطری، و زیر قطری ماتریس  $A$  نگهداری شوند.  
 اکنون دستگاه سه قطری مرتبه  $n$ م زیر را به طور دقیقتری بررسی می‌کنیم:

$$d_1 x_1 + c_1 x_2 = b_1$$

$$a_2 x_1 + d_2 x_2 + c_2 x_3 = b_2$$

$$a_3 x_2 + d_3 x_3 + c_3 x_4 = b_3$$

$$\vdots$$

$$a_{n-1} x_{n-2} + d_{n-1} x_{n-1} + c_{n-1} x_n = b_{n-1}$$

$$a_n x_{n-1} + d_n x_n = b_n$$

با فرض  $d_1 \neq 0$ ،  $x_1$  را از معادله دوم حذف می‌کنیم، معادله دوم جدید به صورت زیر  
 خواهد شد

$$d'_2 x_2 + c_2 x_3 = b'_2$$

که در آن

$$d'_2 = d_2 - \frac{a_2}{d_1} c_1 \quad b'_2 = b_2 - \frac{a_2}{d_1} b_1$$

سپس فرض کنید  $d'_2 \neq 0$ ، معادله فوق را برای حذف  $x_1$  از معادله سوم به کار می‌بریم، معادله سوم جدید به صورت زیر خواهد شد

$$d'_2 x_2 + c_3 x_3 = b'_3$$

که در آن

$$d'_3 = d_3 - \frac{a_3}{d'_2} c_2 \quad b'_3 = b_3 - \frac{a_3}{d'_2} b'_2$$

به همین طریق عمل را ادامه می‌دهیم و طی مرحله  $k$ ام،  $x_k$  را از معادله  $k+1$ ام حذف می‌کنیم (با فرض اینکه  $d'_k \neq 0$ ) و معادله  $k$ ام جدید به صورت زیر خواهد شد

$$d'_{k+1} x_{k+1} + c_{k+2} x_{k+2} = b'_{k+1}$$

که در آن به ازای  $k=1, 2, \dots, n-1$  داریم

$$d'_{k+1} = d_{k+1} - \frac{a_{k+1}}{d'_k} c_k \quad b'_{k+1} = b_{k+1} - \frac{a_{k+1}}{d'_k} b'_k$$

هنگام پس‌جا یگذاری، با فرض اینکه  $d'_n \neq 0$ ، در ابتدا خواهیم داشت

$$x_n = \frac{b'_n}{d'_n}$$

و سپس به ازای  $k=n-1, \dots, 1$  داریم

$$x_k = \frac{b'_k - c_{k+1} x_{k+1}}{d'_k}$$

**الگوریتم ۳.۴ روش حذف برای دستگاه‌های سه قطری.** ضرایب  $a_i, d_i, c_i$  و قسمت سمت راست  $b_i$  از دستگاه سه قطری زیر داده شده‌اند

$$a_i x_{i-1} + d_i x_i + c_i x_{i+1} = b_i \quad i = 1, \dots, n \quad (a_1 = c_n = 0)$$

For  $k=2, \dots, n$ , do:

اگر  $d_{k-1} = 0$ ، ناکامی را اعلام کنید و عملیات را متوقف سازید.

در غیر این صورت، به ازای  $m := \frac{a_k}{d_{k-1}}$  عملیات را

ادامه دهید

$$\begin{cases} d_k := d_k - m * c_{k-1} \\ b_k := b_k - m * b_{k-1} \end{cases}$$

اگر  $d_n = 0$ ، ناکامی را اعلام کنید و عملیات را متوقف سازید

در غیر این صورت به ازای  $x_n := \frac{b_n}{d_n}$  عملیات را

ادامه دهید

For  $k = n-1, \dots, 1$ , do:

$$x_k := \frac{b_k - c_k * x_{k+1}}{d_k}$$

□ مثال ۳.۴: دستگاه خطی زیر را به ازای  $n = 10$  حل کنید

$$2x_1 - x_2 = 1$$

$$-x_{i-1} + 2x_i - x_{i+1} = 0 \quad i = 2, \dots, n-1$$

$$-x_{n-1} + 2x_n = 0$$

این مسئله با برنامه فورترن زیر حل می شود. توجه کنید که الگوریتم ۳.۴ به صورت زیر برنامه زیر ترجمه شده است

TRID(SUB, DIAG, SUP, B, N)

در این زیر برنامه، SUB، DIAG، SUP، و B بردارهای هستند که انتظار می رود ضرایب و قسمت سمت راست دستگاه سه قطری

$$\text{SUB}(i)x_{i-1} + \text{DIAG}(i)x_i + \text{SUP}(i)x_{i+1} = \text{B}(i) \quad i = 1, \dots, N$$

را داشته باشند [که در اینجا از SUB(1) و SUB(N) صرف نظر شده است]. این زیر برنامه محتوای DIAG را تغییر می دهد و بردار جواب را در B ذخیره کرده، برگشت می دهد. جواب دقیق دستگاه به صورت زیر داده شده است

$$x_i = \frac{n+1-i}{n+1} \quad i = 1, \dots, n$$

بنابراین جوابهای محاسبه شده در رقم ششم پس از ممیز، متضمن خطا هستند. این برنامه روی يك کامپیوتر IBM ۳۶۰ اجرا شده است.

```

C  FORTRAN PROGRAM FOR EXAMPLE 4.3
  PARAMETER N=10
  INTEGER I
  REAL A(N),B(N),C(N),D(N)
  DO 10 I=1,N
    A(I) = -1.
    D(I) = 2.
    C(I) = -1.
10  B(I) = 0.
    B(1) = 1.
  CALL TRID ( A, D, C, B, N )
  PRINT 610, (I,B(I),I=1,N)
610 FORMAT('THE SOLUTION IS '/(I5,E15.7))
      STOP
  END
  SUBROUTINE TRID ( SUB, DIAG, SUP, B, N )
  INTEGER N, I
  REAL B(N),DIAG(N),SUB(N),SUP(N)
C  THE TRIDIAGONAL LINEAR SYSTEM
C  SUB(I)*X(I-1) + DIAG(I)*X(I) + SUP(I)*X(I+1) = B(I), I=1,...,N
C  (WITH SUB(1) AND SUP(N) TAKEN TO BE ZERO) IS SOLVED BY FACTORIZATION
C  AND SUBSTITUTION. THE FACTORIZATION IS RETURNED IN SUB , DIAG , SUP
C  AND THE SOLUTION IS RETURNED IN B .
  IF (N .LE. 1) THEN
    B(1) = B(1)/DIAG(1)
      RETURN
  END IF
  DO 11 I=2,N
    SUB(I) = SUB(I)/DIAG(I-1)
    DIAG(I) = DIAG(I) - SUB(I)*SUP(I-1)
11  B(I) = B(I) - SUB(I)*B(I-1)
    B(N) = B(N)/DIAG(N)
  DO 12 I=N-1,1,-1
12  B(I) = (B(I) - SUP(I)*B(I+1))/DIAG(I)
      RETURN
  END

```

### خروجی

جواب برابر است با

1	0.9090915E 00
2	0.8181832E 00
3	0.7272751E 00
4	0.6363666E 00
5	0.5454577E 00
6	0.4545485E 00
7	0.3636391E 00
8	0.2727295E 00
9	0.1818197E 00
10	0.9090990E -01

□

### تمرین

۱-۲۰۴ يك ملاك كارایی يك الگوریتم، تعداد عملیات حسابی لازم برای به دست آوردن جواب است. نشان دهید که اگر الگوریتم ۲۰۴ برای يك دستگاه از مرتبه  $n$ ام به کار رود،  $n(n-1)/2$  عمل تقسیم،  $(n^3-n)/3$  عمل ضرب، و  $(n^3-n)/3$  عمل جمع لازم دارد.

۲-۲۰۴ نشان دهید که الگوریتم بسجایگذاری ۱۰۴، مستلزم  $n$  عمل تقسیم،  $n(n-1)/2$  عمل ضرب، و  $n(n-1)/2$  عمل جمع است.

۳-۲۰۴ در برخی از کامپیوترها عمل تقسیم وقتگیرتر از عمل ضرب است، چگونه ممکن

است الگوریتم ۲.۴، برای چنین کامپیوتری اصلاح شود؟

۲.۴-۴ تعداد عملیات جمع و ضرب لازم برای ضرب یک ماتریس  $n \times n$  در یک  $n$ -بردار را محاسبه کنید.

۲.۴-۵ چند عمل جمع، ضرب و تقسیم برای الگوریتم ۲.۴ لازم است، اگر فقط ماتریس نهایی بالامثلی  $U$  مورد نیاز باشد؟

۲.۴-۶ با استفاده از روش حذف نشان دهید که دستگاه زیر جواب ندارد

$$x_1 + 2x_2 + x_3 = 3$$

$$2x_1 + 3x_2 + x_3 = 5$$

$$3x_1 + 5x_2 + 2x_3 = 1$$

۲.۴-۷ زمان اجرای یک برنامه متضمن الگوریتم ۲.۴، تا حد زیادی با زمان مصروفه در داخلی ترین حلقه معین می شود. بدین دلیل لازم است که تا حد ممکن حلقه مذکور کارا شود. در عین حال در زبان فورترن، ذخیره کردن آرایه ها در حافظه به صورت ستونی انجام می گیرد و در بسیاری از ماشینها کار کردن با یک آرایه به صورت ستون به ستون سریعتر از سطر به سطر انجام می پذیرد.

بدین دلیل الگوریتم ۲.۴ را چنان بازسازی کنید که داخلی ترین حلقه (ها) بر مبنای نمایه سطرها اجرا شود (شوند)، یعنی بازسازی به گونه ای باشد که در هر زمان به عوض سطر، ستون تغییر کند.

۲.۴-۸ دستگاه زیر را با روش حذف حل کنید. تمام محاسبات را تا سه رقم اعشاری گرد کنید.

$$0.21x_1 + 0.32x_2 + 0.12x_3 + 0.31x_4 = 0.96$$

$$0.10x_1 + 0.15x_2 + 0.24x_3 + 0.22x_4 = 0.71$$

$$0.20x_1 + 0.24x_2 + 0.46x_3 + 0.36x_4 = 1.26$$

$$0.61x_1 + 0.40x_2 + 0.32x_3 + 0.20x_4 = 1.53$$

جوابهای خود را با پسجایگذاری در دستگاه اولیه امتحان کنید و میزان دقت عمل را برآورد نمایید. جواب دقیق:  $[1, 1, 1, 1]$  است.

۲.۴-۹ با استفاده از زیر برنامه TRID و به ازای  $n = 30$  و  $h = 0.1$  دستگاه خطی زیر را حل کنید

$$-2(1+h^2)x_1 + x_2 = 1$$

$$x_{i-1} - 2(1+h^2)x_i + x_{i+1} = 0 \quad i = 2, 3, \dots, n-1$$

$$x_{n-1} - 2(1+h^2)x_n = 1$$

۲۰۴-۱۰ با استفاده از قضیه ۶.۴ و فرع لم ۱.۲، ثابت کنید که هر بسجمله‌ای از درجه  $n$  تا بزرگتر از  $n$  به ازای مراکز داده شده  $c_1, \dots, c_n$ ، تنها به یک روش می‌تواند به شکل نیوتنی نوشته شود.

(دانهمایی: دستگاه خطی مربوط به ضرایب آن شکل نیوتنی از بسجمله‌ای را که با تابع داده شده در مراکز  $c_1, \dots, c_n, c_{n+1}$  یکسان است در نظر بگیرد.)

۲۰۴-۱۱ قضیه ۷.۴ را ثابت کنید. (دانهمایی: نخست ثابت کنید که هر جواب معادله  $AX = b$ ، یک جواب  $\tilde{A}X = \tilde{b}$  نیز هست. سپس نشان دهید که هر عمل از نوع مذکور در قضیه، می‌تواند به توسط عملی با همان نوع خنثی شود، لذا نشان دهید که با یک سلسله از عملیات یاد شده،  $AX = b$  می‌تواند از  $\tilde{A}X = \tilde{b}$  به دست آید.)

### ۳.۴ تدبیر لولاگزینی

در الگوریتم حذف ۲.۴، که در بخش قبل ارائه شد، اگر تمامی محاسبات با دقت بسیار زیاد انجام شود هر دستگاه  $AX = b$  با کارایی و با اطمینان حل می‌شود. در صورتی که، طبق معمول، محاسبات با دقت محدود انجام پذیرند، به آسانی می‌توان مثالهایی ارائه داد که برای آنها الگوریتم ۲.۴، جرابهایی کاملاً غلط به دست می‌دهد.

در این بخش فقط یک منشأ ممکن برای چنین خطایی به اختصار ارائه می‌شود و آن تدبیر نادرست لولاگزینی است. در اینجا، منظور ما از تدبیر لولاگزینی<sup>۲</sup> طرح کلی است برای گزینش معادله لولایی (و شاید هم ستون لولایی) در هر مرحله از حذف.

□ مثال ۴.۴: دستگاه

$$0.00003x_1 + 1.0566x_2 = 1.0569$$

$$0.3454x_1 - 2.436x_2 = 1.018$$

دارای جوابهای  $x_1 = 1$  و  $x_2 = 1$  است. برای حل این دستگاه از روش حذف و محاسبات ممیز شناور با ۴ رقم اعشاری استفاده می‌کنیم و اولین معادله را به عنوان معادله لولایی برای مرحله اول (و تنها این مرحله) می‌گزینیم. داریم

$$m_{21} = 0.3454 / 0.00003 = 11510$$

بنابراین

$$\begin{aligned} a_1^{(1)} &= -22436 - (11510)(12566) \\ &= -22436 - 180200 = -180436 \end{aligned}$$

$$\begin{aligned} b_1^{(1)} &= 10018 - (11510)(12569) = 10018 - 180600 \\ &= -180582 \end{aligned}$$

این مقادیر به ما می‌دهند

$$x_2 = -180582 / -180436 = 1001$$

بنابراین از اولین معادله نتیجه می‌شود که

$$\square \quad x_1 = [12569 - (12566)(1001)] / 00003 = 3233$$

يك توضیح «موجه‌نما»<sup>۱</sup> برای این نقص چنین است: دربارهٔ لولای  $00003 = a_{11}$  «خیلی کوچک» است و از آنجا که اگر  $a_{11}$  صفر می‌بود محاسبات با شکست مواجه می‌شد، لذا شکست آور نیست که در محیطی با دقت عمل محدود، این الگوریتم «در نزدیکی صفر» به طریق نامطلوبی عمل کند.

البته در توضیح فوق اصطلاحات تعریف نشده‌ای مانند «خیلی کوچک» و «نزدیک صفر» به کار برده شده‌اند که آن را از درجهٔ اعتبار می‌اندازند. در حقیقت اگر طرفین معادلهٔ اول را در توان مناسبی از ده ضرب کنیم، می‌توانیم  $a_{11}$  را هر قدر که بخواهیم بزرگ کنیم بی‌آنکه تغییر در جواب محاسبه شده ایجاد شود. برای مشاهدهٔ این امر، دوباره دستگاه مثال ۴.۴ را در نظر می‌گیریم که در آن معادلهٔ اول در  $10^m$ ، که  $m$  عددی صحیح است، ضرب شده است

$$00003 \times 10^m x_1 + 12566 \times 10^m x_2 = 12569 \times 10^m$$

$$003454 x_1 - 22436 x_2 = 10018$$

دوباره از معادلهٔ اول به عنوان معادلهٔ لولایی استفاده می‌کنیم و محاسبات ممیز شناور با رقم اعشاری را به کار می‌گیریم، داریم

$$m_{x_1} = \frac{003454}{00003 \times 10^m} = 11510 \times 10^{-m}$$

بنابراین



$$a_{۲۲}^{(۱)} = -۲۰۴۳۶ - (۱۱۵۱۰۰ \times ۱۰^{-۳})(۱۰۵۶۶ \times ۱۰^۳) = -۱۸۰۴۰۰$$

$$b_{۲}^{(۱)} = ۱۰۰۱۸ - (۱۱۵۱۰۰ \times ۱۰^{-۳})(۱۰۵۶۹ \times ۱۰^۳) = -۱۸۰۵۰۰$$

که همان نتایج قبلی است. بنابراین دوباره خواهیم داشت  $x_۲ = ۱۰۰۰۱$  و بالاخره

$$x_۱ = (۰۰۰۰۱ \times ۱۰^۳) / (۰۰۰۰۰۳ \times ۱۰^۳) = ۳۰۳۳۳$$

در واقع عیب این مثال این است که  $|a_{۱۱}|$  در مقایسه با  $|a_{۲۲}|$  کوچک است. لذا خطای نسبتاً کوچک ناشی از گرد کردن در محاسبهٔ  $x_۲$  منجر به اختلاف بزرگی بین جواب محاسبه شدهٔ  $x_۱$  و جواب صحیح  $x_۱$  می‌شود. این امر زمانی تأیید می‌شود که معادلهٔ دوم را به عنوان معادلهٔ لولایی انتخاب کنیم، در این صورت داریم  $۶ \approx |a_{۲۲}/a_{۲۱}|$  که در مقایسه با  $۵۲۲۰ \approx |a_{۲۲}/a_{۱۱}|$  کوچک است. از آنجا

$$m_{۱۱} = \frac{۰۰۰۰۰۳}{۰۰۳۴۵۴} = ۰۰۰۰۰۸۶۸۶$$

و معادلهٔ جدید اولی به صورت زیر درمی‌آید

$$۱۰۵۶۸x_۲ = ۱۰۵۶۸$$

به طوری که  $x_۲ = ۱$  همان جواب صحیح است و سرانجام از معادلهٔ دوم به دست می‌آوریم  $x_۱ = ۱۰$ . اما (در این حالت) حتی اگر خطای گرد کردن موجب شود که داشته باشیم  $x_۲ = ۱۰۰۰۱$  (همان طوری که در مثال ۴.۴ پیش آمد) باز معادلهٔ دوم خواهد داد

$$x_۱ = \frac{۱۰۰۱۸ + ۲۰۴۳۸}{۰۰۳۴۵۴} = ۱۰۰۰۱$$

که نتیجهٔ خوبی است.

تعیین چگونگی تأثیر تدبیرهای مختلف لولایی بر دقت عمل جواب محاسبه شده، بسیار مشکلتر (اگر ناممکن نباشد) است. یک استثنای قابل توجه و حائز اهمیت، دستگاههای خطی با ماتریس ضرایب معین و مثبت، یعنی دستگاههایی هستند که برای ماتریس ضرایب آنها شرایط زیر برقرارند.

$$x^T A x > ۰, \quad x \neq ۰ \quad \text{و به ازای کلیهٔ مقادیر } A = A^T$$

برای چنین دستگاهی، می‌توان نشان داد که خطای جواب محاسبه شده ناشی از خطای گرد کردن طی حذف و پس‌جا‌یگذاری به‌طور قابل قبولی کوچک است [ص ۱۲۷؛ ۴۱]. به شرط آنکه تدبیر پیش با افتادهٔ گزینش لولایی به صورت عوض نکردن جای معادله‌ها اختیار شود. (تمرین ۴.۰۹- الگوریتم کارآمدی برای این حالت است). اما در حال حاضر امکان ارائهٔ «بهترین» تدبیر لولایی برای یک دستگاه خطی کلی امکان‌پذیر نیست، حتی به‌خوبی روشن نیست که عبارت «بهترین» تدبیر لولایی چه معنایی می‌تواند داشته باشد.

به دلیل اقتصادی باید انتخاب معادلهٔ لولایی برای هر مرحله و در آغاز آن مرحله بر اساس وضعیت جاری دستگاه مورد نظر یعنی بی اطلاع قبلی از تأثیر این گزینش بر مراحل بعدی صورت گیرد.

یک تدبیر که در حال حاضر پذیرفته شده، لولایزینی جزئی مدرج است. در این راه کار، ابتدا «اندازه»ی  $d_i$  مربوط به سطر  $i$ ام از ماتریس  $A$  را به ازای  $i = 1, \dots, n$  محاسبه می کنند. اندازه مناسب برای این اندازه عدد زیر است (به بخش ۵.۴ نگاه کنید).

$$d_i = \max_{1 \leq j \leq n} |a_{ij}|$$

پس در آغاز مرحلهٔ کلی یا مرحلهٔ  $k$ ام از الگوریتم حذف ۲.۴، از بین  $(n-k)$  معادله موجود، آن معادله‌ای به عنوان معادلهٔ لولایی انتخاب می شود که مطلقاً بزرگترین ضریب  $p_k$  را نسبت به اندازه معادله داشته باشد. به موجب الگوریتم ۲.۴ این بدان معنی است که عدد صحیح  $j$  بین  $k$  و  $n$  (معمولاً کوچکترین) چنان انتخاب شده است که

$$\frac{|w_{jk}|}{d_j} \geq \frac{|w_{ik}|}{d_i} \quad i = k, \dots, n$$

واضح است که لولایزینی جزئی مدرج در انتخاب تدبیر درست لولایزینی برای دستگاه مثال ۴.۴ دخالت می کند و با دوباره مدرج کردن معادلات، کنار گذاشته نمی شود. این امکان وجود دارد که الگوریتم ۲.۴ چنان اصلاح شود که نه فقط معادلات لولایی، بلکه مجهولاتی هم که باید حذف شوند در معرض انتخاب قرار گیرند. در این تغییر دو جایگشت  $p$  و  $q$  انتخاب می شوند که معادلهٔ  $p_k$ ام را به عنوان معادله‌ای که باید طی مرحلهٔ  $k$ ام حذف  $x_{q_k}$ ،  $k = 1, \dots, n-1$ ، به کار رود تعیین می کند. در لولایزینی کامل  $2$  انتخاب معادلهٔ لولایی و مجهولی که می باید حذف شود با توجه به ضریب مطلقاً بزرگتر هر یک از  $n-k$  مجهول در هر یک از  $n-k$  معادلهٔ مورد نظر انجام می گیرد. بدیهی است که این تدبیر بسیار گرانتر از لولایزینی جزئی مدرج تمام می شود، و بنا بر این چندان به کار گرفته نمی شود، ولو اینکه به طور مسلم این تدبیر برتر از لولایزینی جزئی می باشد.

## تمرین

۳-۳۴ الگوریتم ۲.۴ را به گونه‌ای اصلاح کنید که در آن لولایزینی کامل انجام پذیرد.

۳-۳۴ مثالی از یک دستگاه خطی  $2 \times 2$  ارائه دهید که در آن لولایزینی کامل نتیجه‌ای دقیقتر از لولایزینی جزئی مدرج بدهد. محاسبات در ممیز شناور را با ۴ رقم اعشاری در نظر بگیرید. (داهنمایی:  $a_{11}$  و  $a_{21}$  را در مقایسه با  $a_{12}$  و  $a_{22}$  «کوچک بگیرد».)

۳-۳.۴ دستگاه معادلات خطی زیر را با به کارگیری محاسبات ممیز شناور و چهار رقم اعشاری حل کنید. برای حل آن یکبار معادله اول و یکبار معادله دوم را به عنوان معادله لولایی در نظر بگیرید و نهایتاً حل دستگاه فوق را با روش لولاگزینی کامل نیز انجام دهید و جوابها را در هر سه مورد با جواب دقیق  $x_1 = 1000$  و  $x_2 = 0.2500$  مقایسه کنید.

$$0.1410 \times 10^{-2} x_1 + 0.4004 \times 10^{-1} x_2 = 0.1142 \times 10^{-1}$$

$$0.2000 \times 10^0 x_1 + 0.4912 \times 10^1 x_2 = 0.1428 \times 10^1$$

۴-۳.۴ دستگاه خطی تمرین ۲.۴-۸ را با استفاده از لولاگزینی جزئی مدرج حل کنید و جوابهای خود را با جوابهای تمرین ۲.۴-۸ مقایسه کنید.

### ۴.۴ تجزیه به عوامل مثلثی

روند حذف در الگوریتم ۲.۴ را می توان به عنوان يك تجزیه ماتریس ضرایب  $A$  به سه عامل

$$A = PLU$$

تلقی کرد که در آن  $p$  ماتریس جایگشتی است معرف تعویض سطرها، و  $L$  يك ماتریس واحد پایین مثلثی است که شامل مضارب (در جا لبرترین قسمت آن) به کار برده شده است و بالاخره  $U$  يك ماتریس بالا مثلثی است. در حالتی که ماتریس  $A$  يك ماتریس معین مثبت و متقارن باشد، این دید نسبت به الگوریتم ۴.۴ به يك الگوریتم کارایی (تجزیه به عوامل چولسکی، تمرین ۴.۴-۹ را ببینید) منجر خواهد شد. همچنین لازم است به طرح با ارزش به اصطلاح طرحهای فشرده<sup>۲</sup> (منسوب به دولیتل و کروت<sup>۳</sup>، تمرین ۴.۴-۸ را ببینید) توجه نمود که مزیت آنها در حل دستگاههای خطی توسط ماشینهای حساب رومیزی (یا جیبی) است که تعداد نتایج محاسبات میانی را که می باید ثبت گردند، کاهش می دهد. روشهای مذکور را می توان برای کاهش اثرات ناشی از خطای گرد کردن، به هنگام جمع حاصلزریهای عددی (در بعضی ماشینها) در روش جمع با دقت مضاعف<sup>۴</sup> به کار برد. سرانجام آنکه، توجه به روش حذف از دیدگاه تجزیه به عوامل، امکان تحلیل خطای پسرود را در روش حذف (همان گونه که در بخش ۶.۴ انجام خواهد گرفت) آسان می سازد.

اولاً فرض می کنیم که در طی اجرای این الگوریتم هیچ گونه تعویض سطری اتفاق نیفتد و آنچه را که برای معادله<sup>۵</sup>  $k$  پیش خواهد آمد مورد بررسی قرار می دهیم. به ازای

$k = 1, 2, \dots, i-1$  در طی مرحله<sup>۶</sup>  $k$ ام، معادله از صورت

$$a_{ik}^{(k-1)} x_k + a_{i, k+1}^{(k-1)} x_{k+1} + \dots + a_{in}^{(k-1)} x_n = b_i^{(k-1)}$$

1. Choleski
2. compact schemes
3. Doolittle and Crout
4. double-precision
5. backward error analysis

به صورت

$$a_{i, k+1}^{(k)} x_{k+1} + \dots + a_{in}^{(k)} x_n = b_i^{(k)}$$

درخواهد آمد که در آن

$$a_{ij}^{(k)} = a_{ij}^{(k-1)} - m_{ik} a_{kj}^{(k-1)}$$

$$b_i^{(k)} = b_i^{(k-1)} - m_{ik} b_k^{(k-1)}$$

و ضریب

$$m_{ik} = a_{ik}^{(k-1)} / a_{kk}^{(k-1)}$$

در درایه  $(i, k)$  ام آرایه ای که کار می کنیم ذخیره می گردد. در اینجا  $a_{kj}^{(k-1)}$  و  $b_k^{(k-1)}$  در مرحله  $k$  ام، به ترتیب ضرایب و سمت راست معادله لولایی هستند و از این روش شکل نهایی خود را دارند. این بدان معنی است که بر حسب خروجی از الگوریتم ۲۰۴ یعنی بر حسب ماتریس بالامثلثی  $U$  و بردار  $y$  حاصله در آن الگوریتم، داریم

$$a_{kj}^{(k-1)} = u_{kj}, \quad j = k, \dots, n \quad \text{و} \quad b_k^{(k-1)} = y_k$$

در نتیجه

$$\left. \begin{aligned} m_{ik} &= a_{ik}^{(k-1)} / u_{kk} \\ a_{ij}^{(k)} &= a_{ij}^{(k-1)} - m_{ik} u_{kj}, \quad j = k, \dots, n \\ b_i^{(k)} &= b_i^{(k-1)} - m_{ik} y_k \end{aligned} \right\} k = 1, \dots, i-1$$

چون  $a_{ij}^{(0)} = a_{ij}$  و  $b_i^{(0)} = b_i$  از آنجا نتیجه می شود که

$$m_{ik} = (a_{ik} - m_{i1} u_{1k} - m_{i2} u_{2k} - \dots - m_{i, k-1} u_{k-1, k}) / u_{kk} \quad k = 1, \dots, i-1$$

$$u_{ij} = a_{ij} - m_{i1} u_{1j} - m_{i2} u_{2j} - \dots - m_{i, i-1} u_{i-1, j} \quad j = i, \dots, n \quad (23.4)$$

و

$$y_i = b_i - m_{i1} y_1 - m_{i2} y_2 - \dots - m_{i, i-1} y_{i-1} \quad (24.4)$$

اکنون دوباره این معادلات را چنان می نویسیم که داده های اولیه، یعنی  $A$  و  $b$ ، در سمت راست ظاهر شوند. در این صورت داریم

$$m_{i1} u_{1k} + \dots + m_{ik} u_{kk} = a_{ik} \quad k = 1, \dots, i-1$$

$$m_{i1} u_{1j} + \dots + m_{i, i-1} u_{i-1, j} + u_{ij} = a_{ij} \quad j = i, \dots, n \quad (25.4)$$

و

$$m_{i,1}y_1 + \dots + m_{i,i-1}y_{i-1} + y_i = b_i$$

بنابراین، اگر  $L = (l_{ij})$  را ماتریس واحد پایین مثلثی بگیریم که بر حسب مقادیر نهایی آرایه عملی  $W$  به صورت

$$l_{ij} = \begin{cases} w_{ij} = m_{ij}, & i > j \\ 1, & i = j \\ 0, & i < j \end{cases} \quad (۲۶.۴)$$

داده شده اند، آنگاه می توانیم این معادلات را (به ازای  $i = 1, \dots, n$ ) به شکل ماتریسی ساده زیر بنویسیم

$$LU = A$$

و

$$Ly = b$$

این رابطه تجزیه به عوامل مثلثی را، در صورتی که هیچ تعویضی صورت نگرفته باشد، اثبات می کند. از سوی دیگر، اگر تعویضهایی صورت گرفته بود و ما اول، این تعویضها را انجام داده سپس الگوریتم ۲۰۴ را بدون هیچ گونه تعویضی به کار برده بودیم، در آن صورت محتوای نهایی  $W$  تغییر می کرد. علت آن این است که تمامی عملیات موجود در الگوریتم متضمن کاستن مضربی از یک سطر از سایر سطرهاست تا در سطرهاى دیگر صفری پدید آید و بدین جهت ترتیب نوشتن سطرها امر مهمی به شمار نمی آید. تنها چیزی که می باید به آن توجه داشت این است که وقتی سطری به عنوان یک سطر لولایی به کار می رود، دیگر نباید تغییر کند و به همین دلیل سطرهایی که به عنوان سطرهای لولایی به کار نرفته اند می باید به صورت جداگانه نگهداری شوند.

در نتیجه اگر تعویض سطرها در طی اجرای الگوریتم ۲۰۴ صورت گیرد، آنگاه ماتریسهای  $L$  و  $U$  که به توسط الگوریتم حاصل شده اند، به ازای ماتریس جایگشتی مناسبی مانند  $P$ ، در رابطه زیر صدق می کنند

$$LU = P^{-1}A$$

آنگاه

$$PLU = A \quad (۲۷.۴)$$

و همچنین

$$Ly = P^{-1}b \quad (۲۸.۴)$$

بر حسب بردار  $p$  که در الگوریتم ۲۰۴ برای ثبت تعویض سطرها به کار برده شده،

معادله  $p_k$ ام در طی مرحله  $k$ ام به عنوان معادله لولایی به کار رفته است. از این رو  $p^{-1}$  به ازای جمیع مقادیر  $k$  می باید سطر  $p_k$ ام را به سطر  $k$ ام بدل کند. این امر بدان معنی است که، اگر واقعاً بخواهیم بدانیم، به ازای جمیع مقادیر  $k$  تساوی،  $P\mathbf{i}_k = \mathbf{i}_{p_k}$  برقرار است [بخش (iii) قضیه ۵.۴ را ببینید]. ولی آنچه که برای ما اهمیت دارد، بر حسب خروجی  $\mathbf{p}$  در الگوریتم ۲.۴، عبارت زیر است

$$P^{-1}\mathbf{b} = (b_{p_k})_{k=1}^n$$

به عنوان اولین نمونه از کاربرد دیدگاه تجزیه به عوامل، اکنون در مورد امکان خرد کردن روند حل معادله  $A\mathbf{x} = \mathbf{b}$  به دو مرحله، بررسی به عمل می آوریم، دو مرحله مذکور عبارتند از مرحله تجزیه به عوامل<sup>۱</sup> که در آن عوامل مثلثی  $L$  و  $U$  (و احتمالاً ترتیب متفاوتی از سطرهای  $\mathbf{p}$ ) حاصل می شوند و دیگر مرحله یافتن جواب<sup>۲</sup>، که در طی آن، ابتدا دستگاه مثلثی

$$L\mathbf{y} = P^{-1}\mathbf{b} = (b_{p_k})_{k=1}^n \quad (29.4)$$

نسبت به  $\mathbf{y}$  حل می شود و سپس دستگاه مثلثی  $U\mathbf{x} = \mathbf{y}$  نسبت به  $\mathbf{x}$  با روش پسجایگذاری حل می شود. باید توجه داشت که سمت راست  $\mathbf{b}$ ، تنها در مرحله دوم وارد می شود. از این رو اگر دستگاه معادلات قرار باشد نسبت به مقادیر سمت راست دیگری نیز حل شود. تنها مرحله دوم است که تکرار می شود.

به موجب رابطه (۲۴.۴)، دستگاه مثلثی (۲۹.۴) در الگوریتم ۲.۴ طی مراحل زیر

حل می شود

For  $k = 1, \dots, n$ , do:

$$y_k := b_{p_k} - l_{k1}y_1 - \dots - l_{k,k-1}y_{k-1}$$

در واقع این امر شبیه به الگوریتم پسجایگذاری (۱.۴) است که برای حل  $U\mathbf{x} = \mathbf{y}$  نسبت به  $\mathbf{x}$  به کار می رود با این تفاوت که بررسی معادله‌ها، به علت پایین-مثلثی بودن ماتریس  $L$ ، از اولین معادله تا آخرین معادله صورت می گیرد، تمامی مراحل حل مسئله به گونه زیر ثبت می گردد:

### الگوریتم ۴.۴: پسجایگذاری و پسجایگذاری<sup>۳</sup>

مندرجات نهایی  $n$  ستون اول آرایه مورد عمل  $W$  و  $n$ -برداری  $\mathbf{p}$  از الگوریتم ۲.۴ (که برای دستگاه  $A\mathbf{x} = \mathbf{b}$  به کار گرفته شده) و نیز سمت راست معادله، یعنی  $\mathbf{b}$ ، داده شده‌اند.

1. factorization phase
2. solving phase
3. Forward-and back-substitution

For  $k = 1, \dots, n$ , do:

$$y_k := b_{p_k} - \sum_{j=1}^{k-1} w_{kj} y_j$$

For  $k = n, n-1, \dots, 1$ , do:

$$x_k := \left( y_k - \sum_{j=k+1}^n w_{kj} x_j \right) / w_{kk}$$

اکنون بردار  $\mathbf{x} = (x_i)$  جواب  $A\mathbf{x} = \mathbf{b}$  را در بردارد.

باز توجه کنید که هر دو حاصلجمع گاهی اوقات تهی هستند.

اهمیت عملی بحث قبلی زمانی آشکار می‌شود که تعداد عملیات (با ممیز شناور) را در الگوریتم‌های ۲.۴ و ۴.۴ بشماریم و به موجب تمرین ۲.۲-۲.۴ عمل تقسیم  $n(n-1)/2$  عمل ضرب، و  $n(n-1)/2$  عمل جمع برای انجام حلقه دوم الگوریتم ۴.۴ لازم است. حلقه اول نیز به همین تعداد عملیات نیاز دارد. جز اینکه به هیچ عمل تقسیم احتیاجی ندارد. از این رو الگوریتم ۴.۴ متضمن

$(n-1)n$  عمل جمع و  $n^2$  عمل ضرب یا تقسیم می‌باشد.

با در نظر گرفتن تمرین ۲.۲-۴ این تعداد دقیقاً تعداد عملیاتی است که برای ضرب یک ماتریس  $n \times n$  در یک  $n$ -برداری لازم است.

برعکس، برای محاسبه محتوای نهایی اولین  $n$  ستون ماتریس مورد عمل  $W$  به وسیله الگوریتم ۲.۴،

$n/6 - n^2/2 + n^3/3$  عمل جمع و  $(n^3 - n)/3$  عمل ضرب یا تقسیم لازم است (به تمرین ۲.۴-۵ نگاه کنید). بنابراین عمده کار لازم در حل  $A\mathbf{x} = \mathbf{b}$  از راه روش حذف، به دست آوردن محتوای نهایی ماتریس مورد عمل  $W$ ، یعنی،  $\Theta[(1/3)n^3]$  عمل جمع و به همین تعداد عمل ضرب یا تقسیم به ازای  $n$  های بزرگ است. پیش‌جا‌گذاری و پس‌جا‌گذاری بعدی به همین تعداد عملیات بایک‌درجه کمتر، یعنی  $\Theta(n^2)$  عمل جمع و همین تعداد عمل ضرب به ازای هر مقدار سمت راست معادله لازم دارد. بنابراین، وقتی محتوای نهایی  $W$  را دانستیم می‌توانیم در مدت زمانی که برای محاسبه محتوای نهایی  $W$  لازم است، دستگاه  $A\mathbf{x} = \mathbf{b}$  را برای مقادیر بسیار زیاد سمت راست حل کنیم.

در محاسبه این کارها، روش سنتی را دنبال کرده‌ایم و تنها عملیات با ممیز شناور را به حساب آورده‌ایم به‌ویژه، محاسبه زیرنما‌ی‌ها، تعیین هزینه برای اجرای حلقه‌های DO و سایر هزینه‌های مربوط به ثبت و ضبط یا ساماندهی<sup>۱</sup> را نادیده گرفته‌ایم، این امر به جهت آنست که محاسبات اخیر بسیار سریعتر از عملیات با ممیز شناور انجام می‌گیرند. این نوع حسابرسی روی کامپیوترهای امروزی موردی ندارد و تصویردرستی از مقدار هزینه به دست

نمی‌دهد (تمرین ۲۰۴-۷ را ببینید). از سوی دیگر این مطلب کسه تا چه حد انجام يك کار (از لحاظ محاسبهٔ زمان موردنیاز) به‌ساماندهی در يك برنامه بستگی دارد، در کامپیوترهای مختلف به‌شدت فرق می‌کند و بنابراین ارائهٔ يك بحث کلی در این موضوع در این کتاب دشوار است.

يك زیر برنامهٔ فورترن به‌نام SUBST که متضمن الگوریتم جایگذاری ۴.۴ است در زیر داده شده‌است.

```

SUBROUTINE SUBST( W, IPIVOT, B, N, X )
INTEGER IPIVOT(N), I, IP, J
REAL B(N), W(N,N), X(N), SUM
C***** I N P U T *****
C W, IPIVOT, N ARE AS ON OUTPUT FROM F A C T O R , APPLIED TO THE
C MATRIX A OF ORDER N.
C B IS AN N-VECTOR, GIVING THE RIGHT SIDE OF THE SYSTEM TO BE SOLVED.
C***** O U T P U T *****
C X IS THE N-VECTOR SATISFYING A*X = B.
C***** M E T H O D *****
C ALGORITHM 4.4 IS USED, I.E., THE FACTORIZATION OF A CONTAINED IN
C W AND IPIVOT (AS GENERATED IN FACTOR ) IS USED TO SOLVE A*X = B
C FOR X BY SOLVING TWO TRIANGULAR SYSTEMS.
C
IF (N .LE. 1) THEN
X(1) = B(1)/W(1,1)
RETURN
END IF
IP = IPIVOT(1)
X(1) = B(IP)
DO 15 I=2,N
SUM = 0.
DO 14 J=1,I-1
14 SUM = W(I,J)*X(J) + SUM
IP = IPIVOT(I)
15 X(I) = B(IP) - SUM
C
X(N) = X(N)/W(N,N)
DO 20 I=N-1,1,-1
SUM = 0.
DO 19 J=I+1,N
19 SUM = W(I,J)*X(J) + SUM
20 X(I) = (X(I) - SUM)/W(I,I)
RETURN
END

```

باردیگر، يك زیر برنامهٔ فورترن به‌نام FACTOR می‌دهیم که از الگوریتم حذف ۲.۴ و تدبیرگزینش جزئی لولایی مسدوح به‌کار گرفته شده است و این به‌کارگیری برای محاسبهٔ تجزیه به‌عوامل مثلثی (در صورت امکان) يك ماتریس مفروض  $A$  از مرتبهٔ  $N \times N$ ، ذخیره نمودن عوامل تجزیهٔ موجود در ماتریس  $W$  از مرتبهٔ  $N \times N$ ، و تدبیرگزینش لولایی در يك  $N$ -برداری IPIVOT است که برای استفاده از زیر برنامهٔ SUBST که قبلاً داده شده آماده‌است. استفاده‌کننده باید يك  $N$ -برداری  $D$  به‌عنوان فضای کارا برای ذخیره نمودن «اندازه» سطرهای  $A$ ، نیز فراهم کند. اگر به‌ماتریس  $A$  نیازی نباشد و جا برای ذخیره کم باشد، خود  $A$  را می‌توانیم به‌جای  $W$  در زمرهٔ شناسه‌های دستورالعمل CALL (که این عمل در برخی از گویشهای زبان فورترن غیر قانونی است) به‌کار ببریم. در این صورت عوامل تجزیهٔ به‌دست آمده در آرایهٔ  $A$  ضبط و جایگزین ماتریس اولیه می‌شوند.

## 1. working space



```

SUBROUTINE FACTOR ( W, N, D, IPIVOT, IFLAG )
INTEGER IFLAG, IPIVOT(N), I, ISTAR, J, K
REAL D(N), W(N,N), AWIKOD, COLMAX, RATIO, ROWMAX, TEMP
C***** I N P U T *****
C W ARRAY OF SIZE (N,N) CONTAINING THE MATRIX A OF ORDER N TO BE
C FACTORED.
C N THE ORDER OF THE MATRIX
C***** W O R K A R E A *****
C D A REAL VECTOR OF LENGTH N, TO HOLD ROW SIZES
C***** O U T P U T *****
C W ARRAY OF SIZE (N,N) CONTAINING THE LU FACTORIZATION OF P*A FOR
C SOME PERMUTATION MATRIX P SPECIFIED BY IPIVOT .
C IPIVOT INTEGER VECTOR OF LENGTH N INDICATING THAT ROW IPIVOT(K)
C WAS USED TO ELIMINATE X(K) , K=1,....,N .
C IFLAG AN INTEGER,
C = 1, IF AN EVEN NUMBER OF INTERCHANGES WAS CARRIED OUT,
C = -1, IF AN ODD NUMBER OF INTERCHANGES WAS CARRIED OUT,
C = 0, IF THE UPPER TRIANGULAR FACTOR HAS ONE OR MORE ZERO DIA-
C GONAL ENTRIES.
C THUS, DETERMINANT(A) = IFLAG*W(1,1)*...*W(N,N) .
C IF IFLAG .NE. 0, THEN THE LINEAR SYSTEM A*X = B CAN BE SOLVED FOR
C X BY A
C CALL SUBST (W, IPIVOT, B, N, X )
C***** M E T H O D *****
C THE PROGRAM FOLLOWS ALGORITHM 4.2, USING SCALED PARTIAL PIVOTING.
C
C IFLAG = 1
C INITIALIZE IPIVOT, D
DO 9 I=1,N
IPIVOT(I) = I
ROWMAX = 0.
DO 5 J=1,N
5 ROWMAX = AMAX1(ROWMAX,ABS(W(I,J)))
IF (ROWMAX .EQ. 0.) THEN
IFLAG = 0
ROWMAX = 1.
END IF
9 D(I) = ROWMAX
IF (N .LE. 1) RETURN
C FACTORIZATION
DO 20 K=1,N-1
C DETERMINE PIVOT ROW, THE ROW ISTAR .
COLMAX = ABS(W(K,K))/D(K)
ISTAR = K
DO 13 I=K+1,N
AWIKOD = ABS(W(I,K))/D(I)
IF (AWIKOD .GT. COLMAX) THEN
COLMAX = AWIKOD
ISTAR = I
END IF
13 CONTINUE
IF (COLMAX .EQ. 0.) THEN
IFLAG = 0
ELSE
IF (ISTAR .GT. K) THEN
C MAKE K THE PIVOT ROW BY INTERCHANGING IT WITH
C THE CHOSEN ROW ISTAR .
IFLAG = -IFLAG
I = IPIVOT(ISTAR)
IPIVOT(ISTAR) = IPIVOT(K)
IPIVOT(K) = I
TEMP = D(ISTAR)
D(ISTAR) = D(K)
D(K) = TEMP
DO 15 J=1,N
TEMP = W(ISTAR,J)
W(ISTAR,J) = W(K,J)
15 W(K,J) = TEMP
END IF
C ELIMINATE X(K) FROM ROWS K+1,....,N .
16 DO 19 I=K+1,N
W(I,K) = W(I,K)/W(K,K)
RATIO = W(I,K)
DO 19 J=K+1,N
19 W(I,J) = W(I,J) - RATIO*W(K,J)
CONTINUE
END IF

```

20 CONTINUE  
 IF (W(N,N) .EQ. 0.) IFLAG = 0 RETURN  
 END

بحث فوق، روش بسیار کارآمدی است برای محاسبه وارون يك ماتریس وارونپذیر مفروضی مانند  $A$  از مرتبه  $n$ . همان طوری که در بخش ۱.۴ اشاره شد به ازای  $n, \dots, 1 = z$  ستون  $j$  از ماتریس وارون  $A^{-1}$ ، جواب دستگاه خطی

$$AX = \mathbf{i}_j$$

است. بنابراین، برای محاسبه  $A^{-1}$ ، یکبار از زیر برنامه FACTOR استفاده می شود، سپس هر يك از  $n$  دستگاه  $AX = \mathbf{i}_j, j = 1, \dots, n$ ، به وسیله الگوریتم ۴.۴ یعنی با استفاده از SUBST حل می شود. بنابراین به محض اینکه عملیات حذف انجام پذیرد فقط  $n^2 \times n$  عمل ضرب و تقریباً همین تعداد عمل جمع برای پیدا کردن  $A^{-1}$  لازم است.

در عین حالی که این دستورالعمل ساده را، برای محاسبه وارون يك ماتریس داده ایم بلافاصله باید متذکر شویم که معمولاً همیشه دلیل موجهی برای محاسبه ماتریس وارون وجود ندارد. گاهی اوقات در بعضی مسائل اتفاقاً ممکن است که درایه های  $A^{-1}$  دارای ارزش فیزیکی خاصی باشند. مثلاً در روشهای آماری برای برآوردن يك تابع به داده های تجربی به وسیله روش کوچکترین توانهای دوم، درایه های ماتریس  $A^{-1}$  اطلاعاتی درباره نوع و اندازه خطا در داده ها در بر دارند. اما هر گاه به  $A^{-1}$  فقط به منظور محاسبه بردار  $A^{-1}\mathbf{b}$  (مثلاً در حل  $AX = \mathbf{b}$ ) و یا حاصلضرب ماتریسی  $A^{-1}B$  احتیاج باشد،  $A^{-1}$  نمی بایستی به طور صریح محاسبه شود، بلکه بهتر است الگوریتم جایگزینی ۴.۴ برای محاسبه این حاصلضرب به کار رود. دلیل این توصیه چنین است: محاسبه بردار  $A^{-1}\mathbf{b}$  برای يك  $\mathbf{b}$  داده شده به معنی پیدا کردن جواب معادله خطی  $AX = \mathbf{b}$  است. همان طوری که قبلاً اشاره شد، به محض اینکه ماتریس  $A$  به توسط الگوریتم ۲.۴ به عوامل مثلثی تجزیه شود محاسبه  $A^{-1}\mathbf{b}$  می تواند از راه الگوریتم ۴.۴ به آسانی و دقیقاً با همان تعداد عملیات ضرب و جمعی که برای تشکیل حاصلضرب  $A^{-1}$  و  $\mathbf{b}$  لازم است، انجام پذیرد. بنابراین، همین قدر که تجزیه به عوامل مثلثی برای  $A$  در دست باشند، دانستن  $A^{-1}$  به طور صریح امتیازی در محاسبه  $A^{-1}\mathbf{b}$  ندارد، (زیرا که تشکیل حاصلضرب  $A^{-1}B$  شامل ضرب هر ستون  $B$  در  $A^{-1}$  است، این تذکر برای محاسبه چندین حاصلضرب ماتریسی نیز کاربرد دارد.) از سوی دیگر اولین گام در محاسبه  $A^{-1}$  پیدا کردن عوامل مثلثی برای  $A$  است که سپس با  $n$  بار کاربرد الگوریتم جایگذاری دنبال می شود، بنابراین محاسبه  $A^{-1}$  مستلزم مقدار قابل ملاحظه ای محاسبات اولیه در مقایسه با محاسبه  $A^{-1}\mathbf{b}$  است. علاوه ماتریسی که این چنین محاسبه شده باشد فقط يك وارون تقریبی است و به تعبیری دقت عمل کمتری نسبت به عوامل مثلثی دارد، زیرا که این ماتریس به وسیله محاسبات اضافی دیگر روی عوامل مثلثی به دست می آید.

بنابراین اگر در محاسبه یک حاصلضرب ماتریسی منضم  $A^{-1}$ ، صرفاً  $A^{-1}$  را به کار ببریم، هیچ گونه فایده‌ای نخواهد داشت و حتی فاقد دقت عمل خواهد بود.

در زیر یک برنامه فورترن برای محاسبه وارون یک ماتریس  $N \times N$  مفروض  $A$  داده‌ایم. در این برنامه، از زیر برنامه‌های FACTOR و SUBST که قبلاً داده شده‌اند استفاده شده است. همچنین، نمونه ورودی و نتایج خروجی ذکر شده‌اند. تذکرات زیر برای فهم برنامه ممکن است مفید واقع شوند. مرتبه ماتریس  $A$ ، یعنی  $N$ ، بخشی از ورودی این برنامه است، بنابراین مشخص کردن مرتبه دقیق  $A$  هنگام همگرانی ممکن نیست. از سوی دیگر هر دو زیر برنامه FACTOR و SUBST، خواهان ماتریسهای  $A$  و  $W$  (یا  $W$  با مرتبه دقیق  $N \times N$  هستند. بنابراین در برنامه فورترن زیر، ماتریس  $A$  در یک آرایه یک بعدی ذخیره می‌شود که از این قرارداد زبان فورترن که آرایه  $(I, J)$  ام از یک آرایه دو بعدی  $(N, M)$  با درایه  $((J-1) * N + I)$  ام از آرایه یک بعدی هم‌ارزاست، استفاده می‌شود. قرارداد مشابهی برای ذخیره کردن درایه‌های ستون  $J$  ام  $A^{-1}$  در آرایه یک بعدی  $AINV$  به کار گرفته شده است: درایه  $((J-1) * N + 1)$  از آرایه  $AINV$  به‌عنوان اولین درایه  $N$ -برداری  $X$  از زیر برنامه SUBST به این زیر برنامه داده می‌شود و در این بردار جواب دستگاه  $AX = I$  ذخیره می‌شود.

برنامه فورترن برای محاسبه وارون یک ماتریس داده شده

```

C PROGRAM FOR CALCULATING THE INVERSE OF A GIVEN MATRIX
C CALLS FACTOR, SUBST.
  PARAMETER NMAX=30, NMAXSQ=NMAX*NMAX
  INTEGER I, IBEG, IFLAG, IPIVOT(NMAX), J, N, NSQ
  REAL A(NMAXSQ), AINV(NMAXSQ), B(NMAX)
  1 READ 501, N
501 FORMAT(I2)
  IF (N .LT. 1 .OR. N .GT. NMAX) STOP
C READ IN MATRIX ROW BY ROW
  NSQ = N*N
  DO 10 I=1, N
  10 READ 510, (A(J), J=I, NSQ, N)
  510 FORMAT(5E15.7)
C CALL FACTOR ( A, N, B, IPIVOT, IFLAG )
  IF (IFLAG .EQ. 0) THEN
  PRINT 611
611 FORMAT('MATRIX IS SINGULAR') GO TO 1
  END IF
  DO 21 I=1, N
  21 B(I) = 0.
  IBEG = 1
  DO 30 J=1, N
  B(J) = 1.
  CALL SUBST ( A, IPIVOT, B, N, AINV(IBEG) )
  B(J) = 0.
  30 IBEG = IBEG + N
  PRINT 630
630 FORMAT('THE COMPUTED INVERSE IS '//)
  DO 31 I=1, N

```

ماتریسها و دستگاههای معادلات خطی ۲۱۵

```
31 PRINT 631, I, (AINV(J), J=I, NSQ, N)
631 FORMAT('BROW ', I2, 8E15.7 / (7X, 8E15.7))
GO TO 1
END
```

نمونه‌ای از ورودی

```
3
2. 3. -1.
4. 4. -3.
-2. 3. -1.
```

نتایج به دست آمده

وارون محاسبه شده عبارتست

```
ROW 1 0.2500000E 00 0.0 -0.2499999E 00
ROW 2 0.5000000E 00 -0.1999998E 00 0.9999996E -01
ROW 3 0.1000000E 01 -0.6000000E 00 -0.2000000E 00
```

تمرین

۱-۴.۴ برنامه فورترن برای محاسبه  $A^{-1}$  را که در متن داده شده است طوری اصلاح کنید که به یک برنامه کلیتری برای محاسبه حاصلضرب  $C = A^{-1}B$  منجر شود. یک ماتریس  $n \times n$  (وارونپذیر) و  $B$  یک ماتریس  $n \times m$  است که هر دو داده شده‌اند.

۲-۴.۴ وارون ماتریس ضرایب  $A$  از دستگاه معادلات تمرین ۱-۴.۴ را محاسبه کنید. سپس دقت ماتریس وارون محاسبه شده،  $A_{\text{comp}}^{-1}$ ، را از راه محاسبه  $A_{\text{comp}}^{-1}A$  و  $AA_{\text{comp}}^{-1}$  بررسی کنید.

۳-۴.۴ نشان دهید که ماتریس

$$A = \begin{bmatrix} 2 & 2 & 1 \\ 1 & 1 & 1 \\ 3 & 2 & 1 \end{bmatrix}$$

وارونپذیر است، اما  $A$  را نمی‌توان به صورت حاصلضرب یک ماتریس پایین‌مثلثی در یک ماتریس بالامثلثی نوشت.

۴-۴.۴ ثابت کنید که جمع و ضرب دو ماتریس پایین (بالا) مثلثی یک ماتریس پایین (بالا) مثلثی است و وارون یک ماتریس پایین (بالا) مثلثی یک ماتریس پایین (بالا) مثلثی است.

۵-۴.۴ ثابت کنید که تجزیه به عوامل مثلثی به تعبیر زیر یکتاست: اگر  $A$  وارونپذیر و  $L_1 U_1 = L_2 U_2$ ،  $L_1$  و  $L_2$  ماتریسهای واحد پایین‌مثلثی و  $U_1$  و  $U_2$  ماتریسهای بالامثلثی

باشند، آنگاه  $L_1 = L_2$  و  $U_1 = U_2$  (دوهمینای: تمرین ۱.۴-۸ را برای اثبات وارونپذیری  $U_1$  و  $L_1$  به کار گیرید و سپس نشان دهید که تساوی  $U_1^{-1} L_1^{-1} = U_1 U_1^{-1}$  باید برقرار باشد، که با استفاده از تمرین ۴.۴-۴ نتیجه می شود که  $L_1^{-1} L_1^{-1}$  باید يك ماتریس قطری باشد، بنابراین چون تمام درایه های قطری  $L_1$  و  $L_2$  عدد ۱ می باشند،  $(L_1^{-1} L_1 = I)$ .

۴.۴-۶ با استفاده از نتایج تمرین ۴.۴-۵ نشان دهید که اگر  $A$  متقارن باشد ( $A = A^T$ ) و به عوامل مثلثی  $A = LU$  تجزیه شود، آنگاه داریم  $U = D L^T$  که  $D$  يك ماتریس قطری است که عناصر قطریش همان عناصر قطری  $U$  هستند.

۴.۴-۷ ثابت کنید که اگر ماتریس سه قطری  $A$  بتواند به صورت  $A = LU$  تجزیه شود،  $L$  يك ماتریس پایین-مثلثی و  $U$  يك ماتریس بالا مثلثی است. آنگاه هم  $L$  هم  $U$ ، ماتریس سه قطری خواهند بود. الگوریتم ۳.۴ را به عنوان روشی برای تجزیه ماتریسهای سه قطری تفسیر کنید.

۴.۴-۸ طرحهای فشرده عوامل مثلثی ماتریس  $A$  یعنی  $L$  و  $U$  را با به کار گیری معادله (۲۳.۴) به شکل زیر بنا کنید تا درایه های جالب  $L$  و  $U$  به دست آیند.

$$l_{ij} = (a_{ij} - l_{i1}u_{1j} - \dots - l_{i,j-1}u_{j-1,j}) / u_{jj} \quad i > j \text{ به ازای } j$$

$$u_{ij} = a_{ij} - l_{i1}u_{1j} - \dots - l_{i,i-1}u_{i-1,j} \quad i \leq j \text{ به ازای } j$$

در واقع، محتوای نهایی آرایه کار  $W$  باید پس از انجام همه تغییرات در یک زمان، برای هر درایه، به دست آید و بنابراین مانع نوشتن نتایج مختلف میانی شود. البته این عمل می بایستی به ترتیب منظمی انجام گیرد. مثلا  $l_{ij}$  (به ازای  $i > j$ ) نمی تواند به دست آید مگر آنکه قبلا  $l_{ir}$  به ازای  $r < j$  و  $u_{rj}$  به ازای  $r \leq j$  معین شده باشند، و باز باید قبلا  $l_{ir}$  و  $u_{rj}$  به ازای  $r < j$  معین باشند تا  $u_{ij}$  (برای  $i \leq j$ ) محاسبه گردد.

(الف) برای به دست آوردن  $L$  و  $U$  از  $A$  به روش فشرده، الگوریتمی طرح کنید.  
(ب) الگوریتم مذکور را چنان تغییر دهید که برای گزینش لولایی جزئی مدرج به کار آید.

(پ) اگر الگوریتم مذکور قبلا بر این اساس درست نشده است، آن را طوری تغییر دهید که درونی ترین حلقه بسا همه زیرنمایه ها را در بر گیرد (برای انگیزه آن به تمرین ۲.۴-۷ مراجعه کنید).

۴.۴-۹ روش چولسکی اگر ماتریس  $A$  از مرتبه  $n$ ، حقیقی، متقارن ( $A = A^T$ ) و مثبت معین (یعنی به ازای تمام  $n$ -برداری های غیر صفر  $x$ ،  $x^T A x > 0$ ) باشد، آنگاه امکان پذیر است که  $A$  به صورت  $LDL^T$  تجزیه شود، که در آن  $L$  يك ماتریس واحد پایین-مثلثی و حقیقی و  $D = (d_{ij})$  يك ماتریس قطری (مثبت) است. بنابراین از رابطه (۲۳.۴) خواهیم داشت

$$l_{ij} = (a_{ij} - l_{i1}l_{j1} - \dots - l_{i,j-1}l_{j,j-1}) / d_{jj} \quad i > j \text{ به ازای } j$$

در حالی که داریم

$$d_{jj} = a_{jj} - l_{j1}^2 - \dots - l_{j,j-1}^2$$

بر اساس این معادله‌ها برای تولید (جالبترین قسمت)  $D$  و  $L$ ، يك زیر برنامه فورترن بنویسید و زمانی که  $D$  و  $L$  مشخص شدند، برای حل  $AX = b$  بر حسب  $x$  از راه جایگذاری نیز يك زیر برنامه تهیه کنید.

۴-۱۰ نشان دهید که روش چولسکی، هر وقت که ماتریس  $A$  به صورت  $BB^T$  و  $B$  ماتریس وارونپذیر باشد قابل اجراست.

### ۵.۴ خطا و باقیمانده يك جواب تقریبی؛ نرم‌ها

هر جواب حساب شده يك دستگاه خطی باید، به دلیل خطای گرد کردن، يك جواب تقریبی تلقی شود. در این بخش، درباره مسئله دشوار تعیین خطای يك جواب تقریبی (بدون دانستن آن جواب) بحث می‌شود. در این بحث نرم<sup>۱</sup> را وارد می‌کنیم و از آن به عنوان ابزار مناسبی برای اندازه‌گیری «اندازه»ی بردارها و ماتریسها استفاده می‌کنیم. اگر  $\hat{x}$  يك جواب حساب شده دستگاه خطی  $AX = b$  باشد، آنگاه خطای این جواب برابر با تفاضل زیر است

$$e = x - \hat{x}$$

بدیهی است که این خطا معمولاً بر ما مجهول است (چون در غیر این صورت  $x$  معلوم است و بحثهای دیگر غیر ضروری اند). اما همواره می‌توان باقیمانده<sup>۲</sup> (خطا) یعنی

$$r = Ax - A\hat{x}$$

را محاسبه کرد، زیرا که  $Ax$  دقیقاً سمت راست  $b$  است. سپس این باقیمانده  $r$ ، میزان صدق کردن  $\hat{x}$  را در دستگاه خطی  $AX = b$  معین می‌کند. اگر  $r$  يك بردار صفر باشد، آنگاه  $\hat{x}$  جواب (دقیق) این دستگاه است. یعنی در این صورت  $e$  صفر است. اگر  $\hat{x}$  تقریب خوبی برای جواب  $x$  باشد، می‌توان انتظار داشت که کلیه درایه‌های  $r$ ، حداقل به معنای نسبی، کوچک باشند.

□ مثال ۵.۴: دستگاه خطی ساده زیر را که ریشه یگانه آن،  $x$ ، دارای درایه‌های  $x_1 = x_2 = 1$  در نظر می‌گیریم

$$1.01x_1 + 0.99x_2 = 2$$

$$0.99x_1 + 1.01x_2 = 2$$

جواب تقریبی  $\hat{x} = \begin{bmatrix} 1.01 \\ 1.01 \end{bmatrix}$  خطایی مساوی  $e = \begin{bmatrix} -0.01 \\ -0.01 \end{bmatrix}$  و باقیمانده‌ای برابر

دارد، لذا در این حالت يك باقیمانده «كوجك» (نسبت به سمت راست  $\mathbf{r} = \begin{bmatrix} -0.002 \\ -0.002 \end{bmatrix}$ )

معادله) به يك خطای نسبتاً «كوجك» مربوط می شود. از سوی دیگر جواب تقریبی  $\hat{\mathbf{x}} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$

خطایی مساوی  $\mathbf{e} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$ ، ولی باقیماندهای برابر  $\mathbf{r} = \begin{bmatrix} -0.002 \\ 0.002 \end{bmatrix}$  دارد که باز هم

باقیمانده نسبتاً «كوجك» است در صورتی که خطا اکنون نسبتاً «بزرگ» شده است. اگر برای سمت راست مقادیر دیگری اختیار کنیم می توانیم نتیجه عکس به دست آوریم. دستگاه خطی

$$1001x_1 + 0999x_2 = 2 \quad (30.4)$$

$$0999x_1 + 1001x_2 = -2$$

دارای جواب منحصر به فرد  $x_1 = 100$  و  $x_2 = -100$  است. جواب تقریبی

$\hat{\mathbf{x}} = \begin{bmatrix} 101 \\ -99 \end{bmatrix}$  خطایی مساوی  $\mathbf{e} = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$  ولی باقیماندهای برابر  $\mathbf{r} = \begin{bmatrix} -2 \\ -2 \end{bmatrix}$  دارد

که در این حالت باقیمانده اکنون نسبتاً «بزرگ» است، در حالی که خطا نسبتاً «كوجك» است  $\square$  (فقط يك درصد جواب است).

همان طوری که این مثال نشان می دهد اندازه باقیمانده  $\mathbf{r} = \mathbf{b} - A\hat{\mathbf{x}}$  برای يك جواب تقریبی  $\hat{\mathbf{x}}$ ، همیشه شاخص مطمئنی برای اندازه خطای  $\mathbf{e} = \mathbf{x} - \hat{\mathbf{x}}$  در این جواب تقریبی نیست. اینکه آیا يك باقیمانده «كوجك» مستلزم يك خطای «كوجك» هست یا نیست بستگی به «اندازه»ی ضرایب ماتریس و وارون آن دارد که دقیقاً در زیر توضیح داده می شود. برای این بحث نیاز به وسیله ای برای اندازه گیری «اندازه»ی  $n$ -برداري ها و ماتریسهای  $n \times n$  داریم.

قدر مطلق وسیله ای مناسب برای اندازه گیری «اندازه»ی اعداد حقیقی و حتی اعداد همناقت است. یقین نداریم که بتوانیم يك راه قطعی برای اندازه يك  $n$ -برداري یا يك ماتریس  $n \times n$  تعیین کنیم. مطمئناً يك راه قطعی برای انجام این عمل که در تمام موارد مورد قبول باشد وجود ندارد

برای مثال، معیاری که غالباً برای اندازه يك  $n$ -برداري  $\mathbf{a}$  به کار برده می شود، عدد غیر منفی زیر است

$$\|\mathbf{a}\|_{\infty} = \max_{1 \leq i \leq n} |a_i| \quad (31.4)$$

اکنون فرض کنید که جواب محاسبه شده  $\hat{\mathbf{x}}$  برای  $A\mathbf{x} = \mathbf{b}$ ، که از این راه اندازه گرفته شده است، تا شش رقم اعشاری دقیق باشد، یعنی

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|_{\infty}}{\|\mathbf{x}\|_{\infty}} < 10^{-6} \quad (۳۲.۴)$$

که نشان دهنده يك جواب محاسبه شده بسیار رضایتبخش در موردی است که مجهولات، مثلاً، مقادیر تقریبی يك جواب رضایتبخش يك معادله دیفرانسیل می باشند. اما اگر یکی از مجهولات اتفاقاً درآمد سالانه شما و دیگری تولید ناخالص ملی باشد، نابرابری (۳۲.۴) هیچ گونه اشاره‌ای به رضایتبخش بودن یا نبودن  $\mathbf{x}$  (درحد موردنظر) ندارد، زیرا با برقراری (۳۲.۴)، خطا در محاسبه درآمد سالانه شما (حتی اگر فقط برای یکسال باشد) ممکن است مستقلاً شما را بسیار ثروتمند یا برای تمام عمر بدهکار سازد. معیاری مانند

$$\|\mathbf{a}\| = \max\{10^{10}|a_1|, \max_{2 \leq i \leq n} |a_i|\}$$

(با این فرض که درآمد سالانه مجهول اول باشد) اطلاعات بسیار زیادتری برای اندازه خواهد داد تا معیارهای خاص دیگری که در آنها از چندین عدد (به جای به کار گرفتن فقط يك عدد غیرمنفی) برای بیان «اندازه» يك  $n$ -برداری استفاده می شود. اما در بیشتر موارد، کافی است که اندازه يك  $n$ -برداری با نرم اندازه گیری شود. نرم، بعضی از ویژگیهای قدرمطلق اعداد را حفظ می کند. به ویژه يك نرم به هر  $n$ -برداری  $\mathbf{a}$  يك عدد حقیقی  $\|\mathbf{a}\|$  با شرایط قابل قبول زیر نسبت می دهد، که آن را نرم  $\mathbf{a}$  می نامند، که تابع شرایط مناسب زیر است:

(i) به ازای هر  $n$ -برداری  $\mathbf{a}$ ، نامساوی  $\|\mathbf{a}\| \geq 0$  و تساوی  $\|\mathbf{a}\| = 0$ ،

برقرار است اگر، و فقط اگر،  $\mathbf{a} = \mathbf{0}$ .

(ii) به ازای هر  $n$ -برداری  $\mathbf{a}$  و هر عدد  $\alpha$  داریم،  $\|\alpha\mathbf{a}\| = |\alpha| \|\mathbf{a}\|$  (۳۲.۴)

(iii) به ازای هر دو  $n$ -برداری  $\mathbf{a}$  و  $\mathbf{b}$  داریم،  $\|\mathbf{a} + \mathbf{b}\| \leq \|\mathbf{a}\| + \|\mathbf{b}\|$

اولین شرط، همه  $n$ -برداری‌ها را به استثنای بردار صفر مقید می کند که دارای «طول» مثبت باشند، شرط دوم بیان می کند که، مثلاً،  $\mathbf{a}$  و منفی آن  $-\mathbf{a}$  يك «طول» دارند و طول  $\mathbf{a}$ ، سه برابر طول  $\mathbf{a}$  است. شرط سوم نابرابری مثلثی<sup>۱</sup> نامیده می شود، زیرا که بیان می کند که مجموع طولهای دو ضلع يك مثلث هرگز از طول ضلع سوم آن کوچکتر نیست.

احتمالاً دانشجویان با طول اقلیدسی و یا نرم  $n$ -برداری  $\mathbf{a} = (a_i)$

$$\|\mathbf{a}\|_2 = \sqrt{\mathbf{a}^T \mathbf{a}} = \sqrt{|a_1|^2 + |a_2|^2 + \dots + |a_n|^2}$$



لا اقل در موارد  $n=2$  یا  $n=3$  آشنایی دارند. اما به دلیلی که در زیر توضیح داده می‌شود، ترجیح می‌دهیم که در مثالهای عددی زیر، از نرم ماکسیمم (۳۱.۴) به عنوان راهی برای اندازه‌گیری طول  $n$ -برداری  $\mathbf{a}$  استفاده کنیم. به آسانی می‌توان تحقیق کرد که (۳۱.۴) یک تعریف نرم است، یعنی  $\|\mathbf{a}\| = \|\mathbf{a}\|_\infty$  سه ویژگی نرم را که در (۳۳.۴) آمده است داراست. برای ویژگی (i)،  $\|\mathbf{a}\|_\infty$  عبارت است از ماکسیمم کمیت‌های غیرمنفی و بنابراین خود غیرمنفی است، همچنین  $\|\mathbf{a}\|_\infty = 0$ ، اگر و فقط اگر به ازای کلیه مقادیر  $i$ ،  $|a_i| = 0$ ، عیناً مثل این است که گفته شود  $\mathbf{a} = \mathbf{0}$ ، بعلاوه اگر  $\alpha$  یک عدد باشد، آنگاه داریم

$$\|\alpha \mathbf{a}\|_\infty = \max_i |\alpha a_i| = \max_i |\alpha| |a_i| = |\alpha| \max_i |a_i| = |\alpha| \|\mathbf{a}\|_\infty$$

که اثبات (ii) است. بالاخره

$$\begin{aligned} \|\mathbf{a} + \mathbf{b}\| &= \max_i |a_i + b_i| \leq \max_i (|a_i| + |b_i|) \\ &\leq \max_i |a_i| + \max_i |b_i| = \|\mathbf{a}\|_\infty + \|\mathbf{b}\|_\infty \end{aligned}$$

که اثبات (iii) است.

نرم‌های برداری دیگری که غالباً به کار می‌روند شامل ۱-نرم

$$\|\mathbf{a}\| = \|\mathbf{a}\|_1 = \sum_{i=1}^n |a_i|$$

و موارد متفاوت  $p$ -نرم وزین هستند

$$\|\mathbf{a}\| = \|\mathbf{a}\|_{p, w} = \left( \sum_{i=1}^n |a_i|^p w_i \right)^{1/p}$$

در رابطهٔ بالا  $p$  عددی است بین ۱ و  $\infty$  و اعداد  $w_1, \dots, w_n$  کمیت‌های ثابت و مثبت‌اند. در حالت  $p=2$  و  $w_i=1$ ، (به‌ازای جميع مقادیر  $i$ )، همان نرم آشنای اقلیدسی به دست می‌آید.

وقتی یک نرم برداری انتخاب شد، اندازهٔ متناظر ماتریس  $A$  از مرتبهٔ  $n \times n$  را از مقایسهٔ اندازهٔ  $A\mathbf{x}$  با اندازهٔ  $\mathbf{x}$  تعیین می‌کنیم. نرم ماتریسی و متناظر  $A$  را دقیقاً، چنین تعریف می‌کنیم

$$\|A\| = \max \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} \quad (34.4)$$

که ماکسیمم‌گیری روی تمامی  $n$ -برداریهای (غیر صفر)  $\mathbf{x}$  انجام می‌شود. می‌توان نشان

داد که این ماکسیمم به ازای هر ماتریس  $A$  از مرتبه  $n \times n$  (و انتخاب هر نرم برداری) وجود دارد. نرم ماتریسی  $\|A\|$  با دو ویژگی زیر مشخص می شود

$$\left. \begin{aligned} \|Ax\| &\leq \|A\| \|x\| && \text{به ازای هر } n\text{-برداری } x \\ & && \text{و} \\ \|Ax\| &\geq \|A\| \|x\| && \text{به ازای } n\text{-برداری غیر صفری مانند } x \end{aligned} \right\} (35.4)$$

البته بلافاصله از (35.4) نتیجه می شود که به ازای هر  $x$  با شرط بسز قراری  $\|Ax\| \geq \|A\| \|x\|$ ،  $\|Ax\| = \|A\| \|x\|$  برقرار خواهد بود. سپس می توان نشان داد که ویژگیهای زیر برای نرم ماتریس (34.4) برقرارند:

$$\left. \begin{aligned} (i) & \text{ برای تمامی ماتریسهای نظیر ماتریس } A \text{ از مرتبه } n \times n : \|A\| \geq 0, \\ & \text{ و تساوی } \|A\| = 0 \text{ برقرار است اگر و فقط اگر } A = 0. \\ (ii) & \text{ به ازای جمیع ماتریسهای } A \text{ از مرتبه } n \times n \text{ و جمیع اعداد } \alpha, \\ & \|\alpha A\| = |\alpha| \|A\| \\ (iii) & \text{ به ازای هر دو ماتریس } A \text{ و } B \text{ از مرتبه } n \times n \end{aligned} \right\} (36.4)$$

$$\|A+B\| \leq \|A\| + \|B\|$$

بنابر این اصطلاح «نرم» برای عدد  $\|A\|$  موجه است.

بعلاوه

$$\|AB\| \leq \|A\| \|B\| \text{ داریم } B \text{ از مرتبه } n \times n \text{ و } A \text{ از مرتبه } n \times n \text{ (iv) (37.4)}$$

و بالاخره، اگر ماتریس  $A$  وارون پذیر باشد، آنگاه  $x = A^{-1}(Ax)$ ، بنا بر این  $\|x\| \leq \|A^{-1}\| \|Ax\|$ . از ترکیب این رابطه با (35.4) خواهیم داشت.

$$\frac{\|x\|}{\|A^{-1}\|} \leq \|Ax\| \leq \|A\| \|x\|, \text{ به ازای همه } n\text{-برداریهای } x, (38.4)$$

و هر دو نامساوی نافذند، یعنی با انتخاب بردار مناسب (غیر صفر)  $x$ ، از هر يك می توان يك نامساوی ساخت.

چنانکه مشاهده می شود، نرم ماتریسی

$$\|A\|_r = \max \frac{\|Ax\|_r}{\|x\|_r}$$

که بر اساس نرم اقلیدسی معمولاً بسیار مشکل محاسبه می شود. در حالی که نرم ماتریسی

$$\|A\|_\infty = \max \frac{\|Ax\|_\infty}{\|x\|_\infty}$$

که بر پایهٔ نرم ماکسیمم استوار است، به آسانی قابل محاسبه بوده و با عدد زیر بیان می‌شود

$$\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \quad (۳۹.۴)$$

برای اثبات این، باید نشان دهیم که عدد  $\|A\| = \max_i \sum_j |a_{ij}|$  در دو حکم (۳۵.۴) صدق می‌کند، یعنی

$$\|Ax\|_\infty \leq \left( \max_i \sum_{j=1}^n |a_{ij}| \right) \|x\|_\infty, \quad x, \text{ بردارهای } n\text{-تایی}$$

$$\|Ax\|_\infty \geq \left( \max_i \sum_{j=1}^n |a_{ij}| \right) \|x\|_\infty, \quad x \text{ بردار غیر صفری نظیر } x$$

اما به ازای يك  $x$  دلخواه داریم،

$$\begin{aligned} \|Ax\|_\infty &= \max_{1 \leq i \leq n} \left| \sum_{j=1}^n a_{ij} x_j \right| \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| |x_j| \\ &\leq \max_{1 \leq i \leq n} \left( \left( \max_{1 \leq j \leq n} |x_j| \right) \sum_{j=1}^n |a_{ij}| \right) = \|x\|_\infty \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \end{aligned}$$

که صحت حکم اول را ثابت می‌کند. برای حکم دوم، گیریم  $i_0$  يك عدد صحیح بین ۱ تا  $n$  باشد، به طوری که

$$\sum_{j=1}^n |a_{i_0 j}| = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|$$

فرض می‌کنیم  $x$  يك  $n$ -برداری با نرم ماکسیمم ۱ باشد به طوری که

$$a_{i_0 j} x_j = |a_{i_0 j}| \quad j = 1, \dots, n$$

مثلاً می‌گیریم

$$x_j = \begin{cases} 1 & \text{اگر } a_{i_0 j} \geq 0 \\ -1 & \text{اگر } a_{i_0 j} < 0 \end{cases} \quad j = 1, \dots, n$$

در این صورت برای این بردار  $x$  که غیر صفر بودن آن واضح است، داریم  $\|x\|_\infty = 1$  و

$$\begin{aligned} \|Ax\|_\infty &= \max_{1 \leq i \leq n} \left| \sum_{j=1}^n a_{ij} x_j \right| \geq \left| \sum_{j=1}^n a_{i_0 j} x_j \right| \\ &= \sum_{j=1}^n |a_{i_0 j}| = \|x\|_\infty \left( \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \right) \end{aligned}$$

که اثبات حکم دوم است.

□ مثال ۵.۴ (الف): برای ماتریس ضرایب  $A$  در مثال ۵.۴ به آسانی می توان دریافت که

$$\|A\|_{\infty} = \max\{|1001| + |0999|, |0999| + |1001|\} = 2$$

دیدیم که

$$A \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \quad A \begin{bmatrix} 100 \\ -100 \end{bmatrix} = \begin{bmatrix} 2 \\ -2 \end{bmatrix}$$

و از این رو

$$A \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 100 \\ -100 \end{bmatrix} \right) = \begin{bmatrix} 4 \\ 0 \end{bmatrix}, \quad A \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 100 \\ -100 \end{bmatrix} \right) = \begin{bmatrix} 0 \\ 4 \end{bmatrix}$$

بنابراین

$$A^{-1} = \begin{bmatrix} 25025 & -24975 \\ -24975 & 25025 \end{bmatrix} \quad \text{که نشان می دهد} \quad A \begin{bmatrix} \frac{101}{4} & \frac{-99}{4} \\ -99 & 101 \\ \frac{1}{4} & \frac{1}{4} \end{bmatrix} \\ = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

که در نتیجه

$$\|A^{-1}\|_{\infty} = \max\{|25025| + |-24975|, |-24975| + |25025|\} = 50$$

بنابراین برای این مثال رابطه (۳۸.۴) بیان می کند که

$$0.02\|x\|_{\infty} \leq \|Ax\|_{\infty} \leq 2\|x\|_{\infty}, \quad x \text{ برداری } 2\text{-برداری}$$

با انتخاب  $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$  خواهیم داشت  $Ax = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$  و بنابراین  $\|x\|_{\infty} = 1$  و  $\|Ax\|_{\infty} = 2$

و نامساوی دوم به تساوی تبدیل می شود. با انتخاب  $x = \begin{bmatrix} 100 \\ -100 \end{bmatrix}$  خواهیم داشت

$Ax = \begin{bmatrix} 2 \\ -2 \end{bmatrix}$ ، بنابراین  $\|x\|_{\infty} = 100$  و  $\|Ax\|_{\infty} = 2$  و در این انتخاب  $x$ ، نامساوی

اول به تساوی بدل می شود.

اکنون به بحث در بارهٔ رابطهٔ بین خطای  $\mathbf{e} = \mathbf{x} - \hat{\mathbf{x}}$  در جواب تقریبی  $\hat{\mathbf{x}}$  از  $A\mathbf{x} = \mathbf{b}$  و باقیماندهٔ  $\mathbf{r} = \mathbf{b} - A\hat{\mathbf{x}}$  بازمی‌گردیم داریم

$$\mathbf{r} = A\mathbf{x} - A\hat{\mathbf{x}} = A(\mathbf{x} - \hat{\mathbf{x}}) = A\mathbf{e}$$

بنابراین  $\mathbf{e} = A^{-1}\mathbf{r}$  از این رو، با توجه به  $(A^{-1})^{-1} = A$  از (۳۸.۴) خواهیم داشت

$$\frac{\|\mathbf{r}\|}{\|A\|} \leq \|\mathbf{e}\| = \|A^{-1}\mathbf{r}\| \leq \|A^{-1}\| \|\mathbf{r}\| \quad (۴۰.۴)$$

رابطهٔ فوق در باب خطای نسبی  $\|\mathbf{e}\|/\|\mathbf{x}\|$  بر حسب باقیماندهٔ نسبی  $\|\mathbf{r}\|/\|\mathbf{b}\|$ ، یک کران بالا و یک کران پایین به صورت زیر می‌دهد

$$\frac{\|\mathbf{b}\|}{\|A\| \|\mathbf{x}\|} \cdot \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{x}\|} \leq \frac{\|A^{-1}\| \|\mathbf{b}\|}{\|\mathbf{x}\|} \cdot \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \quad (۴۱.۴)$$

در اینجا، از روی جوابی که برای دستگاه  $A\mathbf{x} = \mathbf{b}$  به دست آمده است؛ می‌توان  $\|\mathbf{x}\|$  را برآورد کرد. و الا، در حالت خاص  $\hat{\mathbf{x}} = \mathbf{0}$  باید (۴۰.۴) را به کار برد، یعنی

$$\frac{\|\mathbf{b}\|}{\|A\|} \leq \|\mathbf{x}\| \leq \|A^{-1}\| \|\mathbf{b}\|$$

و از (۴۱.۴) خواهیم داشت که

$$\frac{1}{\|A\| \|A^{-1}\|} \cdot \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \leq \frac{\|\mathbf{e}\|}{\|\mathbf{x}\|} \leq \|A^{-1}\| \|A\| \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} \quad (۴۲.۴)$$

کرانه‌های موجود در (۴۱.۴) و (۴۲.۴) به تعبیر زیر نافذ هستند.  $A$  و  $\mathbf{x} \neq \mathbf{0}$  هر چه باشند، می‌توان مقادیر غیر صفری برای  $\mathbf{e}$  و  $\mathbf{r}$  چنان انتخاب کرد که یکی از دو نامساوی (۴۱.۴) به تساوی بدل شود. اگر بخواهیم تساوی را از یکی از نامساویهای (۴۲.۴) به دست آوریم، می‌باید  $\mathbf{x}$  خاصی را نیز انتخاب کنیم، ولی همیشه امکان این گونه انتخابها وجود دارد.

به علت اهمیت روابط (۴۱.۴) و (۴۲.۴)، آنها را با کلمات بیان می‌کنیم: اگر برای ماتریس ضرایب وارون پذیر  $A$  از دستگاه خطی  $A\mathbf{x} = \mathbf{b}$  مقدار  $\|A\| \|A^{-1}\|$  را با  $k$  نشان دهیم، آنگاه خطای نسبی یک جواب تقریبی می‌تواند به بزرگی  $k$  برابر، یا دقیقتر بگوییم،  $\|A^{-1}\| \|\mathbf{b}\| / \|\mathbf{x}\|$  برابر باقیماندهٔ نسبی آن باشد، اما می‌تواند به کوچکی  $1/k$  برابر، یا دقیقتر بگوییم  $\|\mathbf{b}\| / (\|A\| \|\mathbf{x}\|)$  برابر باقیماندهٔ نسبی نیز باشد. بنابراین اگر  $k \approx 1$ ، آنگاه خطای نسبی و باقیماندهٔ نسبی همیشه یک اندازه هستند، و باقیماندهٔ نسبی را می‌توان با اطمینان برای برآورد خطای نسبی به کار برد. اما هر چه که  $k$  بزرگتر باشد، به همان نسبت از روی باقیماندهٔ نسبی اطلاعات کمتری می‌توان دربارهٔ خطای نسبی کسب کرد.

عدد  $\|A^{-1}\| \|A\|$ ، عدد شرط  $A$  نامیده می‌شود و گاهی اوقات به صورت زیر خلاصه نویسی می‌شود

$$\text{cond}(A) = \|A\| \|A^{-1}\|$$

باید توجه کرد که عدد شرط یعنی  $\text{cond}(A)$  برای ماتریس  $A$  به نرم ماتریسی که به کار گرفته می‌شود بستگی دارد و می‌تواند برای بعضی ماتریسها با تغییر نرم ماتریسی به طور قابل ملاحظه‌ای تغییر یابد. از سوی دیگر همیشه عدد شرط حداقل برابر ۱ است، زیرا برای ماتریس واحد  $I$  داریم  $\|I\| = \max\|x\|/\|x\| = 1$ ، و طبق رابطه (۳۷.۴) داریم،  $\|I\| = \|AA^{-1}\| \leq \|A\| \|A^{-1}\|$ .

□ **مثال ۶۰۴:** با توجه به محاسبات قبلی می‌بینیم که برای ماتریس ضرایب  $A$  در مثال ۵۰۴ داریم  $\|A\|_{\infty} \|A^{-1}\|_{\infty} = 2 \times 50 = 100$ . بعلاوه مشاهده شد که در مثال ۵۰۴ حقیقتاً خطای نسبی در یک جواب تقریبی می‌تواند به بزرگی ۱۰۰ برابر باقیمانده نسبی باشد، اما می‌تواند فقط  $1/100$  باقیمانده نسبی نیز باشد. □

کرانهای (۴۱.۴) و (۴۲.۴) به عدد  $\|A^{-1}\|$ ، که به آسانی در دسترس نیست، احتیاج دارند. اما در نمونه‌ای از شرایط، عدد  $\|\hat{e}\|$  تخمین خوبی برای  $\|e\|$  است که  $\hat{e}$  جواب محاسبه شده در دستگاه خطی  $Ae = \mathbf{r}$  است. از آنجایی که معمولاً  $\hat{\mathbf{x}}$  از راه روش گاوس به دست می‌آید و یک تجزیه به عوامل  $A$  در دسترس است، بنابراین  $\hat{e}$  با محاسبات بسیار کمتری نسبت به آنچه که برای به دست آوردن  $\hat{\mathbf{x}}$  لازم است، می‌تواند (به وسیلهٔ SUBST) به دست آید. در این امر فرض آن است که  $\mathbf{r}$  بر اساس دقت مضاعف محاسبه شده است. بردار  $\hat{e}$  که این گونه به دست می‌آید نخستین بارست در اصلاح بارستی (الگوریتم ۵۰۴) است که در بخش بعدی مورد بحث قرار خواهد گرفت.

## تمرین

۱-۵۰۴ تحقیق کنید که

$$\|a\| = \|a\|_1 = \sum_{i=1}^n |a_i|$$

می‌تواند برای تمامی  $n$ -برداریهای  $a$  تعریفی برای نرم باشد.

۴-۲۰۵ ثابت کنید که نرم ماتریسی  $\|A\|_1$  متناظر با نرم برداری  $\|a\|_1$  از تمرین ۱-۵۰۴ را می‌توان به صورت زیر محاسبه کرد

$$\|A\|_1 = \max \sum_{i=1}^n |a_{ij}|$$

۳-۵.۴ اگر یک برداری  $a$  را به عنوان نقطه ای در یک صفحه بسا مختصات  $\{a_1$  و  $a_2\}$  تعبیر کنیم، آنگاه ۲-نرم آن،  $\|a\|_2$ ، فاصله اقلیدسی این نقطه از مرکز خواهد شد. بعلاوه، مجموعه تمامی بردارهایی که نرم اقلیدسی آنها برابر ۱ است، دایره ای به شعاع ۱ پیرامون مبدأ مختصات تشکیل می دهند. «دایره ای به شعاع ۱ پیرامون مبدأ» رسم کنید که فاصله «نقطه»ی  $a$  به وسیله نرمهای زیر محاسبه شود (الف)  $\|a\|_1$ ، (ب) با نرم  $\|a\|_{3/2}$ ، (پ) با نرم اقلیدسی  $\|a\|_2$ ، (ت) با نرم  $\|a\|_4$ ، (ث) با نرم ماکسیم  $\|a\|_\infty$ .

۴-۵.۴ با همان تعبیر تمرین ۳-۵.۴، از ۲-بردارها به عنوان نقاط واقع در صفحه نشان دهید که برای هر دو برداری  $a$  و  $b$  سه نقطه  $a$ ،  $b$  و  $a+b$  رؤس مثلثی هستند با اضلاعی به طول (اقلیدسی)  $\|a\|_2$ ،  $\|b\|_2$  و  $\|a+b\|_2$  و اصطلاح «نامساوی مثلثی»، برای ویژگی (iii) مربوط به نرمها [معادله (۳.۴)] را شرح دهید.

۵-۵.۴ نشان دهید، که برای هر ۲-بردار  $a$  و  $b$  و هر نرم برداری داریم

$$|\|a\| - \|b\|| \leq \|a - b\|$$

۶-۵.۴ نشان دهید که برای هر ۲-بردار  $a$  و  $b$  و هر عدد  $\lambda$  بین ۰ و ۱، داریم

$$\|\lambda a + (1 - \lambda)b\| \leq \max(\|a\|, \|b\|)$$

۷-۵.۴ نشان دهید که نرم ماتریسی  $\|A\| = \max(\|Ax\| / \|x\|)$  را نیز می توان به طریق زیر محاسبه کرد

$$\|A\| = \max_{\|x\|=1} \|Ax\|$$

۸-۵.۴ کلیه احکام (۳.۴) را درباره نرمهای ماتریسی ثابت کنید.

۹-۵.۴ با استفاده از تمرین ۷-۵.۴،  $\|A\|_2$  را برای

$$A = \begin{bmatrix} 3 & -5 \\ 6 & 1 \end{bmatrix}$$

محاسبه کنید. (دانهایی: یک ۲-بردار  $x$  دارای ۲-نرم  $\|x\|_2 = 1$  است اگر، و فقط اگر، به ازای مقداری مانند  $\theta$  داشته باشیم،  $x_1 = \cos \theta$  و  $x_2 = \sin \theta$ ).

۱۰-۵.۴ با استفاده از تمرین ۴-۴، عدد شرط ماتریس ضرایب  $A$  از دستگاه خطی تمرین ۸-۲.۴ را محاسبه کنید، سپس خطای نسبی و باقیمانده نسبی جوابهایی را که در تمرین ۸-۲.۴ و ۴-۳.۴ بر حسب این عدد شرط به دست آمده مورد بحث قرار دهید. همچنین تنها

با استفاده از الگوریتم جایگذاری ۴.۲، مقدار  $\hat{\mathbf{e}} = A^{-1}\mathbf{r}$  را به ازای جوابهای فوق (با محاسبه  $\mathbf{r}$  از راه دقت مضاعف) حساب کنید.

## ۶.۴ تحلیل خطای پسر و اصلاح بارستی

در بخش ۵.۲ پیشین عدد شرط ماتریس ضرایب  $A$  از يك دستگاه خطی  $A\mathbf{x} = \mathbf{b}$ ، یعنی

$$\text{cond}(A) = \|A\| \|A^{-1}\| \quad (۴۳.۴)$$

را به عنوان کمیتی مستقل از  $\mathbf{x}$  دربر آورد خطای يك جواب تقریبی شناسایی کردیم. به طور خلاصه: عدد شرط (۴۳.۴) معیاری است برای اینکه ببینیم، باقیمانده نسبی  $\|\mathbf{b} - A\hat{\mathbf{x}}\| / \|\mathbf{b}\|$  از جواب تقریبی  $\hat{\mathbf{x}}$  تا چه اندازه خطای نسبی این ریشه تقریبی، یعنی  $\|\mathbf{x} - \hat{\mathbf{x}}\| / \|\mathbf{x}\|$  را نشان می دهد. بنابراین، عدد شرط معیاری است که به توسط آن می توان امیدوار بود که بتوان به کمک آن يك جواب (تقریبی) «خوب» را از يك جواب تقریبی «بد» از راه نگاه کردن به باقیمانده به خوبی تمیز داد.

واضح است که محاسبه عدد شرط برای يك ماتریس داده شده کاملاً دشوار است، حتی اگر بتوان نرم ماتریس را نسبتاً آسان محاسبه کرد، زیرا (برای محاسبه  $\text{cond}(A)$ ) می باید  $A^{-1}$  را در دست داشت. گاهی اوقات عدد شرط  $\text{cond}(A)$  را می توان به کمک قضیه زیر برآورد کرد، که این امر به شرح بیشتر معنی عدد شرط کمک می کند.

قضیه ۸.۴ برای هر ماتریس وارون پذیر  $A$  از مرتبه  $n \times n$  و هر نرم ماتریسی، عدد شرط  $A$  نشانگر فاصله نسبی  $A$  از نزدیکترین ماتریس وارون ناپذیر  $n \times n$  است. به ویژه داریم:

$$\frac{1}{\text{cond}(A)} = \min \left\{ \frac{\|A - B\|}{\|A\|} \mid B \text{ وارون پذیر نیست} \right\}$$

اثبات کامل این قضیه خسارچ از ظرفیت این کتاب است (ولی به تمرین ۶.۴-۵ نگاه کنید) فقط نشان می دهیم که

$$\frac{1}{\text{cond}(A)} \leq \inf \left\{ \frac{\|A - B\|}{\|A\|} \mid B \text{ وارون ناپذیر است} \right\}$$

یعنی برای هر ماتریس  $n \times n$  و وارون ناپذیر  $B$  داریم

$$\frac{1}{\|A^{-1}\|} \leq \|A - B\| \quad (۴۴.۲)$$

در واقع، اگر  $B$  وارون پذیر نباشد، آنگاه طبق قضیه ۴.۲، يك برداری غیر صفر  $\mathbf{x}$  وجود دارد، به طوری که  $B\mathbf{x} = \mathbf{0}$ ، اما داریم



$$\|A-B\| \|x\| \geq \|(A-B)x\| = \|Ax-Bx\| = \|Ax\| \geq \frac{\|x\|}{\|A^{-1}\|}$$

با به کار گرفتن (۳۸.۴)، و توجه به  $x \neq 0$ ، می‌توانیم رابطهٔ بالا را بر  $\|x\| \neq 0$  تقسیم کنیم تا (۴۴.۴) به دست آید.

برهانی که هم‌اکنون آوردیم، اثبات فرع مهم زیر است

فرض: اگر  $A$  وارونپذیر و  $B$  ماتریسی باشد که

$$\|A-B\| < \frac{1}{\|A^{-1}\|}$$

آنگاه  $B$  نیز وارونپذیر است.

مثلاً در مثال ۵.۴، برای ماتریس

$$A = \begin{bmatrix} ۱۰۰۱ & ۰۰۹۹ \\ ۰۰۹۹ & ۱۰۰۱ \end{bmatrix}$$

پیدامی کنیم که  $\|A^{-1}\|_{\infty} \geq ۱/۰۰۰۲ = ۵۰۰$ ، زیرا که ماتریس  $B = \begin{bmatrix} ۱ & ۱ \\ ۱ & ۱ \end{bmatrix}$  وارونپذیر

نیست و

$$A-B = \begin{bmatrix} ۰۰۰۱ & -۰۰۰۱ \\ -۰۰۰۱ & ۰۰۰۱ \end{bmatrix}$$

دارای نرم ماکسیمم  $\|A-B\|_{\infty} = ۰۰۰۲$  می‌باشد. بنا بر این، چون  $\|A\|_{\infty} = ۲$ ، خواهیم داشت  $\text{cond}(A) \geq ۱۰۰$ . ماتریس مثلثی وارون‌ناپذیر مثال دیگری است از آنها. اگر  $A$  مثلثی باشد با توجه به قضیهٔ ۶.۴ می‌دانیم که درایه‌های قطری  $A$  غیر صفرند، و قرارداد صفر به جای هر یک از درایه‌های قطری  $A$ ، ماتریس  $A$  را وارونپذیر می‌سازد. در نتیجه اگر  $A$  مثلثی باشد، آنگاه

$$\text{cond}(A) \geq \frac{\|A\|_{\infty}}{\min_i |a_{ii}|}$$

عدد شرط در تحلیل مشکلات بعدی حل دستگاه معادلات خطی نیز نقش مهمی دارد اگر دستگاه خطی  $Ax = b$  از یک مسئلهٔ عملی نتیجه شده باشد، باید انتظار داشته باشیم که ضرایب این دستگاه در معرض خطا واقع شده باشند، خواه بدین علت که نتیجهٔ محاسبات دیگر و یا اندازه‌گیری فیزیکی هستند و یا حتی به علت خطای گرد کردن از تبدیل اعداد دهدهی به دودویی طی خواندن اعداد ناشی شده‌اند. بنا بر این موقتاً فرض می‌کنیم که سمت راست

معادله عاری از خطا باشد و در حقیقت می‌خواهیم دستگاه خطی

$$\hat{A}\hat{x} = \mathbf{b} \quad (۴۵.۴)$$

را به جای  $AX = \mathbf{b}$  حل کنیم، در اینجا  $A = \hat{A} + E$  و ماتریس  $E$  شامل خطاهای موجود در ضرایب است. حتی اگر تمامی محاسبات دقیق انجام پذیرند، باز ما فقط جواب  $\hat{x}$  را از رابطه (۴۵.۴) به جای جواب  $x$  از  $AX = \mathbf{b}$  محاسبه کرده‌ایم. اما داریم  $x = A^{-1}\mathbf{b}$  بنابراین با فرض اینکه (۴۵.۴) دارای جواب باشد، داریم

$$\mathbf{x} = A^{-1}\mathbf{b} = A^{-1}\hat{A}\hat{x} = A^{-1}(A + \hat{A} - A)\hat{x} = \hat{x} + A^{-1}(\hat{A} - A)\hat{x}$$

بنابراین با توجه به  $\hat{A} - A = -E$  داریم:

$$\mathbf{x} - \hat{x} = A^{-1}(-E)\hat{x}$$

بنابراین

$$\|\mathbf{x} - \hat{x}\| \leq \|A^{-1}\| \|E\| \|\hat{x}\| = \|A^{-1}\| \|A\| \frac{\|E\|}{\|A\|} \|\hat{x}\|$$

که نتیجه نهایی زیر را به دست می‌دهد

$$\frac{\|\mathbf{x} - \hat{x}\|}{\|\hat{x}\|} \leq \text{cond}(A) \frac{\|E\|}{\|A\|} \quad (۴۶.۴)$$

به بیان لفظی، نسبت تفاضل  $x$  و  $\hat{x}$  بر  $\|\hat{x}\|$  می‌تواند به بزرگی  $\text{cond}(A)$  ضرب‌پذیر تغییر نسبی  $\|E\|/\|A\|$  در ماتریس ضرایب باشد. اگر بدانیم که ضرایب دستگاه خطی  $AX = \mathbf{b}$  فقط تا حدود  $10^{-8}$  (نسبت به اندازه  $A$ ) دقیق‌اند و  $10^8 \approx \text{cond}(A)$ ، آنگاه محاسبه جواب با دقت نسبی بیش از  $10^{4-8}$  مفید فایده‌ای نخواهد بود.

□ مثال ۷.۴: بار دیگر دستگاه خطی (۳۰.۴) در مثال ۵.۴، را بررسی می‌کنیم. قبلاً دیدیم که برای ماتریس ضرایب داریم  $\text{cond}(A) = 100$ . طبق رابطه (۴۶.۴) یک درصد تغییر در ضرایب دستگاه می‌تواند جواب را به‌طور قابل ملاحظه‌ای تغییر دهد. در حقیقت با یک تغییر یک درصد (در جهت صحیح) دستگاه زیر را پدید می‌آوریم

$$x_1 + x_2 = 2$$

$$x_1 + x_2 = -2$$

که اصلاً جواب ندارد زیرا که ماتریس ضرایب در این حال وارون‌پذیر نخواهد بود. □

تحلیل پیشین به کمک تحلیل خطای پسرو در سنجش اثرات ناشی از خطای گرد کردن،

که در طی عملیات حذف و پسجایگذاری در دقت جواب محاسبه شده حاصل می شود، می تواند مفید باشد. در این مورد، از نشانگذاری و اصطلاحات و قراردادهای مربوط به بخش ۳۰۱، استفاده خواهیم نمود.

**قضیه ۹.۴** فرض کنید که برای دستیابی به عوامل تجزیه  $PLU$  متعلق به ماتریس  $A$  از مرتبه  $n \times n$  و از آنجا، رسیدن به حل دستگاه خطی  $AX = b$ ، از الگوریتمهای ۲.۴ و ۴.۴ استفاده کنیم، اما برای به دست آوردن عوامل حساب شده  $P$ ،  $\hat{L}$ ، و  $\hat{U}$  و جواب حساب شده  $\hat{X}$ ، به حساب ممیز شناور همراه با گرد کردن واحد  $u \leq 0.01$  متوسل شویم. در این صورت  $\hat{X}$  دقیقاً در معادله پریشده<sup>۲</sup>

$$(A + PE)\hat{X} = b \quad (47.4)$$

که در آن

$$|E| \leq u_n |P^{-1}A| + u_n(3 + u_n)|\hat{L}| |\hat{U}| \quad (48.4)$$

و

$$u_n := n(1.01)u$$

صدق خواهد کرد.

در اینجا ماتریسی را که از  $B = (b_{ij})$  با قرارداد قدرمطلق هر درایه به جای خود آن درایه، حاصل می شود با  $|B|$  نشان می دهیم، یعنی

$$|B| = (|b_{ij}|)$$

همچنین برای معرفی دو ماتریس  $B$  و  $C$  در حالتی که  $B$  و  $C$  هم مرتبه باشند و داشته باشیم

$$b_{ij} \leq c_{ij} \quad \text{به ازای جميع مقادیر } i \text{ و } j$$

قرارداد زیر را اختیار می کنیم

$$B \leq C$$

قضیه فوق بیان می کند که هر گاه  $n$  «خیلی بزرگ» نباشد و اگر اندازه  $|\hat{L}| |\hat{U}|$  در حدود اندازه  $|A|$  باشد، آنگاه می توانیم دنبال خطاهای حاصل در جواب را از راه تنظیم معادلاتی بگردیم، که مرتبه بزرگی آنها با مرتبه تغییراتی که فقط برای دادن معادلات به ما همین مجبوریم انجام دهیم، یکی باشند. به عبارت دیگر، خطای ناشی از استفاده از حساب با ممیز شناور در جواب احتساب شده، بدتر از خطایی نیست که ما در آغاز به علت اجبار در قبول گرد کردن درایه های  $A$  به صورت اعدادی با ممیز شناور، داشته ایم.

البته، اگر ماتریس  $|\hat{L}||\hat{U}|$ ، خیلی بزرگتر از  $|A|$  باشد، خطاها در ریشه محاسبه شده  $x$ ، ممکن است خیلی بزرگتر از خطاهای ناشی از تبدیل مسئله به اعداد ماشینی با ممیز شناور باشند. توجه کنید که وقتی کران حاصل در ماتریس پریش  $E$  بیش از تحملی باشد که در آن درایه‌های  $A$  دقیق شناخته می‌شوند، می‌توان ماتریس  $|\hat{L}||\hat{U}|$  را (با صرف هزینه) عملاً حساب نمود و به دقت حسابی خیلی بالایی دست یافت. اما مهمتر از همه آنکه چون مرتبه لولا ممکن است تأثیر اساسی بر اندازه  $|\hat{L}||\hat{U}|$  بگذارد، از قضیه ۹.۴ این نتیجه مهم را استخراج می‌کنیم که تدبیر لولاگزی می‌باید معطوف به کوچک نگهداشتن  $|\hat{L}||\hat{U}|$  باشد.

اکنون با استفاده از نمادگذاری و اصطلاحات بخش ۳.۱، برهان ساده قضیه ۹.۴ را ذکر می‌کنیم. ابتدا به تعریفهایی که با کاربرد الگوریتم ۲.۴ (به گونه‌ای که در ماتریس جایگشت  $P$  ذکر شد) بدون تعویض در ماتریس  $A' := P^{-1}A$  صورت می‌گیرد، می‌پردازیم (همان کاری را که در بخش ۴.۴ کردیم). بنا بر این درایه‌های جالب عوامل  $L$  و  $U$  را بر طبق (۲۳.۴) به وسیله روابط زیر محاسبه می‌کنیم

$$l_{ij} = (a'_{ij} - \sum_{k < j} l_{ik} u_{kj}) / u_{jj} \quad i > j$$

$$u_{ij} = a'_{ij} - \sum_{k < i} l_{ik} u_{kj} \quad i \leq j$$

در نتیجه بنا بر بخش ۳.۱، به‌ویژه از مقایسه (۱۲.۱) با (۱۳.۱)، درایه‌های  $\hat{l}_{ij}$  و  $\hat{u}_{ij}$  از عوامل  $\hat{L}$  و  $\hat{U}$  که بر اساس حساب ممیز شناور محاسبه شده‌اند در معادله‌های پریشده زیر صدق می‌کنند

$$\hat{l}_{i1} \hat{u}_{1j} \varepsilon^j + \hat{l}_{i2} \hat{u}_{2j} \varepsilon^{j-1} + \dots + \hat{l}_{ij} \hat{u}_{jj} \varepsilon = a'_{ij} \varepsilon^{j-1} \quad i > j$$

$$\hat{l}_{i1} \hat{u}_{1j} \varepsilon^{i-1} + \hat{l}_{i2} \hat{u}_{2j} \varepsilon^{i-2} + \dots + \hat{l}_{ii} \hat{u}_{ij} = a'_{ij} \varepsilon^{i-1} \quad i \leq j$$

در اینجا هر  $\varepsilon$  بیانگر عددی به فرم  $(1 + \delta)$ ، به‌ازای  $|\delta| \leq u$ ، گرد کردن واحد، است. برای ساده کردن این معادلات می‌بینیم مادامی که نامساوی  $0 \leq u \leq 1$  برقرار است، به‌ازای هر عددی مانند  $\varepsilon$  و برای هر  $r$ ،  $\delta$ ی به‌صورت  $(1 + r\varepsilon)$  وجود دارد، به‌طوری که

$$\varepsilon^r = 1 + r\delta$$

این رابطه نشان می‌دهد که

$$\sum_k \hat{l}_{ik} \hat{u}_{kj} - a'_{ij} = \begin{cases} a'_{ij} (j-1) \delta - \hat{l}_{i1} \hat{u}_{1j} j \delta - \dots - \hat{l}_{ij} \hat{u}_{jj} \times 1 \times \delta & i > j \\ a'_{ij} (i-1) \delta - \hat{l}_{i1} \hat{u}_{1j} (i-1) \delta - \dots - \hat{l}_{ii} \hat{u}_{ij} \times 0 \times \delta & i \leq j \end{cases}$$

و بنا بر این

$$\hat{L}\hat{U} - A' = F \quad (۴۹.۴)$$

که در آن

$$|F| \leq u_n(|A'| + |\hat{L}||\hat{U}|) \quad (۵۰.۴)$$

$$u_n = n(۱۰۱)u$$

این نشان می‌دهد که عوامل محاسبه شده  $\hat{L}$  و  $\hat{U}$  برای  $A'$ ، درست همان عوامل ماتریس پریشیده  $A' + F$  هستند، با ماتریس خطای  $F$  که مرتبه گرد کردن آن همان مرتبه گرد کردن درایه‌های  $A$  است به شرط آنکه ماتریس  $|\hat{L}||\hat{U}|$  خیلی بزرگتر از  $|A|$  نباشد. مراحل محاسبه‌ای که در الگوریتم ۴.۴، یعنی در مرحله حل کردن، مورد استفاده قرار گرفتند تا حدی شبیه، به آن مراحل هستند که در بالا بیان شد. بنابراین می‌توان نشان داد که بردار محاسبه شده  $\hat{y}$  دقیقاً در دستگاه پابین مثلثی پریشیده زیر صدق می‌کند

$$|G| \leq u_n |\hat{L}| \quad \text{با} \quad (\hat{L} + G)\hat{y} = b$$

در حالی که جواب محاسبه شده  $\hat{x}$  دقیقاً در دستگاه خطی پریشیده زیر صدق می‌کند

$$|H| \leq u_n |\hat{U}| \quad \text{با} \quad (\hat{U} + H)\hat{x} = \hat{y}$$

از آنجا نتیجه می‌شود که جواب محاسبه شده در رابطه زیر صدق می‌کند

$$(\hat{L} + G)(\hat{U} + H)\hat{x} = b$$

ولی داریم

$$\begin{aligned} (\hat{L} + G)(\hat{U} + H) &= \hat{L}\hat{U} + \hat{G}\hat{U} + \hat{L}H + GH \\ &= A' + F + \hat{G}\hat{U} + \hat{L}H + GH \\ &= A' + E \end{aligned}$$

که

$$\begin{aligned} |E| &\leq |F| + |G||\hat{U}| + |\hat{L}||H| + |G||H| \\ &\leq u_n(|A'| + |\hat{L}||\hat{U}|) + u_n|\hat{L}||\hat{U}|(۱ + ۱ + u_n) \end{aligned}$$

که قضیه را ثابت می‌کند.

کران رابطه (۴۸.۴) محفوظ نگاه داشته شده است. اگر از لولاکزینی جزئی‌سی استفاده شود، آنگاه کران

$$|E| \leq nu|P^{-1}A| \quad (الف ۵۰.۴)$$

غالباً واقفتر ایانه تر خواهد شد. در هر حال، این گونه کرانها دید روشنی از نتیجهٔ دقتی که در محاسبات جواب حساب شده صورت گرفته، به ما می دهند، زیرا که، برای مثال، با توجه به (۴۶.۲) و (۵۰.۴) درمی یابیم که خطای جواب محاسبه شده نسبت به اندازهٔ این جواب معمولاً در محدودهٔ زیر واقع است:

$$\frac{\|\mathbf{x} - \hat{\mathbf{x}}\|}{\|\hat{\mathbf{x}}\|} \leq \text{cond}(A) \cdot n \cdot u \quad (51.4)$$

اگر دقت را کاملاً کنار بگذاریم، دستگاه خطی  $A\mathbf{x} = \mathbf{b}$  غالباً بد شرط نامیده می شود اگر  $\text{cnod}(A)$  «بزرگ» باشد. تا حدی نزدیکتر به واقعیت، دستگاه خطی را نسبت به دقت عمل مصروفه «بد شرط» گویند اگر  $\text{cond}(A)$  تقریباً  $1/u$  باشد، زیرا بنا بر (۵۱.۴) يك جواب محاسبه شده ممکن است اصلاً شباهتی به جواب (دقیق) نداشته باشد.

□ مثال ۸۰۴: دستگاه خطی زیر را در نظر می گیریم.

$$0.24x_1 + 0.36x_2 + 0.12x_3 = 0.84$$

$$0.12x_1 + 0.16x_2 + 0.24x_3 = 0.52 \quad (52.2)$$

$$0.15x_1 + 0.21x_2 + 0.25x_3 = 0.64$$

سعی می کنیم این دستگاه را با استفاده از الگوریتم حذفی ۲.۲، حساب ممیز شناور با دو رقم اعشاری و لولا گزینی جزئی مدرج، حل کنیم. ترتیب گزینش لولا نتیجه می دهد که داشته باشیم  $\mathbf{p}^T = [1 \ 2 \ 3]$  و محتوای نهایی آرایه کاری به شرح زیر است

$$\begin{bmatrix} 0.24 & 0.36 & 0.12 & 0.84 \\ 0.50 & -0.02 & 0.18 & 0.10 \\ 0.63 & 1.0 & -0.01 & 0.01 \end{bmatrix}$$

اگر محاسبات را ادامه دهیم و از راه پس جایگذاری به جواب تقریبی  $\hat{\mathbf{x}} = \begin{bmatrix} 25 \\ -14 \\ -1 \end{bmatrix}$

باقیمانده  $\mathbf{r} = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.08 \end{bmatrix}$  می رسمیم. زیرا  $\mathbf{x} = \begin{bmatrix} -3 \\ 4 \\ 1 \end{bmatrix}$  جواب دستگاه است، لذا جواب

محاسبه شده  $\hat{\mathbf{x}}$  در اولین رقم با معنی اشتباه است.

نرم ما کسیم برای ماتریس ضرایب  $A$ ی مربوط به این دستگاه عبارت است از

$$\|A\|_{\infty} = ۰.۷۷۲. \text{ به علاوه، ماتریس}$$

$$B = \begin{bmatrix} ۰.۲۵۲ & ۰.۳۶ & ۰.۱۲ \\ ۰.۱۱۲ & ۰.۱۶ & ۰.۲۴ \\ ۰.۱۴۷ & ۰.۲۱ & ۰.۲۵ \end{bmatrix}$$

وارون ناپذیر است (ستون اول از حاصلضرب  $۰.۷$  در ستون دوم به دست آمده است) در حالی که  $\|A - B\|_{\infty} = ۰.۰۱۲$ . بنابراین از قضیه  $۸.۴$  نتیجه می گیریم که

$$\text{cond}(A) \geq \frac{۰.۷۷۲}{۰.۰۱۲} \geq ۶۰$$

بنابراین دستگاه فوق، نسبت به دقت عملی که به کار رفته، بسیار بد شرط است و خطای بسیار بزرگ موجود در جواب محاسبه شده شگفت آور نیست.

سپس، محاسبات را با استفاده از حساب ممیز شناور با سه رقم اعشاری تکرار می کنیم. از آنجا که  $\text{cond}(A) \approx ۶۰$ ، باز انتظار نداریم که جواب حساب شده خیلی دقیق باشد. به موجب الگوریتم  $۲.۴$ ، محتوای ماتریس کار به شرح زیر است

$$\begin{bmatrix} ۰.۲۴ & ۰.۳۶ & ۰.۱۲ & ۰.۸۴ \\ ۰.۵۰ & -۰.۰۲ & ۰.۱۸ & ۰.۱۰ \\ ۰.۶۲۵ & ۰.۷۵ & ۰.۰۴ & ۰.۰۴ \end{bmatrix} \quad (۵۳.۴)$$

و روش پسجایگذاری جواب محاسبه شده  $\hat{\mathbf{x}}^T = [-۳ \quad ۴ \quad ۱]$  را به دست می دهد؛ یعنی، حتی با اینکه دستگاه هنوز نسبت به دقت به کار گرفته شده اندکی بد شرط است، اما جواب

$$\mathbf{b} = \begin{bmatrix} ۰.۸۵۲ \\ ۰.۶۲۰ \\ ۰.۷۴۰ \end{bmatrix} \text{ (دقیق) را می دهد. این امر از تبدیل سمت راست (۵۲.۴) به}$$

می شود. بنا به کار گرفتن تجزیه به عوامل (۵۳.۴) و بنا توجه به الگوریتم (۲.۴)

$$\hat{\mathbf{x}} = \begin{bmatrix} -۳.۳۰ \\ ۴.۰۵ \\ ۱.۵۳ \end{bmatrix} \text{ (تقریبی) جواب را محاسبه می کنیم (باز از حساب با ممیز شناور، با سه رقم}$$

$$\mathbf{r} = \begin{bmatrix} ۰.۰۰۰۲۴ \\ ۰.۰۰۰۰۸ \\ ۰.۰۰۰۲۰ \end{bmatrix} \text{ اعشاری، استفاده می کنیم) که باقیمانده آن است. جواب دقیق}$$

برابر است با  $\mathbf{x} = \begin{bmatrix} -۳۶ \\ ۴۲۵ \\ ۱۵۵ \end{bmatrix}$ ، و جواب محاسبه شده دارای ۱۵ درصد خطاست که با  $\square$  (۵۱.۴) سازگار است.

چنانچه، مثال فوق نشان می‌دهد، يك عدد با شرط بزرگ نسبت به دقتی که به کار رفته ممکن است به خطای نسبتاً بزرگی در جواب محاسبه شده منجر شود، اما همیشه الزاماً چنین نیست.

این مطلب را که آیا يك دستگاه خطی نسبت به دقتی که به کار گرفته شده بدشرط است یا نیست، می‌توان به سادگی [حتی بدون اطلاع از  $\text{cond}(A)$ ] در طی اصلاح بارستی که در زیر از آن بحث می‌شود، تعیین کرد. با توجه به  $\hat{\mathbf{x}}^{(1)}$ ، خطا (مجهول) در جواب تقریبی  $\hat{\mathbf{x}}^{(1)}$  برای  $A\mathbf{x} = \mathbf{b}$  را در بخش ۵.۴ پیدا کردیم و دیدیم که

$$A\mathbf{e} = \mathbf{r} \quad (۵۴.۲)$$

که در آن  $\mathbf{r} = \mathbf{b} - A\hat{\mathbf{x}}^{(1)}$  باقیمانده قابل محاسبه‌ای برای  $\hat{\mathbf{x}}^{(1)}$  است. پس در اینجا يك دستگاه خطی داریم که جواب آن خطای  $\mathbf{e}$  و ماتریس ضرایب آن با ماتریس ضرایب دستگاه اولیه یکی است. اگر  $\hat{\mathbf{x}}^{(1)}$  از راه الگوریتم حذف ۲.۴ به دست آمده باشد، می‌توانیم (۵۴.۴) را با الگوریتم جایگذاری ۴.۴ نسبتاً سریع حل کنیم. گیریم  $\hat{\mathbf{e}}^{(1)}$  يك جواب (تقریبی) برای (۵۴.۴) باشد که بدین طریق محاسبه شده است. در این صورت  $\hat{\mathbf{e}}^{(1)}$  در حالت کلی با  $\mathbf{e}$  یکی نخواهد بود. اما  $\hat{\mathbf{e}}^{(1)}$ ، حداقل هیچ نباشد، می‌بایستی نشانه‌ای از اندازه  $\mathbf{e}$  به دست دهد. اگر  $10^{-s} \approx \|\hat{\mathbf{e}}^{(1)}\| / \|\hat{\mathbf{x}}^{(1)}\|$ ، نتیجه می‌گیریم که احتمالاً اولین  $s$  رقم اعشاری  $\hat{\mathbf{x}}^{(1)}$  و اولین  $s$  رقم اعشاری  $\mathbf{x}$  با هم یکی هستند. پس، انتظار داریم که تقریب  $\hat{\mathbf{e}}^{(1)}$  در همان اندازه دقت  $\mathbf{e}$  باشد. بنابراین انتظار داریم که

$$\hat{\mathbf{x}}^{(2)} = \hat{\mathbf{x}}^{(1)} + \hat{\mathbf{e}}^{(1)}$$

تقریبی بهتر از  $\hat{\mathbf{x}}^{(1)}$  برای  $\mathbf{x}$  باشد. اکنون می‌توانیم، در صورت لزوم، باقیمانده جدید  $\mathbf{r} = \mathbf{b} - A\hat{\mathbf{x}}^{(2)}$  را حساب و (۵۴.۴) را دوباره حل کنیم، و تصحیح جدید  $\hat{\mathbf{e}}^{(2)}$  و تقریب جدید  $\hat{\mathbf{x}}^{(2)} = \hat{\mathbf{x}}^{(1)} + \hat{\mathbf{e}}^{(2)}$  را برای  $\mathbf{x}$  به دست آوریم. تعداد ارقام اعشاری یکسان در تقریبهای متوالی  $\hat{\mathbf{x}}^{(1)}$ ،  $\hat{\mathbf{x}}^{(2)}$ ، ... و همچنین بررسی باقیمانده‌های متوالی بایستی قرینه‌ای از صحت این جواب‌های تقریبی به دست دهند. معمولاً این بارستها تا زمان برقراری شرط  $10^{-s} \approx \|\hat{\mathbf{e}}^{(k)}\| / \|\hat{\mathbf{x}}^{(k)}\|$  که  $s$  تعداد ارقام اعشاری است که در محاسبات به کار گرفته شده، ادامه می‌یابد. می‌توان نشان داد که تعداد مراحل بارستی لازم برای حصول این نتیجه با  $\text{cond}(A)$  افزایش می‌یابد. هر گاه  $\text{cond}(A)$  «خیلی بزرگ» باشد تصحیحهای  $\hat{\mathbf{e}}^{(1)}$

## 1. iterative improvement



$\hat{e}^{(2)}$  و ... ممکن است هرگز کوچک نشوند، و لذا معرف حد اعلاى بدشرطى در دستگاه اوليه باشند.

برای توفیق در روش اصلاح بارستی، شرط اساسی این است که باقیمانده با دقت هرچه بیشتر محاسبه شود. اگر، چنان که متداول است، حساب با ممیز شناور مورد استفاده قرار گیرد، می بایستی باقیمانده از راه حساب با دقت مضاعف، محاسبه شود.

#### الگوریتم ۵.۴ اصلاح بارستی

دستگاه خطی  $Ax = b$  و جواب تقریبی  $\hat{x}$  داده شده اند.

با استفاده از حساب با دقت مضاعف،  $r = b - A\hat{x}$  را محاسبه کنید.

با استفاده از الگوریتم ۲.۴ (یا در صورت امکان فقط از الگوریتم ۴.۴) جواب

تقریبی  $\hat{e}$  از دستگاه خطی  $Ae = r$  را محاسبه کنید.

اگر  $\|\hat{e}\| / \|\hat{x}\|$  «به اندازه کافی کوچک» باشد محاسبات را متوقف کنید و  $\hat{x} + \hat{e}$  را

به عنوان جواب بپذیرید. در غیر این صورت قرار دهید  $\hat{x} := \hat{x} + \hat{e}$  و عملیات بالا را تکرار کنید.

هنگامی می توان از اصلاح بارستی استفاده کرد که يك جواب تقریبی به هر طریق پیدا شده باشد. اصلاح بارستی را همیشه بایستی پس از محاسبه يك جواب تقریبی از راه حذف به کار گرفت زیرا که تصحیحاها را می توان از راه پیشجا یگذاری و پسجا یگذاری نسبتاً ارزان محاسبه کرد. همچنین میزان همگرایی عملیات (در صورت همگرا بودن) نشانه خوبی از شرط دستگاه (نسبت به دقت عمل به کار گرفته شده) است.

□ مثال ۹.۴: برای جواب تقریبی دستگاه (۵۲.۴) که در مثال ۸.۴ محاسبه شد، روش اصلاح بارستی را به کار می گیریم. باقیمانده محاسبه شده، که تا دو رقم اعشاری صحیح و

گرد شده، برابر است با  $r = \begin{bmatrix} 0.0 \\ 0.0 \\ 0.08 \end{bmatrix}$ . با اعمال الگوریتم ۴.۴ بر این سمت راست (یا

استفاده از حساب با ممیز شناور با دو رقم اعشاری) تصحیح  $\hat{e}^{(1)} = \begin{bmatrix} 120 \\ -75 \\ -8 \end{bmatrix}$  به دست

می آید که اندازه آن با اندازه جواب محاسبه شده یکی است. بنا بر این نتیجه می گیریم که دستگاه خطی داده شده نسبت به دقتی که به کار رفته بسیار بدشرط است و دقت عمل بیشتری باید به کار گرفت تا بتوان جواب (۵۲.۴) را محاسبه کرد.

ماتریسها و دستگاههای معادلات خطی ۲۳۷

در مثال ۸.۴ يك جواب تقریبی  $\hat{\mathbf{x}} = \begin{bmatrix} -۳۳۰ \\ +۴۰۵ \\ ۱۵۳ \end{bmatrix}$  برای دستگاه خطی با همان

ماتریس ضرایب، اما با سمت راست متفاوت را با استفاده از حساب با ممیز شناور با سه

رقم اعشاری محاسبه کردیم. باقیمانده محاسبه شده صحیح چنین است  $\mathbf{r} = \begin{bmatrix} ۰.۰۰۰۲۴ \\ ۰.۰۰۰۰۸ \\ ۰.۰۰۰۲۰ \end{bmatrix}$

یا کاربرد الگوریتم ۴.۴ برای این سمت راست  $\mathbf{r}$  (و با به کار بستن همان دقت عمل قبلی)

تصحیح  $\hat{\mathbf{x}}^{(۱)} = \begin{bmatrix} -۰.۰۳ \\ ۰.۰۲ \\ ۰.۰۰۲ \end{bmatrix}$  نتیجه می شود که فقط ۱۰٪ جواب محاسبه شده است و جواب

تصحیح شده  $\hat{\mathbf{x}}^{(۲)} = \begin{bmatrix} -۳.۶ \\ ۴.۲۵ \\ ۱.۵۵ \end{bmatrix}$  را نتیجه می دهد. باقیمانده برای این جواب تقریبی ۰

درمی آید، لذا درست يك مرحله از روش اصلاح بارستی، جواب (دقیق) این مثال را به دست می دهد. □

## تمرین

۱-۶.۴ با استفاده از قضیه ۸.۴ عدد شرط ماتریس زیر را برآورد کنید.

$$A = \begin{bmatrix} ۷ & ۸ & ۹ \\ ۸ & ۹ & ۱۰ \\ ۹ & ۱۰ & ۸ \end{bmatrix}$$

۲-۶.۴ از الگوریتم اصلاح بارستی در جواب محاسبه شده تمرین ۳.۴ استفاده کنید.

۳-۶.۴ ماتریس  $A$  از مرتبه  $n$  را نافذ قطری (سطراً مؤکند) گوئیم اگر به ازای

$i = 1, \dots, n$  نامساوی  $|a_{ii}| > \sum_{j \neq i} |a_{ij}|$  برقرار باشد. با استفاده از فرع قضیه ۸.۴

ثابت کنید که يك ماتریس «نافذ قطری» وارونپذیر است (اذهمایی: فرض کنید  $A = DB$

و  $D$  يك ماتریس قطری باشد که درایه های قطری اش با درایه های قطری ماتریس  $A$  مساوی

است، سپس نشان دهید که  $(\|I - B\|_\infty < 1)$

## 1. diagonally dominant

۴-۶۰۴ عدد شرط ماتریس تمرین ۶-۶۰۲ را از طریق حل دستگاه خطی  $AX = b$  با (اولاً)  $b^T = [۲۴, ۲۷, ۲۷]$ ، (ثانیاً)  $b^T = [۲۴, ۲۶, ۲۶]$  برآورد کنید. از روش اصلاح بارستی استفاده کنید.

۵-۶۰۴ نشان دهید که در نرم ماتریس خاص (۳۹-۴)، برای برقراری رابطه (۴۴-۴) ماتریس وارون ناپذیر  $B$  را می‌توان به گونه‌ی زیر بنا کرد: به موجب رابطه (۳۵-۴)، می‌توان  $X$  با نرم ۱ چنان تعیین نمود که تساوی  $\|A^{-1}X\| = \|A^{-1}\| \|X\|$  برقرار باشد. حال  $B$  را برابر با ماتریس  $A - XX^T$ ، که در آن  $z = y_m^{-1} \mathbf{i}_m$ ،  $y = A^{-1}X$ ، بگیریید و  $m$  را چنان انتخاب کنید که  $\|y\|_\infty = |y_m|$ .

۶-۶۰۴ نشان دهید که مسئله ۶-۶۰۲ را برای یک نرم کلی می‌توان حل کرد، به شرط آنکه بدانیم که به ازای یک  $n$ -برداری غیر صفر مفروض  $y$ ، چگونه یک  $n$ -برداری  $z$  را انتخاب کنیم، تا به ازای جمیع  $n$ -برداریهای  $u$ ، داشته باشیم  $\|u\| \leq z^T u$ ، تساوی زمانی برقرار است که  $u = y$ . درحالی‌که نرم از نوع ۱-نرم باشد،  $z$  را چگونه انتخاب می‌کند؟

### ۷.۴۰ دترمینانها

گرچه فرض بر این است که دانشجویان با مفهوم دترمینانها آشنایی دارند، ولی ما این بخش را به تعریف رسمی دترمینانها اختصاص می‌دهیم و بعضی از خواص اولیه آنها را مورد بحث قرار می‌دهیم.

به هر ماتریس مربعی از اعداد مانند  $A$ ، عددی نسبت داده می‌شود به نام دترمینان  $A$  که با  $\det(A)$  نشان می‌دهند. اگر  $A = (a_{ij})$  یک ماتریس  $n \times n$  باشد، آنگاه دترمینان  $A$  چنین تعریف می‌شود

$$\det(A) = \sum_p \sigma_p a_{۱,p_1} a_{۲,p_2} \cdots a_{n,p_n} \quad (۵۵.۴)$$

در رابطه بالا مجموع روی تمامی  $n!$  جایگشت  $p$  از مرتبه  $n$  در نظر گرفته می‌شود و  $\sigma_p$  مساوی  $+۱$  یا  $-۱$  است بسته به اینکه  $p$  زوج یا فرد باشد (به بخش ۱.۴ نگاه کنید). بنابراین اگر  $n=۱$ ، آنگاه داریم

$$\det(A) = \det[a_{۱,۱}] = a_{۱,۱}$$

درحالی‌که اگر  $n=۲$ ، داریم

$$\det(A) = \det \begin{bmatrix} a_{۱,۱} & a_{۱,۲} \\ a_{۲,۱} & a_{۲,۲} \end{bmatrix} = a_{۱,۱} a_{۲,۲} - a_{۱,۲} a_{۲,۱} \quad (۵۶.۴)$$

به ازای  $n=۳$ ، شش حاصلضرب باید با هم جمع شوند، و به ازای  $n=۱۰$  بیشتر از سه میلیون حاصلضرب، هر یک با ۱۰ عامل، باید محاسبه و با هم جمع شوند تا سمت راست

(۵۵.۲) به دست آید. بنابراین، تعریف (۵۵.۲) برای محاسبهٔ دترمینانها چندان مفید نیست. اما در زیر به ذکر فهرستی از قواعد مربوط به دترمینانها، که به آسانی از تعریف ۵۵.۲ به دست می آید، می پردازیم. سپس نشان می دهیم که چگونه می توان با به کارگیری این قواعد و الگوریتم ۲.۴، مقدار دترمینان را حدوداً با  $O(n^3)$  [به جای  $O(n!)$ ] عمل محاسبه کرد. دترمینان يك ماتریس به دلیل قضیهٔ زیر حائز اهمیت است.

**قضیه ۱۰.۴** گیریم  $A$  يك ماتریس  $n \times n$  باشد در این صورت  $A$  وارونپذیر است اگر و فقط اگر،  $\det(A) \neq 0$ .

ما از این قضیه در بخش بعدی که به محاسبهٔ ویژه مقادارها و ویژه بردارها مربوط است استفاده خواهیم کرد.

در برخی از ماتریسها، دترمینان به آسانی قابل محاسبه است. قاعدهٔ ۱ اگر  $A = (a_{ij})$  يك ماتریس بالا(پایین) مثلثی باشد، آنگاه

$$\det(A) = a_{11} a_{22} \cdots a_{nn}$$

یعنی دترمینان درست برابر است با حاصلضرب درایه های قطری  $A$ . زیرا مثلاً، اگر  $A$  ماتریس بالامثلثی و  $\mathbf{p}$  هر جایگشتی غیر از جایگشت همانی<sup>۱</sup> باشد، به ازای مقداری از  $i$  باید داشته باشیم  $p_i < i$ ، و بنابراین حاصلضرب مربوطه  $a_{1,p_1} a_{2,p_2} \cdots a_{n,p_n}$  باید متضمن درایه های زیر قطری صفر  $a_{i,p_i}$  و در نتیجه، باید صفر باشد. بنابراین اگر  $A$  يك ماتریس بالامثلثی باشد، آنگاه تنها عامل جمع موجود در (۵۵.۲) که صفر بودن آن رانمی توان تضمین کرد، جملهٔ  $a_{11} a_{22} \cdots a_{nn}$  متناظر با جایگشت همانی (زوج)  $\mathbf{p}^T = [1 \ 2 \ \cdots \ n]$  است. به ویژه

$$\det(I) = 1 \quad (57.4)$$

به طور مشابه می توان يك قاعدهٔ دیگر را ثابت کرد. قاعدهٔ ۲. اگر  $P$  ماتریس جایگشتی  $n \times n$  باشد که با جایگشت

$$P\mathbf{i}_j = \mathbf{i}_{p_j} \quad j = 1, \dots, n$$

داده شده است، آنگاه با جایگشتی مانند  $\mathbf{p}$  داریم

$$\det(P) = \begin{cases} 1 & \text{اگر } \mathbf{p} \text{ زوج باشد} \\ -1 & \text{اگر } \mathbf{p} \text{ فرد باشد} \end{cases}$$

قاعده ۳. اگر ماتریس  $B$  از تعویض دو ستون (سطر)  $A$  نتیجه شده باشد، آنگاه

$$\det(B) = -\det(A)$$

□ مثال:

$$\square \quad \det \begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix} = 3 - 4 = -1 \quad , \quad \det \begin{bmatrix} 2 & 1 \\ 3 & 2 \end{bmatrix} = 4 - 3 = 1$$

در نتیجه، هر گاه دو ستون (سطر) ماتریس  $A$  یکی باشند (به طوری که تعویض آنها با هم در  $A$  تغییری به وجود نیاورد)، آنگاه

$$\det(A) = 0$$

قاعده ۴. اگر ماتریس  $B$  از ضرب تمامی درایه‌های یک ستون (سطر) ماتریس  $A$  در عدد  $\alpha$  به دست آمده باشد، آنگاه

$$\det(B) = \alpha \det(A)$$

□ مثال:

$$\square \quad \det \begin{bmatrix} 3 \times 1 & 4 \\ 3 \times 2 & 3 \end{bmatrix} = 9 - 24 = -15 = 3(-5) = 3 \det \begin{bmatrix} 1 & 4 \\ 2 & 3 \end{bmatrix}$$

قاعده ۵. فرض کنید که سه ماتریس  $A_1$  و  $A_2$  و  $A_3$  از مرتبه  $n \times n$  فقط در یک ستون (سطر) مثلا ستون (سطر)  $j$ ام، با هم تفاوت داشته باشند و ستون (سطر)  $j$ ام  $A_3$  حاصل جمع برداری ستون (سطر)  $j$ ام  $A_1$  و ستون (سطر)  $j$ ام  $A_2$  باشد. در این صورت داریم

$$\det(A_1) + \det(A_2) = \det(A_3)$$

□ مثال:

$$\det \begin{bmatrix} 1 & 2 \\ 2 & 2 \end{bmatrix} + \det \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix} = (2 - 4) + (1 - 2) = 3 - 6 = -3 = \det \begin{bmatrix} 1 & 2 \\ 3 & 3 \end{bmatrix}$$

□

قواعد ۱ تا ۵ قضیه‌های ۱۱.۲ و ۱۲.۴ را ایجاب می‌کنند.

قضیه ۱۱.۴ اگر  $A$  و  $B$  ماتریسهای  $n \times n$  باشند، آنگاه

$$\det(AB) = \det(A) \det(B)$$

قضیه ۱۲.۴ هر گاه  $A$  یک ماتریس  $n \times n$  و  $x = (x_i)$  و  $b$  برداری  $n$  باشد به طوری که

$$Ax = b$$

آنگاه به ازای  $j = 1, \dots, n$  داریم

$$\det(A^{(j)}) = x_j \det(A) \quad (58.4)$$

که  $A^{(j)}$  ماتریسی است که از قرار دادن  $b$  به جای ستون  $j$  از  $A$  به دست آمده است. اگر  $A$  وارونپذیر باشد، یعنی اگر  $\det(A) \neq 0$  (بنابر قضیه ۱۰.۴)  $\det(A) \neq 0$ ، آنگاه می توان  $(58.4)$  را نسبت به  $x_j$  حل کرده، به دست آورد:

$$x_j = \frac{\det(A^{(j)})}{\det(A)} \quad j = 1, \dots, n$$

رابطه بالا قاعده کرامر برای درایه های جواب  $x$  از دستگاه خطی  $Ax = b$  است. به دلیل دشواری محاسبه دترمینانها، در حالت کلی قاعده کرامر تنها از لحاظ نظری اهمیت دارد. در حقیقت سریعترین راهی که برای محاسبه  $\det(A)$  یک ماتریس دلخواه  $A$  از مرتبه  $n \times n$  در دست است، استفاده از الگوریتم حذف ۲.۴ برای ماتریس فوق (نادیده گرفتن سمت راست آن) است. در قسمت ۴.۴ دیدیم که این الگوریتم ماتریس  $A$  را به صورت

$$A = PLU$$

تجزیه می کند که  $P$  ماتریسی است جایگشتی که به توسط تدبیر لولاگزینی  $P$  معین می شود،  $L$  یک ماتریس پایین مثلثی است که تمامی درایه های قطری اش مساوی ۱ است، و بالاخره ماتریس بالامثلثی ضرایب  $U = (u_{ij})$  است که تمامی درایه های لولا را روی قطر خود دارد. بنابر قاعده ۱،  $\det(L) = 1$ ، در حالی که بنابر قاعده ۲  $\det(P)$  بسته به اینکه  $P$  زوج یا فرد باشد، یعنی بسته به اینکه تعداد تعویضها در طی حذف زوج یا فرد باشد برابر ۱ یا  $-1$  است. سرانجام باز بنابر قاعده ۱ داریم  $\det(U) = u_{11}u_{22}\dots u_{nn}$ . لذا

$$\det(A) = (-1)^i u_{11}u_{22}\dots u_{nn} \quad (59.4)$$

که  $i$  تعداد تعویضها طی اجرای الگوریتم حذف است. باید توجه داشت که برنامه فورتن FACTOR (در حالتی که وارونپذیری  $A$  معلوم باشد) عدد  $(-1)^i$  را به IFLAG بدل می کند، و از این رو محاسبه  $\det(A)$  را از طریق (59.4)، از روی درایه های قطری آرایه کار  $W$  آسان می سازد.

البته الگوریتم حذف ۲.۴ (حداقل از لحاظ نظری) موقعی نتیجه بخش خواهد بود که  $A$  وارونپذیر باشد. اما اگر  $A$  وارونپذیر نباشد، الگوریتم فوق آن را نشان خواهد داد، در این حال بنابر قضیه ۱۰.۴ خواهیم داشت  $\det(A) = 0$ .

و بالاخره، چنانچه ماتریس  $A$  دارای خصوصیات ویژه‌ای باشد، گاهی استفاده از قاعده زیر مفید واقع می‌شود.

قاعده ۶: بسط دترمینان به توسط کهادها بنا بر تعریف، کهاده  $M_{ij}$  در یک ماتریس  $A = (a_{ij})$  از مرتبه  $n \times n$  عبارت است از دترمینان ماتریس مرتبه  $n-1$  امی، که از حذف سطر  $i$  ام و ستون  $j$  ام ماتریس  $A$  نتیجه می‌شود. به ازای هر  $i$  داریم

$$\det(A) = a_{i1}(-1)^{i+1}M_{i1} + a_{i2}(-1)^{i+2}M_{i2} + \dots + a_{in}(-1)^{i+n}M_{in}$$

و

به ازای هر  $j$  داریم

$$\det(A) = a_{1j}(-1)^{1+j}M_{1j} + a_{2j}(-1)^{2+j}M_{2j} + \dots + a_{nj}(-1)^{n+j}M_{nj}$$

قاعده ۶ به ما این امکان را می‌دهد که یک دترمینان مرتبه  $n$  را به صورت مجموع دترمینانهای از مرتبه  $n-1$  بیان کنیم. با استفاده تراجعی<sup>۲</sup> از این قاعده، می‌توانیم سرانجام  $\det(A)$  را به صورت مجموع دترمینانهای از مرتبه  $1$  بنویسیم. این قاعده به ویژه هنگامی در محاسبه  $\det(A)$  مفید است که  $A$  یک ماتریس تنگ باشد، که در این حالت بیشتر حاصلجمعها از دور خارج می‌شوند. به عنوان مثال، از بسط ماتریس زیر بر حسب کهادها نسبت به سطر اول داریم:

$$\begin{aligned} \det \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} &= 0 \det \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} - 1 \det \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + 0 \det \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \\ &= -1 \det \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = -1 \times 1 = -1 \end{aligned}$$

### تمرین

۱-۷۰۴ با استفاده از قضیه ۱۱.۴ و معادله (۵۷.۴) ثابت کنید که اگر  $A$  وارونپذیر باشد، آنگاه  $\det(A) \neq 0$ .

۲-۷۰۴ با استفاده از قضایای ۱۲.۴ و ۲.۴ ثابت کنید که اگر  $\det(A) \neq 0$ ، آنگاه  $A$  وارونپذیر است.

۳-۷۰۴ تعداد عملیات حسابی لازم برای محاسبه جواب یک دستگاه خطی مرتبه ۲ را در

روشهای زیر تعیین کنید. (الف) در روش حذف و پسجایگذاری، (ب) در روش کرامر.

۴-۷.۴ اگر  $n = 3$ ؛ آنگاه محاسبهٔ مستقیم (۵۵.۴) مستلزم ۱۲ عمل ضرب و ۵ عمل جمع خواهد بود. تعداد عملیات ضرب و جمع لازم برای محاسبهٔ يك دترمینان مرتبهٔ ۳ برحسب بسط به کهادها (قاعدهٔ ۶) چقدر خواهد بود؟ تعداد عملیات ضرب، یا تقسیم و جمع برای همان دترمینان با روش حذف چقدر است؟

۵-۷.۴ ثابت کنید: اگر ماتریس ضرایب دستگاه خطی  $Ax = b$  وارونپذیر باشد، همیشه می‌توانیم (در صورت لزوم) معادلات را طوری مرتب کنیم که تمامی درایه‌های قطری ماتریس ضرایب دستگاه (هم ارز) حاصل، غیر صفر باشند. [دانهمایی: اگر  $A$  وارونپذیر باشد، بنا بر قضیهٔ ۱۵.۴ حداقل یکی از عوامل جمع در (۵۵.۴) باید مخالف صفر باشد].

۶-۷.۴ درستی قواعد ۱ تا ۵ را در حالتی که ماتریسهای مورد بحث از مرتبهٔ ۲ باشند، تحقیق نمایید. سعی کنید قواعد ۴ و ۵ را برای ماتریسهای از هر مرتبهٔ دلخواه ثابت کنید.

۷-۷.۴ قضیهٔ ۱۱.۴ را برای حالتی که  $A$  و  $B$  ماتریسهای مرتبهٔ ۲ باشند، ثابت کنید.

۸-۷.۴ گیریم  $A$  يك ماتریس سه قطری از مرتبهٔ  $n$  باشد. به ازای  $n, \dots, 2, 1, p$  گیریم  $A_p$  يك ماتریس  $p \times p$ ، حاصل از حذف سطورهایی  $1, \dots, p, \dots, n$  و ستونهای  $1, \dots, p, \dots, n$  ماتریس  $A$  باشد. با استفاده از قاعدهٔ ۶ ثابت کنید که هرگاه  $\det(A_0) = 1$  آنگاه داریم

$$\det(A_p) = a_{pp} \det(A_{p-1}) - a_{p, p-1} a_{p-1, p} \det(A_{p-2}), \quad p = 2, 3, \dots, n$$

برنامه‌ای بر اساس این فرمول تراجعی برای محاسبهٔ دترمینان ماتریس سه قطری بنویسید.

### ۸.۴\* مسئله ویژه مقدارها

در بسیاری از مسائل فیزیکی، ویژه مقدارها اهمیت زیادی دارند. برای مثال، پایداری يك هواپیما با جای ویژه مقدارهای يك ماتریس در صفحهٔ هماتافت<sup>۱</sup> مشخص می‌گردد. بسامد<sup>۲</sup> طبیعی ارتعاشات يك باریکهٔ نوری، در واقع ویژه مقدارهای يك ماتریس (متناهی) هستند. ویژه مقدارها به‌طور طبیعی در تحلیل بسیاری از مسائل ریاضی نیز پدید می‌آیند، زیرا این کمیات به‌ویژه وسیلهٔ بسیار مناسب و روشنگری در نمایش يك ماتریس (در شکل قانونی جردن و در شکلهای متشابه) هستند. به‌همین دلیل، هر دستگاه معادلات دیفرانسیل خطی معمولی مرتبهٔ اول با ضرایب ثابت، می‌تواند برحسب ویژه مقادیر ماتریس ضرایبش حل شود، و همچنین رفتار دنبالهٔ توانهای  $A, A^2, A^3, \dots$  در يك ماتریس  $A$ ، برحسب ویژه مقدارهای آن، خیلی آسانتر مورد تحلیل قرار می‌گیرد. این گونه دنباله‌ها در حل بارستی دستگاه



معادلات خطی (و غیر خطی) پدید می آیند.

بدین دلایل و دلایل دیگر در این بخش مقدمه کوتاهی در جایابی<sup>۱</sup> و محاسبه ویژه مقدارها می آوریم. متأسفانه، کیفیت فنی بیان آن خارج از محدوده این کتاب است. کتاب دائرةالمعارفی ج. ۵. ویلکینسن<sup>۲</sup> [۲۴] و کتاب مقدماتیتر ج. و. استیوارت<sup>۳</sup> [۲۳] منابع اطلاعاتی موجودی هستند که روشهای جدیدی مانند روش  $QR$  (با تغییر مکان) و بسیاری از مطالب مشروح حذف شده در صفحات بعد را در بر دارند.

عدد (حقیقی یا همثافت)  $\lambda$  را يك ویژه مقدار ماتریس  $B$  گوئیم، هر گاه به ازای بردار (حقیقی یا همثافت) غیر صفری مانند  $y$  داشته باشیم

$$By = \lambda y \quad (60.4)$$

در این صورت  $n$ -بردار  $y$  يك ویژه بردار متعلق به ویژه مقدار  $\lambda$  نامیده می شود. فرمول (60.4) را می توانیم به شکل زیر بنویسیم

$$(B - \lambda I)y = 0 \quad (61.4)$$

از آنجا که  $y$  قرار است يك بردار غیر صفر باشد، ملاحظه می کنیم که  $\lambda$  يك ویژه مقدار  $B$  است اگر و فقط اگر، دستگاه همگن (61.4) دارای جوابهای غیر صفر باشد. بنابراین لم زیر نتیجه ای است از قضیه 4.4.

لم 4.4 عدد  $\lambda$  يك ویژه مقدار ماتریس  $B$  است اگر و فقط اگر،  $(B - \lambda I)$  وارون پذیر نباشد.

توجه شود که (60.4) و (61.4) ویژه بردار متعلق به  $\lambda$  را تنها بی مضارب عددی (اسکالر)، معین می کند. اگر  $y$  يك ویژه بردار متعلق به  $\lambda$  و  $z$  مضرب بی عددی از  $y$ ، مثلاً  $z = \alpha y$  باشد، آنگاه  $z$  نیز يك ویژه بردار متعلق به  $\lambda$  است، زیرا داریم،

$$Bz = B(\alpha y) = \alpha(By) = \alpha(\lambda y) = \lambda(\alpha y) = \lambda z$$

□ مثال: برای هر بردار  $y$ ، ماتریس همانی  $I$  در معادله زیر صدق می کند

$$Iy = y = 1y$$

بنابراین 1 ویژه مقداری برای  $I$  است و هر بردار غیر صفر، يك ویژه برداری است برای  $I$  متعلق به 1. از آنجا که يك بردار می تواند فقط به يك (یا هیچ) ویژه مقدار متعلق باشد (یا به هیچ ویژه برداری متعلق نباشد)، نتیجه می شود که 1 تنها ویژه مقدار  $I$  است.

ماتریس صفر فقط و فقط يك ویژه مقدار دارد و آن عدد صفر است.  
ماتریس

$$B = \begin{bmatrix} 1 & 2 & 0 \\ 2 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

دارای ویژه مقدار ۱- است، زیرا  $B\mathbf{i}_3 = -\mathbf{i}_3$ ، همچنین  $B(\mathbf{i}_1 + \mathbf{i}_2) = 3(\mathbf{i}_1 + \mathbf{i}_2)$  لذا  $\lambda = 3$  نیز يك ویژه مقدار  $B$  است. و بالاخره  $B(\mathbf{i}_1 - \mathbf{i}_2) = -(\mathbf{i}_1 - \mathbf{i}_2)$  لذا ویژه مقدار ۱- دارای دو ویژه بردار  $\mathbf{i}_3$  و  $(\mathbf{i}_1 - \mathbf{i}_2)$  است که مستقل خطی هستند.  $\square$

اگر ماتریس  $B = (b_{ij})$  بالامتثلی باشد، آنگاه  $\lambda$  يك ویژه مقدار  $B$  خواهد بود اگر و فقط اگر، به ازای مقداری مانند  $\lambda$  داشته باشیم  $\lambda = b_{ii}$ . زیرا که در این صورت ماتریس  $(B - \lambda I)$  نیز بالامتثلی خواهد بود، و بنا بر این به موجب قضیه ۶.۴،  $(B - \lambda I)$  وارونپذیر نیست اگر و فقط اگر، یکی از درایه‌های قطری آن صفر باشد، یعنی اگر فقط اگر، به ازای مقداری مانند  $\lambda$ ، تساوی  $b_{ii} - \lambda = 0$  برقرار باشد. بنا بر این مجموعه ویژه مقدارهای متعلق به يك ماتریس مثلثی همان مجموعه اعدادی است که روی قطر آن قرار دارند.

$\square$  مثال ۱۰۴: بخصوص، تنها ویژه مقدار ماتریس

$$B = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

عدد صفر است و هر دو بردار  $\mathbf{i}_1$  و  $\mathbf{i}_2$  ویژه بردارهای  $B$  و متعلق به این ویژه مقدار هستند. هر ویژه بردار دیگر  $B$  باید ترکیبی خطی از این دو ویژه بردار باشد. چنانچه فرض کنیم  $\mathbf{y} = y_1\mathbf{i}_1 + y_2\mathbf{i}_2 + y_3\mathbf{i}_3$  (که متعلق به تنها ویژه مقدار ۰ است) باشد، آنگاه

$$\begin{aligned} \mathbf{0} &= \mathbf{0} \cdot \mathbf{y} = B\mathbf{y} = y_1 B\mathbf{i}_1 + y_2 B\mathbf{i}_2 + y_3 B\mathbf{i}_3 \\ &= \mathbf{0} + \mathbf{0} + y_3 B\mathbf{i}_3 \end{aligned}$$

از آنجا که  $B\mathbf{i}_3 = \mathbf{i}_3 \neq \mathbf{0}$ ، نتیجه می‌شود که  $y_3 = 0$ ، یعنی  $\mathbf{y} = y_1\mathbf{i}_1 + y_2\mathbf{i}_2$ ، که نشان می‌دهد  $\mathbf{y}$  يك ترکیب خطی از ویژه بردارهای  $\mathbf{i}_1$  و  $\mathbf{i}_2$  است.  $\square$

برای اینکه نشان دهیم که چرا ویژه مقدارها مورد توجه هستند، اکنون به طور مختصر دنباله‌های برداری به شکل

$$\mathbf{z}, B\mathbf{z}, B^2\mathbf{z}, B^3\mathbf{z}, \dots \quad (۶۲.۴)$$

را مورد مطالعه قرار می‌دهیم. دنبالهٔ مذکور در کار بردهای بی که در شروع این بخش نام برده شد، پدید می‌آید. چنین دنباله‌هایی را می‌باید در فصل ۵، در بحث روشهای بارستی برای حل دستگاههای معادلات، مورد بحث قرار دهیم.

فرض کنید که بردار آغازین  $\mathbf{z}^1$  در (۶۲.۴) بتواند به صورت مجموع ویژه بردارهای  $B$  نوشته شود، یعنی

$$\mathbf{z} = \mathbf{y}_1 + \mathbf{y}_2 + \dots + \mathbf{y}_r \quad (۶۳.۴)$$

که داریم

$$B\mathbf{y}_i = \lambda_i \mathbf{y}_i \quad i = 1, \dots, r$$

در این صورت جملهٔ  $m$ ام از دنباله (۶۲.۴)، به شکل ساده

$$B^m \mathbf{z} = \lambda_1^m \mathbf{y}_1 + \lambda_2^m \mathbf{y}_2 + \dots + \lambda_r^m \mathbf{y}_r \quad (۶۴.۴)$$

خواهد بود بنا بر این رفتار دنبالهٔ برداری  $\mathbf{z}^m$  (۶۲.۴) کاملاً با دنبالهٔ عددی ساده

$$\lambda_i^m, \lambda_i^1, \lambda_i^2, \lambda_i^3, \dots \quad i = 1, \dots, r$$

معین می‌شود. برای مثال نتیجه‌ری می‌شود که

$$\lim_{m \rightarrow \infty} B^m \mathbf{z} = \mathbf{0} \quad \text{به‌ازای تمام } \mathbf{z} \text{ها، اگر } |\lambda_i| < 1 \text{، آنگاه} \quad (۶۵.۴)$$

بعلاوه فرض کنید که  $\lambda_i$ ها از لحاظ قدرمطلق مرتب شده باشند، یعنی

$$|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_r|$$

که این امر همیشه با مرتب کردن درست  $\mathbf{y}_i$ ها می‌تواند به دست آید. بعلاوه فرض می‌کنیم که

$$|\lambda_1| > |\lambda_2| \quad (۶۶.۴)$$

در این فرض نه تنها  $\lambda_1$  باید غیر از تمام  $\lambda_i$ های دیگر باشد [که این امر همیشه با جمع کردن تمام  $\mathbf{y}_i$ های متعلق به  $\lambda_1$  در (۶۳.۴) به دست می‌آید، و به موجب آن تنها یک ویژه مقدار متعلق به  $\lambda_1$  حاصل می‌شود] بلکه  $\lambda_i$  دیگری به همان اندازه  $\lambda_1$  وجود ندارد و این قسمت است که موجب می‌شود (۶۶.۴) یک فرض غیر بدیهی باشد.

بنابراین از تقسیم دو طرف (۶۴.۴) بر  $\lambda_1^m$  خواهیم داشت

$$(\lambda_1^{-1} B)^m \mathbf{z} = \mathbf{y}_1 + \sum_{i=2}^r \left( \frac{\lambda_i}{\lambda_1} \right)^m \mathbf{y}_i \quad m = 0, 1, 2, \dots$$

- 
1. starting vector
  2. vector sequence
  3. numerical sequence

بنا بر فرضهای اتخاذ شده، داریم

$$\left| \frac{\lambda_i}{\lambda_1} \right| < 1 \quad i = 2, \dots, r$$

از این رو نتیجه می شود که

$$\lim_{m \rightarrow \infty} (\lambda_1^{-1} B)^m \mathbf{z} = \mathbf{y}_1 \quad (۶۷.۴)$$

می توانیم بگوئیم که اگر  $\mathbf{z}$  بتواند به شکل (۶۳.۴) بر حسب ویژه بردارهای  $B$  چنان نوشته شود که ویژه مقدار  $\lambda_1$  متناظر با  $\mathbf{y}_1$ ، مطلقاً از تمام ویژه مقدارهای دیگر بزرگتر باشد، آنگاه شکل مدرج شده  $B^m \mathbf{z}$  به طرف  $\mathbf{y}_1$  همگرا می شود.

□ مثال ۱۱.۴: قبلاً دیدیم که ماتریس

$$B = \begin{bmatrix} 1 & 2 & 0 \\ 2 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

دارای ویژه بردارهای  $\mathbf{z}_3 = \mathbf{i}_3$ ،  $\mathbf{z}_2 = \mathbf{i}_1 - \mathbf{i}_2$ ،  $\mathbf{z}_1 = \mathbf{i}_1 + \mathbf{i}_2$  و ویژه مقدارهای  $\lambda_1 = \lambda_2 = -1$ ،  $\lambda_3 = 3$  می باشد. این ویژه بردارها مستقل خطی هستند (تمرین ۱۰.۴-۱۰.۵ را ببینید) و از این رو یک پایه برای تمام ۳- بردارها تشکیل می دهند بنا بر این نتیجه می شود که هر ۳- برداری را می توان به صورت حاصلجمعی از ویژه بردارهای  $B$  نوشت. به ویژه بردار  $\mathbf{z}$  که به وسیله  $\mathbf{z}^T = [1 \ 2 \ 3]$  معین می شود، می تواند به صورت

$$\mathbf{z} = \mathbf{y}_1 + \mathbf{y}_2$$

نوشته شود، که در آن

$$\mathbf{y}_1 = 1.05 \mathbf{z}_1 \quad \mathbf{y}_2 = -0.05 \mathbf{z}_2 + 3 \mathbf{z}_3$$

در جدول ۱.۴،  $B^m \mathbf{z}$  و  $\alpha_m B^m \mathbf{z}$  به ازای  $m = 0, \dots, 5$  فهرست شده اند. چنان مدرج شده است که عنصر اول آن برابر با ۱ باشد. به نظر می رسد که  $\mathbf{z}^{(m)}$  به سمت ویژه بردار  $\mathbf{i}_1 + \mathbf{i}_2$  متعلق به  $\lambda_1 = 3$  همگرا می شود. □

**روش توانی<sup>۱</sup>** که برای محاسبه بزرگترین ویژه مقدار (از نظر قدر مطلق) یک ماتریس معین  $B$  به کار می رود مبتنی بر توضیح زیر است. با انتخاب یک  $\mathbf{z}$ ، مثلاً  $\mathbf{z} = \mathbf{i}_1$ ، چند جمله اول دنباله (۶۲.۴) تولید می شود و با ادامه کار نسبتهایی به شکل

## جدول ۱۰۴

$m$	۰	۱	۲	۳	۴	۵
$B^m \mathbf{z}$	۱	۵	۱۳	۲۱	۱۲۱	۳۶۵
	۲	۴	۱۴	۲۰	۱۲۲	۳۶۴
	۳	۳	۳	۳	۳	۳
$\mathbf{z}^{(m)}$	۱	۱۰۰	۱۰۰۰	۱۰۰۰۰	۱۰۰۰۰	۱۰۰۰۰
	۲	۰۰۸	۱۰۰۸	۰۰۹۷۶	۱۰۰۰۸	۰۰۹۹۷
	۳	۰۰۶	۰۰۲۳	۰۰۰۷۳	۰۰۰۲۵	۰۰۰۰۸

$$\mathbf{u}^T B^{m+1} \mathbf{z} / (\mathbf{u}^T B^m \mathbf{z}) \quad (۶۸.۴)$$

محاسبه می‌شود، از (۶۴.۴) نتیجه می‌شود

$$\frac{\mathbf{u}^T B^{m+1} \mathbf{z}}{\mathbf{u}^T B^m \mathbf{z}} = \frac{\lambda_1 \mathbf{u}^T \mathbf{y}_1 + (\lambda_r / \lambda_1)^{m+1} \mathbf{u}^T \mathbf{y}_r + \dots + (\lambda_r / \lambda_1)^{m+1} \mathbf{u}^T \mathbf{y}_r}{\mathbf{u}^T \mathbf{y}_1 + (\lambda_r / \lambda_1)^m \mathbf{u}^T \mathbf{y}_r + \dots + (\lambda_r / \lambda_1)^m \mathbf{u}^T \mathbf{y}_r}$$

بنابراین به شرط  $\mathbf{u}^T \mathbf{y}_1 \neq 0$  و  $|\lambda_1| \geq |\lambda_r| \geq \dots \geq |\lambda_r|$  خواهیم داشت

$$(\mathbf{u}^T B^{m+1} \mathbf{z}) / (\mathbf{u}^T B^m \mathbf{z}) = \lambda_1 + \theta ((\lambda_r / \lambda_1)^m)$$

توجه کنید، در حالتی که  $B = B^T$ ، یعنی  $B$  متقارن باشد در رابطه (۶۸.۴) بردار  $\mathbf{u} = B^m \mathbf{z}$  به کار می‌آید. نسبت حاصله یعنی

$$(\mathbf{u}^T B \mathbf{u}) / (\mathbf{u}^T \mathbf{u})$$

به خارج قسمت ریالی<sup>۱</sup> (برای  $\mathbf{u}$  و  $B$ ) موسوم است که می‌توان به آسانی دید که برابر است با

$$(\mathbf{z}^T B^2 \mathbf{z}) / (\mathbf{z}^T B \mathbf{z})$$

و از این رو با تعداد عملیاتی تا  $\theta ((\lambda_r / \lambda_1)^m)$  با  $\lambda_1$  برابر خواهد شد.

□ مثال ۱۲.۴: از دنباله حاصل در مثال ۱۰.۴، به ازای  $\mathbf{u} = \mathbf{i}$  دنباله نسبتها را به صورت زیر به دست می‌آوریم

$$۵, ۲۰۶, ۳۰۱۵۳۸ \dots, ۲۰۹۵۱۲ \dots, ۳۰۰۱۶۵ \dots$$

در حالی که به ازای  $\mathbf{u} = \mathbf{i}_p$ ، دنبالهٔ زیر به دست می‌آید

$$۲, ۳۰۵, ۲۰۸۵۷۱ \dots, ۳۰۵, ۲۰۹۸۳۶ \dots$$

به نظر می‌رسد که هر دو دنباله به سمت  $۳ = \lambda_1$  همگرا می‌شوند، اما به ازای  $\mathbf{u} = \mathbf{i}_p$ ، دنبالهٔ زیر را خواهیم داشت

$$۱, ۱, ۱, ۱, ۱, \dots$$

که به نظر نمی‌رسد به سمت ۳ همگرا شود.

چون  $B$  متقارن است، دنبالهٔ خارج‌قسمتهای ریالی را نیز محاسبه و نسبت‌های زیر را

پیدا می‌کنیم

$$۱۰۵۷۱۴ \dots, ۲۰۶, ۲۰۹۴۶۵ \dots, ۲۰۹۹۳۹ \dots, ۲۰۹۹۹۳ \dots$$

در این دنباله تقریباً در هر دور یک رقم به دست می‌آوریم که حاکی از این امر است که دنبالهٔ فوق باید به ازای تعداد عملیاتی تا  $O((1/3)^{2m})$  به سمت  $۳ = \lambda_1$  همگرا شود.  $\square$

شق جالب دیگر روش توانسی روش بارست معکوس<sup>۱</sup> است. در اینجا، علاوه بر انتخاب بردار آغازین  $\mathbf{z}$  که در رابطه (۶۳.۴) صدق می‌کند، عددی مانند  $p$  که بایکی از ویژه‌مقدارهای  $B$  برابر نباشد، انتخاب و سپس دنبالهٔ زیر را تشکیل می‌دهند

$$\mathbf{z}, \hat{B}\mathbf{z}, \hat{B}^2\mathbf{z}, \hat{B}^3\mathbf{z}, \dots$$

با

$$\hat{B} = (B - pI)^{-1}$$

باید توجه داشت که در رابطه (۶۳.۴) به ازای هر ویژه بردار  $\mathbf{y}_i$  از  $B$ ، داریم

$$(B - pI)\mathbf{y}_i = B\mathbf{y}_i - p\mathbf{y}_i = (\lambda_i - p)\mathbf{y}_i$$

بنابراین خواهیم داشت

$$(\lambda_i - p)^{-1}\mathbf{y}_i = (B - pI)^{-1}\mathbf{y}_i$$

رابطهٔ فوق نشان می‌دهد که  $\mathbf{z}$ ، مجموع ویژه بردارهای  $(B - pI)^{-1}$  متناظر بسا ویژه‌مقدارهای  $(\lambda_i - p)^{-1}$ ،  $i = 1, \dots, r$ ، نیز هست. اما اگر  $p$  تنها به یکی از ویژه‌مقدارهای

$\lambda_1, \dots, \lambda_p, \lambda_r$ ، مثل  $\lambda_r$ ، کاملاً نزدیک باشد، آنگاه  $(\lambda_j - p)^{-1}$  در مقایسه با سایر ویژه مقادیرهای  $(\lambda_i - p)^{-1}$ ، از لحاظ قدر مطلق بسیار بزرگ خواهد بود و بنابراین بحث پیشین ما در روش توانی، به ما اجازه می‌دهد که چنین نتیجه‌گیری کنیم که شکل مدرج و مناسبی از دنباله  $\dots \hat{B}^2 z, \hat{B} z, z$  به سرعت به سمت ویژه بردار  $y$  متناظر با  $\lambda_r$  همگرا می‌شود، درحین اینکه نسبت‌های متناظر با آن، یعنی

$$u^T \hat{B}^{m+1} z / (u^T \hat{B}^m z)$$

نیز با همان سرعت به سمت عدد  $(\lambda_r - p)^{-1}$  همگرا می‌شود. این امر موجب می‌شود که بارست معکوس، روش بسیار کارایی در وضعیت زیر باشد: تقریب خوبی قبلاً برای ویژه بردار  $B$  به دست آورده‌ایم و می‌خواهیم این تقریب را بهبود بخشیم و (یا) ویژه بردار متناظر با آن را محاسبه کنیم.

همان گونه که ذکر کردیم، بارست معکوس ابتدا به ساختن ماتریس  $\hat{B} = (B - pI)^{-1}$  نیاز دارد. اما همچنان که در بخش ۴.۴ بحث شد، ما چنین عکسی را صراحتاً نمی‌سازیم. بلکه از

$$z^{(m)} := \hat{B}^m z \quad m = 0, 1, 2, \dots$$

به دست می‌آوریم، و توجه می‌کنیم که

$$(B - pI)z^{(m)} = z^{(m-1)}$$

در نتیجه، وقتی که یکبار تجزیه به عوامل  $PLU$  را برای ماتریس  $B - pI$  به دست آوردیم از راه الگوریتم جایگزینی ۴.۴ یعنی قرارداد در  $(n^2)$ ،  $z^{(m)}$  را از  $z^{(m-1)}$  به دست می‌آوریم. این روش، اگر مجبور به انجام آن باشیم، گرانتر از محاسبه صریح حاصلضرب  $\hat{B} z^{(m-1)}$  از  $\hat{B}$  نخواهد بود.

در اینجا یک زیر برنامه فورترن برای انجام بارست عکس آورده‌ایم. در مرحله  $m$ ام،  $u = B^m z$  را انتخاب می‌کنیم، یعنی، در هر مرحله خارج قسمت ریلای را انجام می‌دهیم.

```

SUBROUTINE INVTR ( B, N, EGUSS, VGUSS, W, D, IPIVOT,
                  VALUE, VECTOR, IFLAG )
C   CALLS FACTOR, SUBST.
C   INTEGER IFLAG, IPIVOT(N), I, ITER, ITERMX, J
C   REAL B(N,N), D(N), EGUSS, VALUE, VECTOR(N), VGUSS(N), W(N,N)
C   EPSLON, EVNEW, EVOLD, SCNORM
C***** INPUT *****
C B THE MATRIX OF ORDER N WHOSE EIGENVALUE/VECTOR IS SOUGHT.
C N ORDER OF THE MATRIX B.
C EGUSS A FIRST GUESS FOR THE EIGENVALUE.
C VGUSS N-VECTOR CONTAINING A FIRST GUESS FOR THE EIGENVECTOR.
C***** WORK AREA *****
C W MATRIX OF ORDER N
C D VECTOR OF LENGTH N
C IPIVOT INTEGER VECTOR OF LENGTH N
C***** OUTPUT *****
C VALUE COMPUTED APPROXIMATION TO EIGENVALUE
C VECTOR COMPUTED APPROXIMATION TO EIGENVECTOR
C IFLAG AN INTEGER,

```

## ماتریسها و دستکاهای معادلات خطی ۲۵۱

```

C      = 1 OR -1 (AS SET IN FACTOR), INDICATES THAT ALL IS WELL,
C      = 0 , INDICATES THAT SOMETHING WENT WRONG. SEE PRINTED ERROR
C      MESSAGE .
C***** M E T H O D *****
C      INVERSE ITERATION, AS DESCRIBED IN THE TEXT, IS USED.
C*****
C      THE FOLLOWING T E R M I N A T I O N P A R A M E T E R S A R E S E T
C      H E R E , A T O L E R A N C E E P S L O N O N T H E D I F F E R E N C E B E T W E E N S U C C E S S I V E
C      E I G E N V A L U E I T E R A T E S , A N D A N U P P E R B O U N D I T E R M X O N T H E N U M B E R
C      O F I T E R A T I O N S T E P S .
      DATA EPSLON,ITERMX /.000001,20/
C
C      PUT Z - (EGUESS)*IDENTITY INTO W
      DO 10 J=1,N
      DO 9 I=1,N
        9      W(I,J) = B(I,J)
        10     W(J,J) = W(J,J) - EGRESS
      CALL FACTOR ( W, N, D, IPIVOT, IFLAG )
      IF (IFLAG.EQ. 0) THEN
        PRINT 610
        610    FORMAT(' EIGENVALUE GUESS TOO CLOSE.
          *          , 'NO EIGENVECTOR CALCULATED.')
              RETURN
      END IF
      ITERATION STARTS HERE
C
      PRINT 619
        619    FORMAT(' ITER EIGENVALUE      EIGENVECTOR COMPONENTS'/)
      EVOLD = 0.
      DO 50 ITER=1,ITERMX
C          NORMALIZE CURRENT VECTOR GUESS
      SQNORM = 0.
      DO 20 I=1,N
        20     SQNORM = VGUESS(I)**2 + SQNORM
      SQNORM = SQRT(SQNORM)
      DO 21 I=1,N
        21     VGUESS(I) = VGUESS(I)/SQNORM
C          GET NEXT VECTOR GUESS
      CALL SUBST ( W, IPIVOT, VGUESS, N, VECTOR )
C          CALCULATE RAYLEIGH QUOTIENT
      EVNEW = 2.
      DO 30 I=1,N
        30     EVNEW = VGUESS(I)*VECTOR(I) + EVNEW
      EVALUE = EGRESS + 1./EVNEW
C
      PRINT 630,ITER,EVALUE,VECTOR
        630    FORMAT(I3,E15.7,2X,3E14.7/(20X,3E14.7))
C          STOP ITERATION IF CURRENT GUESS IS CLOSE TO
C          PREVIOUS GUESS FOR EIGENVALUE
      IF ( ABS(EVNEW-EVOLD) .LE. EPSLON*ABS(EVNEW) )
        RETURN
      EVOLD = EVNEW
      DO 50 I=1,N
        50     VGUESS(I) = VECTOR(I)
C
      IFLAG = 0
      PRINT 660,EPSLON,ITERMX
        660    FORMAT(' NO CONVERGENCE TO WITHIN',E10.4,' AFTER',I3,' STEPS. ')
              RETURN
      END

```

□ مثال ۱۳.۴: برای ماتریس  $B$  در مثال ۱۲.۴ به ازای  $\mathbf{z} = [1, 1, 1]^T$  و  $p = 3$ ، که در مثال فوق برای  $\lambda_1 = 3$  از نخستین دنباله نسبتها بهترین حدس است، برنامه INVITR را به کار می بریم

بارست	ویژه مقدار	مؤلفه های ویژه بردار		
1	0.2991801 + 0i	-0.3499093 + 0i	-0.3499093 + 0i	-0.1437446 + 0i
2	0.3000000 + 0i	0.4285478 + 0i	0.4285478 + 0i	0.7232219 - 0i
3	0.3000000 + 0i	-0.4285496 + 0i	-0.4285496 + 0i	-0.2971047 - 0i
4	0.3000000 + 0i	0.4285496 + 0i	0.4285496 + 0i	0.1220522 - 0i
	ویژه مقدار = 0.3000000 + 0i			
	ویژه بردار =			
	0.4285496 + 0i	0.4285496 + 0i	0.1220522 - 0i	



چون  $B$  متقارن است و خارج قسمت ریاضی محاسبه شده بود. خروجی، همگرایی بسیار سریع ویژه بردار (در هر مرحله از بارست، در حدود ۲ رقم اعشاری به دست می آید) و حتی همگرایی بسیار سریع ویژه مقدار را نشان می دهد.

به عنوان مثالی از این امر که برخلاف روش توانی، بارست معکوس برای هر ویژه مقداری می تواند به کار رود، از  $T^3 = [1, 1, 1]$  و  $z = 0$  نیز شروع کرده امیدواریم که بدین وسیله به کوچکترین ویژه مقدار  $B$  (از نظر قدر مطلق) دست یابیم.

بارست	ویژه مقدار	مؤلفه های ویژه بردار		
1	$-0.9000000 + 0i$	$0.1924501 + 0i$	$0.1924501 + 0i$	$-0.5773503 + 0i$
2	$-0.1320000 + 0i$	$0.1005038 + 0i$	$0.1005038 + 0i$	$0.9045340 + 0i$
3	$-0.1033195 + 0i$	$0.3658808 - 0i$	$0.3658809 - 0i$	$-0.9878783 + 0i$
4	$-0.1003661 + 0i$	$0.1232878 - 0i$	$0.1232877 - 0i$	$0.9986311 + 0i$
5	$-0.1000406 + 0i$	$0.4114594 - 0i$	$0.4114604 - 0i$	$-0.9998476 + 0i$
6	$-0.1000045 + 0i$	$0.1371724 - 02$	$0.1371714 - 02$	$0.9999831 + 0i$
7	$-0.1000005 + 0i$	$0.4572417 - 03$	$0.4572513 - 03$	$-0.9999981 + 0i$
8	$-0.1000001 + 0i$	$0.1524206 - 03$	$0.1524109 - 03$	$0.9999998 + 0i$
9	$-0.1000000 + 0i$	$0.5080043 - 04$	$0.5081010 - 04$	$-0.1000000 + 0i$
		$-0.1000000 + 0i$ = ویژه مقدار		
		$=$ ویژه بردار		
		$0.5080043 - 04$	$0.5081010 - 04$	$-0.1000000 + 0i$

در این مورد همگرایی بسیار کند تر است، زیرا  $0$ ، مخصوصاً به ویژه مقدار  $1$  - نزدیک نیست، اما با اختیار ویژه برداری به شکل  $T^3 = [0, 0, 1]$  (به جای شکل کلیتر ممکن  $T^3 = [a, -a, b]$  برای ویژه مقدار  $1$  - از ماتریس  $B$ )، بعد از ۹ بارست به همگرایی می رسیم. □

روش توانی و شکل دیگر آن یعنی بارست معکوس، عموماً کاربرد ی ندارد. به طور کلی اگر قرار باشد که ویژه مقدارها به صورت اعداد هم تافت به دست آیند، در ابتدا حساب اعداد هم تافت باید به کار گرفته شود. راه های ماهرانهٔ ویژه های در دست است که به کمک آنها در حساب اعداد حقیقی می توان از یک ماتریس حقیقی  $B$  به یک جفت ویژه مقدار مزدوج هم تافت دست یافت. مشکل بسیار جدی عبارت از احتمال همگرایی بسیار کند است که زمانی پیش می آید که بزرگترین ویژه مقدار بعدی از لحاظ قدر مطلق به بزرگترین مقدار بسیار نزدیک باشد. در حالی که فرایند  $\Delta^2$ ی اتکین (الگوریتم ۷.۳) می تواند برای شتاب بخشیدن به همگرایی (اگر همگرایی وجود داشته باشد) به کار رود، در حالت نهایی یعنی وقتی که تساوی  $|\lambda_1| = |\lambda_p|$  برقرار باشد و  $\lambda_1 \neq \lambda_p$ ، (برای مثال وقتی که  $\lambda_1 = \lambda_p$ )، همگرایی وجود نخواهد داشت. می توان با تغییر مکان مناسب، راه حلی برای این مسئله پیدا کرده یعنی به جای استفاده از خود ماتریس  $B$  با ماتریس  $B - pI$  کار کرد، به طوری که  $|\lambda_1 - p| > |\lambda_p - p|$  (تمرینات ۶-۸.۴ و ۷-۸.۴ را ببینید).

و بالاخره زمانی که بردار آغازین  $z$  را نتوان به صورت حاصل جمع بردارهای خاص  $B$  نوشت، روش توانی ثابت نظری خود را (آن گونه که در اینجا عنوان کردیم) از دست

می دهد. از آنجایی که ویژه بردارهای  $B$  در دست نیستند، تنها وقتی می توانیم مطمئن باشیم که  $Z$  به صورت حاصلجمعی از ویژه بردارهای ماتریس  $B$  از مرتبه  $n \times n$  می تواند نوشته شود که بدانیم هر  $n$ -برداری را می توان به صورت حاصلجمعی از ویژه بردارهای  $B$  نوشت، ولی در واقع، معنی این امر این است که بخواهیم ویژه بردارهای  $B$  آن قدر باشند که بتوانند یک پایه تشکیل دهند. یک پایه از  $n$ -برداریها را که تماماً ویژه بردارهای ماتریس  $B$  از مرتبه  $n \times n$  باشند، یک مجموعه کامل ویژه بردارهای  $B$  نامند. روشن است که هرگاه  $Z_1, \dots, Z_n$ ، مجموعه کامل ویژه بردارهای ماتریس  $B$  از مرتبه  $n \times n$  و بنا بر این پایه ای برای تمام  $n$ -برداریها باشد، آنگاه هر  $n$ -برداری خاص  $Z$  می تواند به صورت ترکیبی خطی از این ویژه بردارها به ازای ضرایب مناسب  $a_1, \dots, a_n$  به صورت

$$Z = a_1 Z_1 + a_2 Z_2 + \dots + a_n Z_n$$

نوشته شود. اگر  $a_i \neq 0$ ، آنگاه  $y_i = a_i Z_i$  نیز یک ویژه بردار  $B$  است، درحالی که اگر  $a_i = 0$ ، بی آنکه مشکلی پیش آید، جمله  $a_i Z_i$  می تواند از حاصلجمع حذف شود. بدین ترتیب  $Z$  به صورت حاصلجمعی از بردارهای  $B$  به دست می آید، (بجز درحالت نامطلوب  $Z = 0$ ). متأسفانه، همان گونه که قبلاً در مثال ۱۰.۴ دیدیم، معلوم نیست که هر ماتریس مجموعه کاملی از ویژه بردارها داشته باشد.

### تشابه<sup>۱</sup>

این امر که معلوم نیست هر ماتریسی یک مجموعه کاملی از ویژه بردارها داشته باشد، نشانه پیچیدگی ناشی از نظریه ویژه مقدارهاست. این پیچیدگی مربوط به این حکم است که هر ماتریس مربعی را نمی توان به شکل

$$Y \Lambda Y^{-1}$$

که در آن  $\Lambda$  یک ماتریس قطری است، نوشت. زیرا تساوی  $BY = Y \Lambda$  با ماتریس قطری  $\Lambda$  برقرار است، اگر و فقط اگر ستونهای ماتریس  $Y$  شامل ویژه بردارهای  $B$  باشند. درحالی که چنین ماتریس  $Y$  ای از مرتبه  $n \times n$  وارون پذیر خواهد بود، اگر و فقط اگر، ستونهای آن تشکیل یک پایه برای  $n$ -برداریها دهند.

دو ماتریس  $A$  و  $B$  را مشابه<sup>۲</sup> گوئیم، هرگاه به ازای ماتریسی (وارون پذیر) مانند  $C$  داشته باشیم

$$A = C^{-1} B C$$

در ماتریسهای مشابه، ویژه مقدارها و ویژه بردارهای متناظر به آنها یکی هستند. در واقع اگر به ازای بردار غیر صفری مانند  $x$ ، داشته باشیم  $(B - \lambda I)x = 0$  و  $A = C B C^{-1}$

آنگاه  $Cx$  نیز غیر صفر بوده و داریم  $AC = CB$ ، از این رو خواهیم داشت

$$(A - \lambda I)Cx = C(B - \lambda I)x = C \cdot 0 = 0$$

به طور خلاصه، متناظر با هر زوج ویژه مقدار-ویژه بردار  $(\lambda, x)$  از  $B$ ، یک زوج ویژه مقدار-ویژه بردار  $(\lambda, Cx)$  از  $A$  وجود دارد.

مطالب فوق به عنوان اولین مرحله در محاسبه ویژه مقادیرهای  $B$  بیانگر یک تبدیل تشابهی از  $B$  به ماتریس  $A = C^{-1}BC$  است که برای آن ویژه مقادیرها را تا حدی آسانتر می توان محاسبه کرد.

برای مثال، اگر بتوان یک ماتریس بالامثلثی مشابه با  $B$  یافت، آنگاه تمام ویژه مقادیرهای  $B$  معین می شوند، زیرا این ویژه مقادیرها تماماً روی قطر  $T$  قرار دارند. در حقیقت می توان ثابت کرد که:

**قضیه ۱۳.۴:** قضیه شورا هر ماتریس مربعی  $B$  را می توان به صورت  $U^{-1}TU$  نوشت که در آن،  $T$  ماتریس بالامثلثی و  $U$  ماتریس یکانسی است، یعنی  $U^H U = I$ .

این حقیقت که  $U$  ماتریس یکانسی است، این نتیجه جالب را در بردارد که به ازای هر  $x$  تساوی  $\|x\|_2 = \|Ux\|_2$  برقرار است. از این رو داریم  $\|B\|_2 = \|T\|_2$ ، لذا ماتریس بالامثلثی  $T$  مشابه  $B$  بوده و حتی اندازه آن با اندازه  $B$  نیز یکی می باشد. با وجود این، متأسفانه برای ساختن این ترکیبهای  $U$  و  $T$  معمولاً از روند بارستی استفاده می کنند.

اما همواره ممکن است که با تعداد  $\Theta(n^3/3)$  عملیات ممیز شناور،  $B$  را از راه مشابهت به ماتریس  $H = (h_{ij})$ ، که تقریباً مثلثی یا هسنبرگ<sup>۳</sup> یعنی به صورت

$$h_{ij} = 0, \quad i > j - 1$$

است تبدیل کرد. بنابر این قسمت پایین-مثلثی  $H$  بجز احتمالاً اولین قطر فرعی پایین قطر اصلی، برابر صفر است.  $H$  از روی  $B$  با  $n - 2$  تبدیل تشابهی ساده به دست می آید که هر کدام از آنها یک ستون صفر دیگر در زیر اولین قطر پایینی ایجاد می کنند.

مثلاً ممکن است که از تقارنهای محوری هاوس هولدر<sup>۴</sup>، یعنی ماتریسهایی به شکل

$$\alpha y^T y = \gamma \quad \text{با} \quad R(y) = I - \alpha y y^T \quad (69.4)$$

به صورت زیر استفاده کرد: فرض کنید که در ستونهای  $1, 2, \dots, k - 1$  در زیر قطر فرعی صفر آمده باشد، مثل حالتی که در آن به ازای  $k = 1$  با  $H^{(1)} = B$  پیش می آید. پس، می خواهیم رابطه

$$H^{(k+1)} = (R(y))^{-1} H^{(k)} R(y)$$

را چنان تشکیل دهیم که نخستین  $k - 1$  ستون تغییر نیافته بماند، در حالی که اکنون در ستون

$k$ ام در زیر اولین قطر فرعی نیز صفر داریم. بدین دلیل قبل از همه توجه داریم که وارون  $R(y)$ ، خود  $R(y)$  است، زیرا داریم

$$(I - \alpha y y^T)(I - \alpha y y^T) = I - \alpha y y^T - \alpha y y^T + \alpha^2 y (y^T y) y^T$$

و  $\alpha y^T y = 2$ . از این رو خواهیم داشت،  $H^{(k+1)} = R(y)H^{(k)}R(y)$ ، به همین ترتیب، محاسبه می شود که

$$\|R(y)x\|_2 = \|x\|_2 \quad (70.4)$$

یادآوری می کنیم که  $\|z\|_2^2 = z^T z$ . این مطلب، اتفاقاً نام «تقارن محوری» را روشن می کند. بعد باید بدانیم که کوتاهترین راه برای به دست آوردن حاصلضرب ماتریسی  $AR(y)$  این است که به جای هر سطر  $x^T$  از  $A$  بردار سطر  $x^T - \alpha(x^T y)y^T$  را بگذاریم. از این رو با انتخاب

$$y_1 = \dots = y_k = 0 \quad (71.4)$$

$k$ ستون اول از ماتریس  $R(y)H^{(k)}$  مشابه با  $k$ ستون  $H^{(k)}$  خواهد بود. بعد، باید بدانیم کوتاهترین راه به دست آوردن حاصلضرب ماتریسی  $R(y)A$  این است که به جای هر ستون  $x$  از  $A$  بردار ستونی  $x - \alpha(y^T x)y$  را بگذاریم. چون  $H^{(k)}$  در ستونهای ۱ تا  $k-1$  و در پایین سطر  $(k-1)$  ام خود صفر دارد، این امر نشان می دهد که انتخاب (71.4) ما را مطمئن می سازد که  $(k-1)$  ستون اول از  $R(y)H^{(k)}R(y)$  مشابه با  $(k-1)$  ستون  $H^{(k)}$  است. این مطلب ما را با مسئله انتخاب  $y_{k+1}, \dots, y_n$  به گونه ای که ستون  $k$ ام از  $R(y)H^{(k)}$  در سطرهای  $k+2, \dots, n$  صفر داشته باشد مواجه می سازد. با توجه به (71.4)، معنی این مطلب این است که می باید  $R(y)$  بردار

$$\hat{x} = [0, \dots, 0, h_{k+1, k}^{(k)}, \dots, h_{n, k}^{(k)}]^T$$

را به ازای مقدار اسکالر (عدد دوار) مانند  $-\beta$ ، بر روی بردار  $-\beta i_{k+1}$  بنگارد [در اینجا عنصر  $(i, j)$  ام  $H^{(k)}$  با  $h_{ij}^k$  نشان داده شده است]. با توجه به رابطه (70.4). این امر به معنای آن است که

$$\beta = \pm \|\hat{x}\|_2 = ((h_{k+1, k}^{(k)})^2 + \dots + (h_{n, k}^{(k)})^2)^{1/2}$$

بعلاوه

$$\hat{x} - (-\beta) i_{k+1} = \hat{x} - R(y)\hat{x} = \hat{x} - (\hat{x} - \alpha(y^T \hat{x})y) = \alpha(y^T \hat{x})y$$

نشان می دهد که  $y$  می باید یک مضرب عددی (اسکالر) از بردار  $\hat{x} + \beta i_{k+1}$  باشد. این امر مشخص می کند که انتخاب زیر برای  $y$  کار را به تمام خواهد رسانید:

$$y_i = \begin{cases} \hat{x}_{k+1} + \text{signum}(\hat{x}_{k+1}) \|\hat{\mathbf{x}}\|_2 & i = k+1 \\ \hat{x}_i & i > k+1 \end{cases} \quad (۷۲.۲)$$

یعنی  $\beta = \text{signum}(\hat{x}_{k+1}) \|\hat{\mathbf{x}}\|_2$ . در اینجا علامت  $\beta$  چنان انتخاب می‌شود که در محاسبهٔ  $y_{k+1}$ ، مانع از بین رفتن ارقام با معنی شود.  $\alpha$  ی متناظر به آن می‌تواند به سادگی به شکل

$$\alpha = 1/(\beta y_{k+1}) \quad (۷۳.۴)$$

نوشته‌شود. بدین ترتیب، بعد از  $n-2$  مرحله از این قبیل، ماتریس

$$H = R^{-1}BR$$

به دست می‌آید که در آن  $H$ ، ماتریس هسبرگ است و

$$R = R(y^{(1)})R(y^{(2)}) \dots R(y^{(n-2)})$$

حاصل ضرب تقارنهای محوری هاوس هولدر است، از این رو داریم

$$R^{-1} = R(y^{(n-2)}) \dots R(y^{(2)})R(y^{(1)})$$

روشن است که تقارن هاوس هولدر، ماتریس متقارن حقیقی است (اگر  $y$  حقیقی باشد)، بنا بر این در صورتی که  $B$  ماتریس متقارن و حقیقی باشد،  $H$  نیز متقارن حقیقی خواهد بود. بنا بر این در صورتی که  $B$  متقارن حقیقی باشد،  $H$  سه قطری و متقارن است. اکنون، برای راحتی، مطالب بالا را به صورت رسمی درمی‌آوریم.

**الگوریتم ۶.۴:** تبدیل تشابهی به شکل ماتریس بالامثلثی هسبرگ با استفاده از تقارنهای محوری هاوس هولدر

ماتریس  $A$  از مرتبهٔ  $n$  که در  $n$  ستون اول از آرایهٔ کاری  $H$  از مرتبهٔ  $(n+2) \times n$  حفظ شده، مفروض است.

For  $k = 1, \dots, n-2$ , do:

$$\beta := \text{signum}(h_{k+1,k}) \left( \sum_{j=k+1}^n (h_{jk})^2 \right)^{1/2}$$

$$h_{k+1,k} := h_{k+1,k} + \beta$$

$$\alpha^{-1} := \beta h_{k+1,k}$$

For  $j = 1, \dots, n$ , do:

$$\gamma := \left( \sum_{i=k+1}^n h_{ik} h_{ji} \right) / \alpha^{-1}$$

```

For  $i = k + 1, \dots, n$ , do:
   $h_{ji} := h_{ji} - \gamma h_{ik}$ 
For  $j = k + 1, \dots, n$ , do:
   $\gamma := \left( \sum_{i=k+1}^n h_{ik} h_{ij} \right) / \alpha^{-1}$ 
For  $i = k + 1, \dots, n$ , do:
   $h_{ij} := h_{ij} - \gamma h_{ik}$ 
 $h_{k+1, n+1} := h_{k+1, k}$    $h_{k+1, n+2} := \alpha^{-1}$ 
 $h_{k+1, k} := -\beta$ 

```

بنابراین  $H$ ، متضمن جزء جالب از يك ماتریس بالامثلثی هسبرگ است که باماتریس ورودی  $A$  در قسمت بالامثلثی هسبرگ در اولین  $n$  ستون و  $n$  سطرش مشابه است. بعلاوه ماتریس فوق شامل اطلاعات کاملی است دربارهٔ بردارهای  $\mathbf{y}$  و عددوارهای  $\alpha$  که تقارنهای محوری هاوس هولدر گونا گونی را که به کار رفته اند به مامی دهند. زمانی که ویژه بردارهای ماتریس بالامثلثی هسبرگ می باید مجدداً به ویژه بردارهای ماتریس اصلی  $A$  تبدیل شوند، این اطلاعات لازم خواهند بود.

روشی که معمولاً برای یافتن تمام ویژه مقدارهای يك ماتریس کلی  $B$  توصیه می شود، روش **QR** است. این کار با تبدیل ماتریس فوق به شکل ماتریس هسبرگ  $H$ ، که هم اکنون ذکر کردیم، آغاز می شود. زمانی که این کار انجام گیرد، ماتریس  $H$  نخستین عضو در دنبالهٔ زیر خواهد بود

$$A^{(0)} = H, A^{(1)}, A^{(2)}, \dots$$

که در آن  $A^{(k+1)}$  به گونهٔ زیر از  $A^{(k)}$  به دست می آید: به ماتریس یکانی  $Q$  و ماتریس بالامثلثی (یسا راست-مثلثی)  $R$  به صورت  $A^{(k)} = QR$  تجزیه و سپس به شکل زیر تشکیل می شود

$$A^{(k+1)} = RQ = Q^{-1}(QR)Q$$

بنابراین  $A^{(k+1)}$  با  $A^{(k)}$  مشابه است. بعلاوه چون  $A^{(k)}$  يك ماتریس هسبرگ است،  $A^{(k+1)}$  نیز به شکل يك ماتریس هسبرگ است. این امر تعداد عملیات لازم برای به دست آوردن عوامل تجزیهٔ ماتریس فوق را به مقدار زیاد کاهش می دهد. اما در بسیاری از موارد، به ازای مقدار بسیار زیاد  $k$ ،  $A^{(k)}$  به سمت يك ماتریس بالامثلثی که درایه های قطری آن لزوماً تمام ویژه مقدارهای  $B$  هستند، همگرا می شود.

بیان چیزئیات و نظریه‌ای که در پشت این محاسبه وجود دارد ظریف و دشوار است، به‌ویژه اینکه، برای تسریع همگرایی به‌جای خود  $A^{(k)}$ ، ماتریس  $(A^{(k)} - s_k I)$  با تغییر مکان  $s_k$  تجزیه می‌شود. اما خواننده می‌باید از این حقیقت آگاه باشد که این روش و سایر روشها، به‌ویژه آنهایی که برای دسته‌های خاصی از ماتریسهای  $B$  مناسب‌اند، به‌صورت بسته‌های پیش‌ساخته‌ای<sup>۱</sup> از زیر برنامه‌های فورترن موسوم به EISPACK ترجمه شده‌اند که از سوی آزمایشگاه ملی آرگون<sup>۲</sup> یا بسیاری از مراکز محاسباتی علمی مسقیماً در دسترس علاقه‌مندان قرار داده می‌شود، توضیح کامل در مورد بسته‌های پیش‌ساختهٔ فوق، که شامل فهرست برنامه‌هاست، در کتاب اسمیت و دیگران<sup>۳</sup> [۳۲] یافت می‌شود.

### جایابی<sup>۴</sup>

گاهی تنها یک برآورد تقریبی از یک یا تمام ویژه‌مقدارهای  $B$  مورد نظر است. حتی اگر شخص مایل به محاسبهٔ ویژه‌مقدارها باشد، ممکن است که لازم آید از اطلاعاتی در مورد جای تقریبی ویژه‌مقدارها آغاز کند. چنین اطلاعاتی را قضیهٔ جایابی به دست می‌دهد که در صفحهٔ همتافت ناحیه‌هایی را که ویژه‌مقدارهای  $B$  در آنها قرار دارند معین می‌کند.

اگر  $Bx = \lambda x$ ، آنگاه داریم  $\|Bx\| = \|\lambda x\| = |\lambda| \|x\|$ ، که ایجاب می‌کند که در حالت  $x \neq 0$  داشته باشیم،  $|\lambda| \leq \|B\|$ . این مطلب، ثابت می‌کند که رابطهٔ

$$|\lambda| \leq \|B\| \quad \text{به‌ازای هر ویژه‌مقدار } \lambda \text{ در } B \quad (۷۴.۴)$$

و به‌ازای هر نرم ماتریسی برقرار باشد. یک بیان دقیقتر آن در زیر می‌آید:

قضیهٔ ۱۴.۴: **قرصهای گرشگورین<sup>۵</sup>** هر ویژه‌مقدار  $\lambda$  در ماتریس  $B = (b_{ij})$  از مرتبهٔ  $n \times n$  در رابطهٔ زیر صدق می‌کند

$$|b_{ii} - \lambda| \leq \sum_{\substack{j=1 \\ j \neq i}}^n |b_{ij}| \quad \text{به‌ازای مقداری از } i$$

به‌عبارت دیگر، همهٔ ویژه‌مقدارهای  $B$  در ناحیه‌ای که از اجتماع<sup>۶</sup> برخی قرصها در صفحهٔ همتافت پدید می‌آید، قرار دارند. در واقع اگر داشته باشیم

$$|b_{ii} - \lambda| > \sum_{\substack{j=1 \\ j \neq i}}^n |b_{ij}| \quad i = 1, \dots, n$$

- 
- |                 |                                |                |
|-----------------|--------------------------------|----------------|
| 1. package      | 2. Argonne National Laboratory | 3. Smith et al |
| 4. Localization | 5. Gershgorin's disks          | 6. union       |

آنگاه ماتریس  $B - \lambda I$  نافذ قطری (سطراً) مؤکد خواهد بود، از این رو با توجه به تمرین ۳-۶.۴، پس  $B - \lambda I$  وارونپذیر است یعنی  $\lambda$  یک ویژه مقدار  $B$  نیست.

□ مثال ۱۴.۴: به موجب رابطه (۷۴.۴)، قدرمطلق هر ویژه مقدار ماتریس

$$B = \begin{bmatrix} 1 & 2 & 0 \\ 2 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

در مثال ۱۱.۴، می بایسد بزرگتر از  $\|B\|_\infty = 3$  نباشد. قرصهای گرشگورین اطلاعات مفصلتری در اختیار ما قرار می دهد، بدین صورت که هر ویژه مقدار  $\lambda$  از ماتریس  $B$  می باید

$$|1 - \lambda| \leq 2 \quad \text{یا در} \quad |-1 - \lambda| \leq 0$$

□

صدق کند.

در یک ماتریس هرمیتی، به ویژه یک ماتریس متقارن حقیقی، همه ویژه مقدارها حقیقی هستند. این شبیه به یک ماتریس قطری است، یعنی دارای یک مجموعه کامل از ویژه بردارهاست. این یک نتیجه ای است که بد آسانی از قضیه شور به دست می آید، تمرین ۸.۴-۱۵ را ببینید برای یک ماتریس هرمیتی  $B$ ، هم

$$\max_x \frac{x^H B x}{(x^H x)} \quad \text{و} \quad \min_x \frac{x^H B x}{(x^H x)}$$

ویژه مقدارهای  $B$  هستند و سایر ویژه مقدارهای  $B$  بین این دو قرار دارند. به یاد بیاورید که این خارج قسمت های ریالی قبلاً در این قسمت در هنگام بحث از روش توانی آمده اند. از ترکیب لم ۴.۴ و قضیه ۱۰.۴ قضیه جابجایی دقیق زیر به دست می آید.

قضیه ۱۵.۴  $\lambda$  یک ویژه مقدار ماتریس  $B$  است اگر و فقط اگر  $\lambda$  یک ریشه معادله مشخصه زیر باشد

$$\det(B - \lambda I) = 0$$

ماتریس  $(B - \lambda I)$  با ماتریس  $B$  فقط این تفاوت را دارد که از هر یک از درایه های قطری  $B$  مقدار  $\lambda$  کسب شده است. اگر نماد کروکتور<sup>۲</sup>،  $\delta_{ij}$  را برای نشان دادن آرایه  $(j, i)$  ماتریس همانی به کار گیریم، یعنی

$$\delta_{ij} = \begin{cases} 1 & \text{اگر } i = j \\ 0 & \text{اگر } i \neq j \end{cases}$$

آنگاه داریم:



$$(B - \lambda I) \text{ ماتریس } a(i, j) = b_{ij} - \lambda \delta_{ij}$$

بنابراین

$$\begin{aligned} \det(B - \lambda I) &= \sum_p \sigma_p(b_{1p_1} - \lambda \delta_{1p_1})(b_{2p_2} - \lambda \delta_{2p_2}) \dots (b_{np_n} - \lambda \delta_{np_n}) \\ &= \sum_p \sigma_p \left( \prod_{i \neq p_i} b_{ip_i} \right) \left( \prod_{i=p_i} (b_{ip_i} - \lambda) \right) \end{aligned}$$

رابطهٔ بالا نشان می‌دهد که  $\det(B - \lambda I)$  مجموع بسجمله‌ای‌هایی با متغیر  $\lambda$  است. از آنجا که هر حاصلجمع دارای  $n$  عامل است، در نتیجه هر حاصلجمع یک بسجمله‌ای است برحسب  $\lambda$  از درجهٔ حداکثر  $n$ ، در حالی که جمعیتهٔ متناظر با جایگشت همانی  $p^T = [1 \ 2 \ \dots \ n]$  به صورت

$$(b_{11} - \lambda)(b_{22} - \lambda) \dots (b_{nn} - \lambda)$$

و از این رو نسبت به  $\lambda$  دقیقاً از درجهٔ  $n$  است. در نتیجه  $\det(B - \lambda I)$  که به عنوان تابعی از  $\lambda$  مورد بررسی قرار گرفته است، یک بسجمله‌ای است دقیقاً از درجهٔ  $n$  برحسب  $\lambda$

$$p(\lambda) = \det(B - \lambda I) = (-\lambda)^n + n \text{ از درجهٔ کمتر از } (-\lambda)^n$$

این بسجمله‌ای، بسجمله‌ای مشخصهٔ  $B$  نامیده می‌شود.

□ مثال ۱۵.۴: اگر

$$B = \begin{bmatrix} 1 & 2 & 0 \\ 2 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

آنگاه

$$p(\lambda) = \det(B - \lambda I) = \det \begin{bmatrix} 1-\lambda & 2 & 0 \\ 2 & 1-\lambda & 0 \\ 0 & 0 & -1-\lambda \end{bmatrix}$$

و از بسط این دترمینان برحسب آخرین سطر یا آخرین ستون داریم

$$\begin{aligned} p(\lambda) &= (-1-\lambda) \det \begin{bmatrix} 1-\lambda & 2 \\ 2 & 1-\lambda \end{bmatrix} = -(1+\lambda)[(1-\lambda)^2 - 4] \\ &= -(1+\lambda)(1-\lambda+2)(1-\lambda-2) \end{aligned}$$

بنابراین ویژه‌مقدارهای  $A$  یعنی ریشه‌های  $p(\lambda)$  اعداد  $۱ -$  و  $۳$  هستند که قبلاً با روشهای متفاوت دیگری در آغازین بخش محاسبه شدند.  $\square$

از آنجا که  $p(\lambda)$  بسجمله‌ای از درجه  $n$  می‌تواند حداکثر  $n$  ریشه متمایز داشته باشد (به بخش ۱۰.۲ نگاه کنید)، لذا نتیجه می‌شود که  $p(\lambda)$  حداکثر  $n \times n$  ماتریس می‌تواند  $n$  ویژه‌مقدار داشته باشد. از سوی دیگر به موجب قضیه بنیادی جبر، هر بسجمله‌ای از درجه مثبت حداقل یک ریشه دارد (به قضیه ۱۰.۱ نگاه کنید)؛ از این رو برای هر ماتریس مربعی حداقل یک ویژه‌مقدار وجود دارد. این ویژه‌مقدارها، اگر  $B$  یک ماتریس حقیقی باشد، حتی ممکن است هم‌تاف نیز باشند. به موجب قضیه ۱۵.۴، تکنیکهایی که برای یافتن ریشه‌های معادلات (مخصوصاً معادلات بسجمله‌ای که در فصل ۳ بحث شد) به کار می‌روند برای یافتن ویژه‌مقدارها نیز به کار خواهند رفت.

برای مثال، روش درونیایی درجه دوم (روش مولر) را که در بخش ۷.۳ مورد بحث قرار گرفت، می‌توان برای پیدا کردن یک یا چند ویژه‌مقدار حقیقی یا هم‌تاف از یک ماتریس داده شده به کار برد. برای استفاده از این روش فقط باید بتوانیم بسجمله‌ای  $p(\lambda)$  را به ازای هر مقدار  $\lambda$  محاسبه کنیم. از آنجا که برای یک مقدار داده شده  $\lambda$ ،  $p(\lambda)$  اصلاً دترمینانی از مرتبه  $n$  است، از هر روشی که برای محاسبه دترمینان به کار می‌رود می‌توان استفاده کرد. به ویژه، این کار را می‌توان با روش حذفی که در بخش ۷.۴ شرح داده شد، انجام داد. اما بهتر است در ابتدا ماتریس فوق را به ماتریس هسبرگک و یا احتمالاً به شکل ماتریس سه‌قطری که قبلاً مورد بحث قرار گرفته تبدیل کرد. زیرا این کار هزینه محاسبه یک دترمینان را از تعداد عملیات  $\mathcal{O}(n^3)$  به  $\mathcal{O}(n^2)$  کاهش می‌دهد. در هر صورت با اعمال درونیایی مربعی برای یافتن یک ریشه  $\xi_1 = \lambda$  از بسجمله‌ای مشخصه  $p(\lambda) = \det(A - \lambda I)$  به ترتیب زیر عمل می‌کنیم:

- گیریم  $\lambda_0 = -1$ ،  $\lambda_1 = 1$  و  $\lambda_2 = 0$ .

۲. دترمینان زیر را محاسبه کنید

$$p(\lambda_i) = \det \begin{bmatrix} a_{11} - \lambda_i & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} - \lambda_i & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} - \lambda_i \end{bmatrix} \quad i = 0, 1, 2$$

۳. الگوریتم ۱۱.۳ را به کار برید تا همگرایی به سمت ریشه  $\xi_1$  حاصل شود.

۴. برای پیدا کردن ریشه بعدی، این روند را تکرار کنید اما به جای  $p(\lambda)$  تابع اندازه‌

1. deflated function

تقلیل یافته

$$\frac{p(\lambda)}{\lambda - \xi_1}$$

را به کار گیرید

۰۵. به طریقی که در قسمت ۷.۳ شرح داده شد کار را ادامه دهید.

روش درونیابی درجهٔ دوم در مقایسه با روشهای پیشرفته تر از نظر کارایی محاسباتی چندان قابل رقابت نیست. ولی به کار گرفتن این روش بسیار آسان و کاملاً کلی و تقریباً همیشه همگراست و در بیشتر موارد دقت عمل کافی به دست می‌دهد. این روش را می‌توان برای حل مسئلهٔ کلیتر ویژه مقدار زیر نیز به کار گرفت

$$\det(A - \lambda B) = 0$$

در رابطهٔ بالا  $A$  و  $B$  ماتریسهایی از مرتبهٔ  $n$  هستند.

□ مثال ۱۶.۴: ارتعاشهای آزاد ساختارهای ساده. در مهندسی راه و ساختمان مسئله‌ای که غالباً با آن مواجه می‌شویم، مسئلهٔ تعیین فراوانی طبیعی  $\lambda$ ی ارتعاشهای آزاد یک ساختار غیر میرا است که در قبال چندین جرم و چندین درجهٔ آزادی<sup>۳</sup> حاصل می‌شود. این مسئله را می‌توان به شکل

$$Ax = \lambda Mx \quad (۷۵.۴)$$

بیان کرد که در آن

$M$  = ماتریس جرم دستگاه

$A$  = سختی<sup>۴</sup> دستگاه

$x$  = تغییر مکان طبیعی

از آنجا که (۷۵.۴) معرف یک دستگاه معادلات همگن است، اگر دترمینان ضرایب آن صفر باشد، یعنی اگر تساوی

$$\det(A - \lambda M) = 0 \quad (۷۵.۴ \text{ الف})$$

برقرار باشد، آنگاه یک جواب غیر بدیهی<sup>۵</sup>  $x$  خواهد داشت. بنا بر این اگر ماتریسهای  $A$  و  $M$  داده شده باشند مقادیری از  $\lambda$  که در (۷۵.۴ الف) صدق می‌کنند فراوانیهای طبیعی هستند. روش مولر را می‌توان مستقیماً برای پیدا کردن این ویژه مقادیرها به کار برد. مثلاً در یک دستگاه معین با  $M = I$  و ماتریس سختی  $A$  به صورت

1. structure

2. undamped

3. degrees of freedom

4. stiffness

$$A = \begin{bmatrix} 4 & -1 & -1 & -1 \\ -1 & 4 & -1 & -1 \\ -1 & -1 & 4 & -1 \\ -1 & -1 & -1 & 4 \end{bmatrix}$$

فراوانی طبیعی  $\lambda_i$  این دستگاه را پیدا کنید.

يك برنامه کامپیوتری با استفاده از حذف گاوسی ۲۰۴ برای محاسبه دترمینانها و الگوریتم ۱۱.۳ مولر به عنوان ریشه یاب، ویژه مقادارها را به شرح زیر برآورد نموده است.

$$\lambda_1 = 100000000$$

$$\lambda_2 = 500000234$$

$$\lambda_3 = 499999954$$

$$\lambda_4 = 499999973$$

به آسانی دیده می شود که مقدار دقیق ویژه مقادارها به ترتیب برابر ۱، ۵، ۵، ۵ هستند. با این مثال کارایی روش مولر به عنوان يك ریشه یاب، نشان داده شده است که در اینجا يك ریشه سه گانه<sup>۱</sup> با دقت عمل نسبتاً خوبی محاسبه شده است.

□ مثال ۱۷.۴: درایه های يك ماتریس سه قطری  $B$  به صورت زیر تولید شده اند

$$b_{i,i+1} = 0.5 \quad i = 1, 2, \dots, n-1$$

$$b_{i+1,i} = 0.5 \quad i = 1, 2, \dots, n-1$$

$$b_{ij} = 0 \quad \text{برای تمامی } i \text{ و } j \text{ دیگر}$$

يك برنامه کامپیوتری بنویسید که ویژه مقادارهای  $B$  را به ازای  $n = 20$  پیدا کند. به ازای  $n = 20$ ، روش مولر در کامپیوتر IBM ۷۰۹۴ نتایج زیر را تولید کرده است:

$$\pm 0.074730093 \quad \pm 0.73305187$$

$$\pm 0.22252094 \quad \pm 0.82623877$$

$$\pm 0.36534102 \quad \pm 0.90096886$$

$$\pm 0.500000001 \quad \pm 0.95557281$$

$$\pm 0.62348980 \quad \pm 0.98883083$$



بیان گردید (تمرین ۷.۴-۱۱ را نگاه کنید)، (۷۶.۴) ثابت می‌شود. خصوصیت بازگشتی<sup>۱</sup> (۷۶.۴) این امکان را فراهم می‌سازد که  $p_n(\lambda)$  با تعداد عملیاتی در حدود  $3n$  محاسبه شود. بعلاوه به راحتی می‌توان از رابطه بازگشتی فوق نسبت به  $\lambda$  مشتق گرفت، و این عمل، محاسبه  $p'_n(\lambda)$  را به صورت بازگشتی امکانپذیر و بنابراین کاربرد روش نیوتن را میسر می‌سازد.

اگر به ازای مقداری مانند  $i$  داشته باشیم  $b_i = 0$  آنگاه از رابطه بازگشتی (۷۶.۴) می‌فهمیم که بسجمله‌ای  $p_{i-1}(\lambda)$  عاملی از  $p_n(\lambda) = \det(B - \lambda I)$  است. بنابراین ریشه‌های  $p_n(\lambda)$ ، ریشه‌های دوسجمله‌ای  $p_{i-1}(\lambda)$  و  $p_n(\lambda)/p_{i-1}(\lambda)$  با درجه کمتر هستند و می‌توان کار را بر این دوسجمله‌ای متمرکز کرد. در غیر این صورت اگر به ازای جمیع مقادیر  $i$ ،  $b_i \neq 0$ ، آنگاه  $B$  دارای  $n$  ویژه‌مقدار متمایز است. همچنین، دنباله  $p_0(\lambda), p_1(\lambda), \dots, p_n(\lambda)$  از مقادیر محاسبه شده در طی محاسبه  $p_n(\lambda)$ ، اطلاعات اضافی زیر را به همراه دارد: تعداد تغییر علامتها (بسا قاطعیت) در دنباله مذکور با تعداد ویژه‌مقادیر  $B$ ، که از  $\lambda$  کمتر است برابر است، این امر به علت این واقعیت است که بسجمله‌ایهای  $p_0(\lambda), p_1(\lambda), \dots, p_n(\lambda)$  يك دنباله استورم<sup>۲</sup> تشکیل می‌دهند که به سرعت بازه‌هایی متضمن تنها يك ویژه‌مقدار را پدید می‌آورند.

□ مثال ۱۸.۴: برای ماتریس مثال ۱۷.۴، رابطه بازگشتی (۷۶.۴) به صورت ساده زیر درمی‌آید

$$p_0(\lambda) = 1, p_1(\lambda) = -\lambda, p_j(\lambda) = -(\lambda p_{j-1}(\lambda) + p_{j-2}(\lambda)/4) \quad j > 1$$

با انتخاب  $n = 10$  و  $\lambda = 0$  دنباله زیر حاصل می‌شود که پنج تغییر علامت دارد

$$(p_j(\lambda)) = 1, 0, -\frac{1}{4}, 0, \frac{1}{16}, 0, -4^{-3}, 0, 4^{-4}, 0, -4^{-5} \\ = -0.000976563$$

به ازای  $\lambda = 0.2$ ، در عوض دنباله زیر حاصل می‌شود که شش تغییر علامت را نشان می‌دهد.

$$(p_j(\lambda)) = 1, -0.2, -0.21, 0.092, 0.034, -0.029, \\ -0.0025, 0.0079, -0.00095, -0.00018, 0.000059857$$

[در اینجا جز برای مقدار  $p_{10}(\lambda)$ ، اعداد تا دو رقم اول با معنی داده شده‌اند.] در نتیجه در بازه  $[0, 0.2]$  دقیقاً يك ویژه‌مقدار  $B$  وجود دارد. الگوریتم قاعده تصحیح خطای

اصلاح شده (الگوریتم ۳.۳) که با این بازه آغاز می‌شود در چهار مرحله (دریک هیولت-پاکارد ۶۷) ویژه مقدار  $142314837$  را به دست می‌دهد که با ویژه مقدار درست  $\cos(5\pi/11) = 142314838$  (تمرین ۸.۴-۲ را ببینید). □

### تمرین

۱-۸.۴ گیریم  $a, b$  عدد دوار (اسکالر) و  $A$  یک ماتریس مربعی باشد. ثابت کنید که اگر  $\lambda$  یک ویژه مقدار  $A$  باشد، آنگاه  $a\lambda + b$  یک ویژه مقدار ماتریس  $aA + bI$  خواهد بود. [داهنمایی:  $(aA + bI)x$  را که در آن  $x$  یک ویژه بردار متعلق به  $\lambda$  است، در نظر بگیرید.]

۲-۸.۴ ثابت کنید که اگر  $\lambda$  یک ویژه مقدار ماتریس مربعی  $A$  و  $p(x)$  یک بسجمله‌ای غیر مشخصی باشد، آنگاه  $p(\lambda)$  یک ویژه مقدار  $p(A)$  خواهد بود (تمرین ۱۰.۴-۱۲ را ببینید).

۳-۸.۴ گیریم  $A$  یک ماتریس سه قطری از مرتبه  $n$  با درایه‌های قطری صفر باشد و به ازای  $i = 1, \dots, n-1$ ,  $j = 1, \dots, n$ ,  $a_{i, i+1} = a_{i+1, i} = a_i$ . به ازای  $i = 1, \dots, n$ ,  $x^{(j)}$  برداری باشد که درایه  $i$ ام آن به ازای  $i = 1, \dots, n$  به صورت

$$x_i^{(j)} = \sin[ij\pi/(n+1)]$$

است. ثابت کنید که

$$Ax^{(j)} = 2 \cos\left(\frac{j\pi}{n+1}\right) x^{(j)} \quad j = 1, \dots, n$$

۴-۸.۴ با استفاده از تمرینهای ۱-۸.۴ و ۳-۸.۴ ثابت کنید که اگر  $A$  یک ماتریس سه قطری باشد که در آن به ازای جميع مقادیر  $i$ ,  $a_{ii} = d$  و  $a_{i, i+1} = a_{i+1, i} = e$ ، آنگاه همه ویژه مقدارهای  $A$  اعدادی به شکل زیر خواهند بود

$$d + 2e \cos \frac{j\pi}{n+1} \quad j = 1, \dots, n$$

۵-۸.۴ با استفاده از روش توانی ویژه مقدار با قدر مطلق ماکسیمم و ویژه بردار متناظر با آن را برای ماتریس سه قطری  $A$  از مرتبه ۲۵، با  $a_{ii} = 4$  و  $a_{i, i+1} = a_{i+1, i} = -1$  به ازای جميع مقادیر  $i$ ، بر آورد کرده، بسا جواب دقیق حاصله از تمرین ۲-۸.۴ مقایسه کنید.

۶-۸.۴ سعی کنید که ویژه مقدار با قدرمطلق ماکسیمم ماتریس  $A$  در مثال ۳-۸.۴ (بافرض  $n=21$ ) را با استفاده از روش توانی برآورد کنید. در مورد مشکلاتی که با آن مواجه می شود، توضیح دهید.

۷-۸.۴ اگر ماتریسی دویا چندین ویژه مقدار با قدرمطلق ماکسیمم برابر داشته باشد، روش توانی کارا نیست. در مورد اینکه چگونه تمرین ۱-۸.۴ می تواند برای این مشکل فائق آید، بحث و سعی کنید که این راه رفع مشکل را در مورد مسئله تمرین ۶-۸.۴ به کار برید.

۸-۸.۴ نشان دهید که ماتریس  $B = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$  دارای یک مجموعه کامل از ویژه بردارها

نیست.

۹-۸.۴ گیریم  $x$  و  $y$  دو ویژه بردار ماتریس  $A$ ، به ترتیب متعلق به ویژه مقدارهای  $\lambda$  و  $\mu$  باشند. نشان دهید که اگر  $\lambda \neq \mu$ ، آنگاه  $x$  و  $y$  مستقل خطی هستند.

۱۰-۸.۴ با استفاده از تمرین ۹-۸.۴ نشان دهید که ماتریس  $\begin{bmatrix} 1 & 2 \\ 0 & 2 \end{bmatrix}$  می یابد مجموعه

کاملی از ویژه بردارها داشته باشد.

۱۱-۸.۴ کلیه ویژه مقدارهای ماتریس

$$A = \begin{bmatrix} 4 & -1 & -1 & -1 \\ -1 & 4 & -1 & -1 \\ -1 & -1 & 4 & -1 \\ -1 & -1 & -1 & 4 \end{bmatrix}$$

را با تعیین صریح بسجمله ای مشخصه آن محاسبه کنید، و سپس ریشه این بسجمله ای را به دست آورید.

۱۲-۸.۴ ماتریس  $A$  در تمرین ۱۱-۸.۴ را از راه تقارنهای محوری هاوس هولدر به شکل سه قطری  $B$  تبدیل کنید (چون بسجمله ای مشخصه  $A$  یک ریشه سه گانه دارد، بر طبق مثال ۱۶.۴، حداقل دوتا از  $b_i$ ها صفرند)، سپس ویژه مقدارهای  $B$  را محاسبه کنید.

۱۳-۸.۴ کلیه ویژه مقدارهای ماتریس سه قطری  $B$  در مثال ۱۷.۴ را با استفاده از رابطه بازگشتی (۷۶.۴)، از ویژگی دنباله استورم برای جدانمودن ویژه مقادیرها، و سپس از روش نیوتن برای دستیابی به تک تک ویژه مقدارها، حساب کنید.



۸۰۴-۱۴ پس از حل تمرین ۸۰۴-۱۳، روش توانی معکوس<sup>۱</sup> را برای تعیین ویژه بردارهای متناظر به کاربرد.

۸۰۴-۱۵ ثابت کنید که یک ماتریس هرمیتی با یک ماتریس قطری متشابه است و همه ویژه مقادیر آن حقیقی هستند. (داهنمایی: نشان دهید که اگر  $B$  یک ماتریس هرمیتی باشد، آنگاه ماتریس بالامثلثی که در قضیه شور به دست آمد، لزوماً هرمیتی است).

۸۰۴-۱۶ از روش مولر برای تعیین فراوانیهای طبیعی مثال ۱۶۰۴، دمورد

$$A = \begin{bmatrix} -2 & 0 & 1 & 0 \\ 0 & -2 & 0 & 1 \\ 1 & 0 & -2 & 0 \\ 0 & 1 & 0 & -2 \end{bmatrix} \quad M = \begin{bmatrix} 2 & 1 & 0 & 0 \\ 1 & 2 & 1 & 0 \\ 0 & 1 & 2 & 1 \\ 0 & 0 & 1 & 2 \end{bmatrix}$$

استفاده کنید.

۸۰۴-۱۷ فرض کنید که ماتریس  $A$  از مرتبه  $n$ ، دارای یک مجموعه کامل از ویژه بردارهای  $x^{(1)}, \dots, x^{(n)}$  باشد. ثابت کنید که  $A$  با یک ماتریس قطری (که درایه‌های قطری آن باید ویژه مقادیر  $A$  باشند) مشابه است. [داهنمایی: ماتریس  $C^{-1}AC$  را که در آن به ازای  $1, \dots, n$  تساوی  $Ci_j = x^{(j)}$  برقرار است، در نظر بگیرید. چرا  $C$  می‌باید وارون پذیر باشد؟].

۸۰۴-۱۸ تقلیل اندازه برای روش توانی فرض کنید که با روش توانی یسا روشهای دیگر، یک ویژه مقدار  $\lambda$  از ماتریس  $A$  از مرتبه  $n$  و ویژه بردار  $x$  متناظر با آن را محاسبه کرده‌ایم، و فرض کنید که  $x_n \neq 0$  گیریم.  $B$  ماتریسی از مرتبه  $n-1$  باشد که از ماتریس  $C^{-1}AC$  با حذف آخرین سطر و آخرین ستونش به دست آمده و در آن، تساویهای  $Ci_j = i_j$  و  $Ci_n = x$ ، به ازای  $1, \dots, n-1$ ، برقرارند. ثابت کنید که کلیه ویژه مقادیر  $A$ ، احتمالاً به استثنای ویژه مقدار  $\lambda$ ، ویژه مقادیر  $B$  نیز هستند.



## دستگاه معادلات و بهینه‌سازی نامقید

یک دستگاه کلی  $n$  معادله و  $n$  مجهول بر حسب  $x_1, x_2, \dots, x_n$  را می‌توان همواره به شکل

$$f_i(x_1, x_2, \dots, x_n) = 0 \quad i = 1, \dots, n \quad (1.5)$$

نوشت که در آن  $f_1, \dots, f_n$  توابع  $n$  متغیره هستند. نمادگذاری برداری مذکور در فصل ۲ را همچنان ادامه می‌دهیم و معادله (۱.۵) را به صورت فشرده‌تر زیر می‌نویسیم

$$\mathbf{f}(\mathbf{x}) = \mathbf{0} \quad (2.5)$$

بنا بر این  $\mathbf{f}$  تابعی است برداری مقداراً. مقدار این تابع به ازای  $n$ -برداری

$$\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_n(\mathbf{x})]^T \quad \mathbf{x} = [x_1, x_2, \dots, x_n]^T$$

است.

با این نمادگذاری نه تنها در نوشتن صرفه‌جویی شده است بلکه بیانگر این حقیقت است که روشهای بارستی که برای حل یک معادله یک مجهولی که در فصل ۳ مورد بحث قرار گرفت، در اینجا هم از جهتی باید به کار آیند. بخصوص روش بارست نقطه ثابت و روش نیوتن و بعضی از شکل‌های مختلف این روش مورد بحث قرار می‌گیرند. ولی نمی‌توانیم تجزیه و تحلیل ریاضی این روشها را عمیقاً مورد بحث قرار دهیم. بحث جامع از کلیه مطالب

موجود در این زمینه را می توان در شرح تفصیلی ارتکا و رایبولت [۳۳] برای این موضوع، به دست آورد. حل دستگاه معادلات همچنان یک زمینه فعال تحقیقاتی، بخصوص در ساختن الگوریتمهای کاراست.

یک مثال خاص از دستگاه (۱.۵) دستگاه خطی

$$Ax - b = 0$$

است که حل مستقیم آن مفصلاً در فصل ۴ مورد بحث قرار گرفت، ولی در اینجا دستگاه کلی (۲.۵) معمولاً به صورت بارستی می باید حل شود؛ یعنی، بسا حل یک رشته از دستگاههای خطی که معمولاً با روشهای مستقیم بحث شده در فصل ۴، انجام می گیرد. اما در بعضی از روشهای بارستی بخصوص روشهای واهلش<sup>۲</sup>، آن رشته از دستگاههای خطی که می باید حل شوند به اندازه ای ساده هستند که صرفه در این است که این روشها، برای سیستمهایی که خود خطی می باشند به کار گرفته شوند (و به کار گرفته شده اند). در حل دستگاههای خطی، ما به این گونه روشهای بارستی توجه خاصی خواهیم کرد.

وسر انجام، رابطه نزدیک بین حل دستگاههای معادلات و جستجو برای فرینه های یک تابع حقیقی  $n$ -متغیره را که در قسمت اول این فصل بیشتر توضیح داده خواهد شد، مطمح نظر قرار می دهیم.

### ۱.۵\* بهینه سازی و تندترین کاهش<sup>۳</sup>

بهینه سازی منبعی است همیشگی برای حل دستگاههای معادلات، و پاره ای از روشهای حل دستگاهها مستقیماً تحت تأثیر آن قرار دارند. برای یادآوری، اگر بخواهیم مینیمم (بسا ما کسیمم) یک تابع حقیقی  $n$ -متغیره  $F(x) = F(x_1, \dots, x_n)$  را پیدا کنیم، کافی است که به مقادیر آن تابع در نقاط بحرانی<sup>۴</sup> یعنی نقاطی از  $x$  که به ازای آنها رابطه

$$\nabla F(x) = 0$$

برقرار است، نگاه کنیم. در اینجا  $\nabla F$  گرادیان<sup>۵</sup>  $F$  یعنی بردار

$$\nabla F = \left[ \frac{\partial F}{\partial x_1} \quad \dots \quad \frac{\partial F}{\partial x_n} \right]^T$$

است که درایه های آن اولین مشتقهای جزئی تابع  $F$  هستند

$$f_i := \frac{\partial F}{\partial x_i} \quad i = 1, \dots, n$$

- 
1. Ortega and Rheinholdt      2. relaxation      3. steepest descent  
4. critical points                5. gradient

بدین علت  $\nabla F(\mathbf{x})$  نوشته می‌شود که می‌خواهیم به نقطه  $\mathbf{x}$  که گرادیان در آن نقطه محاسبه شده تأکید کنیم.

به یاد آورید که گرادیان  $\nabla F$  همان «مشتق اول» تابع  $n$ -متغیره  $F(\mathbf{x})$  است: بنابراین قضیه ۸.۱ مشتق تابع یک متغیره

$$g(t) = F(\mathbf{x} + t\mathbf{u})$$

در  $t = 0$  به صورت زیر معین می‌شود

$$g'(0) = \nabla F(\mathbf{x})^T \mathbf{u} = f_1(\mathbf{x})u_1 + \dots + f_n(\mathbf{x})u_n$$

عدد بالا اطلاعاتی درباره رفتار تابع  $F$ ، هنگامی که بخواهیم از نقطه  $\mathbf{x}$  در جهت  $\mathbf{u}$  حرکت کنیم، به دست می‌دهد. بنابراین  $F$  در کلیه جهتهای  $\mathbf{u}$  که زاویه اش نسبت به بردار گرادیان  $\nabla F(\mathbf{x})$  کمتر از  $90^\circ$  باشد با بیشترین سرعت در جهت گرادیان، افزایش می‌یابد. دلیل این امر وجود رابطه زیر است

$$(\nabla F)^T \mathbf{u} = \|\nabla F\|_2 \|\mathbf{u}\|_2 \cos \theta. \quad (3.5)$$

که در آن  $\theta$  زاویه بین دو بردار است. درحقیقت این نوعی اتحاد منطقی است زیرا اگر از ما خواسته شود که زاویه  $\theta$  بین دو بردار  $\mathbf{u}$  و  $\mathbf{v}$  را تعیین کنیم، معمولاً به جواب زیر متوسل می‌شویم

$$\theta = \cos^{-1} \left( \frac{\mathbf{v}^T \mathbf{u}}{\|\mathbf{v}\|_2 \|\mathbf{u}\|_2} \right)$$

اما نکته حائز اهمیت این است که، به ازای  $\nabla F \neq 0$ ، تقریباً نیمی از جهتهای ممکن  $\mathbf{u}$ ، مثلاً برای آنهایی که داریم  $(\nabla F)^T \mathbf{u} > 0$ ، موجب افزایش مقدار  $F$  می‌شوند و حداکثر افزایش فقط و فقط زمانی حاصل می‌شود که  $\mathbf{u}$  با  $\nabla F$  موازی باشد. با بحنی مشابه، تقریباً نیمی از تمام جهتهای ممکن  $\mathbf{u}$  موجب کاهش مقدار  $F$  می‌شوند و بیشترین کاهش زمانی است که  $\mathbf{u}$  با  $-\nabla F$  موازی باشد. در نتیجه  $\mathbf{x}$  نمی‌تواند جز در حالت  $\nabla F = 0$ ، برای  $F$  مینیمم یا ماکسیمم باشد.

□ مثال ۱.۵: تابع  $F(\mathbf{x}) = x_1^2 + x_2^2 - 2x_1^2 + 3x_2^2 - 8$  دارای گرادیان

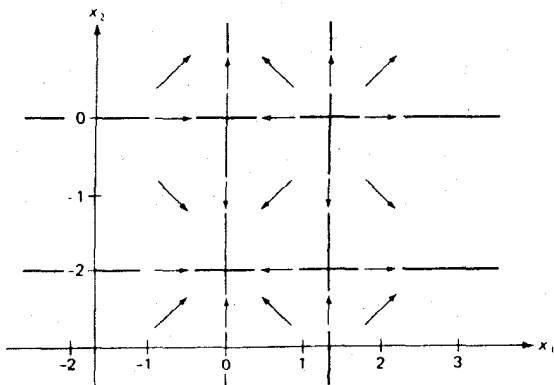
$$\nabla F = [3x_1^2 - 4x_1 \quad 3x_2^2 + 6x_2]^T$$

است. بنابراین معادله  $\nabla F(\mathbf{x}) = 0$  دارای چهار جواب  $(0, 0)$ ،  $(0, -2)$ ،  $(4/3, 0)$  و  $(4/3, -2)$  است. برای درک ماهیت این نقاط بحرانی، و برای آنکه تمرینی بسا گرادیانها داشته باشیم، اکنون نواحی مختلفی را که در آنها صفحه  $(x_1, x_2)$  به وسیله منحنیهای زیر قطع می‌شوند مورد بررسی قرار می‌دهیم

$$\left(\frac{\partial F}{\partial x_1}\right)(\mathbf{x}) = 0, \quad \left(\frac{\partial F}{\partial x_2}\right)(\mathbf{x}) = 0$$

ملاحظه می‌کنیم که  $f_1 = \partial F / \partial x_1$  روی دوخط مستقیم  $x_1 = 0$  و  $x_1 = 4/3$  صفر می‌شود و بین این خطوط مقدار آن منفی و در جاهای دیگر مثبت است. بنا بر این اولین مؤلفهٔ گرادیان  $\nabla F$  بین این دوخط منفی و در جاهای دیگر مثبت است. همچنین  $f_2 = \partial F / \partial x_2$  که روی دوخط مستقیم  $x_2 = -2$  و  $x_2 = 0$  صفر می‌شود، بین این خطوط منفی و در جاهای دیگر مثبت است. بحث فوق تصویر کیفی شکل ۱۰۵ را برای جهت گرادیان در نواحی مختلف که به وسیلهٔ خطوط  $f_1 = 0$  و  $f_2 = 0$  تعریف می‌شوند به دست می‌دهد. از روی این شکل مشخص است که نقطهٔ بحرانی  $(0, -2)$  یک ماکسیمم موضعی است (زیرا کلیهٔ گرادیانهای مجاور آن، به این نقطه متوجه‌اند) در حالی که نقطهٔ بحرانی  $(4/3, 0)$  یک مینیمم موضعی است (زیرا کلیهٔ گرادیانهای مجاور آن از این نقطه رو برمی‌گردانند). دو نقطهٔ بحرانی دیگر، نقاط فرین<sup>۲</sup> نیستند بلکه نقاط زینی<sup>۳</sup> هستند، زیرا در مجاورت آنها گرادیانهای قرار دارند که هم به آنها متوجه‌اند و هم از آنها رو برمی‌گردانند. □

یک روش اساسی برای پیدا کردن نقطهٔ فرین، روش تندترین کاهش (یا افزایش) است. این روش به روش کوشی<sup>۴</sup> برمی‌گردد و در آن سعی بر آن است که مسئله پیدا کردن مینیمم یک تابع حقیقی  $n$  متغیره، از راه چندبار پیدا کردن مینیمم یک تابع یک متغیره حل شود. فکر اساسی به شرح زیر است: یک تقریب  $\mathbf{x}$  برای مقدار مینیمم  $\mathbf{x}^*$  از تابع  $F$  داده شده است. سپس در طول خط مستقیم ماربر  $\mathbf{x}$  و در جهت  $-\nabla F(\mathbf{x})$  به دنبال آن مینیممی از



شکل ۱۰۵ نمایش جهت‌های گرادیان برای تابع  $F(x_1, x_2) = x_1^2 + x_2^3 - 2x_1^2 + 3x_2^2 - 8$

1. component
2. extrema
3. saddle
4. Cauchy

از  $F$  که به  $\mathbf{x}$  از همه نزدیکتر است می‌گردیم. یعنی مینیمم  $t^* > 0$  را که نزدیکترین نقطه به صفر تابع

$$g(t) = F(\mathbf{x} - t \nabla F(\mathbf{x}))$$

که تابعی است از یک متغیر تصادفی<sup>۱</sup> پیدا کنیم، و پس از پیدا کردن آن، تقریب بعدی برای مینیمم  $\mathbf{x}^*$  را نقطه  $\mathbf{x} - t^* \nabla F(\mathbf{x})$  می‌گیریم.

**الگوریتم ۱۰۵:** تندترین کاهش تابع هموار  $F(\mathbf{x})$  از  $n$ -بردار  $\mathbf{x}$  و یک تقریب  $\mathbf{x}^{(0)}$  برای مینیمم (موضعی)  $\mathbf{x}^*$  از  $F$ ، داده شده‌اند.

For  $m = 0, 1, 2, \dots$ , do until satisfied :

$$\mathbf{u} := \Delta F(\mathbf{x}^{(m)})$$

اگر  $\mathbf{u}$  مساوی ۰ شد، آنگاه الگوریتم را متوقف سازید.

در غیر این صورت مینیمم  $t^* > 0$  را که به ۰ تابع

$g(t) = F(\mathbf{x}^{(m)} - t\mathbf{u})$  از همه نزدیکتر باشد تعیین کنید

$$\mathbf{x}^{(m+1)} := \mathbf{x}^{(m)} - t^* \mathbf{u}$$

□ **مثال ۲۰۵:** برآورد  $\mathbf{x}^{(0)} = [1, -1]^T$  برای مینیمم موضعی  $(0, 4/3)$  از تابع

$F(x_1, x_2) = x_1^2 + x_2^2 - 2x_1 + 3x_2^2 - 8$  در مثال ۱۰۵ داده شده‌است، داریم

$$\nabla F(\mathbf{x}^{(0)}) = [-1, -3]^T$$

بنابراین در اولین مرحله از روش تندترین کاهش، به دنبال یافتن مینیمم تابع

$$g(t) = F(1+t, -1+3t)$$

$$= (1+t)^2 + (-1+3t)^2 - 2(1+t) + 3(-1+3t)^2 - 8$$

می‌رویم. با قراردادن  $g'(t) = 0$  خواهیم داشت

$$0 = 3(1+t)^2 + 3(3t-1)^2 - 2(1+t) + 3 \times 2(3t-1)$$

$$= 84t^2 + 2t - 10$$

که دارای دو جواب  $5/14$  - یا  $1/3$   $= (-2 \pm \sqrt{4 + 3360}) / 168$  است  $t^*$  است (با

استفاده از فرمول حل معادله درجه دوم). ریشه مثبت  $t^* = 1/3$  را انتخاب می‌کنیم زیرا

می‌خواهیم از  $\mathbf{x}^{(0)}$  در جهت گرادیان  $-\nabla F(\mathbf{x}^{(0)})$  حرکت کنیم. این انتخاب، خود مینیمم

□ یعنی  $\mathbf{x}^{(1)} = [4/3, 0]^T$  را به دست می‌دهد.

واضح است که روش تندترین کاهش، ضامن کاهش مقدار تابع در هر مرحله است، یعنی برای اثبات همگرایی (با فرض  $\|x^{(m)}\| \geq \epsilon$  يك عدد ثابت، به ازای جميع مقادیر  $m$ ) روش فوق قرارداد. ولی به سادگی می توان مثالهایی ارائه داد که نشان دهند روش فوق ممکن است خیلی به کندی همگرا شود.

□ مثال ۳۰۵: تابع  $F(x) = x_1^2 + \alpha x_2^2$  با  $\alpha > 0$ ، دارای يك مینیمم کلی در  $x = 0$  است. گرادیان این تابع، خطی و به صورت

$$\nabla F = [2x_1, 2\alpha x_2]^T$$

است بنا بر این بلافاصله می توان نقطهٔ بحرانی یگانهٔ آن را از دستگاه

$$2x_1 = 0$$

$$2\alpha x_2 = 0$$

به دست آورد. اما برای اینکسه نکته ای را متذکر شویم به جای استفاده از روش فوق، از روش تندترین کاهش استفاده می کنیم، این امر مستلزم تعیین مینیمم تابع

$$g(t) = F(x - t\nabla F(x))$$

$$= F(x_1(1-2t), x_2(1-2\alpha t))$$

است. با قراردادن  $g'(t) = 0$ ، معادلهٔ زیر را به دست خواهیم آورد

$$0 = 2(x_1(1-2t))(-2) + \alpha 2(x_2(1-2\alpha t))(-2\alpha)$$

که جواب آن عبارت است از  $t^* = (1/2)(x_1^2 + \alpha^2 x_2^2) / (x_1^2 + \alpha^2 x_2^2)$  بنا بر این اگر حدس معمولی ما  $x = [x_1, x_2]^T$  باشد، آنگاه

$$\frac{x_1 x_2 (\alpha - 1)}{x_1^2 + \alpha^2 x_2^2} [\alpha^2 x_2 - x_1]^T$$

حدس بعدی ما خواهد بود.

اکنون  $x$  را به شکل خاص  $c[\alpha, \pm 1]^T$  می گیریم. در این حال حدس بعدی ما

به صورت

$$c \frac{\alpha - 1}{\alpha + 1} [\alpha, \mp 1]^T$$

خواهد شد. یعنی خطا به نسبت  $(\alpha - 1)/(\alpha + 1)$  کاهش یافته است. برای مثال اگر  $\alpha = 100$  و  $\mathbf{x}^{(0)} = [1, 0.001]^T$ ، آنگاه بعد از ۱۰۰ مرحله استفاده از روش تندترین کاهش، نقطهٔ زیر را خواهیم داشت

$$\mathbf{x}^{(100)} = \left(\frac{\alpha - 1}{\alpha + 1}\right)^{100} [1, 0.001]^T = [0.0135\dots, 0.000135\dots]^T$$

□ که نسبت به حدس اولیه، جواب باز به مقدار  $7/8$  کمتر است.

در شکل ۲.۵ قسمتی از بارست تندترین کاهش برای مثال ۲.۵ نشان داده شده است. برای فهمیدن این شکل می‌باید به دو نکتهٔ زیر توجه کرد:

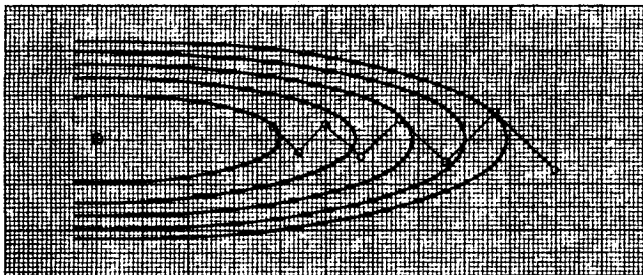
(i) از آنجا که  $(d/dt)F(\mathbf{x} + t\mathbf{u}) = (\nabla F(\mathbf{x} + t\mathbf{u}))^T \mathbf{u}$  در نقطهٔ مینیمم  $\mathbf{x} - t^* \nabla F(\mathbf{x})$  از تابع  $F$ ، در جهت گرادیان منفی، برامتداد آن عموداً است، یعنی

$$\nabla F(\mathbf{x}^{(m+1)})^T \nabla F(\mathbf{x}^{(m)}) = 0$$

(ii) یک تابع دو متغیرهٔ  $F(x_1, x_2)$  در صفحهٔ  $(x_1, x_2)$  غالباً به وسیلهٔ خطوط تراز یا خطوط مرزی آن یعنی با منحنیهای

$$F(x_1, x_2) = \text{ثابت}$$

توصیف می‌شود. این خطوط در شکل ۲.۵ نشان داده شده‌اند. از آنجا که الزاماً گرادیان در یک نقطه بر خط تراز ماربر آن نقطه عمود است (تمرین ۱.۵-۳) بنا بر این خطوط مزبور اطلاعاتی دربارهٔ جهت گرادیان به دست می‌دهند. همان طوری که این مثال نشان می‌دهد انتخاب جهت تندترین کاهش ممکن است تدبیر



شکل ۲.۵ زمانی که به دنبال یک مینیمم در یک محدودهٔ باریک باشیم، روش تندترین کاهش ممکن است به نحو غیر مؤثری موجب پس‌و‌پیش شدن شود.



جوابی به حساب آید، اما غالباً طرح دقیق و مناسبی نیست. امروزه روشهای کاهش پیچیده تری به کار می رود که در آن  $\mathbf{x}^{(m+1)}$  از  $\mathbf{x}^{(m)}$  به گونه زیر به دست می آید

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} + t_m \mathbf{u}^{(m)} \quad (۲.۵)$$

در اینجا،  $\mathbf{u}^{(m)}$  يك جهت کاهش است، یعنی  $\nabla F(\mathbf{x}^{(m)})^T \mathbf{u}^{(m)} < 0$  و  $t_m$  از راه خط جویبی<sup>۲</sup>، یعنی با مینیمم کردن تقریبی تابع

$$g(t) = F(\mathbf{x}^{(m)} + t\mathbf{u}^{(m)})$$

پیدا می شود. اگر گرادیان  $F$  در دسترس باشد، آنگاه عمل خط جویبی به پیدا کردن ریشه مناسب تابع

$$g'(t) = \nabla F(\mathbf{x}^{(m)} + t\mathbf{u}^{(m)})^T \mathbf{u}^{(m)}$$

تبدیل می شود، و برای این منظور روشهای فصل سه را می توان به کار گرفت. گرچه باید به یاد داشته باشیم که دقتی که همراه این ریشه است به درجه نزدیکی ریشه مزبور به مینیمم  $F(\mathbf{x})$  بستگی دارد.

اگر گرادیان  $F$  در دسترس نباشد (یا تصور شود که محاسبه آن پرهزینه است) آنگاه پیشنهاد می شود که به شکلی از درونیابی<sup>۳</sup> درجه دوم استفاده به عمل آید. در زیر نمونه ای از آن آمده است.

**الگوریتم ۲.۵:** خط جویبی از راه درونیابی درجه دوم تابع  $g(t)$  با  $g'(0) < 0$ ، و يك عدد مثبت  $t_{\max}$  و تحمل (مساهاله<sup>۴</sup>) مثبت  $\epsilon$  داده شده اند

$$s_1 := 0 \quad ۱.$$

$$s_2 := s_1 \quad ۲ \text{ اگر } s_2 \leq t_{\max} \text{ را طوری انتخاب کنید که } s_1 \leq s_2 \leq t_{\max} \text{ و } 0 \leq g[s_1, s_2] \leq 0 \quad ۲.$$

۳ اگر  $s_2 = s_1 = t_{\max}$ ، آنگاه قرار دهید  $t_m := t_{\max}$  و ازالگوریتم خارج کنید<sup>۵</sup> و نمونه<sup>۶</sup> سهمی  $p_2(t)$  را که در نقاط  $s_1, s_2, s_3$  بر  $g(t)$  منطبق می باشد بررسی کنید.

۴ اگر  $g[s_1, s_2, s_3] \leq 0$ ، بنابراین  $p_2(t)$  مینیمم ندارد، و سپس قرار دهید  $s_3 := t_{\max}$  و به مرحله سوم بروید<sup>۷</sup> و نمونه مینیمم  $s$  را در  $p_2(t)$  محاسبه کنید، یعنی

$$s := (s_1 + s_2 - g[s_1, s_2] / g[s_1, s_2, s_3]) / 2$$

- |                      |                |              |
|----------------------|----------------|--------------|
| 1. descent direction | 2. line search | 3. tolerance |
| 4. IF                | 5. EXIT        | 6. ELSE      |
|                      |                | 7. GO TO     |

۵. اگر  $t_{\max} > s$ ، آنگاه قرارداد دهید  $(s_1, s_2, s_3) := (s_2, s_3, t_{\max})$  و به مرحله سوم بروید

۱۰.۵ اگر  $\varepsilon \leq |g(s) - \min_i g(s_i)|$  یا  $|g(s) - p_\nu(s)| \leq \varepsilon$  آنگاه قرار دهید  $s := t_m$  و از الگوریتم خارج شوید.

وگرنه یک سری سه نقطه‌ای مرتب‌جدید  $(s_1, s_2, s_3)$  از مجموعه چهار نقطه‌ای  $\{s, s_1, s_2, s_3\}$  چنان انتخاب کنید که یارابطه  $g[s_2, s_3] < 0 < g[s_1, s_3]$  برقرار باشد یا، اگر برقراری رابطه فوق امکان‌پذیر نباشد،  $\max_i g(s_i)$  تا حد ممکن کوچک شود و سپس به مرحله ۴ بروید.

در خروج از الگوریتم،  $t_m$  به‌عنوان تقریبی از مینیمم  $g(t)$  در  $[0, t_{\max}]$  محسوب می‌شود.

توجه کنید که خروج از مرحله ۱۰.۵ تضمین نمی‌کند که  $t_m$  محاسبه شده به مینیمم  $g(t)$  «نزدیک» باشد، تمرین ۱۰.۵-۵ را ببینید. وقتی که الگوریتم ۲.۰۵ به‌عنوان قسمتی از الگوریتم مینیمم‌سازی چند متغیره<sup>۲</sup> به کار گرفته می‌شود، معمولاً این الگوریتم با  $s_3 = s_2 = 0$  شروع می‌شود [زیرا معمولاً  $\nabla F(\mathbf{x}^{(m)})^T \mathbf{u}^{(m)} = g'(0)$  در دست است] و  $s_3 = t_{\max} = 1$ ، و مرحله ۱۰.۵ به شکل « $s := t_m$ : قرارداد و از الگوریتم خارج شوید» ساده می‌شود. می‌توان نشان داد که این عمل اشکالی ندارد، به شرطی که جهت جستجو یعنی  $\mathbf{u}^{(m)}$  طوری انتخاب شود که  $\mathbf{x}^{(m)} + \mathbf{u}^{(m)}$  یک مینیمم موضعی از یک تابع درجه دوم، که تقریبی از تابع  $F$  در حوالی  $\mathbf{x}^{(m)}$  است، باشد.

اشاره کردیم که بهینه‌سازی منجر به پیدایش دستگاه‌های معادلات، یعنی، دستگاه‌هایی به صورت خاص

$$\nabla F(\mathbf{x}) = \mathbf{0}$$

می‌شود. برعکس دستگاه دلخواه از  $n$  معادله  $n$  مجهولی

$$\mathbf{f}(\mathbf{x}) = \mathbf{0}$$

را می‌توان اساساً به وسیله بهینه‌سازی حل کرد، زیرا مثلاً هر مینیمم تابع

$$F(\mathbf{x}) := \|\mathbf{f}(\mathbf{x})\|^2 = (f_1(\mathbf{x}))^2 + \dots + (f_n(\mathbf{x}))^2 \quad (5.5)$$

یک جواب معادله  $\mathbf{f}(\mathbf{x}) = \mathbf{0}$  است و بالعکس. برای این تابع خاص  $F$  داریم

$$\frac{\partial F}{\partial x_i} = \sum_{j=1}^n 2f_j(\mathbf{x}) \left( \frac{\partial f_j}{\partial x_i} \right)$$

1. new ordered three-point sequence
2. multivariate minimization algorithm

یا

$$\nabla F = \mathbf{v}(\mathbf{f}')^T \mathbf{f} \quad (۶.۵)$$

که در آن

$$\mathbf{f}' := \left( \frac{\partial f_i}{\partial x_j} \right)_{i,j=1}^n$$

ماتریس ژاکوبی<sup>۱</sup> تابع بردار-مقداری  $\mathbf{f}$  است.

## تمرین

۱-۱۰۵ نقاط بحرانی تابع

$$F(x_1, x_2) = x_1^2/3 + x_2^2 x_1 + 3$$

را با ترسیم منحنیهای  $\partial F/\partial x_1 = 0$  و  $\partial F/\partial x_2 = 0$  پیدا کنید. سپس با استفاده از جهت‌های گرد ادیان مجاور آنها، این نقاط را به نقاط ما کسیم، مینیم و نقاط زینی دسته‌بندی کنید.

۲-۱۰۵ با استفاده از روش تندترین کاهش و افزایش، مینیم و ماکزیمم تابع مذکور در تمرین ۱-۱۰۵ را تا  $10^{-6}$  رقم صحیح پیدا کنید.

۳-۱۰۵ گیریم  $\mathbf{u}$  جهت تماس بر یک خط تراز، ثابت  $F(x_1, x_2) =$  در نقطه  $\mathbf{x} = [x_1, x_2]^T$  باشد. با استفاده از قضیه ۸.۱ ثابت کنید که  $(\nabla F)(\mathbf{x})^T \mathbf{u} = 0$ .

۴-۱۰۵ یک زیر برنامه فورتون برای اجرای الگوریتم ۲.۵ بنویسید، سپس از آن برای حل تمرین ۲-۱۰۵ استفاده کنید (توجه: پیدا کردن ما کسیم تابع  $F$  - مثل پیدا کردن مینیمم تابع  $F$  - است).

۵-۱۰۵ (س. ر. را ببینن<sup>۲</sup> [۳۴]): گیریم  $h(t)$  یک تابع هموار در بازه  $[a, b]$  باشد، و  $h'(t) > 0$  و  $h(a) = h(b)$ .

(الف) ثابت کنید که  $h(t)$  دارای یک مینیمم منحصر به فرد  $t^*$  در  $[a, b]$  است.  
 (ب) پیدا کردن  $t^*$  را با انتخاب بازه  $[\alpha, \beta]$  که متضمن  $t^*$  باشد، و سپس با به کارگیری الگوریتم ۲.۵ برای ورودی  $t_{max} = \beta - \alpha$  و  $t_{max} = \beta - \alpha$  و  $g(t) = h(t - \alpha)$ ، به ازای مقداری مانند  $\epsilon > 0$  و انتخاب اولیه  $t_{max}/2$  و  $t_{max}$  (و به ازای  $(s_1, s_2, s_3)$ ، در نظر بگیرید. در این صورت بر آورد حاصله  $t_m^*$  برای  $t^*$  بستگی به مقادیر  $\alpha$  و  $\beta$ ،  $\epsilon$  پیدا می‌کند. ثابت کنید: اگر، به ازای کلیه این گونه مقادیر  $\alpha$  و  $\beta$ ، داشته باشیم  $\lim_{\epsilon \rightarrow 0} t_m = t^*$ ، آنگاه  $h(t)$  می‌باید

یک سهمی باشد. [دانه‌مایی:  $\alpha$  و  $\beta$  را طوری انتخاب کنید که تساوی  $h(\alpha) = h(\beta)$  برقرار باشد.]

(ج) نتیجه‌گیری کنید که الگوریتم ۲.۵ ممکن است کلاً برآورد خوبی برای مینیمم  $g$  به دست ندهد (حتی اگر  $\varepsilon$  بسیار کوچک باشد)، مگر اینکه  $g$  به یک سهمی نزدیک باشد.

۶-۱۰.۵: تقریب با روش کوچکترین توانهای دوم یک کار محاسباتی معمولی مستلزم تعیین پارامترهای  $a_1, \dots, a_k$  است به طوری که مدل  $y = R(x; a_1, \dots, a_k)$  با داده‌های  $(x_i, y_i)$  به ازای  $i = 1, \dots, N$ ، به خوبی بپردازد، یعنی به طوری که

$$R(x_i; a_1, \dots, a_k) = y_i + \varepsilon_i \quad i = 1, \dots, N$$

که  $N$ -برداری  $\varepsilon = [\varepsilon_1 \ \varepsilon_2 \ \dots \ \varepsilon_N]^T$  تا حد ممکن کوچک باشد.

(الف) با فرض اینکه  $R$  به طور هموار<sup>۲</sup> به بردار پارامتری  $\mathbf{a} = [a_1 \ a_2 \ \dots \ a_k]^T$  بستگی داشته باشد، نشان دهید که انتخاب  $\mathbf{a}^*$  که موجب مینیمم شدن  $\|\varepsilon\|_2$  می‌شود، می‌باید در معادلات به اصطلاح فرمال<sup>۳</sup> زیر صدق کند

$$A[R(x_1; \mathbf{a}^*), R(x_2; \mathbf{a}^*), \dots, R(x_N; \mathbf{a}^*)]^T = A\mathbf{y}$$

که در رابطه بالا، ماتریس  $A$  از مرتبه  $k \times N$  با رابطه زیر معین می‌شود

$$\left( \frac{\partial R}{\partial a_i}(x_j; \mathbf{a}^*) \right)$$

(ب) اعداد خاص  $a_1, a_2$  را در مدل

$$y = a_1 e^{a_2 x}$$

طوری پیدا کنید که با تغییر فوق، به بهترین وجه با مشاهدات زیر برازنده باشد.

$x_i$	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰
$y_i$	۱۲۴۸	۱۲۱۰	۰۸۱	۰۶۱	۰۴۵	۰۳۳	۰۲۴	۰۱۸	۰۱۳	۰۱۰

\* ۲.۵ روش نیوتن

در فصل ۳، هنگام حل یک معادله یک مجهولی

$$f(\xi) = 0$$

بر حسب  $\xi$ ، روش نیوتن را این گونه به دست آوردیم که (اولاً) از بسط تیلر

$$f(x+h) = f(x) + f'(x)h + \mathcal{O}(h^2)$$

تابع  $f$  در نقطهٔ  $x$  استفاده کردیم و سپس (ثانیساً) از جملهٔ بالاترین رتبهٔ  $\mathcal{O}(h^2)$  صرف نظر کردیم، معادلهٔ «خطی شده»

$$0 = f(x) + f'(x)h$$

را به جای معادلهٔ اصلی  $0 = f(x+h)$ ، نسبت به  $h$  حل کردیم و در نتیجه

$$h = -f(x)/f'(x)$$

و تقریب «بهیود یافته»

$$x - f(x)/f'(x)$$

را به دست آوردیم. اکنون که سعی داریم يك  $n$ -بردار  $\xi$  را طوری معین کنیم که در دستگاه معادلهٔ

$$f(\xi) = 0$$

صدق کند، عیناً به همان روش عمل می کنیم. بنا بر قضیهٔ ۹.۱ می دانیم که مؤلفهٔ  $i$ ام تابع  $f_i$  از تابع بردار-مقداری  $f$  در رابطهٔ زیر صدق می کند

$$f_i(x+h) = f_i(x) + (\nabla f_i(x))^T h + \mathcal{O}(\|h\|^2)$$

و آن در صورتی است که  $f_i$  مشتق جزئی و پیوستهٔ مراتب اول و دوم داشته باشد. بنابراین داریم

$$f(x+h) = f(x) + f'(x)h + \mathcal{O}(\|h\|^2) \quad (7.5)$$

که در آن ماتریس  $f'$ ، ماتریس ژاکوبی  $f$  در  $x$  نامیده می شود و به صورت زیر داده می شود

$$f'(x) = (\partial f_i / \partial x_j)_{i,j=1}^n$$

مجدداً از جملهٔ بالاترین مرتبهٔ  $\mathcal{O}(\|h\|^2)$  صرف نظر و معادلهٔ «خطی شده»

$$0 = f(x) + f'(x)h$$

را به جای معادلهٔ اصلی  $0 = f(x+h)$ ، نسبت به مقدار تصحیحی  $h$  حل می کنیم و در نتیجه جواب

$$h = -f'(x)^{-1} f(x)$$

را به شرط آنکه ماتریس ژاکوبی  $f'(x)$  وارون پذیر باشد، به دست می آوریم. بدین طریق برای

ξ تقریب جدید

$$\mathbf{x} - \mathbf{f}'(\mathbf{x})^{-1} \mathbf{f}(\mathbf{x})$$

را به دست می‌آوریم. این مرحله، مرحله‌ی اساسی روش نیوتن برای يك دستگاه است. البته معادله‌ی نیوتن

$$\mathbf{f}'(\mathbf{x})\mathbf{h} = -\mathbf{f}(\mathbf{x})$$

نسبت به مقدار تصحیحی  $\mathbf{h}$  بر حسب  $\mathbf{x}$ ، يك دستگاه خطی است و باروشهای مستقیم مذکور در فصل ۴ حل می‌شود.

الگوریتم ۳.۵: روش نیوتن برای يك دستگاه معادلات. دستگاه  $n$  معادله‌ی  $n$  مجهولی

$$\mathbf{f}(\xi) = \mathbf{0}$$

که در آن  $\mathbf{f}$  تابعی است مقدار-برداری و دارای مؤلفه‌های هموار همراه با حدس اولیه‌ی  $\mathbf{x}^{(0)}$  برای جواب ξ از دستگاه فوق، داده شده است

For  $m = 0, 1, 2, \dots$ , until satisfied, do:

$$\lfloor \mathbf{x}^{(m+1)} := \mathbf{x}^{(m)} - \mathbf{f}'(\mathbf{x}^{(m)})^{-1} \mathbf{f}(\mathbf{x}^{(m)})$$

می‌توان نشان داد که روش نیوتن به سمت ξ همگرا می‌شود، به شرط آنکه  $\mathbf{x}^{(0)}$  به اندازه‌ی کافی به ξ نزدیک باشد، و به شرطی که ماتریس  $\mathbf{f}'$ ، ژاکوبی  $\mathbf{f}$ ، پیوسته و  $\mathbf{f}'(\xi)$  وارونپذیر باشد. بعلاوه اگر مشتقات جزئی مرتبه‌ی دوم مؤلفه‌های تابع  $\mathbf{f}$  نیز پیوسته باشند، آنگاه رابطه‌ی

$$\|\xi - \mathbf{x}^{(m+1)}\| \leq c \|\xi - \mathbf{x}^{(m)}\|^2$$

به ازای مقدار ثابتی مانند  $c$  و جمیع مقادیر  $m$  به اندازه‌ی کافی بزرگ برقرار است. به عبارت دیگر، روش نیوتن به صورت عبارت درجه‌ی دوم همگرا می‌شود (مثال ۶.۵ را نگاه کنید).

□ مثال ۳.۵: اعداد  $0 < \xi_0 < \xi_1 < \xi_2 < \dots < \xi_n < 1$  را طوری تعیین کنید که به ازای  $\xi_0 = 0$  و  $\xi_{n+1} = 1$  و  $G(x) = x^3$  داشته باشیم

$$G'(\xi_i) = G[\xi_{i-1}, \xi_{i+1}] \quad i = 1, \dots, n$$

این کار مستلزم حل دستگاه

$$3\xi_i^2 = (\xi_{i+1}^3 - \xi_{i-1}^3) / (\xi_{i+1} - \xi_{i-1}) \quad i = 1, \dots, n$$

یا

$$\mathbf{f}(\xi) = \mathbf{0}$$



همگرایی درجه دوم، هم در کاهش مقدار خطای باقیمانده  $F(X^{(x)})$  و هم در کاهش مقدار تصحیحی نیوتن یعنی  $h$  برای  $x^{(m)}$  مشخص است. محاسبات لازم با دقت مضاعف (تقریباً ۱۷ رقم اعشاری) روی کامپیوتر UNIVAC ۱۱۱۰ انجام گرفته است.

استفاده از روش نیوتن، اشکالاتی به همراه دارد که در این مثال ساده آشکار نیست. عموماً دو اشکال عمده وجود دارد: (۱) عدم همگرایی به دلیل حدس نامناسب اولیه و (۲) هزینه ساختن معادله نیوتن به طور صحیح و حل آن نسبت به مقدار تصحیحی  $h$ . در زیر، هر دو اشکال را به طور جداگانه مورد بحث قرار می‌دهیم.

دو نظر که در امر تقویت، یا حداقل در تسریع به همگرایی تا حدی موفق بوده‌اند عبارت‌اند از تداوم<sup>۲</sup> یا نشانیدن<sup>۳</sup> و میراندن<sup>۴</sup>. در تداوم، مسئله حل  $f(\xi) = 0$  را به نحو مناسبی به عنوان آخرین مسئله از یک خانواده مسائل یک پارامتری پیوسته

$$g(\xi, t) = 0$$

با شرط

$$g(x, 1) = f(x)$$

مورد توجه قرار می‌دهند، که  $g(x, 0)$  تابعی است که در حل

$$g(\xi, 0) = 0$$

مشکلی ایجاد نخواهد کرد.

پس از اینکه  $\xi^{(0)}$  طوری پیدا شد که در تساوی  $g(\xi^{(0)}, 0) = 0$  صادق باشد، دنباله

$$0 = t_0 < t_1 < \dots < t_N = 1$$

$$g(\xi^{(i)}, t_i) = 0$$

را به ازای  $i = 1, 2, \dots, N$  و با روش نیوتن و بسا استفاده از بردار  $\xi^{(i-1)}$  به عنوان حدس اولیه و یا شاید، حتی بردار برونمایی شده

$$\xi^{(i-1)} + (t_i - t_{i-2}) \frac{\Delta \xi^{(i-2)}}{\Delta t_{i-2}}$$

(اگر  $i > 1$ ) حل می‌کنند. امید این است که مسائل مجاور هم یعنی  $g(\xi, t_i) = 0$  و  $g(\xi, t_{i-1}) = 0$  آنقدر بهم نزدیک باشند که یک جواب خوب از یکی، یک حدس اولیه مناسبی برای جواب دیگری باشد. انتخاب معمولی برای تابع  $g$  عبارت است از

$$g(x, t) = tf(x) + (1-t)(f(x) - f(x^{(0)}))$$



$$g(x, t) = tf(x) + (1-t)(x - x^{(0)})$$

در روش میرای<sup>۱</sup> نیوتن مقدار  $x$  در بارست بعدی یعنی  $x^{(m+1)} = x^{(m)} + h$  قابل قبول واقع نمی‌شود، اگر این امر به افزایش خطای مانده منجر شود، یعنی اگر  $\|f(x^{(m+1)})\|_2 > \|f(x^{(m)})\|_2$  در چنین موردی بردارهای  $h/2^i$  به ازای  $i = 1, 2, \dots$  بررسی می‌شوند و  $x^{(m+1)}$  به عنوان اولین برداری که به ازای آن خطای مانده کمتر از  $\|f(x^{(m)})\|_2$  است به کار گرفته می‌شود.

**الگوریتم ۴۰۵:** روش میرای نیوتن برای يك دستگاه دستگاه  $n$  معادله  $n$  مجهولی  $f(x) = 0$  که در آن  $f$  تابعی است بردارمقداری و دارای مؤلفه‌های هموار تابعی، به همراه يك حدس اولیه  $x^{(0)}$  برای يك جواب  $\bar{x}$  از دستگاه فوق داده شده است.

For  $m = 0, 1, 2, \dots$  until satisfied, do:

$$\begin{cases} h := -f'(x^{(m)})^{-1}f(x^{(m)}) \\ * i := \min\{j : 0 \leq j, \|f(x^{(m)} + h/2^j)\|_2 < \|f(x^{(m)})\|_2\} \\ x^{(m+1)} := x^{(m)} + h/2^i \end{cases}$$

بدون مطالعه روشن نیست که آیا مرحله \* همیشه می‌تواند اجرا شود یا نه. برای اینکه  $i$  تعریف شود، لازم و کافی است که جهت نیوتن  $h$  در  $x = x^{(m)}$  برای تابع

$$F(x) = \|f(x)\|_2^2$$

يك جهت کاهشی باشد. چون داریم

$$\nabla F(x) = 2f'(x)^T f(x)$$

و بنا بر (۵.۵) و (۶.۵)،  $h$  يك جهت کاهشی برای  $F$  در  $x$  است اگر و فقط اگر رابطه

$$(f'(x)^T f(x))^T h < 0$$

برقرار باشد. از سوی دیگر داریم  $h = -f'(x)^{-1}f(x)$  و بنا بر این

$$(\nabla F(x))^T h = 2(f'(x)^T f(x))^T (-f'(x)^{-1}f(x))$$

$$= -2f(x)^T f(x) = -2\|f(x)\|_2^2 < 0$$

این رابطه نشان می‌دهد که جهت نیوتن برای  $F(x) = \|f(x)\|_2^2$  حقیقتاً يك جهت کاهشی است، بنا بر این عدد صحیح  $i$  در مرحله \* به خوبی تعریف شده است. گرچه در عمل به جای مرحله \* مرحله زیر را می‌گذارند

$$i := \min\{j : 0 \leq j \leq j_{\max}, \|f(x^{(m)} + h/2^j)\|_2 < \|f(x^{(m)})\|_2\}$$

اگر  $i$  تعریف نشده باشد، آنگاه خروج از الگوریتم به جهت وجود اشکال و غیره یا  $z_{\max}$  که از قبل تعیین و مثلاً برابر ۱۰ انتخاب شده شروع کنید.

□ مثال ۵.۵: دستگاه  $f(x) = 0$

$$f_1(x) = x_1 + 3 \ln|x_1| - x_1^2 \quad f_2(x) = 2x_2^2 - x_1x_2 - 5x_1 + 1$$

دارای چندین جواب است. به این دلیل، به منظور تأمین همگرایی به سمت یک جواب خاص یا به طور کلی برای تأمین همگرایی، تخمین اولیه باید با دقت انتخاب شود. معادله‌های نیوتن به صورت زیر هستند

$$\begin{bmatrix} 1 + 3/x_1 & -2x_2 \\ 2x_2 - x_1 & -5 \end{bmatrix} h = -f(x)$$

اگر با حدس اولیه  $x^{(0)} = [2 \ 2]^T$  شروع کنیم، بارستهای مندرج در جدول زیر به دست می‌آیند.

$m$	$x^{(m)}$		$\ f(x^{(m)})\ _2$	$\ h\ _2$
0	2.	2.	0.500 + 1	0.238 + 2
1	-18.1588	-10.5794	0.572 + 3	0.112 + 2
2	-8.3710	-5.2287	0.142 + 3	0.543 + 1
3	-3.5525	-2.7191	0.351 + 2	0.266 + 1
4	-1.2015	-1.4728	0.860 + 1	0.198 + 1
5	-0.0004	0.0945	0.234 + 2	0.187 + 4
6	0.0451	-1866.2415	0.348 + 7	0.933 + 3
7	0.0233	-933.1179	0.871 + 6	0.467 + 3
8	0.0108	-466.5520	0.218 + 6	0.233 + 3
		etc.		

روشن است که بارستها اساساً همگرا نیستند. اما اکنون با به کار گرفتن روش هیوای نیوتن و شروع با همان حدس اولیه، جدول زیر به دست می‌آید.

$m$	$x^{(m)}$		$\ f(x^{(m)})\ _2$	$\ h\ _2$	$i$
0	2.	2.	0.500 + 1	0.238 + 2	4
1	0.7400698	1.2137849	0.299 + 1	0.160 + 1	1
2	0.5310238	0.4415855	0.205 + 1	0.217 + 1	2
3	0.5178341	-0.1001096	0.178 + 1	0.372 + 1	3
4	0.5584838	-0.5637875	0.173 + 1	0.832 + 1	6
5	0.5847026	-0.6910621	0.172 + 1	0.967 + 1	6
6	0.6215780	-0.8376443	0.171 + 1	0.937 + 1	6
7	0.6657612	-0.9772562	0.171 + 1	0.684 + 1	5
8	0.7448782	-1.1760004	0.169 + 1	0.328 + 1	3
9	0.9489394	-1.5313175	0.163 + 1	0.676	0
10	1.5501608	-1.8410875	0.105 + 1	0.315	0
11	1.3892191	-1.5703845	0.132	0.473 - 1	0
12	1.3735386	-1.5257440	0.249 - 2	0.781 - 3	0
13	1.3734783	-1.5249650	0.608 - 6	0.156 - 6	0
14	1.3734783	-1.5249648	0.843 - 7		

در اینجا برای هر بارست، عدد صحیح  $k$  که در مرحله  $k$  از الگوریتم  $۴.۵$  معین شده نیز آورده شده است. در ابتدا مقدار  $h$  که در نظر گرفته شده نسبتاً بزرگ است، اما این مقدار تا اندازه  $۱/۶۴ = ۱/۲^6$  میرا می شود. همچنین مقدار  $\|f(x^{(m)})\|$  از خطای مانده نیز به زحمت از یک مرحله تا مرحله دیگر کاهش می یابد. اما در نهایت تمام مراحل روش نیوتن اجرا شده و بارستها، همان طوری که باید، به صورت درجه دوم همگرا می شوند. (حقیقتاً مشاهده این بارستها روی پایانه کامپیوتری تجربه مهیجی است. سرانجام هنگامی که همگرایی درجه دوم برقرار می شود، احساس شادی به شخص دست می دهد.)  
محاسبات فوق با دقت ساده روی کامپیوتر UNIVAC ۱۱۱۰ انجام گرفته است.  
بنابر این خطای  $\|f(x^{(۴)})\|$  در سطح خطاهای عادی است.  $\square$

اشکال دوم در استفاده از روش نیوتن مربوط به تشکیل معادله نیوتن وحل آن به ازای مقدار تصحیحی  $h$  است. اگر  $f$  تابع نسبتاً پیچیده ای باشد، ساختن ماتریس ژاکوبی خود کار دشواری است، زیرا در این صورت امکان اشتباهات زیادی از جمله در مشتقگیری و یا رمز گذاری عناصر  $f'$  وجود دارد. معمولاً نتیجه این نوع اشتباهات سبب از دست رفتن همگرایی درجه دوم، و یا در موارد حاد موجب از دست رفتن همگرایی به طور کلی می باشد. اکنون برخی مراکز محاسباتی، برنامه هایی برای مشتقگیری نمادی از عبارتها و یا حتی از توابعی که به وسیله یک زیر برنامه تولید شده اند، فراهم کرده اند. استفاده از این گونه برنامه ها کمک زیادی در ساختن ماتریس ژاکوبی به ما می کند. اگر چنین برنامه هایی در دسترس نباشند، می توانیم ماتریس ژاکوبی رمز شده  $f'(x)$  خود را از مقایسه آن در یک نقطه  $x$  با تقریبهایی عددی نه چندان دقیق آرایه های آن به شکل

$$\frac{\partial f_i}{\partial x_j}(x) \approx [f_i(x + \varepsilon i_j) - f_i(x)] / \varepsilon \quad (۸.۵)$$

یا

$$\frac{\partial f_i}{\partial x_j}(x) \approx [f_i(x + \varepsilon i_j) - f_i(x - \varepsilon i_j)] / (2\varepsilon) \quad (۹.۵)$$

که با هر دو رابطه فوق، در مبحث حسابان (به فصل ۷ نگاه کنید) آشنا شده ایم آزمایش کنیم. یک راه دیگر این است که فقط، به رمز گذاری توابع  $f_1, \dots, f_n$  و سپس به کار گرفتن فرمولهای (۸.۵) یا (۹.۵) برای ساختن یک تقوید مناسب  $J$  برای  $f'(x)$  اکتفا کنیم. این امر مستلزم انتخاب مناسبی برای اندازه  $\varepsilon$  است (به بخش ۱.۷ نگاه کنید). گیریم  $J_m$  ماتریس ژاکوبی  $f'(x^{(m)})$  یا یک تقریب مناسبی برای آن باشد. وقتی

1. terminal

2. single precision

3. coding

4. step size

که  $J_m$  ساخته شد، می‌باید معادله

$$J_m \mathbf{h} = -\mathbf{f}(\mathbf{x}^{(m)})$$

را نسبت به مقدار تصحیحی  $\mathbf{h}$  حل کرد. در حالت کلی  $J_m$  يك ماتریس کامل از مرتبه  $n$  است، لذا تعداد  $\Theta(n^3)$  عمل برای به دست آوردن  $\mathbf{h}$  لازم است. از سوی دیگر، اگر همگرایی وجود داشته باشد و  $\mathbf{f}'(\mathbf{x})$  به طور پیوسته به  $\mathbf{x}$  بستگی داشته باشد،  $\mathbf{f}'(\mathbf{x}^{(m)})$  و  $\mathbf{f}'(\mathbf{x}^{(m+k)})$  تفاوت بسیار کمی با هم خواهند بود که به جای  $\mathbf{f}'(\mathbf{x}^{(m+k)})$  از  $\mathbf{f}'(\mathbf{x}^{(m)})$  استفاده کنیم. زیرا با یکبار تجزیه کردن  $\mathbf{f}'(\mathbf{x}^{(m)})$  می‌توان دستگاه فوق را به ازای مقادیر متعدد سمت راست با هزینه‌ای معادل با  $\Theta(n^2)$  عملیات، حل کرد. ایس روش اصلاح شده نیوتن است، که در آن تساوی  $J_{m+k} = \mathbf{f}'(\mathbf{x}^{(m)})$  به ازای  $k = 0, 1, 2, \dots$  برقرار است، تا وقتی که (پامگر آنکه) کندی سرعت همگرایی اعلام کند که  $J_{m+k}$  به عنوان آخرین ماتریس ژاکوبی به کار گرفته خواهد شد.

يك انحراف خیلی زیادتر از روش نیوتن، در روشهای به اصطلاح باب روز کردن ماتریس<sup>۱</sup> مطرح شده است، در این روش  $J_{m+1}$  از  $J_m$  با افزودن يك ماتریس از مرتبه<sup>۲</sup> يك یا دو که بستگی به  $J_m$ ،  $\mathbf{h}$ ،  $\mathbf{x}^{(m)}$ ،  $\mathbf{f}(\mathbf{x}^{(m)})$  و  $\mathbf{f}(\mathbf{x}^{(m+1)})$  دارد به دست می‌آید. هدف اصلی این است که  $J_{m+1}$  طوری انتخاب شود که با توجه به روابط

$$\delta \mathbf{x} = \mathbf{x}^{(m+1)} - \mathbf{x}^{(m)} \quad \text{و} \quad \delta \mathbf{f} = \mathbf{f}(\mathbf{x}^{(m+1)}) - \mathbf{f}(\mathbf{x}^{(m)})$$

رابطه

$$J_{m+1}(\delta \mathbf{x}) = \delta \mathbf{f}$$

به دست آید این، مطلب دور از ذهنی نیست، زیرا در مواردی که به ازای  $\mathbf{x}$  نزدیک به  $\mathbf{x}^{(m)}$  تساوی  $J_{m+1} = \mathbf{f}'(\mathbf{x})$  برقرار باشد، يك تساوی تقریبی موجود است. اگر ماتریسی که به  $J_m$  افزوده می‌شود از مرتبه<sup>۱</sup> يك یا دو باشد، آنگاه این امکان وجود دارد که تغییر حاصله در  $K_m = J_m^{-1}$  به صورت افزودن ماتریسی که به آسانی قابل محاسبه است بیان شود. بنابراین با حفظ رد  $K_m$  به جای رد  $J_m$ ، می‌توان از تجزیه<sup>۳</sup>  $J_m$  اجتناب کرد. یکی از روشهای معمول از این نوع، روش پرویدن<sup>۳</sup> است. در این روش، ابتدا  $K_0 = \mathbf{f}'(\mathbf{x}^{(0)})^{-1}$  محاسبه و سپس  $K_{m+1}$  از روی  $K_m$ ، به توسط رابطه

$$K_{m+1} = K_m - \frac{[K_m(\delta \mathbf{f}) - \delta \mathbf{x}](\delta \mathbf{x})^T K_m}{(\delta \mathbf{x})^T K_m(\delta \mathbf{f})} \quad (10.5)$$

و رابطه

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} - K_m \mathbf{f}(\mathbf{x}^{(m)}) \quad (11.5)$$

محاسبه می‌شود.  $J_m = K_m^{-1}$  متناظر به آن در رابطه

$$J_{m+1}(\delta \mathbf{x}) = \delta \mathbf{f}$$

صدق می‌کند در حالی که رابطه  $J_{m+1} \mathbf{z} = J_m \mathbf{z}$  به ازای جميع مقادیر  $\mathbf{z}$  عمود بر  $\delta \mathbf{x}$  برقرار است. در عمل، در این روش بارستی، از روش میرا کردن نیز استفاده می‌شود.

## تمرین

۱-۲۰۵ با استفاده از روش نیوتن جوابهای دستگاه  $\nabla F(\mathbf{x}) = \mathbf{0}$  را که در آن  $F$  تابع داده شده در تمرین ۱-۱۰۵ است، پیدا کنید. نتیجه کار خود را با آنچه که در تمرین ۲-۱۰۵ خواسته شده است مقایسه کنید.

۲-۲۰۵ ثابت کنید: که اگر  $a < c < b$  و  $G'(c) = G[a, b]$  و  $G'''(x)$  و  $G''(x)$  هر دو در بازه  $[a, b]$  مثبت باشند، آنگاه  $c > (a+b)/2$ . [دانه‌مایی: گیریم  $c = (a+b)/2$  و با بسط کلیه متغیرها به صورت سری تیلر در حول  $c$  نشان دهید که  $G'(c) < G[a, b]$ . در غیر این صورت مستقیماً از (۸.۷) استفاده کنید].  
نتیجه گیری کنید که ماتریس ژاکوبی  $f'(x)$  در مثال ۴.۵ به طور قطری اکیداً نافذ و بنابراین وارون پذیر است.

۳-۲۰۵ با استفاده از روش نیوتن جواب دستگاه نسبتاً پیچیده زیر در فاصله  $1 \leq x, y \leq 1$  پیدا کنید.

$$\cos \left[ \frac{x^2 - \sqrt{\sin(xy) + 3}}{4 + (xy)^2} \right] + \sin(3xy - 1) = 0.934$$

$$\exp \{ \cos [(xy)^3 - 3] \} + \tan \left[ \frac{x}{y} (0.08 + \cos x) \right] = 1.079$$

(البته شناسه‌های توابع مثلثاتی بر حسب رادیان اندازه گیری می‌شوند).

اگر همگرایی درجه دوم به دست نیامورد، با استفاده از (۸.۵) یا (۹.۵) رمزگذاری ماتریس ژاکوبی را آزمون کنید.

۴-۲۰۵ با به کار گیری روش نیوتنی میرا، جواب مسئله‌ای را که در مثال ۵.۵ مورد بحث قرار گرفت با نقطه شروع  $\mathbf{x}^{(0)} = [2 \ 1]^T$  پیدا کنید.

۵-۲۰۵ سعی کنید مسئله مثال ۵.۵ را با روش تداوم و نقطه شروع  $\mathbf{x}^{(0)} = [2 \ 1]^T$  و استفاده از  $t_0, \dots, t_N = 0, 0.1, 0.3, 0.6, 1$  (در مراحل اول، بارستها را

آنقدر ادامه دهید که همگرایی درجهٔ دوم پدیدار شود.

۶-۲۰۵ مسئلهٔ مثال ۴.۵ را به‌ازای  $n = 10$  و  $G(x) = x^5$  حل کنید.

### ۵-۳ بارست نقطهٔ ثابت و روشهای واهلش

روش نیوتن و بعضی از صورت‌های متفاوت آن، که در قسمت ۲.۵ مورد بحث قرار گرفتند، همگی مثالهایی از بارست نقطهٔ ثابت هستند. در اینجا می‌توان معادلهٔ

$$f(\xi) = 0$$

را به‌صورتی هم‌ارز با آن:

$$\xi = g(\xi)$$

نوشت و سپس با شروع از حدس اولیهٔ  $x^{(0)}$  دنبالهٔ

$$x^{(m+1)} = g(x^{(m)}) \quad m = 0, 1, 2, \dots$$

را تولید کرد با این امید که دنبالهٔ فوق به‌سمت نقطهٔ ثابت  $\xi$  از  $g$  همگرا شود.

برای مثال، روش نیوتن یک روش بارستی نقطهٔ ثابت از این گونه‌است، بسا تابع

بارستی  $g$ ، که به‌صورت رابطهٔ زیر داده شده‌است

$$g(x) = x - f'(x)^{-1}f(x)$$

به‌صورت کلیتر، در روشهای شبه-نیوتنی، یک تابع بارستی به‌صورت زیر به‌کار برده می‌شود.

$$g(x) = x - Cf(x) \quad (12.5)$$

که در آن  $C = C(x)$  یک ماتریس است. روش واهلش، که بعداً در این قسمت مورد بحث قرار خواهد گرفت، دید دیگری است برای تشکیل تابعهای بارستی برای حل  $f(\xi) = 0$  به‌وسیلهٔ بارست نقطهٔ ثابت.

تجزیه و تحلیل روش بارست نقطهٔ ثابت برای دستگاه معادلات با آنچه که برای یک

معادله در فصل ۳ ارائه شد، یک تفاوت جزئی دارد و آن اینکه در این حال دربارست  $m$ ، اندازهٔ خطای  $x^{(m)} - \xi$  به‌جای مقادیر قدرمطلق، توسط نرم معین می‌شود.

قضیهٔ ۱۰۵ فرض کنید که تابع بارست  $g$ ، مجموعهٔ بسته‌ای مانند  $S$  را بر خود بنگارد. یعنی اگر  $x$  جزء مجموعهٔ  $S$  باشد،  $g(x)$  نیز به  $S$  متعلق است و بعلاوه فرض کنید که  $g$  روی  $S$

انقباضی<sup>۱</sup> باشد، یعنی رابطه

$$\|g(x) - g(y)\| \leq K \|x - y\|$$

به ازای جميع مقادير  $x$  و  $y$  در  $S$  و مقداری مانند  $K < 1$  برقرار باشد، در این صورت  $g$  (i) يك نقطه ثابت در  $S$  دارد.

(ii) اگر  $\xi$  يك نقطه ثابت  $g$  در  $S$  باشد، آنگاه روش بارست نقطه ثابت که با  $x^{(0)}$  در

$S$  شروع می شود به سمت  $\xi$  همگرا خواهد شد، یعنی برای دنباله ای نظیر

$$\lim_{m \rightarrow \infty} \|\xi - x^{(m)}\| = 0 \text{ داریم، به ازای } m = 0, 1, 2, \dots, x^{(m+1)} = g(x^{(m)})$$

به بیان صریحتر

$$\|\xi - x^{(m)}\| \leq \frac{K}{1-K} \|x^{(m)} - x^{(m-1)}\| \quad (13.5)$$

بنابراین

$$\|\xi - x^{(m)}\| \leq \frac{K^m}{1-K} \|x^{(1)} - x^{(0)}\| \quad (14.5)$$

فرضیات فوق ناظر بر این هستند، که می توانیم با هر  $x^{(0)}$  در  $S$  شروع کنیم و به ازای

$m = 0, 1, 2, \dots$  بارست  $x^{(m+1)} = g(x^{(m)})$  را تا بینهایت به ازای هر  $x^{(m)}$  در  $S$

ادامه دهیم. وانگهی با استدلالی که خارج از ظرفیت این کتاب است (مثلا با استفاده از

تمامیت فضای  $n$ -بعدی)، قسمت (i) نتیجه می شود، و سرانجام برای به دست آوردن برآورد

(۱۳.۵) و از آنجا (۱۴.۵)، ملاحظه می کنیم که

$$\begin{aligned} \|\xi - x^{(m)}\| &= \|g(\xi) - g(x^{(m-1)})\| \\ &\leq K \|\xi - x^{(m-1)}\| \end{aligned} \quad (15.5)$$

و از این رو چون  $g$  انقباضی است، بنا بر نامساوی مثلثی (۴-۳۳، iii) داریم

$$\begin{aligned} \|\xi - x^{(m-1)}\| &\leq \|\xi - x^{(m)}\| + \|x^{(m)} - x^{(m-1)}\| \\ &\leq K \|\xi - x^{(m-1)}\| + \|x^{(m)} - x^{(m-1)}\| \end{aligned}$$

یا

$$(1-K) \|\xi - x^{(m-1)}\| \leq \|x^{(m)} - x^{(m-1)}\|$$

حال از ترکیب نامساوی فوق با (۱۵.۵)، (۱۳.۵) به دست می آید.

□ مثال ۶۰۵: روش نیوتن يك روش بارست نقطه ثابت با تابع بارست زیر است

$$g(x) = x - f'(x)^{-1} f(x)$$

بنابراین

$$f'(x)[g(x) - x] = -f(x)$$

در حالی که بنا بر رابطه (۷۰۵) و با فرض آنکه  $f$  مشتقات جزئی و پیوسته اول و دوم داشته باشد، داریم

$$\begin{aligned} 0 &= f(\xi) = f(x + \xi - x) \\ &= f(x) + f'(x)(\xi - x) + \theta(\|\xi - x\|^2) \end{aligned}$$

بنابراین از قراردادن  $-f'(x)[g(x) - x]$  به جای  $f(x)$  داریم

$$0 = f'(x)[- (g(x) - x) + (\xi - x)] + \theta(\|\xi - x\|^2)$$

یا

$$f'(x)[g(x) - \xi] = \theta(\|\xi - x\|^2)$$

عبارت بالا بیانگر این است که، به ازای مقداری ثابت مانند  $c$  داریم

$$\|g(x) - \xi\| \leq \|f'(x)^{-1}\| \cdot c \cdot \|\xi - x\|^2$$

اکنون اگر  $f'(\xi)$  وارونپذیر باشد، آنگاه چون بنا بر فرض  $f'(x)$  پیوسته است، می‌توانیم مقدار مثبتی مانند  $\delta$  و یک  $M$  طوری پیدا کنیم که  $f'(x)^{-1}$  به ازای جميع مقادیر  $x$  در محدوده  $\delta$  از  $\xi$  وجود داشته و یک نرم ماتریسی نا بیشتر از  $M$  داشته باشد. اما در این صورت، اگر  $\varepsilon$  کوچکتر از  $\delta$  و  $(Mc)^{-1}$  انتخاب شود، به ازای کلیه مقادیر  $x$  در مجموعه بسته

$$S = \{x : \|\xi - x\| \leq \varepsilon\}$$

$f'(x)^{-1}$  وجود دارد (بنابراین  $g(x)$  تعریف شده است) و

$$\|g(x) - \xi\| \leq Mc \|\xi - x\|^2 \leq (Mc\varepsilon) \|\xi - x\| \leq \|\xi - x\|$$

لذا  $g$  مجموعه بسته  $S$  را بر خودش می‌نگارد. بعلاوه اگر  $\varepsilon < \|\xi - x\|$ ، آنگاه  $\|\xi - x\| < K = Mc \|\xi - x\|$ ، بنابراین  $\xi$  يك نقطه ثابت جاذب از  $g$  است، و بارستی که با هر  $x^{(0)}$  کمتر از  $\varepsilon$  در محدوده  $\xi$  شروع شود، به سمت  $\xi$  همگرا خواهد بود. □



حال، به عنوان يك مثال ديگر حل دستگاه خطی

$$A\xi = \mathbf{b} \quad (۱۶.۵)$$

را به وسیلهٔ بارست نقطهٔ ثابت در نظر می گیریم. این گونه روشهای بارستی تماماً می توانند بر پایهٔ فکر وارون تقریبی<sup>۱</sup> استوار باشند. منظور ما این است که به ازای هر ماتریس  $C$  رابطهٔ

$$\|I - CA\| < ۱ \quad (۱۷.۵)$$

به ازای مقداری از نرم ماتریسی برقرار باشد.

لم ۱۰۵ اگر  $C$  يك وارون تقریبی ماتریس  $A$  باشد، یعنی اگر به ازای يك نرم ماتریسی، نامساوی  $\|I - CA\| < ۱$  برقرار باشد، آنگاه  $C$  و  $A$  هر دو وارون پذیرند.

در حقیقت، اگر  $C$  یا  $A$  وارون پذیر نباشند، آنگاه ماتریس  $CA$  نیز وارون پذیر نخواهد بود (به تمرین ۱۰۴-۸ نگاه کنید). در این حال بنا بر قضیهٔ ۴.۴، می توانیم  $\mathbf{0} \neq \mathbf{x}$  پیدا کنیم که داشته باشیم  $C\mathbf{A}\mathbf{x} = \mathbf{0}$ . اما در این صورت رابطهٔ

$$\mathbf{0} \neq \|\mathbf{x}\| = \|(I - CA)\mathbf{x}\| \leq \|I - CA\| \|\mathbf{x}\| < \|\mathbf{x}\|$$

به دست می آید که بی معنی است.

بخصوص، اگر  $A$  يك وارون تقریبی داشته باشد، (۱۶.۵) دقیقاً دارای يك جواب

است.

در رابطه با يك وارون تقریبی  $C$  از  $A$ ، تابع بارست زیر را در نظر می گیریم

$$\begin{aligned} \mathbf{g}(\mathbf{x}) &= C\mathbf{b} + (I - CA)\mathbf{x} \\ &= \mathbf{x} + C(\mathbf{b} - A\mathbf{x}) \end{aligned}$$

توجه کنید که تابع بارست فوق از نوع تابع شبه-نیوتنی، یعنی به شکل  $\mathbf{g}(\mathbf{x}) = \mathbf{x} - C\mathbf{f}(\mathbf{x})$  است، اگر  $\mathbf{f}(\mathbf{x}) = A\mathbf{x} - \mathbf{b}$  همچنین

$$\begin{aligned} \mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y}) &= C\mathbf{b} + (I - CA)\mathbf{x} - [C\mathbf{b} + (I - CA)\mathbf{y}] \\ &= (I - CA)(\mathbf{x} - \mathbf{y}) \end{aligned}$$

در نتیجه داریم

$$\|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y})\| \leq \|I - CA\| \|\mathbf{x} - \mathbf{y}\| \quad (۱۸.۵)$$

که نشان می دهد،  $\mathbf{g}$  انقباضی بوده و رابطهٔ زیر برقرار است

$$K = \|I - CA\| < ۱$$

بنابراین بارست نقطه ثابت

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} + C(\mathbf{b} - A\mathbf{x}^{(m)}) \quad m = 0, 1, 2, \dots$$

با شروع از هر نقطه  $\mathbf{x}^{(0)}$ ، به سمت جواب منحصر به فرد  $\mathbf{x}$  از معادله (۱۶.۵) همگرا خواهد بود، و خطا در هر مرحله از بارست، حداقل به اندازه یک عامل  $K = \|I - CA\|$  کاهش می‌یابد.

□ مثال ۷.۵: فرض کنید ماتریس  $A$  نافذ قطری سطرأ مؤکد باشد یعنی داشته باشیم

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}| \quad i = 1, \dots, n$$

گیریم  $D = \text{diag}(a_{11}, a_{22}, \dots, a_{nn})$  قطر  $A$  باشد، در این صورت رابطه

$$\|I - D^{-1}A\|_{\infty} = \max_i \{1 - \sum_j |a_{ij}| / |a_{ii}|\} = \max_i \sum_{j \neq i} |a_{ij}| / |a_{ii}| < 1$$

جدول ۱.۵

Jacobi				Gauss-Seidel		
$x_1^{(m)}$	$x_2^{(m)}$	$x_3^{(m)}$	$m$	$x_1^{(m)}$	$x_2^{(m)}$	$x_3^{(m)}$
0	0	0	0	0	0	0
1.2	1.2	1.2	1	1.2	1.08	0.972
0.96	0.96	0.96	2	0.9948	1.0033	1.00019
1.008	1.008	1.008	3	0.99965	1.000016	1.000033
0.9984	0.9984	0.9984	4			
1.00032	1.00032	1.00032	5			
0.999936	0.999936	0.999936	6			

نشان می‌دهد که  $D$  یک وارون تقریبی  $A$  است. شکل بارستی متناظر با آن یعنی

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} + D^{-1}(\mathbf{b} - A\mathbf{x}^{(m)}) \quad m = 0, 1, 2, \dots \quad (19.5)$$

همان بارست ژاکوبی است. توجه کنید که می‌توان  $\mathbf{x}^{(m+1)}$  را از روی  $\mathbf{x}^{(m)}$  با حل معادله  $i$ ام نسبت به مجهول  $i$ ام به ازای هر  $i$  و به دنبال آن مقادیر فعلی سایر مجهول‌ها را از فرمول:

$$x_i^{(m+1)} = (b_i - \sum_{j \neq i} a_{ij}x_j^{(m)}) / a_{ii} \quad i = 1, \dots, n$$

به دست آورد. برای دستگاه خطی خاص

$$10x_1 + x_2 + x_3 = 12$$

$$x_1 + 10x_2 + x_3 = 12$$

$$x_1 + x_2 + 10x_3 = 12$$

روش بارست ژاکوبی که از نقطه  $\mathbf{x}^{(0)} = \mathbf{0}$  شروع می‌شود بردارهای  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(p)}$  واکه در جدول ۱۰۵ آورده شده‌اند تولید می‌کنند. به نظر می‌رسد که دنبالهٔ فوق به شکل مطلوبی به سمت جواب دستگاه یعنی  $[1 \ 1]^T$  همگراست. برای این مثال داریم

$$\|I - D^{-1}A\|_{\infty} = \max_i \left\{ \frac{1}{10} + \frac{1}{10} \right\} = 0.2$$

به طوری که می‌توانیم انتظار داشته باشیم که خطا در هر مرحله حداقل به اندازهٔ عامل ۰.۲ کاهش می‌یابد، این امر نیز از اعداد جدول ۱۰۵ آشکاراست. □

البته پیدا کردن وارون تقریبی  $C$  برای  $A$  اساساً آسان است. برای مثال رابطهٔ  $C = A^{-1}$  این کار را خواهد کرد و بارست مربوطه به آن در یک مرحله همگرا خواهد شد. اما نکته‌ای که در به کار گرفتن روش بارستی برای حل  $A\xi = \mathbf{b}$  حائز اهمیت است، این است که در وهلهٔ اول با روش بارستی ممکن است یک جواب تقریبی با دقت قابل قبولی بسیار سریعتر از حل مستقیم  $A\xi = \mathbf{b}$  به دست آورد. از این رو، مسئلهٔ مهم، انتخاب  $C$  به طریقی که بتوانیم به ازای هر بردار  $\mathbf{r}$  بردار  $C\mathbf{r}$  را نسبت به بردار  $A^{-1}\mathbf{r}$  با محاسبات کمتری به دست آوریم. به عنوان نمونه، می‌توان  $C$  را وارون یک ماتریس قطری (مانند بارست ژاکوبی) و یا وارون یک ماتریس مثلثی (مانند بارست گاوس-زایدل که در زیر بحث می‌شود) و یا وارون حاصلضرب دو ماتریس مثلثی (مانند اصلاح بارستی الگوریتم ۵.۴) و یا حتی وارون یک ماتریس سه قطری و غیره گرفت.

**الگوریتم ۵.۵:** بارست نقطه ثابت برای دستگاه‌های خطی دستگاه خطی  $A\xi = \mathbf{b}$  از مرتبهٔ  $n$  داده شده است

یک ماتریس  $C$  از مرتبهٔ  $n$  طوری انتخاب کنید که

- (i) به ازای یک بردار مفروض  $\mathbf{r}$ ، بردار  $C\mathbf{r}$  «به آسانی» محاسبه شود
- (ii) به ازای حداقل یک نرم ماتریسی داشته باشیم،  $\|I - CA\| < 1$ ، یک  $n$ -بردار  $\mathbf{x}^{(0)}$ ، مثلاً  $\mathbf{x}^{(0)} = \mathbf{0}$  را انتخاب کنید

For  $m = 0, 1, 2, \dots$ , until satisfied, do:

$$\mathbf{x}^{(m+1)} := \mathbf{x}^{(m)} + C(\mathbf{b} - A\mathbf{x}^{(m)})$$

در صورت نبودن خطای گرد کردن، دنبالهٔ حاصله  $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$  به سمت جواب دستگاه خطی داده شده همگرا می‌شود.

مانند فصل ۳، در اینجا عبارت «until satisfied» برای تأکید روی ناتمام بودن توضیح مفروض به کار رفته است. برای کامل کردن الگوریتم، می‌باید معیارهای دقیقی برای توقف بارستها تعریف شود. معیارهای نمونه به صورت زیرند:

$$\| \mathbf{x}^{(m)} - \mathbf{x}^{(m-1)} \| \leq \varepsilon \quad \text{الگوریتم را متوقف سازید اگر (الف):}$$

$$\frac{\| \mathbf{x}^{(m)} - \mathbf{x}^{(m-1)} \|}{\| \mathbf{x}^{(m)} \|} < \varepsilon \quad \text{یا اگر (ب):}$$

$$m > M \quad \text{یا اگر (ج): به ازای يك } M \text{ داده شده}$$

معیار آخر می‌باید همیشه در هر برنامه‌ای که الگوریتم فوق را اجرا می‌کند گنجانیده شود. تذکری را که برای اولین بار در بخش ۶.۱ داده شد، تکرار می‌کنیم: این حقیقت که رابطه

$$\| \mathbf{x}^{(m)} - \mathbf{x}^{(m-1)} \| < \varepsilon$$

برقرار است بدان معنی نیست که داشته باشیم

$$\| \mathbf{x}^{(m)} - \xi \| < \varepsilon$$

اما با توجه به (۱۳.۵) و (۱۸.۵) می‌دانیم که

$$\| \mathbf{x}^{(m)} - \xi \| \leq \frac{K}{1-K} \| \mathbf{x}^{(m)} - \mathbf{x}^{(m-1)} \| \quad (20.5)$$

$$K = \| I - CA \|$$

مثلا در مثال ۷.۵ برای بارست ژاکوبی به دست آوردیم که  $\| I - CA \|_{\infty} \leq 0.2$  و  $\| \mathbf{x}^{(6)} - \mathbf{x}^{(5)} \|_{\infty} = 0.000384$ . بنا بر این رابطه (۲۰.۵)، برآورد زیر را به دست می‌دهد

$$\| \xi - \mathbf{x}^{(6)} \|_{\infty} \leq \frac{0.2}{1-0.2} 0.000384 = 0.000096$$

و در واقع  $\| \xi - \mathbf{x}^{(6)} \| = 0.000064$ ، به طوری که خطا فقط به اندازه ۵۰ درصد زیادی برآورد شده است. متأسفانه، برآورد مطلوبی از  $\| I - CA \|$  معمولاً مشکل به دست می‌آید، یا آنکه این برآورد آنقدر به ۱ نزدیک است که مخارج  $1 - K$  در (۲۰.۵) را بیش از حد کوچک می‌سازد و بنا بر این کران حاصله در  $\| \xi - \mathbf{x}^{(m)} \|$  قابل استفاده نیست.

لازم است متذکر شویم که  $C$  ممکن است وارون تقریبی  $A$  باشد حتی اگر به ازای يك نرم ماتریسی داشته باشیم

$$\| I - CA \| > 1$$

تمام آنچه‌کسه برای يك وارون تقریبی  $C$  برای  $A$  (و برای همگرایی بارست نقطه ثابت

متناظر به آن) لازم است، این است که  $I - CA$  یک نرم ماتریسی کمتر از ۱ داشته باشد.

$$\text{مثلا ماتریس } B = \begin{bmatrix} 0.9 & 0.9 \\ 0 & 0 \end{bmatrix} \text{ در رابطه زیر صدق می کند}$$

$$\|B\|_{\infty} = 1.8 > 1$$

با این همه، در برخی نرم‌های ماتریسی داریم  $\|B\| < 1$ ، مثلا  $\|B\|_1 = 0.9$ . بنابراین، مهم پیدا کردن روشهایی است برای بیان اینکه به ازای یک نرم ماتریسی، نامساوی  $\|B\| < 1$  برقرار باشد (بی آنکه کلیهٔ نرم‌های ماتریسی ممکن را بررسی کنیم). قضیهٔ زیر (در اصل) چنین روشی را در اختیار ما قرار می‌دهد.

**قضیهٔ ۲.۵** گیریم  $\rho(B)$  شعاع طیفی ماتریس  $B$  باشد، یعنی

$$\rho(B) = \max \{ |\lambda| : \lambda \text{ يك ویژه مقدار ماتریس } B \text{ است} \}$$

در این حال به ازای جميع مقادیر  $\varepsilon > 0$  یک نرم برداری وجود دارد به طوری که نرم ماتریسی مربوط به  $B$  در رابطه  $\|B\| \leq \rho(B) + \varepsilon$  صدق می‌کند.

در نتیجه  $C$  یک وارون تقریبی برای ماتریس  $A$  است اگر و فقط اگر  $\rho(I - CA) < 1$ . علاوه بر چه شعاع طیفی ماتریس  $I - CA$  کوچکتر باشد همگرایی در روش بارست نقطهٔ ثابت، در پایان سریعتر خواهد بود.

مطلب فوق را می‌توان این گونه نشان داد، که در حال دستگاه  $A\xi = b$  خطای

$$e^{(m)} = \xi - x^{(m)}$$

در مرحلهٔ  $m$  بارست نقطهٔ ثابت

$$x^{(m+1)} = x^{(m)} + C(b - Ax^{(m)}), \quad m = 0, 1, 2, \dots$$

به ازای کلیهٔ مقادیر  $m$  در رابطه زیر صدق می‌کند

$$e^{(m+1)} = (I - CA)e^{(m)}$$

بنابراین به ازای کلیهٔ مقادیر  $n$  و با داشتن  $B = I - CA$ ، خواهیم داشت  $e^{(m)} = B^m e^{(0)}$ . این امر نشان می‌دهد که دنبالهٔ خطاها یعنی  $e^{(0)}, e^{(1)}, e^{(2)}, \dots$  به شکلی است که در فصل ۴ (به ۲.۴ تا ۶۷.۴ نگاه کنید) در رابطه با روش توانی مورد بحث قرار گرفت. در آنجا بیان کردیم که دنبالهٔ نرمال شدهٔ

$$e^{(m)} / \|e^{(m)}\| \quad m = 0, 1, 2, \dots$$

معمولاً به سمت يك ویژه بردار  $B = I - CA$  متعلق به بزرگترین ویژه مقدار (از نظر قدرمطلق)  $B$  همگرا می‌شود، یعنی

$$e^{(m+1)} = Be^{(m)} \approx \lambda e^{(m)}$$

که در آن  $|\lambda| = \rho(B)$ . بنابراین در نهایت به‌طور کلی مقدار خطا در هر مرحله از بارست به اندازه عامل  $\rho(B)$ ، و نه سریعتر، کاهش می‌یابد.

اکنون مثالهای خاصی از بارست نقطه ثابت برای دستگاههای خطی را مورد بحث قرار می‌دهیم. یکی از این مثالها روش اصلاح بارستی است که در فصل قبل مورد بحث قرار گرفت. برای یادآوری، در این روش ابتدا مانده  $\mathbf{r}^{(m)} = \mathbf{b} - A\mathbf{x}^{(m)}$  برای  $m$  امین جواب تقریبی  $\mathbf{x}^{(m)}$  محاسبه می‌شود، سپس با استفاده از تجزیه مثلثی  $A$  که هنگام حذف محاسبه می‌شود، جواب (تقریبی) برای  $\mathbf{y}^{(m)}$  دستگاه خطی  $A\mathbf{y} = \mathbf{r}^{(m)}$  به دست می‌آید و با افزودن  $\mathbf{y}^{(m)}$  به  $\mathbf{x}^{(m)}$  جواب تقریبی مناسب (چنین امید است)  $\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} + \mathbf{y}^{(m)}$  حاصل می‌شود. در حالت کلی بردار  $\mathbf{y}^{(m)}$  جواب (دقیق)  $A\mathbf{y} = \mathbf{r}^{(m)}$  نیست. علت این امر تا حدی ناشی از خطای گرد کردن هنگام پیشجایگذاری و پسجایگذاری است. می‌توان نشان داد که سهم اصلی در ایجاد خطا در  $\mathbf{y}^{(m)}$ ، معمولاً ناشی از بی‌دقتی در تجزیه  $A$  به صورت مثلثی  $PLU$ ، یعنی ناشی از این امر است که  $PLU$  تنها تقریبی است از  $A$ . اگر از خطای گرد کردن به هنگام پیشجایگذاری و پسجایگذاری صرف‌نظر کنیم، خواهیم داشت

$$\mathbf{y}^{(m)} = (PLU)^{-1} \mathbf{r}^{(m)} = (PLU)^{-1} (\mathbf{b} - A\mathbf{x}^{(m)})$$

بنابراین

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} + (PLU)^{-1} (\mathbf{b} - A\mathbf{x}^{(m)})$$

مطالب فوق نشان می‌دهند که روش اصلاح بادستی، حالت خاصی از روش بارست نقطه ثابت است و  $C$  همان تجزیه  $A$  به صورت مثلثی  $PLU$  است. اما برای رده‌هایی از ماتریسهای  $A$ ، ماتریسی مانند  $C$  می‌توان پیدا کرد که در حالت‌های (i) و (ii) الگوریتم ۵.۵ صادق باشد و حجم محاسباتی آن نسبت به محاسباتی که برای تجزیه مثلثی  $A$  لازم است، بسیار کمتر باشد. برای يك دستگاه خطی با چنین ماتریس ضرایبی، مفیدتر است که روش حذفی را کنار بگذاریم و جواب دستگاه را مستقیماً با الگوریتم ۵.۵ محاسبه کنیم.

برای بحث در مورد دو روش معمول برای انتخاب  $C$ ، ماتریس ضرایب  $A = (a_{ij})$  را به صورت حاصلجمع يك ماتریس اکیداً پایین مثلثی  $\hat{L} = (\hat{l}_{ij})$  و يك ماتریس قطری  $\hat{D} = (\hat{d}_{ij})$  و يك ماتریس اکیداً بالامثلی  $\hat{U} = (\hat{u}_{ij})$  می‌نویسیم،

$$A = \hat{L} + \hat{D} + \hat{U}$$

که در آن

$$\hat{l}_{ij} = \begin{cases} a_{ij} & i > j \\ 0 & i \leq j \end{cases} \quad \hat{d}_{ij} = \begin{cases} a_{ij} & i = j \\ 0 & i \neq j \end{cases} \quad \hat{u}_{ij} = \begin{cases} 0 & i \geq j \\ a_{ij} & i < j \end{cases}$$

بعلاوه، فرض می‌کنیم که تمام درایه‌های قطری  $A$  غیر صفر، یعنی  $D$  وارونپذیر باشد. اگر این شرط برقرار نباشد، ابتدا معادلات را طوری مرتب می‌کنیم که شرط فوق برقرار شود، که این کار هنگامی که  $A$  وارونپذیر باشد همیشه شدنی است. (به تمرین ۷.۴-۵ نگاه کنید).

در بارست ژاکوبی یا روش تغییر مکانهای همزمان<sup>۱</sup>، همان گونه که در مثال ۷.۵-۵ بحث شد، می‌توان  $C = D^{-1}$  انتخاب کرد.

اگر بارست ژاکوبی همگرا باشد، قسمت قطری  $A$  یعنی  $\hat{D}$ ، تقریب نسبتاً خوبی برای  $A$  است که در نتیجه داریم

$$\|B\| = \|I - \hat{D}^{-1}A\| < 1$$

اما در این حالت انتظار می‌رود که قسمت پایین-مثلثی  $A$  یعنی  $\hat{L} + \hat{D}$  حتی تقریب بهتری برای  $A$  باشد، یعنی انتظار می‌رود که داشته باشیم

$$\|I - (\hat{L} + \hat{D})^{-1}A\| \leq \|I - \hat{D}^{-1}A\| < 1$$

لذا به نظر می‌رسد که روش بارست نقطه ثابت با  $\hat{L} + \hat{D} = C^{-1}$ ، نسبت به روش ژاکوبی، بارستی با همگرایی سریعتر است. گرچه در حالت کلی مطلب فوق صحیح نیست، اما برای رده‌های مختلفی از ماتریس  $A$  صادق است؛ مثلاً هرگاه  $A$  نافذ قطری سطرأ مؤکد باشد یا هرگاه  $A$  سه قطری (یا به‌طور کلیتر هرگاه  $A$  یک بلوک سه قطری با بلوکهای قطری قطری باشد) یا هرگاه  $A$  دارای عناصر قطری مثبت و عناصر غیر قطری منفی باشد.

بارست نقطه ثابت با  $\hat{L} + \hat{D} = C^{-1}$  به بارست گاوس-زایدل<sup>۲</sup> یا روش تغییر مکانهای پیاپی<sup>۳</sup> موسوم است. در این روش داریم

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} + (\hat{L} + \hat{D})^{-1}(\mathbf{b} - A\mathbf{x}^{(m)})$$

یا

$$(\hat{L} + \hat{D})\mathbf{x}^{(m+1)} = (\hat{L} + \hat{D} - A)\mathbf{x}^{(m)} + \mathbf{b}$$

یا

$$\hat{D}\mathbf{x}^{(m+1)} = -\hat{L}\mathbf{x}^{(m+1)} - \hat{U}\mathbf{x}^{(m)} + \mathbf{b}$$

که نتیجه آن فرمول زیر است

$$x_i^{(m+1)} = \frac{-\sum_{j < i} a_{ij}x_j^{(m+1)} - \sum_{j > i} a_{ij}x_j^{(m)} + b_i}{a_{ii}} \quad i = 1, \dots, n$$

1. method of simultaneous displacements
2. Gauss-Seidel iteration
3. method of successive displacements

ظاهراً با دانستن  $x_1^{(m+1)}, \dots, x_{i-1}^{(m+1)}$  می‌توانیم درایه  $i$  ام  $\mathbf{x}^{(m+1)}$  را محاسبه کنیم.

**الگوریتم ۶.۵:** بارست گاوس-زایدل، دستگاه خطی  $A\mathbf{x} = \mathbf{b}$  از مرتبه  $n$ ، که همه درایه‌های قطری ماتریس ضرایب آن  $A = (a_{ij})$  غیر صفرند داده شده‌است.

درایه‌های  $B = (b_{ij})$  و  $\mathbf{c} = (c_i)$  را به وسیله روابط زیر محاسبه کنید

$$b_{ij} = \begin{cases} -a_{ij}/a_{ii} & i \neq j \\ c & i = j \end{cases}$$

$$c_i = \frac{b_i}{a_{ii}} \quad \text{به ازای جمیع مقادیر } i \text{ و } j$$

یک  $\mathbf{x}^{(0)}$ ، مثلاً  $\mathbf{x}^{(0)} = \mathbf{0}$  را انتخاب کنید

For  $m = 1, 2, \dots$ , until satisfied, do:

For  $i = 1, \dots, n$ , do:

$$x_i^{(m)} := \sum_{j=1}^{i-1} b_{ij} x_j^{(m)} + \sum_{j=i+1}^n b_{ij} x_j^{(m-1)} + c_i$$

اگر یک نرم ماتریسی از  $(\hat{L} + \hat{D})^{-1} \hat{U}$  کمتر از یک باشد، آنگاه دنباله  $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots$

که به طریق فوق تولید می‌شود به سمت جواب دستگاه داده شده همگراست.

بردارهای  $\mathbf{x}^{(1)}$  و  $\mathbf{x}^{(2)}$  که از به کار بردن روش بارستی گاوس-زایدل برای

دستگاه خطی مثال ۷.۵ نتیجه شده‌اند، در جدول ۱۰.۵ آورده شده‌اند. توجه کنید که در این

مثال، بارست گاوس-زایدل خیلی سریعتر از بارست ژاکوبی همگرا می‌شود. دقت عمل بعد

از سه مرحله بهتر از دقت عملی است که بعد از شش مرحله در بارست ژاکوبی به دست آمده

است.

در بارست ژاکوبی از درایه‌های  $\mathbf{x}^{(m)}$  تنها برای محاسبه بارست بعدی  $\mathbf{x}^{(m+1)}$

استفاده می‌شود؛ در حالی که در روش بارستی گاوس-زایدل درایه  $\mathbf{x}^{(m)}$  برای محاسبه

درایه‌های بعدی  $\mathbf{x}^{(m)}$  به کار برده می‌شود، و بدین دلیل است که نام تغییر مکانهای همزمان و تغییر

مکانهای پیاپی بر آن گذارده شده است. بخصوص بارست ژاکوبی مستلزم به خاطر سپردن

مقادیر دو بارست است، در حالی که بارست گاوس-زایدل فقط به یک بردار نیاز دارد.

می‌توان نشان داد که اگر ماتریس ضرایب  $A$  نافذ قطری (سطراً) مؤکد باشد، بارست

گاوس-زایدل همگراست. این روش همچنین همگراست اگر  $A$  معین و مثبت باشد، یعنی

اگر  $A$  حقیقی و متقارن باشد و به ازای جمیع بردارهای  $\mathbf{y}$  غیر صفر نامساوی زیر برقرار

باشد

$$\mathbf{y}^T A \mathbf{y} > 0$$



بالاخره از بین تکنیکهای متعدد موجود برای تسریع تقارب بارست نقطهٔ ثابت، ما به بیش واهلش پیمایی<sup>۱</sup> (با علامت اختصاری بوپ) اشاره می‌کنیم که در آن از بارست پیشنهادهی گاوس-زایدل دایر بر تغییر از  $x^{(m)}$  به  $x^{(m+1)}$  فراتر می‌روند. مثلاً به جای گرفتن رابطهٔ

$$x_i^{(m+1)} = x_i^{(m)} - \left( \sum_{j < i} a_{ij} x_j^{(m+1)} + \sum_{j > i} a_{ij} x_j^{(m)} - b_i \right) / a_{ii} \quad i$$

به‌روال گاوس-زایدل دایر بر تغییر از  $x^{(m)}$  به  $x^{(m+1)}$  گامی فراتر می‌نهند و از رابطهٔ

$$x_i^{(m+1)} = x_i^{(m)} - \omega \left( \sum_{j < i} a_{ij} x_j^{(m+1)} + \sum_{j > i} a_{ij} x_j^{(m)} - b_i \right) / a_{ii} \quad i$$

که در آن  $\omega$  (بزرگتر از ۱) پارامتر بیش واهلش نامیده می‌شود، استفاده می‌کنند. بارست حاصل را می‌توان به روشنی به صورت رابطهٔ زیر، که ولی چندان روشن نیست، نوشت

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} + C_\omega (\mathbf{b} - A\mathbf{x}^{(m)})$$

ماتریس بارست متناظر، به صورت  $(\hat{D} + \omega \hat{L})^{-1} (\hat{D} - \omega \hat{U})$  است. از لحاظ نظری، پارامتر بیش واهلش  $\omega$  می‌باید طوری انتخاب شود که  $\rho(I - C_\omega A)$  تا حد امکان کوچک باشد. البته در وهلهٔ اول، این کار در مقایسه با حل دستگاه معادلات  $A\mathbf{x} = \mathbf{b}$ ، مشکلتیر است. اما ممکن است بخواهیم چنین دستگاه خطی را برای سمت راستهای متفاوت و فقط تا حدی دقیق حل کنیم، که در این صورت صرفه در این است که از راه آزمایش مقدار «مطلوبی» برای  $\omega$  تعیین نماییم. همچنین برای ماتریسهای خاص  $A$  که در حل عددی معادلات استاندهٔ دیفرانسیل با مشتقات جزئی با آنها برخورد می‌کنیم، می‌توان  $\rho(I - C_\omega A)$  را بر حسب شعاع طیفی ماتریس بارست  $(\hat{D} + \hat{U})^{-1} \hat{L}$  - متعلق به بارست ژاکوبی نوشت و بنا بر این می‌توان اظهار نظری کیفی دربارهٔ بهترین انتخاب بهینهٔ  $\omega$  نمود. یک نمونهٔ انتخاب برای  $\omega$  مقداری بین ۱٫۲ و ۱٫۶ است.

چنانچه قبلاً اشاره شد، روشهای بارستی معمولاً برای دستگاههای خطی بزرگ که دارای ماتریس ضرایب تنگ می‌باشند به کار گرفته می‌شوند. در ماتریسهای تنگ تعداد عناصر غیر صفر زیاد نیست و بنابراین در هر مرحله تعداد عملیات حسابی که باید انجام گیرد کم است. علاوه روشهای بارستی نسبت به روشهای خطای گرد کردن حساسیت کمتری دارند و فقط اندازهٔ خطای گرد کردن حاصل در هر مرحله مهم است. از طرف دیگر روشهای بارستی همیشه همگرا نخواهند بود و حتی وقتی هم که همگرا می‌شوند ممکن است همگرای آنقدر کند باشد که رسیدن به جواب مستلزم بارستهای زیاد و غیر عملی باشد. در دستگاههای بزرگ، برای رسیدن به جوابی با چهار یا پنج رقم اعشار، تعداد کل بارستهای لازم ممکن است به چند صد تا برسد.

زیربنایی فکری که روش بارست ژاکوبی و گاوس-زایدل بر آن قرار دارد فکر واهلش است که این فکر در زمینه دستگاه کلی زیر که يك دستگاه  $m$  معادله  $m$  مجهولی غیرخطی است نیز به‌خوبی مصداق دارد

$$f(\xi) = 0$$

ساده‌ترین صورت دستگاه فوق این است که فرض شود معادلات طوری مرتب شده‌اند که بتوان با حل معادله  $i$ ام یعنی

$$f_i(\xi_1, \dots, \xi_n) = 0$$

نسبت به مجهول  $i$ ام، معادله هم‌ارز زیر را به‌دست آورد

$$\xi_i = g_i(\xi_1, \dots, \xi_{i-1}, \xi_{i+1}, \dots, \xi_n)$$

در این صورت برای يك تقریب مفروض  $X$  برای  $\xi$ ، سعی می‌شود که با تغییر مؤلفه  $i$ ام آن به‌صورت

$$g_i(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$$

این مؤلفه بهبود یابد.

واژه «واهلش» برای فرایند فوق منسوب به سوتول<sup>۱</sup> است. درحقیقت تخمین جاری  $X$  همان جواب دقیق برای دستگاه وابسته<sup>۲</sup>

$$x_i = g_i(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) + r_i \quad i = 1, \dots, n$$

است که در آن جمله‌های  $r_i$ ، که مربوط به خطا هستند، بدین جهت به حساب آمده‌اند که دستگاه مجبور شود جواب  $X$  را بپذیرد. در روش واهلش، درقبال برداشتن جمله  $r_i$  از معادله  $i$ ام، برای مؤلفه  $i$ ام مقدار جدیدی به‌دست می‌آید و بنا بر این مؤلفه فوق اصلاح می‌شود. واهلش معمولاً به‌اسلوب گاوس-زایدل انجام می‌گیرد، یعنی مقدار جدید مؤلفه  $i$ ام بلافاصله برای اصلاح مؤلفه دیگر به‌کار برده می‌شود. بعداً به‌اسلوبی منظم، بررسی تمام معادله‌ها تا به آخر انجام می‌گیرد و کلیه مؤلفه‌های  $X$  تغییر می‌یابد. اتمام يك بررسی کامل يك زده‌ایش<sup>۳</sup> نامیده می‌شود.

شکلهای متفاوت و مفیدی برای فکر واهلش وجود دارد. مثلاً ممکن است در بعضی موارد سهلتر باشد که به‌جای معادله  $i$ ام شکل هم‌ارز آن

$$\xi_i = g_i(\xi)$$

گذاشته شود، که در این فرمول هم‌ارز نیز، قسمت سمت راست آشکارا به  $\xi$  بستگی دارد.

به عنوان مثال دیگر، ممکن است با تغییر چندین مؤلفه از حدس جاری، معادله  $\mathbf{x}$  را فوراً واجد شرط کرد. به عبارت دیگر، ممکن است تخمین جدید  $\mathbf{x} + \alpha \mathbf{y}^{(i)}$  را طوری تعیین کرد که تساوی

$$f_i(\mathbf{x} + \alpha \mathbf{y}^{(i)}) = 0$$

که در آن برداری است ثابت و بستگی به  $i$  دارد، برقرار شود. البته برای واهلش معمولی داریم  $\mathbf{y}^{(i)} = \mathbf{i}_i$ .

□ مثال ۸.۵: سعی می کنیم دستگاه معادلات خطی مثال ۴.۵ یعنی

$$3\xi_i^2 = \xi_{i-1}^2 + \xi_{i-1}\xi_{i+1} + \xi_{i+1}^2 \quad i = 1, \dots, n$$

را با توجه به  $1 = \xi_{n+1} < \dots < \xi_1 < \xi_0 = 0$  با روش بارست گاوس-زایدل حل کنیم. بنابراین مانند مثال ۴.۵ با شروع از تخمین اولیه  $\mathbf{x} = [1 \ 2 \ \dots \ n]^T / (n+1)$  و  $n = 3$  بارست زیر را اجرا می کنیم

For  $m = 0, 1, 2, \dots$ , until satisfied, do:

For  $i = 1, \dots, n$ , do:

$$x_i = \sqrt{(x_{i-1}^2 + x_{i-1}x_{i+1} + x_{i+1}^2)/3}$$

در جدول زیر فهرستی از چند بارست اول که بعد از هر زدایش به دست آمده ارائه شده است. همگرایی خطی است (بنابراین با همگرایی روش نیوتن قابل مقایسه نیست)، اما همگرایی آنقدر منظم است که شتاب همگرایی را می توان پیدا کرد. استفاده از روش بیش-واهلش متوالی به ازای  $102 = \omega$  درست در ۱۰ زدایش، ۲۱ بارست فوق را به دست می دهد.

$m$	$x^{(m)}$	$\ x^{(m)} - x^{(m-1)}\ _1$
0	0.2500000	0.5000000
1	0.2886751	0.5361404
2	0.3095408	0.5612525
3	0.3240393	0.5734927
4	0.3311062	0.5794904
5	0.3345689	0.5824365
	...	...
10	0.3378204	0.5852071
11	0.3378695	0.5852490
	...	...
20	0.3379170	0.5852896
21	0.3379171	0.5852896

## تمرین

۱-۳۰۵ دستگاه معادلات

$$x - \sinh y = 0$$

$$2y - \cosh x = 0$$

را با استفاده از روش بارست نقطه ثابت حل کنید. يك جواب نزدیک به  $[0.6 \quad 0.6]^T$  وجود دارد.

۲-۳۰۵ با آزمایش، انتخاب مطلوبی برای پارامتر بیش و اهلس تعیین کنید، و آن را برای اهلس پایایی در مثال ۸.۵ به کار ببرید.

عملیات فوق را برای  $n = 10$  و سپس برای مسئله ۲۰۵-۶ مزبوط به آن نیز انجام دهید.

۳-۳۰۵ سعی کنید دستگاه

$$x^2 + xy^3 = 9 \quad 3x^2y - y^3 = 4$$

را با روش بارست نقطه ثابت حل کنید.

۴-۳۰۵ نشان دهید که روش بارست نقطه ثابت با ماتریس بارست

$$B = \begin{bmatrix} 0.9 & 1000 \\ 0 & 0.9 \end{bmatrix}$$

همگرا خواهد بود، حتی اگر  $\|B\|_\infty = \|B\|_1 > 1000$ .

۵-۳۰۵ با استفاده از قضیه شور ثابت کنید که برای هر ماتریس مربعی  $B$  و هر  $\varepsilon > 0$ ، يك نرم برداری وجود دارد که نرم ماتریسی متناظر آن در رابطه  $\|B\| \leq \rho(B) + \varepsilon$  صدق می‌کند. (داهنمایی: نرم برداری را به شکل  $\|x\| := \|DUX\|_\infty$  تشکیل دهید که در آن  $U$  به وسیله قضیه شور چنان معین می‌شود که  $A = U^{-1}BU$  يك ماتریس بالامتلی باشد و  $D = \text{diag}[1, \delta, \delta^2, \dots, \delta^{n-1}]$  چنان انتخاب می‌شود که کلیه درایه‌های غیرقطری  $D^{-1}AD$  از لحاظ قدرمطلق مقداری کمتر از  $\varepsilon/n$  داشته باشند.)

۶-۳۰۵ نشان دهید که روشهای بارستی ژاکوبی و گاوس-زایدل که برای حل دستگاه معادلات خطی  $A\xi = b$ ، که در آن  $A$  ماتریس بالامتلی و ارونبدیری است، به کار می‌روند در مراحل متناهی زیادی همگرا می‌شوند.

## ۳-۳-۵ دستگاه معادلات

$$\begin{bmatrix} 4 & -1 & 0 & 0 \\ -1 & 4 & -1 & 0 \\ 0 & -1 & 4 & -1 \\ 0 & 0 & -1 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

را با استفاده از روش بارست ژاکوبی و بارست گاوس-زایدل حل کنید. همچنین با استفاده از الگوریتم ۳-۴؛ تجزیه به عوامل ماتریس ضرایب دستگاه فوق را به دست آورید. آنگاه با به کارگیری روش اصلاح بارستی و همان حدس اولیه، دستگاه فوق را حل کنید. حجم محاسبات (عملیات ممیز شناور) لازم برای به دست آوردن جواب تقریبی بسا دقت مطلق کمتر از  $10^{-6}$  را در مورد سه روش فوق، برآورد کنید.

۳-۳-۵ ثابت کنید که روش بارست ژاکوبی همگراست، اگر ماتریس ضرایب  $A$ ی دستگاه نافذ قطری ستوناً مؤکد باشد یعنی

$$|a_{jj}| > \sum_{i \neq j} |a_{ij}| \quad j = 1, \dots, n$$

(دادهمایی: از نرم ماتریسی متناظر با نرم برداری،  $\|x\| = \sum_{i=1}^n |a_{ii}x_i|$  استفاده کنید.)

## تقریب

در این فصل مسئله تقریب زدن يك تابع کلی را به كمك ردهای از توابع ساده تر بررسی می کنیم. برای این توابع تقریب دو کاربرد وجود دارد. نخست گذاردن توابع ساده تر به جای توابع پیچیده تا عملیات معمولی زیادی نظیر مشتقگیری و انتگرالگیری یا حتی محاسبه توابع بتواند بسیار راحت تر انجام گیرد. کاربرد مهم دیگر آن بازیابی يك تابع از روی اطلاعات جزئی در بساب آن، مثلاً از يك جدول تقریبی (احتمالاً فقط تقریبی) است. متداولترین رده توابع تقریب کننده، بسجمله‌ایهای جبری، بسجمله‌ایهای مثلثاتی و اخیراً هم توابع بسجمله‌ای-تکه‌ای بوده‌اند. با بهترین تقریب و تقریب مطلوب به وسیله هر يك از این رده‌ها را بررسی می کنیم.

## ۱.۶ تقریب یکنواخت به وسیله بسجمله‌ایها

در این بخش مسئله ساختن يك بسجمله‌ای  $p(x)$  از درجه نایبتر از  $n$  را که تابع داده شده  $f(x)$  را به طور یکنواخت در بازه  $a \leq x \leq b$  تقریب می زند مورد بررسی قرار می دهیم. منظور این است که خطا در تقریب  $P(x)$  برای تابع  $f(x)$  را، به وسیله عدد یا نرم زیر اندازه می گیریم

$$\|f - p\|_{\infty} = \max_{a \leq x \leq b} |f(x) - p(x)| \quad (1.6)$$

آرمانی صحبت کنیم، می خواهیم بهترین تقریب یکنواخت از  $\Pi_n$  را به دست آوریم،

یعنی می‌خواهیم يك بسجمله‌ای  $p_n^*(x)$  از درجهٔ نایبشتر از  $n$  به‌دست آوریم که برای آن داشته باشیم

$$\|f - p_n^*\|_\infty = \min_{p \in \Pi_n} \|f - p\|_\infty \quad (۲.۶)$$

در اینجا نماد  $p \in \Pi_n$  را به‌عنوان تلخیصی از عبارت « $p$  يك بسجمله‌ای از درجهٔ نایبشتر از  $n$  است» به‌کار برده‌ایم؛ به‌عبارت دیگر،  $p_n^*$  بسجمله‌ای خاصی است از درجهٔ نایبشتر از  $n$  که همانقدر به‌تابع  $f$  نزدیک است که يك بسجمله‌ای از درجهٔ نایبشتر از  $n$  ممکن است به آن نزدیک باشد. عدد  $\|f - p_n^*\|_\infty$  را با

$$\text{dist}_\infty(f, \Pi_n)$$

نشان می‌دهیم و آن را فاصلهٔ یکنواخت بسجمله‌ای  $f$  از درجهٔ نایبشتر از  $n$  در بازهٔ  $a \leq x \leq b$  می‌نامیم.

پیش از آنکه بحثی دربارهٔ ساختن يك بسجمله‌ای مطلوب یا بهترین بسجمله‌ای تقریب‌زننده انجام دهیم، اندکی به بررسی راه‌های برآورد  $\text{dist}_\infty(f, \Pi_n)$  می‌پردازیم. مثلاً اگر چنین برآوردی نشان دهد که  $\text{dist}_\infty(f, \Pi_n) > ۱۰$ ، و ما در پی تقریبی باشیم که تا دورقم آخر بعد از ممیز مطلوب باشد، آنگاه بیهوده وقت و نیرو صرف ساختن  $p_n^*$  نخواهیم کرد. برای چنین منظوری به‌ویژه بسیار مهم است که گرانه‌های پایینی برای  $\text{dist}_\infty(f, \Pi_n)$  به‌دست آوریم و یکی از راه‌های به‌دست آوردن این گرانه‌ها به‌صورت زیر است.

با توجه به فصل ۲ خاطر نشان می‌کنیم که

$$g[x_0, \dots, x_{n+1}] = \sum_{i=0}^{n+1} g(x_i) / w'(x_i)$$

که در آن داریم

$$w(x) = (x - x_0) \cdots (x - x_{n+1})$$

(به‌تمرین ۲.۲-۱ نگاه کنید) موگفتیم که اگر  $g(x)$  يك بسجمله‌ای از درجهٔ نایبشتر از  $n$  باشد، این  $(n+1)$  امین عامل تفاضل برابر یا صفر خواهد بود (به‌تمرین ۲.۲-۵ نگاه کنید). بنابراین برای هر بسجمله‌ای خاص  $p \in \Pi_n$ ، داریم

$$\begin{aligned} f[x_0, \dots, x_{n+1}] &= f[x_0, \dots, x_{n+1}] - p[x_0, \dots, x_{n+1}] \\ &= (f - p)[x_0, \dots, x_{n+1}] \\ &= \sum_{i=0}^{n+1} (f(x_i) - p(x_i)) / w'(x_i) \end{aligned}$$

در نتیجه اگر  $x_0, \dots, x_{n+1}$  همگی در بازه  $a \leq x \leq b$  باشند، آنگاه

$$|f[x_0, \dots, x_{n+1}]| \leq \|f - p\|_\infty \cdot W(x_0, \dots, x_{n+1})$$

که در آن عدد مثبت  $W(x_0, \dots, x_{n+1})$  از رابطه

$$W(x_0, \dots, x_{n+1}) = \sum_{i=0}^{n+1} 1/|w'(x_i)| \quad (۳.۶)$$

معین می‌شود. اکنون  $p$  را چنان انتخاب می‌کنیم که خصوصیت  $p_n^*$  را دارا باشد. بنابراین  $\|f - p\|_\infty = \text{dist}_\infty(f, \pi_n)$  و کران پایین زیر را به دست می‌آوریم

$$|f[x_0, \dots, x_{n+1}]|/W(x_0, \dots, x_{n+1}) \leq \text{dist}_\infty(f, \pi_n) \quad (۴.۶)$$

□ مثال ۱۰۶: به ازای  $n=1$  و  $x_0 = -1, x_1 = 0, x_2 = 1$  داریم

$$\begin{aligned} f[x_0, x_1, x_2] &= \frac{f(x_0)}{(x_1 - x_0)(x_2 - x_0)} + \frac{f(x_1)}{(x_0 - x_1)(x_2 - x_1)} \\ &\quad + \frac{f(x_2)}{(x_0 - x_2)(x_1 - x_2)} \\ &= f(-1)/2 - f(0) + f(1)/2 \end{aligned}$$

از این رو  $W(-1, 0, 1) = 2$  و بنابراین برای  $a \leq -1$  و  $b \leq 1$  خواهیم داشت

$$|f[-1, 0, 1]|/2 \leq \text{dist}_\infty(f, \pi_1)$$

مثلا، برای  $f(x) = e^x$  داریم  $0.5843 \leq \text{dist}_\infty(e^x, \pi_1) \leq 0.5854$  و در نتیجه،  $\text{dist}_\infty(e^x, \pi_1) \geq 0.5843$  □

استفاده از کران پایین (۴.۶) مستلزم محاسبه اعداد  $w'(x_i) = \prod_{j \neq i} (x_i - x_j)$  برای تشکیل  $W(x_0, \dots, x_{n+1})$  است به منظور ملاحظه يك روش کارا برای انجام این امر ۱۴-۱۰۶ نگاه کنید). برای انتخاب بعضی مقادیر  $x_i$ ، این اعداد شکل ساده خاصی پیدا می‌کنند. مثلا اگر

$$x_i = \cos \frac{i}{n+1} \pi \quad i = 0, \dots, n+1 \quad (۵.۶)$$

آنگاه

$$\frac{1}{w'(x_i)} = \frac{2^{n-1}}{n+1} (-1)^i \begin{cases} 1 & \text{اگر } i \text{ مساوی صفر یا } n+1 \text{ باشد} \\ 2 & \text{در غیر این صورت} \end{cases} \quad (۶.۶)$$



از این رو،  $W(x_0, \dots, x_{n+1}) = 2^n$  (به تمرین ۱.۶-۵ نگاه کنید) و بنابراین اگر بازه  $a \leq x \leq b$  شامل هر دو مقدار ۱ و -۱ باشد، خواهیم داشت

$$\frac{1}{2(n+1)} \left| f(1) - 2f\left(\cos \frac{\pi}{n+1}\right) + 2f\left(\cos \frac{2\pi}{n+1}\right) - \dots - (-1)^n f(-1) \right| \leq \text{dist}_\infty(f, \pi_n) \quad (7.6)$$

برای استفاده از این کران پایین در بازه‌های دیگر، اول باید با یک تغییر خطی متغیرها بازه مورد نظر را به بازه  $-1 \leq x \leq 1$  منتقل ساخت.

□ مثال ۷.۶: در بازه استاندارد  $-1 \leq x \leq 1$ ، تابع تقریب  $\pi_n$  را برای تابع

$$f(x) = \tan \pi / 4x$$

در نظر بگیرید. این یک تابع فرد است، یعنی  $f(-x) = f(x)$ ، بنابراین به ازای مقادیر فرد  $n$ ، کران پایین برابر با صفر می‌شود که این امر هیچ کمکی به ما نمی‌کند. در عوض تابع تقریب  $\pi_n$  را در نظر می‌گیریم. در این صورت با توجه به رابطه (۷.۶) داریم

$$\frac{1}{10} \left| f(1) - 2f\left(\cos \frac{\pi}{5}\right) + 2f\left(\cos \frac{2\pi}{5}\right) - 2f\left(\cos \frac{3\pi}{5}\right) + 2f\left(\cos \frac{4\pi}{5}\right) - f(-1) \right| \leq \text{dist}_\infty(f, \pi_4)$$

و یا  $0.000203 \leq \text{dist}_\infty(f, \pi_4) \leq 0.00041 \dots$  در حقیقت می‌توان نشان داد که

$$\text{dist}_\infty(f, \pi_4) = 0.00041 \dots$$

□ و بنابراین کران پایین کاملاً خوب است.

در رابطه با این کرانه‌های پایین، قضیهٔ زیر وجود دارد که به «دولا واله-پوسن»<sup>۲</sup> نسبت داده می‌شود که ما را از محاسبهٔ  $w'(x_i)$  برحذر می‌دارد، اما مستلزم ساختن تقریب کننده برای  $p \in \pi_n$  است.

قضیه ۱.۶ فرض کنید که خطای  $f(x) - p(x)$  در تقریب بسجمله‌ای  $p \in \pi_n$  برای  $f$  در نقاط  $x_0 < x_1 < \dots < x_{n+1}$  تغییر علامت دهد، یعنی

$$(-1)^i [f(x_i) - p(x_i)] \varepsilon \geq 0 \quad i = 0, \dots, n+1$$

که در آن  $\varepsilon = \text{signum}[f(x_0) - p(x_0)]$ . اگر به ازای جمیع مقادیر  $i$  داشته باشیم  
 آنگاه  $a \leq x_i \leq b$

$$\text{dist}_\infty(f, \pi_n) \geq \min_i |f(x_i) - p(x_i)|$$

در واقع، همان گونه که در قضیه فرض شده است، اگر نقاط  $x_i$  مرتب شده باشند،  
 آنگاه خواهیم داشت

$$(-1)^{n+1-i} w'(x_i) > 0 \quad i = 0, \dots, n+1$$

و بنابراین تمام جمعوندها در حاصلجمع

$$(f-p)[x_0, \dots, x_{n+1}] = \sum_{i=0}^{n+1} [f(x_i) - p(x_i)] / w'(x_i)$$

دارای یک علامت هستند. اما این بدان معناست که

$$|f[x_0, \dots, x_{n+1}]| = \sum_{i=0}^{n+1} |(f(x_i) - p(x_i)) / w'(x_i)|$$

$$\geq \min_i |f(x_i) - p(x_i)| \cdot W(x_0, \dots, x_{n+1})$$

و این رابطه همراه با رابطه (۴.۶)، قضیه را ثابت می کند.

اکنون فرض کنید ترتیبی اتخاذ کنیم که در قضیه ۱.۶، به ازای  $i = 0, 1, \dots, n+1$   
 تساوی  $\|f - p\|_\infty = |f(x_i) - p(x_i)|$  نیز برقرار باشد. در این صورت خواهیم داشت

$$\|f - p\|_\infty \geq \text{dist}_\infty(f, \pi_n) \geq \min_i |f(x_i) - p(x_i)| = \|f - p\|_\infty$$

و چون اولین و آخرین عبارت در این رشته نامساویها با هم یکی هستند، بنابراین کلا باید  
 تساوی برقرار باشد. بخصوص بسجمله ای  $p$  می باید بهترین تقریب یکنواخت برای  $f$   
 از نوع  $\pi_n$  باشد. این نکته اثبات نیمه آسان قضیه زیر منسوب به چیشف است.

**قضیه ۲.۰۶** تابع  $f$  که در بازه  $a \leq x \leq b$  پیوسته است دقیقاً بهترین تقریب یکنواخت از  
 نوع  $\pi_n$  را در  $a \leq x \leq b$  دارد. بسجمله ای  $p \in \pi_n$  بهترین تقریب یکنواخت برای  $f$  در  
 بازه  $a \leq x \leq b$  است اگر، و فقط اگر،  $n+2$  نقطه  $a \leq x_0 < \dots < x_{n+1} \leq b$   
 چنان وجود داشته باشند که

$$(-1)^i [f(x_i) - p(x_i)] = \varepsilon \|f - p\|_\infty \quad i = 0, \dots, n+1 \quad (۸.۶)$$

که در آن  $\varepsilon = \text{signum}[f(x_0) - p(x_0)]$ . در اینجا درحالتی که  $f^{(n+1)}(x)$  در بازه

$a \leq x \leq b$  تغییر علامت ندهد، خواهیم داشت  $a = x_0$  و  $b = x_{n+1}$ .

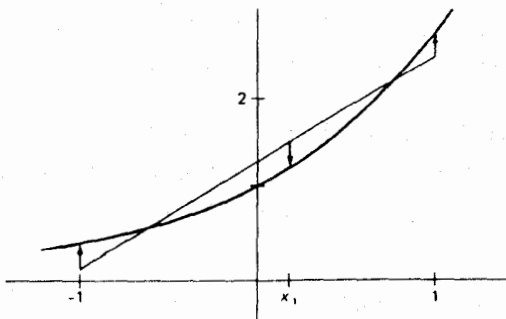
يك برهان این قضیهٔ اساسی را می‌توان در هر کتاب درسی در باب نظریهٔ تقریب، مثلاً در کتاب رایس<sup>۱</sup> [۱۷] یا ریولین<sup>۲</sup> [۳۵] پیدا کرد.

□ **مثال ۳۰۶:** باردیگر به دست آوردن تقریب برای  $f(x) = e^x$  را در بازهٔ استاندارد  $-1 \leq x \leq 1$  بررسی می‌کنیم. در مثال ۱۰۶ مشاهده کردید که  $\text{dist}_\infty(f, \pi_1) \geq 0.27$ . اکنون  $p(x)$  را به صورت  $p(x) = a + bx$ ، با  $b = (e^1 - e^{-1})/2$  و  $a = (e - bx_1)/2$  که در آن  $e^{x_1} = f'(x_1) = p'(x_1) = b$  یا  $x_1 = \ln b$ ، انتخاب می‌کنیم؛ به شکل ۱۰۶ نگاه کنید. در این صورت می‌توان تحقیق کرد که  $f(x) - p(x)$  در شرط تناوب (۸۰۶) با  $n=1$  و  $x_0 = -1$ ،  $x_1 = 1$  صدق می‌کند، یعنی

$$\begin{aligned} f(-1) - p(-1) &= -[f(x_1) - p(x_1)] = f(1) - p(1) \\ &= (e^1 + e^{-1})/2 - a = 0.27880\dots \end{aligned}$$

لذا این خط مستقیم خاص می‌باید بهترین تقریب یکنواخت برای  $e^x$  در بازهٔ  $-1 \leq x \leq 1$  از  $\pi_1$  باشد و  $\text{dist}(e^x, \pi_1) = 0.27880\dots$ . این رابطه نشان می‌دهد که کران پایین به دست آمده در مثال ۱۰۶ کاملاً دقیق است. □

يك مثال به‌ویژه مهم با بهترین تقریب یکنواخت در بازهٔ  $-1 \leq x \leq 1$  از  $\pi_n$  برای تابع  $f(x) = x^{n+1}$  فراهم شده است. خطا در این تقریب، به طوری که مشاهده خواهد شد، مضربی است از  $T_{n+1}(x)$ ، که بسجمله‌ای چبیشف از درجهٔ  $n+1$  است. بنابراین تعریف، بسجمله‌ای چبیشف از درجهٔ  $k$  (در بازهٔ  $-1 \leq x \leq 1$ ) با قاعدهٔ زیر معین می‌شود



شکل ۱۰۶. بهترین تقریب یکنواخت خط راست برای  $e^x$  در  $-1 \leq x \leq 1$

$$T_k(\cos \theta) = \cos k\theta \quad (۹.۶)$$

بنا بر این

$$T_0(x) = 1 \quad T_1(x) = x \quad (۱۰.۶)$$

و به موجب فرمول جمع برای توابع مثلثاتی داریم

$$T_{k+1}(x) = 2xT_k(x) - T_{k-1}(x) \quad k = 1, 2, \dots \quad (۱۱.۶)$$

از رابطه فوق خواهیم داشت

$$T_2(x) = 2xT_1(x) - T_0(x) = 2x^2 - 1$$

$$T_3(x) = 2xT_2(x) - T_1(x) = 2x(2x^2 - 1) - x = 4x^3 - 3x$$

وهشت بسجمله‌ای اول از این بسجمله‌ایها در جدول ۱.۶ و نمودار پنج بسجمله‌ای اول آنها در شکل‌های ۲.۶ و ۳.۶ داده شده‌اند.

رابطه بازگشتی (۱۱.۶) روشن می‌سازد که  $T_k(x)$  که با (۹.۶) تعریف شده در واقع یک بسجمله‌ای است دقیقاً از درجه  $k$  با ضریب بزرگترین درجه  $2^{k-1}$ . بعلاوه از تعریف (۹.۶) آشکار است که

$$|T_k(x)| \leq 1 \quad -1 \leq x \leq 1 \quad \text{به‌ازای جمیع مقادیر} \quad (۱۲.۶)$$

و  $T_k(x)$  متناوباً در  $k+1$  نقطه

### جدول ۱.۶

---


$$T_0(x) = 1$$

$$T_1(x) = x$$

$$T_2(x) = 2x^2 - 1$$

$$T_3(x) = 4x^3 - 3x$$

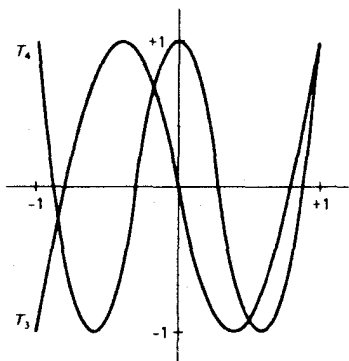
$$T_4(x) = 8x^4 - 8x^2 + 1$$

$$T_5(x) = 16x^5 - 20x^3 + 5x$$

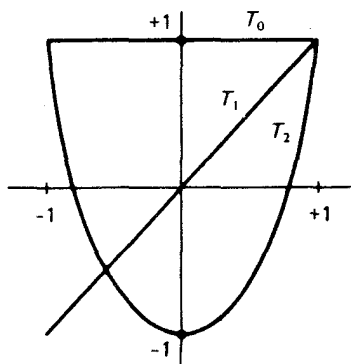
$$T_6(x) = 32x^6 - 48x^4 + 18x^2 - 1$$

$$T_7(x) = 64x^7 - 112x^5 + 56x^3 - 7x$$


---



شکل ۳.۶



شکل ۳.۶

$$x_j = \cos \frac{k-j}{k} \pi \quad j = 0, \dots, k$$

به این کران  $\pm 1$  می‌رسد. یعنی از قاعدهٔ (۹.۶) به دست می‌آوریم

$$T_k(x_j) = T_k\left(\cos \frac{k-j}{k} \pi\right) = \cos k \frac{k-j}{k} \pi = (-1)^{k-j} \quad j = 0, \dots, k$$

اما این رابطه نشان می‌دهد که مخصوصاً برای بسجمله‌ای  $p_n(x)$  از درجهٔ نایبتر از  $n$  تساوی

$$2^{-n} T_{n+1}(x) = x^{n+1} - p_n(x)$$

برقرار است و این بسجمله‌ای، بنا بر قضیهٔ ۳.۶ بهترین تقریب یکنواخت برای  $x^{n+1}$  در بازهٔ  $-1 \leq x \leq 1$  است و نیز

$$\text{dist}_\infty(x^{n+1}, \pi_n) = 1/2^n \quad (13.6)$$

در حالت کلی به دست آوردن بهترین تقریب یکنواخت از  $\pi_n$  کار ساده‌ای نیست. با فرض مشتق‌پذیر بودن تابع  $f(x)$ ، می‌توان بر اساس قضیهٔ ۳.۶ و با قید

$$a \leq x_0 < \dots < x_{n+1} \leq b$$

دستگاه معادلات غیر خطی

$$f(x_i) - p_n^*(x_i) = (-1)^i d \quad i = 0, \dots, n+1$$

$$\phi(x_i)[f'(x_i) - p_n^{*'}(x_i)] = 0 \quad i = 0, \dots, n+1 \quad (14.6)$$

را برای نقاط  $x_0, \dots, x_{n+1}$ ، و برای  $n+1$  ضریب از  $p_n^*(x)$  و عدد (مثبت یا منفی)  $d = \pm \|f - p_n^*\|_\infty$  حل نمود. در اینجا

$$\phi(x) = \begin{cases} 0 & \text{اگر } x = a \text{ یا } x = b \\ 1 & \text{در غیر این صورت} \end{cases}$$

تابع  $\phi(x)$  عاملی است برای تمیز بین یک فرینۀ داخلی خطای  $f(x) - p_n^*(x)$  که در آن مشتق اول می‌بایسد صفر شود، و یک فرینۀ کرانی که در آن نیازی به صفر شدن مشتق نیست (گرچه می‌باید در نامساوی دیگری که در اینجا ذکر نشده است صدق کند). در الگوریتم رمه<sup>۲</sup> و در شکل دیگر آن یعنی شکل مورناهان-رنج<sup>۳</sup> (به کتاب رایس [۱۷] نگاه کنید) سعی بر آن است که این دستگاه معادلات را با روش نیوتن که در فصل ۵ بحث شد حل کنند. اما برای انجام این کار دستگاه معادلات را بسا ساختار خاصی از دستگاه غیر خطی (۱۴.۶) تطابق می‌دهند. یک تخمین اولیه با یک بسجمله‌ای درونیاب مناسب  $p_n \in \pi_n$  برای  $f(x)$  و بسا به کارگیری ضرایب  $p_n(x)$  و فرینۀ موضعی  $f(x) - p_n(x)$  به آسانی به دست می‌آید. در اینجا وقت خوانندگان را برای بحث دربارهٔ جزئیات ایجاد بهترین بسجمله‌ای تقریب یکنواخت نمی‌گیریم. زیرا تقریباً بهترین تقریب را می‌توان با سعی کمتر و بسا درونیابی مناسب به دست آورد.

در حقیقت با توجه به قضیۀ ۲.۶ می‌دانیم که خطای  $f(x) - p_n^*(x)$  در بهترین تقریب یکنواخت برای تابع پیوستۀ  $f(x)$  در بازۀ  $a \leq x \leq b$  می‌باید  $n+1$  مرتبه تغییر<sup>۴</sup> کند، یعنی می‌باید در رابطۀ

$$(-1)^i [f(x_i) - p_n^*(x_i)] = \varepsilon \|f - p_n^*\|_\infty \quad i = 0, \dots, n+1$$

با شرایط  $\varepsilon = \text{signum}[f(x_0) - p_n^*(x_0)]$  و  $a \leq x_0 < \dots < x_{n+1} \leq b$  صدق نماید. در این صورت به موجب قضیۀ ارزش متوسط برای توابع پیوسته (قضیۀ ۱.۳)، باید نقاطی مانند  $\xi_0 < \dots < \xi_n$  با شرط  $x_i < \xi_i < x_{i+1}$ ، به ازای جمیع مقادیرند، طوری وجود داشته باشند که در این نقاط، خطای  $f(x) - p_n^*(x)$  برابر صفر شود، یعنی در این نقاط، در درونیابی  $f(x)$  مقدار  $p_n^*(x)$  بهترین تقریب باشد. بنابراین در اصل تنها اگر بدانیم که در کدام نقاط باید درونیابی را انجام دهیم آنگاه می‌توانیم حتی بهترین تقریب را از راه درونیابی بنا کنیم.

اما خاطر نشان می‌سازیم که بهترین تقریب برای  $x^{n+1}$  از  $\pi_n$  در بازۀ استاندارد  $1 \leq x \leq -1$  مضر بی است از بسجمله‌ای چیشف از درجه  $n+1$  یعنی  $T_{n+1}(x)$ ، که بنا بر خود تعریف (۹.۶) در رابطۀ  $n+1$  نقطه

1. extremum
2. Remez
3. Murnaghan-Wrench
4. alternate

$$\xi_{k, n+1} = \cos \frac{2k+1}{2n+2} \pi \quad k = 0, \dots, n \quad (15.6)$$

صفر می‌شود. این بدان معنی است که برای تابع خاص  $f(x) = x^{n+1}$  می‌توانیم بهترین تقریب‌سازیکنوخت را از  $\pi_n$ ، از راه درونیایی در نقاط (۱۵.۶)، به اصطلاح نقاط چبیشف برای بازه استاندارد  $-1 \leq x \leq 1$ ، به دست آوریم. چنانکه معلوم خواهد شد، این شیوه عمل، تقریبهای خوبی (اگر بهترین نباشد) برای توابع پیوسته تولید می‌کند.

برای اینکه ببینیم چرا باید چنین باشد با توجه به روابط (۱۶.۲) یا (۳۷.۲) یادآوری می‌کنیم که خطای  $f(x) - p_n(x)$  در بسجمله‌ای درونیاب برای  $f(x)$  در نقاط  $x_0, \dots, x_n$  در رابطه زیر صدق می‌کند

$$f(x) - p_n(x) = f[x_0, \dots, x_n](x - x_0) \dots (x - x_n)$$

در نتیجه، به موجب رابطه (۴.۶) داریم

$$|f(x) - p_n(x)| \leq |x - x_0| \dots |x - x_n| \cdot W(x_0, \dots, x_n, x) \text{dist}_\infty(f, \pi_n)$$

به شرط آنکه  $x_0, \dots, x_n$  و  $x$  در بازه مورد نظر قرار گرفته باشند. اکنون قرار می‌دهیم  $x = x_{n+1}$ . لذا با توجه به رابطه (۳.۶) داریم

$$W(x_0, \dots, x_n, x) = \sum_{i=0}^{n+1} \left| \prod_{\substack{j=0 \\ j \neq i}}^{n+1} |x_j - x_i| \right|$$

و بنا بر این

$$\begin{aligned} & |x - x_0| \dots |x - x_n| W(x_0, \dots, x_n, x) \\ &= |x_{n+1} - x_0| \dots |x_{n+1} - x_n| W(x_0, \dots, x_{n+1}) \\ &= \sum_{i=0}^{n+1} \prod_{\substack{j=0 \\ j \neq i}}^{n+1} \left| \frac{x_{n+1} - x_j}{x_i - x_j} \right| \\ &= 1 + \sum_{i=0}^n |l_i(x_{n+1})| \end{aligned}$$

که در آن

$$l_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j} \quad (16.6)$$

تأمین بسجمله‌ای لاگرانژ است [به (۵.۲) و (۶.۲) نگاه کنید]. این مطلب اثبات قضیه زیر است

قضیه ۳.۶ گیریم  $p_n(x)$  یک بسجمله‌ای از درجهٔ نا بیشتر از  $n$  باشد که  $f(x)$  را در نقاط  $x_0 < x_1 < \dots < x_n$  در بازهٔ مطلوب  $a \leq x \leq b$  درونیابی کند. در این صورت داریم

$$\text{dist}_\infty(f, \pi_\varphi) \leq \|f - p_n\|_\infty \leq (1 + \|\Lambda_n\|_\infty) \text{dist}_\infty(f, \pi_n) \quad (۱۷.۶)$$

که در آن  $\Lambda_n$  با تساوی

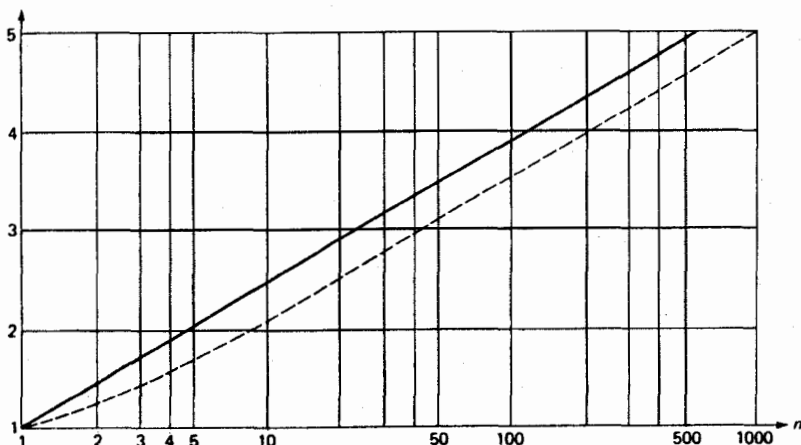
$$\Lambda_n(x) = \sum_{i=0}^n |l_i(x)|$$

داده شده است و بسجمله‌ای لاگرانژ  $l_i(x)$  با رابطهٔ (۱۶.۶).

این امر ایجاب می‌کند که نقاط درونیابی  $x_0, \dots, x_n$  در بازهٔ  $a \leq x \leq b$  چنان انتخاب شود که نرم یکنواخت  $\|\Lambda_n\|_\infty$  از تابع لیگند  $\Lambda_n(x)$  تا حد ممکن کوچک باشد. انفاقاً این امر تقریباً به وسیلهٔ نقاط چبیشف (۱۵.۶) که در بازهٔ مورد نظر  $a \leq x \leq b$  تنظیم شده‌اند، یعنی با نقاط

$$x_i = \left[ a + b + (a - b) \cos \frac{\varphi(i+1)\pi}{\varphi(n+1)} \right] / 2 \quad i = 0, \dots, n \quad (۱۸.۶)$$

امکانپذیر می‌شود. در شکل ۴.۶،  $\|\Lambda_n\|_\infty$  بر حسب این نقاط به عنوان تابعی از  $n$  ترسیم شده است. اعداد  $\|\Lambda_n^e\|_\infty$  که متناظر با نقاط به اصطلاح مبسوط چبیشف<sup>۲</sup>



شکل ۴.۶ عدد  $\|\Lambda_n\|_\infty$  برای نقاط چبیشف (خط پر) و برای نقاط مبسوط چبیشف (خط چین) به عنوان تابعی از  $n$ .



$$x_i^e = \left[ a + b + (a - b) \left( \cos \frac{2i+1}{2n+2} \pi \right) / \left( \cos \frac{\pi}{2n+2} \right) \right] / 2$$

$$i = 0, \dots, n \quad (18.6 \text{ الف})$$

هستند نیز مشخص شده‌اند. می‌توان نشان داد که  $\|\Lambda_n^e\|_\infty$  به‌ازای جمیع مقادیر  $n$  در محدودهٔ  $0.5 \leq n$  کوچکترین مقدار ممکنه برای  $\|\Lambda_n\|_\infty$  است.

از شکل ۴.۶ و قضیهٔ ۳.۶ چنین بر می‌آید که به‌ازای  $n \leq 47$ ، خطا در بسجمله‌ای درونیاب  $f(x)$  در نقاط مسوط چیبیشف (۱۸.۶ الف)، هیچگاه از چهار برابر بهترین خطای ممکنه بزرگتر نیست و معمولاً کوچکتر از آن است. اگر مثلاً در تمام نقاط بازهٔ  $a \leq x \leq b$ ،  $p_n^*(x)$  بهترین تقریب یکنواخت در محدودهٔ  $10^{-5}$  از  $f(x)$  باشد، آنگاه بسجمله‌ای درونیاب، در بدترین حالت، در محدودهٔ  $10^{-5} \times 4$  از  $f(x)$  خواهد بود، که با از دست دادن دقتی به‌اندازهٔ کمتر از نصف یک رقم اعشاری معادل است. از دست دادن این مقدار دقت معمولاً با درونیایی از راه یک بسجمله‌ای از یک یا دو درجه بالاتر امکانپذیر است. بعکس، اگر  $\Lambda_n^u$  معرف تابع لبگ برای فواصل یکنواخت از نقاط درونیایی، به‌نحوی که هنگام درونیایی از جدول پیش می‌آید باشد، آنگاه خواهیم داشت:

$$\|\Lambda_n^u\|_\infty \geq e^{n/2} \quad (19.6)$$

که با افزایش  $n$  به‌سرعت افزایش می‌یابد. (برای به‌دست آوردن نتیجه‌ای از این نوع به کتاب ریولین [۳۵؛ صفحهٔ ۹۹] مراجعه کنید).

□ مثال ۴.۶: در مثال ۲.۶، برای  $f(x) = \tan(\pi/4x)$  در بازهٔ متعارف  $-1 \leq x \leq 1$ ، کران پایین  $\text{dist}_\infty(f, \pi_4) = 0.0004100$  است. اگر برای این  $f(x)$  در پنج نقطهٔ مسوط چیبیشف (۱۸.۶ الف) درونیایی انجام گیرد، بسجمله‌ای  $p(x)$  (که به‌دلیل تقارن آرمانی از درجهٔ ۳ است) به‌دست خواهد آمد که برای آن داریم  $\|f - p\|_\infty = 0.000576\dots$  که بزرگی آن تنها به‌اندازهٔ ۱۲ برابر کوچکترین مقدار خطای ممکن است. با افزودن تنها یک نقطهٔ درونیایی [که از لحاظ محاسباتی ارزانتراز ساختن  $p_4^*(x)$  است] یک بسجمله‌ای از درجهٔ ۵ تولید می‌شود که فاصلهٔ آن از  $f(x)$  برابر است با  $0.000068\dots$  که بهبود قابل توجهی است نسبت به  $\text{dist}_\infty(f, \pi_4) = 0.0004100$ .

### تمرین

۱-۱۰۶ با استفاده از (۷.۶) کران پایین  $\text{dist}_\infty(e^x, \pi_4)$  را در بازهٔ  $-1 \leq x \leq 1$  برآورد کنید. و آن را با فاصلهٔ تابع  $e^x$  از بسجمله‌ای  $p_4(x)$  از درجهٔ نایبتر از ۳، که با تابع  $e^x$  در چهار نقطهٔ مسوط چیبیشف هماهنگ است مقایسه کنید [به‌رابطهٔ (۱۸.۶ الف) به‌ازای  $n = 3$  نگاه کنید].

۲-۱.۶ تمرین ۱-۱.۶ را برای بازه  $0 \leq x \leq 1$  تکرار کنید. (داهنمایی: در عوض، تابع  $e^{(x+1)/2}$  را در بازه  $-1 \leq x \leq 1$  بررسی کنید).

۳-۱.۶ در تمرینهای ۱-۱.۶ و ۲-۱.۶، با استفاده از بسجمله‌ای درونیاب  $p_p(x)$  و قضیه ۱.۶، کران پایین دیگری برای  $\text{dist}_\infty(e^x, \pi_p)$  به دست آورید (توجه: برای تعیین بزرگترین کران پایین، باید فرینت  $e^x - p_p(x)$  را، مثلاً با روش نیوتن، محاسبه کرد).

۴-۱.۶  $p_p^*(x)$  را برای  $e^x$  در بازه متعارف  $-1 \leq x \leq 1$  محاسبه کنید. [داهنمایی: در این حالت از روش نیوتن برای حل دستگاه معادلات (۱۳.۶) استفاده کنید، به عنوان حدس اولیه برای  $p_p^*(x)$  بسجمله‌ای درونیاب  $p_p(x)$  را، که در تمرین ۱-۱.۶ ساخته شد بگیریید و به عنوان حدس اولیه برای نقاط تساوب  $x_0 < x_1 < \dots < x_p$  از فرینت موضعی  $x_0^{(0)} < x_1^{(0)} < \dots < x_p^{(0)} = 1 - x_0^{(0)}$  مربوط به خطای  $e^x - p_p(x)$  شروع کنید. توجه نمایید که بنا بر قضیه ۲.۶، داریم  $x_0 = 1$  و  $x_p = -1$ ].

۵-۱.۶ رابطه (۶.۶) را ثابت کنید [داهنمایی: تحقیق کنید که، چون  $x_j$ ها فرینتهای موضعی  $T_{n+1}(x)$  هستند، به ازای مقداری از  $x_j$  که با رابطه (۵.۶) معین می شود و به ازای یک ثابت مناسب  $c_n$ ، تساوی  $w(x) = c_n(1-x^2)T'_{n+1}(x)$  برقرار است. معادله دیفرانسیل  $xT'_k(x) = kT'_k(x) - k^2T_k(x)$  را با مشتقگیری از رابطه (۹.۶) نسبت به  $\theta$  به دست آورید و از آن برای حذف  $T'_{n+1}(x)$  از عبارت  $w'(x)$  حاصل استفاده کنید. و نیز از این معادله دیفرانسیل برای اثبات تساوی  $xT'_{n+1}(x) = (n+1)T_{n+1}(x)$  به ازای  $x_0 = x_{n+1} = x$ ، کمک بگیرید. سرانجام به برقراری تساویهای  $T'_{n+1}(x_j) = (-1)^j$  و  $T'_{n+1}(x_j) = 0$  به ازای  $j = 1, \dots, n$  احتیاج پیدا می کنید].

۶-۱.۶ ثابت کنید که بهترین تقریب یکنواخت  $p_n^*(x)$  برای تابع کوژا  $f(x)$  در بازه  $a \leq x \leq b$ ، تقریب  $p_n^*(x) = p_1(x) + (1/2) \min_{a \leq y \leq b} \{f(y) - p_1(y)\}$  است که در آن  $p_1(x)$  خط مستقیمی است که در نقاط  $a$  و  $b$  با  $f(x)$  منطبق است.

۷-۱.۶ گیریم  $p_n^*(x)$  بهترین تقریب یکنواست برای  $f(x)$  در بازه متعارف  $-1 \leq x \leq 1$  باشد. با استفاده از یکنوایی  $p_n^*(x)$ ، ثابت کنید که  $p_n^*(x)$  تابعی فرد (زوج) است، در صورتی که  $f(x)$  فرد (زوج) باشد یعنی، به ازای جمیع مقادیر  $x$  تساوی

$$f(-x) = -f(x) \quad (f(-x) = f(x))$$

برقرار است.

از آنجا نتیجه گیری کنید که کران پایین حاصل در مثال ۲.۶، برای

$$\text{dist}_\infty \left( \tan \frac{\pi}{4} x, \pi_p \right)$$

کران پایین برای  $\text{dist}_\infty(\tan(\pi/4)x, \pi_3)$  نیز هست.

۸-۱۰۶ فرض کنید که تابع  $\psi(x)$  بر بسجمله‌ایهای از درجهٔ نایبشتر از  $n$  در بازهٔ  $a \leq x \leq b$  متعامد است یعنی، به ازای جمیع مقادیر  $p \in \pi_n$ ، تساوی  $\int_a^b \psi(x)P(x)dx = 0$  برقرار است. در این صورت ثابت کنید که به ازای هر تابع خاص پیوستهٔ  $f(x)$  رابطهٔ

$$\left| \int_a^b \psi(x)f(x)dx \right| / \int_a^b |\psi(x)|dx \leq \text{dist}_\infty(f, \pi_n)$$

برقرار است.

۹-۱۰۶ از فرمول جمع کوسینوسها برای اثبات، رابطهٔ (۱۰-۶) استفاده کنید.

۱۰-۱۰۶ برای تابع  $f(x) = \sqrt{x}$  یک بسجمله‌ای تقریب مناسب از درجهٔ  $n$ ، به ازای  $n = 1, 2, 3, \dots, 10$  در بازهٔ  $0 \leq x \leq 1$  پیدا و تحقیق کنید که

$$\text{dist}_\infty(\sqrt{x}, \pi_n) \approx \text{const } n^{-1}$$

با توجه به نتیجهٔ فوق درجهٔ  $n$  لازم برای برقراری  $\text{dist}_\infty(\sqrt{x}, \pi_n) \leq 10^{-6}$  را برآورد کنید.

۱۱-۱۰۶ تمرین ۱۰-۱۰۶ را در بازهٔ  $-1 \leq x \leq 1$  تکرار و فرض کنید که

$$\text{dist}_\infty(\sqrt{|x|}, \pi_n) \approx \text{const } n^{-\alpha}$$

مقدار  $\alpha$  را تخمین بزنید.

۱۲-۱۰۶ محاسبات مثال ۴.۲ را تکرار کنید ولی این بار برای نقاط درونیایی، به جای استفاده از نقاط درونیایی با فواصل مساوی، از نقاط مبسوط چیبشف (۱۸.۶ الف) استفاده کنید. نتایج به دست آمده را با نتایج مثال ۴.۲ مقایسه و سعی کنید که آنها را بر حسب شکل‌های ۴.۶ و رابطهٔ (۱۹.۶) بیان کنید.

۱۳-۱۰۶ تمرین ۱۲-۱۰۶، را منتهی برای تابع  $f(x) = |x|$  تکرار کنید. (این مثال خوبی برای این امر است که در تقریب با بسجمله‌ایها، رفتار نامناسب تابع در یک ناحیه به تقریب بد در همه جا منجر می‌شود. به کار گرفتن بسجمله‌ای-تکه‌ای تقریب ساز روش خوبی برای اجتناب از این صورت ناخوش آیند تقریب با بسجمله‌ایهاست.)

۱۴-۱۰۶ ثابت کنید که کران پایینی که از رابطهٔ (۴.۶) به دست آمد می‌تواند به صورت  $|f[x_0, \dots, x_{n+1}] / g[x_0, \dots, x_{n+1}]|$  محاسبه شود با شرط  $x_0 < x_1 < \dots < x_{n+1}$  و در این رابطه  $g(x)$  هر تابعی است که برای آن به ازای جمیع مقادیر  $i$ ، تساوی

\* کلمهٔ انگلیسی  $\text{const}$  در رابطهٔ فوق، به مفهوم مقدار ثابت آمده است. م.

$g(x_i) = (-1)^i$  برقرار است. آنگاه الگوریتم (۳.۲) را چنان تغییر دهید که محاسبه  $g[x_0, \dots, x_{n+1}]$  همزمان با محاسبه  $f[x_0, \dots, x_{n+1}]$  انجام گیرد.

## ۲.۶ برازاندن داده‌ها

تاکنون دربارهٔ تقریب‌زدن يك تابع از راه درونیایی در بعضی نقاط بحث کرده‌ایم. در این شیوهٔ عمل قبلاً فرض می‌شود که مقادیر  $f(x)$  در این نقاط معلوم‌اند. از این رو درونیایی در شرایط معمول زیر کمتر به کار می‌رود (اگر نگوئیم که بکلی خطرناک است): تابع  $f(x)$  معرف رابطهٔ بین دو کمیت فیزیکی  $x$  و  $y = f(x)$  است و اعداد  $f_n$  از طریق اندازه‌گیری و یا آزمایش دیگر به دست آمده‌اند که تنها مقدار  $f(x)$  را در  $x_n$  تقریب می‌کنند، یعنی

$$f(x_n) = f_n + \varepsilon_n \quad n = 1, \dots, N$$

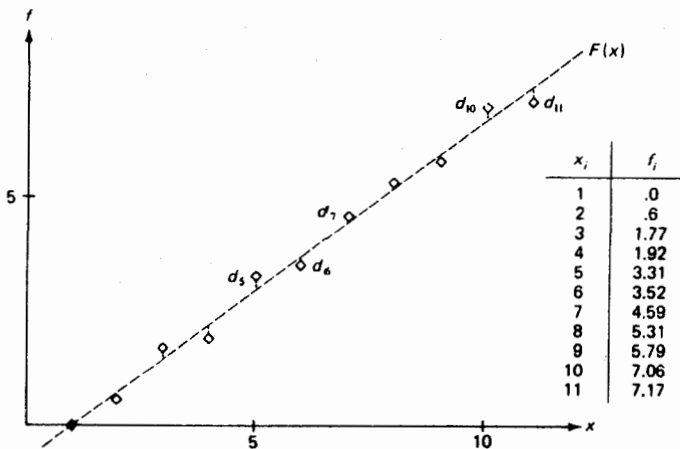
که در اینجا خطاهای تجربی  $\varepsilon_n$  معلوم نیستند. مسئلهٔ برازاندن داده‌ها عبارات است از بازایی تابع  $f(x)$  از روی داده‌های (تقریبی) مفروض  $f_n$  به ازای  $n = 1, \dots, N$ . دقیق بگوئیم، هرگز نمی‌دانیم که در اعداد  $f_n$  خطا وجود دارد. بلکه بر اساس اطلاعات موجود از تابع  $f(x)$  یا حتی تنها از روی احساس در مورد چگونگی تابع، اظهار نظر می‌کنیم که  $f(x)$  تابعی به پیچیدگی  $f_n$  نیست و یا با آن سرعتی که اعداد  $f_n$  نشان می‌دهند، تغییر نمی‌کند و بنابراین به این باور می‌رسیم که می‌باید در اعداد  $f_n$  خطا وجود داشته باشد.

برای مثال، داده‌هایی را که در شکل ۵.۶ مشخص شده‌اند در نظر می‌گیریم. در اینجا

$$x_n = x \quad n = 1, \dots, 11$$

اگر دلیلی برای این باور داشته باشیم که تابع  $f(x)$  يك خط راست است، مطمئناً در داده‌های مفروض خطا وجود دارد. اگر فقط بدانیم که  $f(x)$  يك تابع کوژ است، باز می‌توانیم نتیجه بگیریم که داده‌ها دارای خطا هستند. حتی اگر چیزی دربارهٔ  $f(x)$  ندانیم، باز ممکن است وسوسه شویم و از روی شکل ۵.۶ نتیجه بگیریم که  $f(x)$  يك خط راست است، ولو اینکه اکنون زمینهٔ محکمی برای آن نداشته باشیم. اما خواه در بازهٔ  $f(x)$  چیزی بدانیم یا ندانیم می‌توانیم از روی داده‌های رسم شده نتیجه‌گیری کنیم که بسیاری از اطلاعات مربوط به  $f(x)$  که در داده‌های  $f_n$  وجود دارند، می‌توانند به خوبی با يك خط راست نشان داده شوند.

خلاصه بگوئیم، برازاندن داده‌ها بر اساس این باور نهاده شده‌است که داده‌های مفروض  $f_n$  شامل مؤلفه‌ای است که به کندی تغییر می‌کند و به  $f(x)$  نزدیک می‌شود و یا اطلاعاتی دربارهٔ آن می‌دهد و يك مؤلفهٔ دیگری با دامنهٔ نسبتاً کوتاه‌است که در مقایسه با



شکل ۵.۶ تقریب خط راست با روش کوچکترین توانهای دوم.

مؤلفهٔ اول به سرعت تغییر می‌کند و معرف خطا یا نوفه<sup>۱</sup> در داده‌هاست. کاری که می‌باید انجام گیرد این است که به وسیلهٔ تابعی مانند  $F^*(x)$  داده‌هایی تقریب زده یا برازنده شود که  $F^*(x)$  شامل یا معرف بیشترین (اگر نه همه) اطلاعات موجود در داده‌ها در مورد  $f(x)$  و متضمن اندکی خطا یا نوفه (اگر اصلاً هست) باشد. برای برخی داده‌ها این کار عملاً با انتخاب تابعی مانند

$$F(x) = F(x; c_1, \dots, c_k) \quad (۲۰۰۶)$$

که به پارامترهای معین  $c_1, \dots, c_k$  بستگی دارد، انجام می‌گیرد. معمولاً سعی می‌شود که تابع  $F(x)$  چنان انتخاب شود که به طور خطی به پارامترها بستگی داشته باشد، به طوری که  $F(x)$  بتواند به صورت

$$F(x) = c_1 \phi_1(x) + c_2 \phi_2(x) + \dots + c_k \phi_k(x) \quad (۲۱۰۶)$$

درآید که در آن  $\{\phi_i\}$ ، مجموعه‌ای است از توابع که از قبل انتخاب شده است و  $\{c_i\}$  پارامترهایی هستند که می‌باید تعیین شوند. برای مثال،  $\{\phi_i\}$  ممکن است به صورت مجموعه‌ای از تکجمله‌ایهای  $x^{i-1}$  یا مجموعه‌ای از توابع مثلثاتی  $\{\sin \pi i x\}$  باشد. معمولاً در مقایسه با تعداد نقاط داده‌ها یعنی  $N$ ، عدد  $k$  کوچک است. انتظار ما این است که از یک طرف  $k$ ، به قدر کافی بزرگ باشد تا اطلاعات مربوط به  $f(x)$ ، موجود در داده‌ها، با انتخاب مناسبی برای پارامترهای  $c_1, \dots, c_k$ ، به خوبی نشان داده شود، در حالی که از طرف دیگر،

$k$  آنقدر کوچک باشد که اجازه تولید مجدد خطا یا نوفه را ندهد.

زمانی که افراد مجرب در فن برآزاندن داده‌ها تصمیم بگیرند که شکل صحیح (۲۰.۶) را برای تابع تقریب به کار برند، می‌باید برای پارامترهای  $c_i$  مقادیر خاصی مانند  $c_1^*, \dots, c_k^*$  چنان تعیین کنند که به تقریب «مطلوب»  $F(x; c_1^*, \dots, c_k^*) = F^*(x)$  دست یابند. فکر کلی این است که  $\{c_i\}$  طوری انتخاب شود که مقادیر انحراف

$$d_n = f_n - F(x_n; c_1, \dots, c_k) \quad n = 1, \dots, N$$

به طور همزمان تا حد امکان کوچک شوند. (برای چنین انحرافی در یک مثال نمونه، به شکل ۵.۶ نگاه کنید). بر حسب اصطلاحات مذکور در فصل ۴، سعی می‌شود که یک نرم از  $N$ -برداری  $\mathbf{d} = [d_1 \ d_2 \ \dots \ d_N]^T$  که به عنوان تابعی از  $c_1, \dots, c_k$ ، حداقل شود. انتخاب متداول برای نرم به شرح زیر است

(i) -۱ نرمی

$$\|\mathbf{d}\|_1 = \sum_{n=1}^N |d_n|$$

این نرم موقعی انتخاب می‌شود که خواسته باشیم انحراف میانگین تا حد امکان کوچک باشد، یا

(ii)  $-\infty$  نرمی

$$\|\mathbf{d}\|_\infty = \max_{1 \leq n \leq N} |d_n|$$

وقتی که بخواهیم کلیه انحرافات به طور یکنواخت کوچک باشند.

اما، اگر سعی کنیم که این مسئله حداقل سازی را در قالب فصل ۵ یا به روشهای دیگر حل کنیم، فوراً در خواهیم یافت که تعیین مینیمم‌های  $c_1^*, \dots, c_k^*$  به حل یک دستگاه معادلات غیر خطی منجر می‌شود [مثلاً، برای مسئله مربوط به تقریب یکنواخت در یک بازه، به دستگاه معادلات (۱۴.۶) نگاه کنید]. بنابراین رسم چنین است که آن نرمی را که باید مینیمم شود، ۲-نرمی انتخاب می‌کنند.

$$\|\mathbf{d}\|_2 = \left( \sum_{n=1}^N |d_n|^2 \right)^{1/2}$$

زیرا این نرم در پیدا کردن مینیمم مقادیر  $c_i^*$  به دستگاه معادلات خطی منجر می‌شود. در این حال تقریب نتیجه شده  $F(x; c_1^*, \dots, c_k^*)$  را تقریب با روش کوچکترین مربعات برای داده‌های مفروض نامند.

اکنون دستگاه معادلات برای استخراج  $c_i^*$ ها را به دست می‌آوریم. از آنجا که تابع

جذرها، یکنواست، مینیمم سازی  $\|d\|_2$  مثل مینیمم سازی  $\|d\|_2^2$  خواهد بود. البته برای اینکه  $\mathbf{c}^* = [c_1 \dots c_k]^T$  مینیمم تابع

$$E(\mathbf{c}) = E(c_1, \dots, c_k) = \|d\|_2^2 = \sum_{n=1}^N [f_n - F(x_n; \mathbf{c})]^2$$

باشد، لازم است که گرادیان  $E$  در  $\mathbf{c}^*$  صفر شود، یعنی

$$\nabla E(\mathbf{c}^*) = \mathbf{0}$$

(به بخش ۱۰۵ نگاه کنید)، بنابراین چون بنا بر (۲۱.۶) داریم

$$(\partial/\partial c_i)[f_n - F(x_n; \mathbf{c})] = -\phi_i(x_n)$$

$\mathbf{c}$  می باید در به اصطلاح معادلات نرمال

$$-\sum_{n=1}^N [f_n - F(x_n; \mathbf{c}^*)] \phi_i(x_n) = 0 \quad i = 1, \dots, k \quad (22.6)$$

صدق کند.

صفت «نرمال» بدین دلیل به این معادلات داده شده است که آنها تصریح می کنند که بردار خطای  $\mathbf{e} = [e_1 \ e_2 \ \dots \ e_N]^T$  با  $\mathbf{e} = [e_1 \ e_2 \ \dots \ e_N]^T$  به ازای کلیهٔ مقادیر  $n$ ، باید بر هر کدام از  $k$ -برداریهای

$$\boldsymbol{\phi}_i = [\phi_i(x_1) \ \phi_i(x_2) \ \dots \ \phi_i(x_N)]^T \quad i = 1, \dots, k$$

نرمال یا عمود باشد. زیرا برحسب این  $N$ -برداریهای معادلات (۲۲.۶) به صورت

$$-\sum_{n=1}^N e_n \phi_i(x_n) = 0 \quad i = 1, \dots, k$$

در می آیند. از آنجا که تابع تقریب کلی به شکل

$$F(x) = c_1 \phi_1(x) + c_2 \phi_2(x) + \dots + c_k \phi_k(x)$$

است، این رابطه می گوید که بردار خطا باید (با این تعبیر) بر همهٔ توابع تقریبی ممکن عمود باشد، یعنی

$$\mathbf{e}^T (c_1 \boldsymbol{\phi}_1 + c_2 \boldsymbol{\phi}_2 + \dots + c_k \boldsymbol{\phi}_k) = 0 \quad c_k, \dots, c_1 \text{ مقادیر } c_1, \dots, c_k$$

که این رابطه بردار  $c_1^* \boldsymbol{\phi}_1 + c_2^* \boldsymbol{\phi}_2 + \dots + c_k^* \boldsymbol{\phi}_k$  را به عنوان تصویر قائم<sup>۲</sup> بردار داده‌ها، روی ابرصفحه‌ای<sup>۳</sup> که بر اثر  $\boldsymbol{\phi}_1, \boldsymbol{\phi}_2, \dots, \boldsymbol{\phi}_k$  تنیده شده معرفی می کند.

1. square-root function

2. orthogonal projection

3. hyperplane

معادلات نرمال را مجدداً به شکل

$$\sum_{j=1}^k c_j^* \phi_j^T \phi_i = f^T \phi_i \quad i = 1, \dots, k \quad (۲۳.۶)$$

می‌نویسیم تا صریحاً بیان کنیم که این معادلات دستگاهی از  $k$  معادله خطی با  $k$  مجهول  $c_1^*, \dots, c_k^*$  را تشکیل می‌دهند. چنانچه خواهیم دید، این دستگاه همواره دست کم دارای یک جواب است  $[ \phi_i(x) ]$  هرچه باشد. بعلاوه هر جوابی از (۲۳.۶) مقدار  $E(c_1, \dots, c_k)$  را مینیمم می‌کند.

برای ارائه یک مثال، اکنون به یافتن تقریب با روش کوچکترین مربعات برای داده‌هایی که در شکل ۵.۶ با یک خط راست مشخص شده‌اند می‌پردازیم. در این مثال داریم  $n = 1, \dots, 11$ ،  $x_n = n$  و

$$F(x; c_1, c_2) = c_1 + c_2 x$$

لذا  $k = 2$  و  $\phi_1(x) = 1$  و  $\phi_2(x) = x$ . دستگاه خطی (۲۳.۶) به شکل

$$11c_1^* + 66c_2^* = 41704$$

$$66c_1^* + 506c_2^* = 328705$$

در می‌آید که وقتی با روش حذف گاوس حل شود، جواب منحصر زیر را به دست خواهد داد:

$$c_1^* = -0.7314\dots \quad c_2^* = 0.7437\dots$$

خط راست حاصله نیز در شکل ۵.۶ رسم شده است.

در اینجا باید به این واقعیت ناخوشایند اشاره شود که اگر ماتریس ضرایب (۲۳.۶) در اکثر موارد چنان بدشرط نمی‌شود که کاربرد مستقیم الگوریتم حذفی ۲.۴ را به نتایج غیر قابل اعتمادی منجر سازد، آنگاه همه چیز طبق دلخواه روبه‌راه می‌شود. این مطلب با مثال ساده زیر نشان داده شده است.

□ مثال ۵.۶: مقادیر تقریبی  $f_n \approx f(x_n)$  به ازای

$$x_n = 10 + \frac{n-1}{5} \quad n = 1, \dots, 6$$

داده شده‌اند و دلایلی داریم که معتقد باشیم این داده‌ها می‌توانند به نحو مناسبی با یک سهمی نشان داده شوند. بنابراین مقادیر زیر را انتخاب می‌کنیم

$$\phi_1(x) = 1 \quad \phi_2(x) = x \quad \phi_3(x) = x^2$$



در این حالت ماتریس ضرایب  $A$ ، حاصل از دستگاه (۲۳.۶) برابر است با

$$A = \begin{bmatrix} 6 & 63 & 66272 \\ 63 & 66272 & 69678 \\ 66272 & 69678 & 7339385664 \end{bmatrix}$$

و نتیجه می شود که  $\|A\|_{\infty} \approx 8 \times 10^4$  از طرف دیگر به ازای

$$Ax = \begin{bmatrix} -0.002 \\ -0.018 \\ -1.28 \end{bmatrix} \quad x = \begin{bmatrix} 10007 \\ -200 \\ 0.0099 \end{bmatrix} \quad \text{خواهیم داشت}$$

از این رو با توجه به نامساوی (۳۸.۴) یعنی

$$\|Ax\| \geq \frac{\|x\|}{\|A^{-1}\|}$$

خواهیم داشت

$$\|A^{-1}\|_{\infty} \geq 78 \quad \text{یا} \quad 1.28 = \|Ax\|_{\infty} \geq \frac{10007}{\|A^{-1}\|_{\infty}}$$

بنابراین عدد شرط  $A$  برابر است با

$$\text{cond}(A) = \|A\|_{\infty} \|A^{-1}\|_{\infty} \geq 10^5$$

درواقع همان گونه که نتایج خاص زیر نشان می دهند، عدد شرط  $A$  بسیار بزرگتر از  $10^5$  است.  $f(x)$  را به صورت زیر انتخاب می کنیم

$$f(x) = 10 - 2x + \frac{x^2}{10}$$

و مقادیر دقیق داده ها یعنی  $f_n = f(x_n)$ ،  $n = 1, \dots, 6$ ، را به کار می بریم. در این صورت چون  $f(x)$  یک بسجمله ای از درجه ۲ است،  $F^*(x)$  باید خود  $f(x)$  باشد، بنابراین باید داشته باشیم

$$c_0^* = 10 \quad c_1^* = -2 \quad c_2^* = 0.1$$

با به کار بردن الگوریتم حذفی ۲.۴، برای حل دستگاه (۲۳.۶) روی کامپیوتر CDC۶۵۰۰، نتایج زیر به دست می آید

$$c_1^* = 9999999997437 \dots \quad c_2^* = -1999999999511 \dots$$

$$c_3^* = 099999999976 \dots$$

لذا محاسبه در حساب ممیز شناور با ۱۴ رقم اعشاری برای این دستگاه  $3 \times 3$ ، نتیجه‌ای به دست می‌دهد که تنها در حدود ۸ رقم اعشاری آن صحیح است. اگر درایه  $(3, 3)$  مربوط به  $A$  را به صورت  $7339396$  گرد نموده و محاسبه را تکرار کنیم، نتایج به دست آمده حیرت‌انگیز است

$$c_1^* = 69035 \dots \quad c_2^* = -19243 \dots \quad c_3^* = 090639 \dots$$

به‌طور مشابه اگر تمام محاسبات به صورت ممیز شناور با هفت رقم اعشاری صورت پذیرد، نتایج برابرند با

$$\square \quad c_1^* = 89492 \dots \quad c_2^* = -19712 \dots \quad c_3^* = 090863$$

مثال فوق به وضوح نشان می‌دهد که بدون انجام کارهای مقدماتی، اقدام به حل معادلات نرمال می‌تواند خطرناک باشد. اقدام مذکور متضمن انتخاب دقیق  $\phi_i(x)$  است. راه ساده‌ای که برای اجتناب از مسئله بدشرطی به نظر می‌رسد، انتخاب  $\phi_i(x)$  است به گونه‌ای که در نقاط  $x_1, \dots, x_N$  متعامد باشند، یعنی

$$\phi_j^T \phi_i = \sum_{n=1}^N \phi_i(x_n) \phi_j(x_n) = 0 \quad \text{به‌ازای } j \neq i \quad (24.6)$$

اگر رابطه (۲۴.۶) برقرار باشد، معادلات (۲۳.۶) به صورت زیر تبدیل می‌شوند

$$c_i^T \phi_i^T \phi_i = f^T \phi_i \quad i = 1, \dots, k \quad (25.6)$$

که بدیهی است که حل آنها، دیگر مشکل زیادی ندارد.

البته این روش زیبا، جدا از شرط مسئله، تنها مسئله‌ای را جایگزین مسئله دیگر می‌کند و اکنون می‌باید توابع متعامد را پیدا کنیم. اگر بخواهیم که  $\phi_i$ ها نیز به صورت بسجمله‌ای باشند، با استفاده از یک رابطه بازگشتی سه‌جمله‌ای که برای دنباله‌های بسجمله‌ایهای متعامد معتبر است، می‌توان چنین توابع بسجمله‌ای متعامدی را به‌نحو کارا درست کرد. این مطلب را در بخشهای ۳.۶ و ۴.۶ مورد بحث قرار می‌دهیم. اگر، همان‌گونه که در عمل غالباً اتفاق می‌افتد، نتوانیم  $f(x)$  را به صورت یک بسجمله‌ای فرض کنیم، بایسد از راههای دیگری برای ساختن توابع متعامد مناسب استفاده کنیم. یک نمونه از چنین تکنیکی الگوریتم اصلاح شده گرام-اشمیت<sup>۱</sup> است که در برخی از متون درسی آمده است (برای مثال به کتاب رایس [۱۷] مراجعه کنید). به صورت دیگر، باید به انتخاب  $\phi_1(x), \dots, \phi_k(x)$  که «تقریباً»

## 1. Gram-Schmidt

متعامد باشند قناعت کرد. معنی این عبارت مهم توضیح این حقیقت است که ماتریس ضرایب دستگاه (۲۳.۶) برای این گونه  $\phi_i(x)$ ها تقریباً «قطری» مثلاً نافذ قطری است. اگر نقاط  $x_1, \dots, x_n$  در يك بازه  $(a, b)$  تقریباً به طور یکنواخت توزیع شده باشند، آنگاه  $\phi_1(x), \dots, \phi_n(x)$  تقریباً به سمت عمود بودن میل خواهند کرد، اگر هر  $\phi_i(x)$  در بازه  $(a, b)$  یکبار بیشتر از  $\phi_{i-1}(x)$  تغییر علامت دهد (به تمرین ۲.۶-۳ رجوع کنید).

### تمرین

۲.۶-۱. تقریب باروش کوچکترین مربعات را برای داده‌های مرسوم در شکل ۵.۶، به وسیلهٔ توابعی به صورت

$$F(x) = c_1 + c_2 x + c_3 \sin [123(x-1)]$$

از حل معادلات نرمال مناسب، محاسبه کنید. آیا حس می‌کنید که این تقریب معرف تمام اطلاعات مربوط به  $f(x)$  موجود داده‌ها هست؟ چرا؟

۲.۶-۲. با دنبال کردن برهان مذکور در متن، معادله‌های نرمال برای بهترین  $c_1^*$ ،  $c_2^*$  را در حالت

$$F(x) = F(x; c_1, c_2) = c_1 e^{c_2 x}$$

به دست آورید. آیا این معادله‌های نرمال هنوز خطی هستند؟

۲.۶-۳. با استفاده از توابع زیر، تمام محاسبات مثال ۵.۶ را تکرار کنید

$$\phi_1(x) = 1 \quad \phi_2(x) = x - 10.5 \quad \phi_3(x) = (x - 10.3)(x - 10.7)$$

با توجه به آخرین بند این بخش، این معادلات نرمال می‌باید حالا خیلی به شرط تر باشند، آیا همینطور است؟

### ۳.۶\* بسجمله‌ایهای متعامد

در این بخش، بعضی از ویژگیهای مربوط به دنباله‌های بسجمله‌ایهای متعامد و مثالهای خاص آنها را به اختصار مورد بحث قرار می‌دهیم. اگرچه انگیزهٔ مستقیم این بحث از مسئلهٔ تقریب به روش کوچکترین توانهای دوم با بسجمله‌ایها (که در بخش بعد باید بحث شود) ناشی شده است، ولی ما از آن برای بسجمله‌ایهای متعامد بعداً در زمینه‌های مختلف، مثلاً در بخش ۳.۷ استفاده خواهیم کرد. برای فراهم کردن زمینه برای آن بخش، در اینجا ما از يك مفهوم تعامد توابع استفاده می‌کنیم که کلیتر از آن چیزی است که در بخش ۲.۶ آمده است.

در آنچه که ذیلا می آید، گیریم  $(a, b)$  یک بازه مفروض و  $w(x)$  یک تابع مفروضی تعریف شده (حداقل) در بازه  $(a, b)$  و مثبت در آن بازه باشد. بعلاوه حاصلضرب داخلی هر دو تابع  $g(x)$  و  $h(x)$  [معین در بازه  $(a, b)$ ] یعنی

$$\langle g, h \rangle$$

را به یکی از دو طریق زیر تعریف می کنیم:

$$\langle g, h \rangle = \int_a^b g(x)h(x)w(x) dx \quad (۲۶.۶)$$

یا

$$\langle g, h \rangle = \sum_{n=1}^N g(x_n)h(x_n)w(x_n) \quad (۲۷.۶)$$

در حالت اول فرض می کنیم که به ازای جمیع توابع  $g(x)$  و  $h(x)$  مورد نظر، این انتگرال (حداقل یک انتگرال ناسره<sup>۱)</sup> وجود دارد، در حالت دوم فرض می کنیم که  $N$  نقطه  $x_1, \dots, x_N$ ، همگی در بازه  $(a, b)$  مفروض باشند و در طی بحث، ثابت دو نظر گرفته شوند. توجه کنید که به ازای  $w(x) \equiv 1$ ، رابطه (۲۷.۶) به حاصلضرب داخلی دو تابع  $g^T h = h^T g$  تبدیل می شود که در بحث تقریب با روش کوچکترین مربعات در بخش ۲.۶ ظاهر شد.

با معین بودن حاصلضرب داخلی دو تابع، دو تابع  $g(x)$  و  $h(x)$  را متعامد (بریکدیگر) گویند هر گاه

$$\langle g, h \rangle = 0$$

برای مثال به آسانی تحقیق می شود که توابع  $g(x) \equiv 1$ ،  $h(x) = x$  متعامدند اگر حاصلضرب داخلی آنها برابر عبارت زیر باشد

$$\langle g, h \rangle = \int_{-1}^1 g(x)h(x) dx$$

و نیز متعامدند اگر حاصلضرب داخلی آنها برابر

$$\langle g, h \rangle = \sum_{n=-10}^{10} g(n)h(n)$$

یا برابر با

$$\langle g, h \rangle = \int_{-1}^1 \frac{g(x)h(x)}{(1-x^2)^{1/2}} dx$$

باشد. توابع  $g(x) = \sin nx$  و  $h(x) = \sin mx$  به ازای اعداد صحیح  $m$  و  $n$ ،  $n \neq m$  متعامد هستند اگر تساوی

$$\langle g, h \rangle = \int_0^{2\pi} g(x)h(x) dx$$

برقرار باشد، چنانکه این مطلب برای  $g(x) = \sin nx$ ،  $h(x) = \cos mx$  صادق است. به علاوه  $P_0(x)$ ،  $P_1(x)$ ،  $P_2(x)$ ، ... را یک دنباله متناهی یا نامتناهی از بسجمله‌ایهای متعامد نامند به شرط آنکه همهٔ  $P_i(x)$ ها بر یکدیگر عمود باشند و هر  $P_i(x)$  یک بسجمله‌ای دقیقاً از درجهٔ  $i$  باشد. به عبارت دیگر

(i) به ازای هر  $i$ ،  $i$  یک بسجمله‌ای از درجهٔ کمتر از  $i$   $P_i(x) = (\alpha_i x^i + i$  با  $\alpha_i \neq 0$  باشد.

(ii) هر گاه  $j \neq i$ ، آنگاه  $\langle P_i, P_j \rangle = 0$ .

برای مثال توابع

$$P_0(x) \equiv 1 \quad P_1(x) = x \quad P_2(x) = 3x^2 - 1$$

یک دنباله از سه بسجمله‌ای متعامد را تشکیل می‌دهند اگر

$$\langle g, h \rangle = \int_{-1}^1 g(x)h(x) dx$$

قبلاً اشاره کردیم که  $\langle P_0, P_1 \rangle = 0$ ، و نیز داریم

$$\langle P_0, P_2 \rangle = \int_{-1}^1 1(3x^2 - 1) dx = x^3 - x \Big|_{-1}^1 = 0$$

و

$$\langle P_1, P_2 \rangle = \int_{-1}^1 x(3x^2 - 1) dx = \frac{3}{4}x^4 - \frac{1}{2}x^2 \Big|_{-1}^1 = 0$$

گیریم  $P_0(x)$ ،  $P_1(x)$ ، ...،  $P_k(x)$  یک دنبالهٔ متناهی از بسجمله‌ایهای متعامد باشند. آنگاه ویژگیهای زیر می‌توانند ثابت شوند:

و ویژگی ۱ اگر  $P(x)$  یک بسجمله‌ای غیر مشخص از درجهٔ نایبتر از  $k$  باشد، آنگاه به ازای ضرایب  $d_0, \dots, d_k$  که به وسیلهٔ  $p(x)$  به‌طور منحصر به‌فرد تعیین شده‌اند، می‌توان  $p(x)$  را به صورت

$$p(x) = d_0 P_0(x) + d_1 P_1(x) + \dots + d_k P_k(x) \quad (28.6)$$

نوشت. به ویژه اگر

$$p(x) = (a_k x^k + k \text{ از درجه کمتر از } k)$$

و اگر ضریب بزرگترین درجه  $P_k(x)$  برابر با  $\alpha_k$  باشد، آنگاه

$$d_k = \frac{a_k}{\alpha_k}$$

این ویژگی از قسمت (i) مذکور در بالا، و با استقرار روی  $k$  نتیجه می‌شود. برای مثال بالا، می‌توان بسجمله‌ای کلی از درجه نایبتر از ۲

$$p_2(x) = a_0 + a_1 x + a_2 x^2$$

را به صورت زیر نوشت

$$p_2(x) = \left(a_0 + \frac{a_2}{3}\right) P_0(x) + a_1 P_1(x) + \frac{a_2}{3} P_2(x)$$

با ترکیب ویژگی ۱ با (ii)، ویژگی ۲ به دست می‌آید. ویژگی ۲ اگر  $p(x)$  یک بسجمله‌ای از درجه کمتر از  $k$  باشد، آنگاه  $p(x)$  و  $P_k(x)$  متعامد خواهند بود، یعنی

$$\langle p, P_k \rangle = 0$$

اگر در مثال بالا، بگیریم  $p(x) = 1 + x$  خواهیم داشت

$$\langle p, P_2 \rangle = \int_{-1}^1 (1+x)(3x^2-1) dx = \frac{3}{4} x^4 + x^3 - \frac{1}{2} x^2 - x \Big|_{-1}^1 = 0$$

این ویژگی به اصطلاح «بی‌گزند» دارای چندین نتیجه مهم است. ویژگی ۳ اگر حاصلضرب داخلی با رابطه (۲۶.۶) معین شود، آنگاه  $P_k(x)$  دارای کاربش حقیقی ساده است که تمامی آنها در بازه  $(a, b)$  قرار دارند، یعنی  $P_k(x)$  به ازای  $k$  نقطه متمایز  $\xi_{1,k}, \xi_{2,k}, \dots, \xi_{k,k}$  در بازه  $(a, b)$ ، به شکل زیر است

$$P_k(x) = \alpha_k (x - \xi_{1,k})(x - \xi_{2,k}) \dots (x - \xi_{k,k}) \quad (29.6)$$

در مثال خودمان داریم

$$P_0(x) \equiv \alpha_0 \equiv 1 \quad P_1(x) = \alpha_1 (x - \xi_{1,1}) = 1 \times (x - 0)$$

$$P_2(x) = \alpha_2 (x - \xi_{1,2})(x - \xi_{2,2}) = 3 \left(x + \frac{1}{\sqrt{3}}\right) \left(x - \frac{1}{\sqrt{3}}\right)$$

یک دلیل ساده و یژگی ۳ به قرار زیر است: گیریم  $k > 0$  و  $\xi_{1,k}, \dots, \xi_{r,k}$  همه نقاتی از بازه  $(a, b)$  باشند که در آنها  $P_k(x)$  تغییر علامت می‌دهد. آنگاه می‌گوییم

$$r \geq k$$

زیرا اگر  $r$  کمتر از  $k$  باشد، آنگاه به ازای  $\bar{x} \in (\max_i \xi_{i,k}, b)$  بسجمله‌ای

$$p(x) = P_k(\bar{x})(x - \xi_{1,k})(x - \xi_{2,k}) \dots (x - \xi_{r,k})$$

یک بسجمله‌ای از درجه کمتر از  $k$  خواهد بود که در هر نقطه از بازه  $(a, b)$ ، علامتی مشابه با  $P_k(x)$  دارد. از این‌رو، از یک طرف بنا بر یژگی ۳ داریم

$$\int_a^b p(x) P_k(x) w(x) dx = \langle p, P_k \rangle = 0$$

در حالی که از طرف دیگر، به ازای جمیع مقادیر  $x \in (a, b)$  جز  $\xi_{1,k}, \dots, \xi_{r,k}$  خواهیم داشت  $\langle p(x) P_k(x) w(x) \rangle > 0$  مطمئناً این دو نتیجه متناقض یکدیگرند. در نتیجه باید داشته باشیم  $r \geq k$ ، یعنی  $P_k(x)$  می‌باید حداقل  $k$  مرتبه در بازه  $(a, b)$  تغییر علامت دهد. اما چون  $P_k(x)$  یک بسجمله‌ای از درجه  $k$  و هر  $\xi_{i,k}$  یک ریشه از  $P_k(x)$  است، لذا  $r$  نمی‌تواند از  $k$  بزرگتر باشد (بخش ۱۰۲ را ببینید)، بنا بر این  $r$  می‌باید با  $k$  برابر باشد یعنی  $k$  نقطه متمایز  $\xi_{i,k}$ ،  $i = 1, \dots, k$ ، دقیقاً ریشه‌های  $P_k(x)$  هستند.

به‌طور مشابه ثابت می‌شود که وقتی حاصلضرب داخلی با رابطه (۲۷.۶) معین شده باشد رابطه (۲۹.۶) برقرار خواهد بود به شرط آنکه بین  $x_n$ ها حداقل  $k$  نقطه متمایز وجود داشته باشند.

دیژگی ۴ بسجمله‌ایهای متعامد در یک رابطه بازگشتی سه جمله‌ای صدق می‌کنند. اگر قرار دهیم

$$A_i = \frac{\alpha_i + 1}{\alpha_i} \quad \text{به‌ازای جمیع مقادیر } i$$

$$P_{-1}(x) \equiv 0$$

و اگر

$$S_i = \langle P_i, P_i \rangle$$

به‌ازای  $i = 0, \dots, k-1$ ، مخالف صفر باشد، آنگاه این رابطه بازگشتی می‌تواند به‌صورت

$$P_{i+1}(x) = A_i(x - B_i)P_i(x) - C_i P_{i-1}(x) \quad i = 0, 1, \dots, k-1 \quad (30.6)$$

نوشته شود، که در آن مقادیر  $B_i$  و  $C_i$  عبارتند از

$$B_i = \frac{\langle xP_i(x), P_i(x) \rangle}{S_i} \quad i = 1, \dots, k-1$$

و

$$C_i = \begin{cases} \text{اختیاری است} & i = 0 \\ \frac{A_i S_i}{A_{i-1} S_{i-1}} & i > 0 \end{cases}$$

این ویژگی می‌تواند برای تولید دنباله‌های کلی بسجمله‌ایهای متعامد (به شرط آنکه اعداد  $S_i$  و  $B_i$  بتوانند محاسبه شوند و  $S_i$ ها صفر نباشند) به کار رود. در این روال معمولاً ضرایب پیشرو  $\alpha_i$ ، یا هم‌ارز با آنها،  $A_i$ ، چنان انتخاب می‌شوند که به تعبیری دنباله حاصل بالاخص ساده باشد.

□ مثال ۶.۶: بسجمله‌ای لژاندر. اگر حاصلضرب داخلی با

$$\langle g, h \rangle = \int_{-1}^1 g(x)h(x) dx$$

داده شده باشد. بسجمله‌ایهای متعامد حاصل، بسجمله‌ایهای لژاندر نامیده می‌شوند با شروع از

$$P_0(x) \equiv 1$$

به دست می‌آید

$$S_0 = \int_{-1}^1 1 dx = 2 \quad S_0 B_0 = \int_{-1}^1 x \cdot 1 dx = 0$$

از این رو، با توجه به ویژگی ۴، با انتخاب  $A_1 = 1$ ، به ازای جمیع مقادیر  $i$ ، خواهیم داشت

$$P_1(x) = x$$

بعلاوه

$$S_1 = \int_{-1}^1 x^2 dx = \frac{2}{3} \quad S_1 B_1 = \int_{-1}^1 x \cdot x^2 dx = 0 \quad C_1 = \frac{S_1}{S_0} = \frac{1}{3}$$

لذا باز بنا بر ویژگی ۴



$$P_2(x) = x^2 - \frac{1}{3}$$

مجدداً

$$S_2 = \int_{-1}^1 \left(x^2 - \frac{1}{3}\right)^2 dx = \frac{8}{45} \quad S_2 B_2 = \int_{-1}^1 x \left(x^2 - \frac{1}{3}\right)^2 dx = 0$$

$$C_2 = \frac{4}{15}$$

بنابراین

$$P_3(x) = x^3 - \frac{3}{5}x$$

معمولاً چنین است که بسجمله‌ایهای لژاندر را به صورت نرمال درمی آورند تا داشته باشیم

$$P_k(1) = 1 \quad k \text{ به ازای جميع مقادير}$$

با این نرمال‌سازی، ضرایب در رابطهٔ بازگشتی به صورت زیر درمی آیند

$$A_k = \frac{2k+1}{k+1} \quad B_k = 0 \quad C_k = \frac{k}{k+1} \quad k = 0, 1, 2, \dots$$

لذا

$$P_{k+1}(x) = \frac{(2k+1)xP_k(x) - kP_{k-1}(x)}{k+1}$$

□

جدول ۲.۶، چندتا از نخستین بسجمله‌ایهای لژاندر را به ما می‌دهد.

جدول ۲.۶

$k$	$P_k(x)$
۰	۱
۱	$x$
۲	$(3/2)(x^2 - 1/3)$
۳	$(5/2)[x^3 - (3/5)x]$
۴	$(35/8)[x^4 - (6/7)x^2 + 3/35]$

□ مثال ۷.۶: بسجمله‌ایهای چیبیشف. اگر حاصلضرب داخلی با

$$\langle g, h \rangle = \int_{-1}^1 \frac{g(x)h(x)}{(1-x^2)^{1/2}} dx$$

داده شده باشد، بسجمله‌ایهای چیبیشف  $T_k(x)$  که در بخش ۱.۶، معرفی شدند به دست خواهند آمد. ما قبلاً رابطهٔ بازگشتی آنها را

$$T_{k+1}(x) = 2xT_k(x) - T_{k-1}(x) \quad k = 1, 2, \dots$$

از روی رابطهٔ معرف آنها

$$T_k(\cos \theta) = \cos k\theta$$

□ به دست آوردیم.

□ مثال ۸.۶: بسجمله‌ایهای هرمیت. زمانی که حاصلضرب داخلی

$$\langle g, h \rangle = \int_{-\infty}^{\infty} g(x)h(x)e^{-x^2} dx$$

به کار برده شده باشد،  $H_k(x)$  به دست می‌آید. با اجرای نرمال‌سازی معمولی این بسجمله‌ایها در رابطهٔ بازگشتی زیر صدق می‌کنند

$$H_{k+1}(x) = 2xH_k(x) - 2kH_{k-1}(x) \quad k = 0, 1, 2, \dots$$

چندتا از نخستین بسجمله‌ایهای هرمیت در جدول ۳.۶ داده شده‌اند.

جدول ۳.۶

$k$	$H(x)$
۰	۱
۱	$2x$
۲	$4x^2 - 2$
۳	$8x^3 - 12x$
۴	$16x^4 - 48x^2 + 12$

□

□ مثال ۹.۶: بسجمله‌ایهای تعمیم‌یافته لاگرا.  $L_k^\alpha(x)$  ها در ارتباط با حاصلضرب داخلی زیر

$$\langle g, h \rangle = \int_0^\infty g(x)h(x)x^\alpha e^{-x} dx$$

می‌باشند. ضرایب برای رابطه بازگشتی چنین‌اند

$$A_k = -\frac{1}{k+1} \quad B_k = 2k + \alpha + 1 \quad C_k = \frac{k + \alpha}{k + 1}$$

تولید پنج نخستین بسجمله‌ای لاگرا (به ازای  $\alpha = 0$ ) به‌عهده دانشجویان گذاشته شده است (تمرین ۳.۶-۱ را ببینید). □

دو مثال آخر در انتگرالگیری عددی بر بازه‌های نیمه‌نامتناهی یا نامتناهی اهمیت خاصی دارند (بند ۳.۷ را ببینید).

با بحث از الگوریتمی برای ارزیابی بسجمله‌ای مفروضی برحسب بسجمله‌ایهای متعامد، این قسمت را به پایان می‌رسانیم. فرض کنید که  $P_0(x), \dots, P_k(x)$  يك دنباله متناهی از بسجمله‌ایهای متعامد باشند، و فرض کنید که يك بسجمله‌ای  $P(x)$  از درجه نا کمتر از  $k$  برحسب  $P_i(x)$  به‌ما داده شده است، یعنی ضرایبی مانند  $d_0, \dots, d_k$  داریم به‌طوری که تساوی

$$p(x) = d_0 P_0(x) + d_1 P_1(x) + \dots + d_k P_k(x) \quad (31.6)$$

برقرار است. در محاسبه  $p(x)$  در يك نقطه خاص  $\bar{x}$ ، می‌توانیم از رابطه بازگشتی سه‌جمله‌ای (۳۰.۶) برای  $P_i(x)$  به‌صورت زیر استفاده کنیم: به‌موجب رابطه (۳۰.۶) داریم

$$P_k(\bar{x}) = A_{k-1}(\bar{x} - B_{k-1})P_{k-1}(\bar{x}) - C_{k-1}P_{k-2}(\bar{x})$$

و بنا بر این

$$p(\bar{x}) = d_0 P_0(\bar{x}) + \dots + d_{k-2} P_{k-2}(\bar{x}) + (d_{k-1} - d_k C_{k-1}) P_{k-1}(\bar{x}) \\ + [d_{k-1} + d_k A_{k-1}(\bar{x} - B_{k-1})] P_{k-1}(\bar{x})$$

یا با استفاده از قرارداد اختصار

$$\bar{d}_k = d_k \quad \bar{d}_{k-1} = d_{k-1} + \bar{d}_k A_{k-1}(\bar{x} - B_{k-1})$$

داریم

$$p(\bar{x}) = d_0 P_0(\bar{x}) + \dots + d_{k-2} P_{k-2}(\bar{x}) + (d_{k-2} - \bar{d}_k C_{k-1}) P_{k-2}(\bar{x}) + \bar{d}_{k-1} P_{k-1}(\bar{x}) \quad (۳۲.۶)$$

باز به موجب رابطه (۳۰.۶)

$$P_{k-1}(\bar{x}) = A_{k-2}(\bar{x} - B_{k-2}) P_{k-2}(\bar{x}) - C_{k-2} P_{k-2}(\bar{x})$$

و از قراردادن این مقدار در رابطه (۳۲.۶)، خواهیم داشت

$$p(\bar{x}) = d_0 P_0(\bar{x}) + \dots + (d_{k-2} - \bar{d}_{k-1} C_{k-2}) P_{k-2}(\bar{x}) + \bar{d}_{k-2} P_{k-2}(\bar{x})$$

که در اینجا از قرارداد اختصاری زیر استفاده کرده ایم

$$\bar{d}_{k-2} = d_{k-2} + \bar{d}_{k-1} A_{k-2}(\bar{x} - B_{k-2}) - \bar{d}_k C_{k-1}$$

اگر به همین صورت ادامه دهیم، متوالیاً روابط

$$\bar{d}_j = d_j + \bar{d}_{j+1} A_j(\bar{x} - B_j) - \bar{d}_{j+2} C_{j+1} \quad j = k-2, \dots, 0$$

را محاسبه می‌کنیم که سرانجام به تساوی  $p(\bar{x}) = \bar{d}_0 P_0(\bar{x}) = \bar{d}_0 \alpha_0$  خواهیم رسید.

**الگوریتم ۱.۶ ضرب تودرتو برای بسجمله‌ایهای متعامد ضرایب**

$$A_j, B_j, C_j, j = 0, \dots, k-1$$

برای رابطه بازگشتی سه‌جمله‌ای (۳۰.۶)، که بسجمله‌ایهای متعامد  $P_k(x), \dots, P_0(x)$  در آن صدق می‌کنند، و نیز مقدار ثابت  $\alpha_0 = P_0(x)$  و ضرایب  $d_0, \dots, d_k$  مربوط به  $p(x)$  در رابطه (۳۱.۶) و نقطه  $\bar{x}$  داده شده‌اند.

$$\bar{d}_k := d_k$$

اگر  $k=0$ ، آنگاه از الگوریتم خارج شوید

$$\bar{d}_{k-1} := d_{k-1} + \bar{d}_k A_{k-1}(\bar{x} - B_{k-1})$$

اگر  $k=1$ ، آنگاه از الگوریتم خارج شوید

For  $j = k-2, k-3, \dots, 0$ , do:

$$\bar{d}_j := d_j + \bar{d}_{j+1} A_j(\bar{x} - B_j) - \bar{d}_{j+2} C_{j+1}$$

هنگام خروج از الگوریتم،  $p(\bar{x})$  به صورت زیر تعیین می‌شود

$$p(\bar{x}) = \bar{d}_0 P_0(x)$$

به کار بردن زبان فورترن برای این الگوریتم با این اشکال کوچک همراه است که

برخی از گویشهای زبان فورترن اجازه استفاده از زیرنمایشه صفر را به ما نمی دهد. و نیز مقدار حافظه مورد نیاز و تعداد محاسبات لازم از يك مجموعه از بسجمله آریهای متعامد تا مجموعه دیگر تغییر می کند.

□ مثال ۱۰۰۶: در حالتی که بسجمله ایهای متعامد، بسجمله ایهای چبیش باشند، برنامه ای به زبان فورترن برای الگوریتم ۱۰۰۶، بنویسید.

در این حالت نیازی به نگهداری  $A_i$ ،  $B_i$ ،  $C_i$  در آرایه ها نیست، زیرا به  $i$  بستگی ندارند، و نیز محاسبه  $\bar{d}_j$  تنها به محاسبه  $\bar{d}_{j+1}$  و  $\bar{d}_{j+2}$  نیاز دارد، از این رو به نگهداری تمام آرایه  $\bar{d}_i$ ،  $k, \dots, 0, i$ ، نیازی نیست.

تابع چبیشف به زبان فورترن که در زیر آمده است، مسئله داده شده را حل می کند. NTERMS عبارت از تعداد جملات در  $p(x)$  است، یعنی  $p(x)$  از درجه نایبتر از NTERMS-۱ است. قرار بر این است که هم NTERMS و هم ضرایب  $D(i) = d_{i-1}$ ،  $i = 1, \dots, \text{NTERMS}$ ، با COMMON POLY نامگذاری شوند.

```

REAL FUNCTION CHEB (X)
C RETURNS THE VALUE OF THE POLYNOMIAL OF DEGREE .LT. NTERMS WHOSE
C CHEBYSHEV COEFFICIENTS ARE CONTAINED IN D .
      INTEGER NTERMS, K
      REAL D,X, PREV,PREV2,TWOX
      COMMON /POLY/ NTERMS,D(30)
      IF ( NTERMS .EC. 1) THEN
        CHEB = D(1)
                                RETURN
      END IF
      TWOX = 2.*X
      PREV2 = 0.
      PREV = D(NTERMS)
      IF (NTERMS .GT. 2) THEN
        DO 10 K=NTERMS-1,2,-1
          CHEB = D(K) + TWOX*PREV - PREV2
          PREV2 = PREV
          PREV = CHEB
10      CONTINUE
      END IF
      CHEB = D(1) + X*PREV - PREV2
                                RETURN
END

```

□

## تمرین

۱-۳۰۶ با استفاده از رابطه بازگشتی مناسب، پنج نخستین بسجمله ای لاگر را (به ازای  $\alpha = 0$ ) بسازید

۲-۳۰۶ ریشه های بسجمله ایهای  $P_4(x)$ ،  $P_3(x)$ ،  $P_2(x)$ ،  $P_1(x)$  را بیابید.

۳-۳۰۶ ریشه های بسجمله ایهای هر میت  $H_4(x)$  و  $H_3(x)$  و  $H_2(x)$  را به دست آورید.

۴-۳۰۶ بسجمله ای  $p(x) = x^4 + 2x^3 + x^2 + 2x + 1$  را به صورت حاصلجمع از بسجمله ایهای  $P_4(x)$  بیان کنید.

۵-۳۰۶ مستقیماً تحقیق کنید که بسجمله ای  $P_4(x)$  از  $P_3(x)$  بر هر بسجمله ای از درجه ۲، متعامد است.

۳-۳۰۶ ثابت کنید که اگر  $P_k(x)$  بسجمله‌ای لژاندری از درجه  $k$  باشد، آنگاه

$$\int_{-1}^1 [P_k(x)]^2 dx = \frac{2}{2k+1}$$

از رابطه بازگشتی سه‌جمله‌ای که بسجمله‌ایهای لژاندر در آن صدق می‌کنند استفاده کنید.

۷-۳۰۶ گوییم  $P_0(x), P_1(x), \dots$  يك دنباله از بسجمله‌ایهای متعامد باشند و  $x_0, \dots, x_k, \dots$  ریشه‌ی متمایز از  $P_{k+1}(x)$  ثابت کنید که به‌ازای این نقاط، بسجمله‌ایهای لژاندر  $k, \dots, 0, i$ ،  $l_i(x) = \prod_{j \neq i} (x - x_j) / (x_i - x_j)$  بر یکدیگر عمودند. [داهنمایی: نشان دهید که به‌ازای  $i \neq j$ ، رابطه  $l_i(x)l_j(x) = P_{k+1}(x)g(x)$  که در آن  $g(x)$  يك بسجمله‌ای است از درجه‌ی نایبتر از  $k$ ، برقرار است.]

#### ۴.۶ تقریب با روش کوچکترین مربعات به‌وسیله بسجمله‌ایها

در این بخش استفاده از دنباله‌های بسجمله‌ایهای متعامد برای محاسبه تقریبهای بسجمله‌ایهای (وزین) از راه کوچکترین مربعات مورد بحث قرار می‌گیرد.

گوییم  $f(x)$  تابعی تعریف شده بر بازه  $(a, b)$  باشد، و فرض کنید که بخواهیم  $f(x)$  را بر بازه  $(a, b)$  با يك بسجمله‌ای از درجه‌ی نایبتر از  $k$  تقریب کنیم. اگر اختلاف بین  $f(x)$  و  $p(x)$  به‌وسیله

$$\langle f(x) - p(x), f(x) - p(x) \rangle = \begin{cases} \int_a^b [f(x) - p(x)]^2 w(x) dx \\ \text{یا} \\ \sum_{n=1}^N [f(x_n) - p(x_n)]^2 w(x_n) \end{cases}$$

(۳۳.۶)

که در آن حاصلضرب داخلی توسط رابطه (۲۶.۶)، (۲۷.۶) داده شده، اندازه‌گیری شود، آنگاه طبیعی است که باید در جستجوی يك بسجمله‌ای از درجه‌ی نایبتر از  $k$  باشیم که به‌ازای آن (۳۳.۶) تا حد ممکن کوچک باشد. این بسجمله‌ای را تقریب  $f(x)$  به‌وسیله بسجمله‌ایهایی از درجه‌ی نایبتر از  $k$ ، با روش کوچکترین مربعات (وزین) نامند.

مسئله یافتن يك چنین بسجمله‌ای را در بخش ۲.۶ برای حالت خاصی که حاصلضرب داخلی با رابطه (۲۷.۶)، و تابع وزن با  $w(x) \equiv 1$  داده شده بود، حل کردیم. در حالت کلی، مراحل انجام کار به شرح زیر است: فرض کنید که برای حاصلضرب داخلی، بتوانیم يك دنباله از بسجمله‌ایهای متعامد  $P_0(x), \dots, P_k(x)$  به‌دست آوریم. به‌موجب ویژگی ۱ مربوط به چنین دنباله‌هایی (قسمت ۳.۶ را ببینید)، به‌ازای ضرایب مناسب  $d_0, \dots, d_k$

هر بسجمله‌ای  $p(x)$  از درجهٔ نایبتر از  $k$  می‌تواند به شکل

$$p(x) = d_0 P_0(x) + \dots + d_k P_k(x)$$

نوشته شود. از قرارداد این  $p(x)$  در (۳۳.۶)، نتیجه می‌گیریم که می‌خواهیم به‌ازای تمام انتخابهای ممکنه از  $d_0, \dots, d_k$

$$E(d_0, \dots, d_k) = \langle f(x) - d_0 P_0(x) - \dots - d_k P_k(x), f(x) - d_0 P_0(x) - \dots - d_k P_k(x) \rangle$$

به حداقل برسد. اگر ماتند بخش ۲.۶ عمل کنیم، نشان داده می‌شود که «بهترین» ضرایب  $d_0^*, \dots, d_k^*$  باید در معادلات نرمال

$$d_0^* \langle P_0, P_i \rangle + d_1^* \langle P_1, P_i \rangle + \dots + d_k^* \langle P_k, P_i \rangle = \langle f, P_i \rangle, \quad i = 0, \dots, k$$

که به‌علت تعامد  $P_j(x)$  ها، به

$$d_i^* \langle P_i, P_i \rangle = \langle f, P_i \rangle \quad i = 0, \dots, k$$

بدل می‌شوند صدق کنند. از این‌رو، اگر

$$S_i = \langle P_i, P_i \rangle \quad i = 0, \dots, k$$

همه غیر صفر باشند، آنگاه بهترین ضرایب اصلاً به‌صورت زیر معین می‌شوند

$$d_i^* = \frac{\langle f, P_i \rangle}{S_i} \quad i = 0, \dots, k \quad (34.6)$$

□ مثال ۱۰.۶: از مجموعهٔ تمام بسجمله‌ایهای  $p(x)$  از درجهٔ نایبتر از سه آن يك را که انتگرال

$$\int_{-1}^1 [e^x - p(x)]^2 dx$$

را مینیمم می‌سازد محاسبه کنید.

در این حالت داریم  $f(x) = e^x$  و حاصلضرب داخلی با

$$\langle g, h \rangle = \int_{-1}^1 g(x)h(x) dx$$

معین می‌شود. از مثال ۶.۶ چنین برمی‌آید که بسجمله‌ایهای متعامد برای این حاصلضرب داخلی، بسجمله‌ایهای لژاندر هستند. با استفاده از جدول ۲.۶ مربوط به این بسجمله‌ایها،

محاسبات زیر را انجام می‌دهیم

$$\langle f, P_0 \rangle = \int_{-1}^1 e^x dx = e - \frac{1}{e}$$

$$\langle f, P_1 \rangle = \int_{-1}^1 e^x x dx = \frac{2}{e}$$

$$\langle f, P_2 \rangle = \frac{2}{3} \int_{-1}^1 e^x \left( x^2 - \frac{1}{3} \right) dx = e - \frac{2}{e}$$

$$\langle f, P_3 \rangle = \frac{5}{4} \int_{-1}^1 e^x \left( x^3 - \frac{3}{5}x \right) dx = -5e + \frac{37}{e}$$

می‌توان نشان داد که برای بسجمله‌ایهای لژاندر (تمرین ۳.۶-۶ را ببینید)، داریم

$$S_i = \langle P_i, P_i \rangle = \frac{2}{2i+1} \quad \text{به‌ازای جمیع مقادیر } i$$

لذا  $S_0 = 2$ ،  $S_1 = 2/3$ ،  $S_2 = 2/5$ ،  $S_3 = 2/7$ ،  $S_4 = 2/9$  با استفاده از (۳.۶) برای محاسبه  $d_i^*$  و استفاده از  $e = 2.718281828$ ، ملاحظه می‌کنیم که تقریب با روش کوچکترین توانهای دوم برای  $e^x$  با بسجمله‌ایهای مکعبی بر بازه  $(-1, 1)$  برابراست با

$$p^*(x) = 1.175201192P_0(x) + 1.103638322P_1(x) \\ + 0.3578143506P_2(x) + 0.07045563367P_3(x)$$

اگر به‌جای  $P_i(x)$ ها مقادیر هم‌ارز آنها را برحسب توانهای  $x$  بگذاریم و از جدول ۲.۶ استفاده و حاصل را مرتب کنیم، خواهیم داشت

$$p^*(x) = 0.9962940183 + 0.9979548730x \\ + 0.5267215260x^2 + 0.1761390842x^3$$

□ ماکسیمم انحراف این بسجمله‌ایها از  $e^x$  بر بازه  $(-1, 1)$  در حدود ۰.۰۱۱ است.

اگر نتوان بسجمله‌ایهای متعامد مناسبی درجدولها به‌دست آورد، باید آنها را پدید آورد. این کار به‌کمک رابطه بازگشتی سه‌جمله‌ای (۳.۶) می‌تواند انجام گیرد. اکنون یک بیان الگوریتمی از این تکنیک برای حالت عملاً مهمی که در آن حاصلضرب داخلی برابراست با

$$\langle g, h \rangle = \sum_{n=0}^N g(x_n)h(x_n)w(x_n) \quad (35.6)$$



می‌دهیم، در اینجا  $x_1, \dots, x_N$  نقاط ثابتی بر بازهٔ  $(a, b)$  هستند.

**الگوریتم ۲.۶:** تولید بسجمله‌ایهای متعامد، برای سادگی، به دست آوردن بسجمله‌ایهای متعامدی را در نظر می‌گیریم که ضریب جملهٔ پیشرو در آنها ۱ باشد. لذا

$$A_i = \alpha_i = 1 \quad i \text{ به ازای جميع مقادیر } i$$

مرحلهٔ ۰ قرار دهید  $P_0(x) \equiv 1$ . بعلاوه عبارت زیر را محاسبه کنید

$$S_0 = \langle P_0, P_0 \rangle = \sum_{n=1}^N w(x_n)$$

اگر  $N \geq 1$  و  $w(x) > 0$ ، آنگاه  $S_0$  برابر با صفر نیست و می‌توان محاسبهٔ

$$P_1(x) = (x - B_0)P_0(x) = x - B_0$$

را ادامه داد که بنا بر ویژگی ۴، مربوط به بسجمله‌ایهای متعامد (بخش ۳.۶ را ببینید)، داریم

$$B_0 = \frac{\langle xP_0(x), P_0(x) \rangle}{S_0} = \frac{\sum_{n=1}^N x_n w(x_n)}{S_0}$$

وقتی که  $P_1(x), \dots, P_0(x)$  را حساب شده تلقی کنیم، مرحلهٔ عمومی یا مرحلهٔ  $j$ ام به گونهٔ زیر انجام می‌گیرد:  
مرحلهٔ  $j$ ام عبارت

$$S_j = \langle P_j, P_j \rangle = \sum_{n=1}^N [P_j(x_n)]^2 w(x_n)$$

را محاسبه می‌کنیم. چون  $P_j(x)$  يك بسجمله‌ای است دقیقاً از درجهٔ  $j$ ،  $S_j$  تنها زمانی می‌تواند صفر باشد که بیش از  $j$  نقطه از نقاط  $x_1, \dots, x_N$  متمایز نباشند. از این رو اگر بیش از  $j$  نقطهٔ متمایز بین  $x_n$ ها وجود داشته باشد، می‌توانیم

$$B_j = \frac{\langle xP_j(x), P_j(x) \rangle}{S_j} = \frac{\sum_{n=1}^N x_n [P_j(x_n)]^2 w(x_n)}{S_j}$$

و

$$C_j = \frac{S_j}{S_{j-1}}$$

را محاسبه کنیم و بسجمله‌ای متعامد بعدی را به صورت زیر به دست آوریم

$$P_{j+1}(x) = (x - B_j)P_j(x) - C_j P_{j-1}(x) \quad (۳.۶.۶)$$

□ مثال ۱۲.۶: با استفاده از بسجمله‌ایهای متعامد، مسئله تقریب‌زدن با روش کوچکترین مربعات در مثال ۵.۶ حل کنید.

برای این مثال داریم  $f(x) = 10 - 2x + x^2/10$

$$x_n = 10 + \frac{n-1}{5} f_n = f(x_n) \quad n = 1, \dots, 6$$

و در پی آن بسجمله‌ای از درجه نایبتر از ۲ هستیم که مقدار

$$\sum_{n=1}^6 [f_n - P(x_n)]^2$$

را مینیمم کند، یعنی حاصلضرب داخلی (۲۷.۶) را به‌ازای  $w(x) \equiv 1$  باید مورد بحث قرار دهیم. با متابعت از الگوریتم ۲.۶، مقادیر زیر را محاسبه می‌کنیم

$$S_0 = \sum_{n=1}^6 1 = 6 \quad \text{از این رو} \quad P_0(x) \equiv 1$$

$$B_0 = \sum_{n=1}^6 \left[ 10 + \frac{n-1}{5} \right] / S_0 = \frac{63}{6} = 10.5$$

بنابراین

$$P_1(x) = x - 10.5 \quad S_1 = \sum_{n=1}^6 \left[ \frac{n-1}{5} - 10.5 \right]^2 = 0.7$$

و چون  $S_1 \neq 0$ ، به محاسبه  $P_2(x)$  می‌پردازیم. اگر از هفت رقم اعشاری استفاده و نتیجه را گرد کنیم، خواهیم داشت

$$B_1 = \sum_{n=1}^6 \frac{[10 + (n-1)/5][(n-1)/5 - 10.5]^2}{S_1} = \frac{735}{0.7} = 105$$

$$C_1 = \frac{S_1}{S_0} = \frac{0.7}{6} = 0.1166667$$

از این نتایج چنین به‌دست خواهیم آورد

$$P_2(x) = (x - 10.5)^2 - 0.1166667 \quad S_2 = 0.05973332$$

در مرحله بعد با استفاده از رابطه (۳۴.۶) و حساب ممیزشناور با هفت رقم اعشار، بهترین ضرایب  $d_0^*$ ،  $d_1^*$  و  $d_2^*$  را برای تقریب با روش کوچکترین مربعات محاسبه می‌کنیم،

که نتایج این محاسبات چنین‌اند  $p^*(x) = d_0^* p_0(x) + d_1^* p_1(x) + d_2^* p_2(x)$

$$d_0^* = \sum_{n=1}^9 \frac{f_n}{9} = 0.036666667$$

$$d_1^* = \sum_{n=1}^9 \frac{f_n P_1(x_n)}{0.9} = 0.1$$

$$\square \quad d_2^* = \sum_{n=1}^9 \frac{f_n P_2(x_n)}{0.054973337} = 0.099999999$$

برای مقایسه این نتایج با آنچه که در مثال ۵.۶ محاسبه شد،  $p^*(x)$  را بر حسب  $x$ ، محاسبه می‌نویسیم و خواهیم داشت

$$\begin{aligned} p^*(x) &= 0.036666667 + 0.1(x - 10.5) \\ &\quad + 0.099999999[(x - 10.5)^2 - 0.11666667] \\ &= 0.036666667 - 1.05 + 0.099999999(110.25 - 0.11666667) \\ &\quad + [0.1 + 0.099999999(-21)]x + 0.099999999x^2 \end{aligned}$$

لذا با این روش محاسبه، مقادیر  $c_i^*$  های مثال ۵.۶ برابر خواهند بود با

$$c_0^* = 9.999998 \dots \quad c_1^* = -1.9999998 \dots$$

$$c_2^* = 0.9999999 \dots$$

برعکس، در مثال ۵.۶، وقتی معادلات نرمال (۲۳.۶) را بر حسب  $c_i^*$  ها با استفاده از حساب ممیز شناور با هفت رقم اعشار به‌طور مستقیم حل کردیم، مقادیر زیر به دست آمدند

$$c_0^* = 8.9992 \dots \quad c_1^* = -1.9712 \dots \quad c_2^* = 0.90863 \dots$$

بنابراین، با استفاده از بسجمله‌ایهای متعامد، نتایج به‌سود مؤثری را در این مثال نشان می‌دهند.

در ضمن، معمولاً زحمت محاسبه  $p^*(x)$  بر حسب  $x$  را، به‌خود نمی‌دهند، بلکه هر زمان که  $p^*(x)$  باید محاسبه شود، الگوریتم ۱.۶ را، همراه با  $d_i^*$ ، که محاسبه شده، به‌کار می‌برند، زیرا که ضرایب  $B_i$  و  $C_i$  رابطه بازگشتی در دست هستند.

تولید بسجمله‌ایهای متعامد و محاسبه بهترین ضرایب  $d_i^*$ ، برای صرفه‌جویی در حافظه، با اجرا به‌زبان فورترن بهتر ترکیب و به‌یک عمل بدل می‌شود. برای محاسبه  $d_j^*$  و  $P_{j+1}(x)$ ، تنها اعداد

$$n = 1, \dots, N \quad P_{j-1}(x_n), \quad P_j(x_n)$$

مورد نیازند. بنابراین، اگر به‌محض در دسترس بودن  $P_j(x_n)$ ،  $n = 1, \dots, N$ ، محاسبه  $d_j^*$

شود آنگاه، همین که  $P_{j+1}(x)$  و  $P_{j+2}(x)$  محاسبه شدند  $P_j(x)$  ها،  $n=1, \dots, N$  می توانند بدون ایجاد هیچ گونه مسئله ای نادیده گرفته شوند. باز احتیاجی نیست که  $P_j(x)$  را مثلاً، بر حسب توانهای  $x$  بسازیم، زیرا که، فقط بمقادیر آنها در  $x_n$ ،  $n=1, \dots, N$  نیاز داریم.

```

SUBROUTINE ORTPOL ( X, F, W, NPOINT, PJM1, PJ, ERROR )
C CONSTRUCTS THE DISCRETE WEIGHTED LEAST SQUARES APPROXIMATION BY POLY-
C NOMIALS OF DEGREE .LT. NTERMS TO GIVEN DATA.
C***** I N P U T *****
C (X(I), F(I)), I=1,...,NPOINT GIVES THE ABSISSAE AND ORDINATES OF
C THE GIVEN DATA POINTS TO BE FITTED.
C W NPOINT-VECTOR CONTAINING THE POSITIVE WEIGHTS TO BE USED.
C NPOINT NUMBER OF DATA POINTS.
C***** I N P U T VIA COMMON BLOCK P O L Y *****
C NTERMS GIVES THE ORDEK (= DEGREE + 1) OF THE POLYNOMIAL APPROXIMANT.
C***** W O R K A R E A S *****
C PJM1, PJ ARKAYS OF LENGTH NPOINT TO CONTAIN THE VALUES AT THE X'S
C OF THE TWO MOST RECENT ORTHOGONAL POLYNOMIALS.
C***** O U T P U T *****
C ERROR NPOINT-VECTOR CONTAINING THE ERROR AT THE X'S OF THE POLYNOM-
C IAL APPROXIMANT TO THE GIVEN DATA.
C***** O U T P U T VIA COMMON BLOCK P O L Y *****
C B, C ARRAYS CONTAINING THE COEFFICIENTS FOR THE THREE-TERM RECUR-
C RENCE WHICH GENERATES THE ORTHOGONAL POLYNOMIALS.
C D COEFFICIENTS OF THE POLYNOMIAL APPROXIMANT TO THE GIVEN DATA WITH
C RESPECT TO THE SEQUENCE OF ORTHOGONAL POLYNOMIALS.
C THE VALUE OF THE APPROXIMANT AT A POINT Y MAY BE OBTAINED BY A
C REFERENCE TO ORTVAL(Y) .
C***** M E T H O D *****
C THE SEQUENCE P0, P1, ..., PNTERMS-1 OF ORTHOGONAL POLYNOMIALS WITH
C RESPECT TO THE DISCRETE INNER PRODUCT
C (P,Q) = SUM ( P(X(I))*Q(X(I))*W(I), I=1,...,NPOINT)
C IS GENERATED IN TERMS OF THEIR THREE-TERM RECURRENCE
C PJP1(X) = (X - B(J+1))*PJ(X) - C(J+1)*PJM1(X) ,
C AND THE COEFFICIENT D(J) OF THE WEIGHTED LEAST SQUARES APPROXIMAT-
C ION TO THE GIVEN DATA IS OBTAINED CONCURRENTLY AS
C D(J+1) = (F,PJ)/(PJ,PJ), J=0,...,NTERMS-1 .
C ACTUALLY, IN ORDER TO REDUCE CANCELLATION, (F,PJ) IS CALCULATED AS
C (ERROR,PJ), WITH ERROR = F INITIALLY, AND, FOR EACH J, ERROR RE-
C DUCED BY D(J+1)*PJ AS SOON AS D(J+1) BECOMES AVAILABLE.
C
C INTEGER NPOINT,NTERMS, I,J
C REAL E,C,D,ERROR(NPOINT),F(NPOINT),PJ(NPOINT),PJM1(NPOINT),
C * D(NPOINT),X(NPOINT), P,S(20)
C COMMON /POLY/ NTERMS,B(20),C(20),D(20)
C
C DO 9 J=1,NTERMS
C E(J) = 0.
C D(J) = 0.
9 S(J) = 0.
C(1) = 0.
C DO 10 I=1,NPOINT
C D(1) = D(1) + F(I)*W(I)
C B(1) = B(1) + X(I)*W(I)
10 S(1) = S(1) + W(I)
C D(1) = D(1)/S(1)
C DO 11 I=1,NPOINT
11 ERROR(I) = F(I) - D(1)
C IF (NTERMS .EQ. 1) RETURN
C B(1) = B(1)/S(1)
C DO 12 I=1,NPOINT
C PJM1(I) = 1.
12 PJ(I) = X(I) - B(1)
C
C DO 30 J=2,NTERMS
C DO 21 I=1,NPOINT
C P = PJ(I)*W(I)
C D(J) = D(J) + ERROR(I)*P
C P = P*PJ(I)
C B(J) = B(J) + X(I)*P

```

```

21      S(J) = S(J) + P
      D(J) = D(J)/S(J)
      DO 22 I=1,NPOINT
22      ERROR(I) = ERROR(I) - D(J)*PJ(I)
      IF (J .EQ. NTERMS) RETURN
      B(J) = B(J)/S(J)
      C(J) = S(J)/S(J-1)
      DO 27 I=1,NPOINT
      P = PJ(I)
      PJ(I) = (X(I) - D(J))*PJ(I) - C(J)*PJM1(I)
27      PJM1(I) = P
30 CONTINUE
      END
      RETURN

```

محاسبهٔ  $D(j)$ ، به گونه‌ای که در این زیر برنامه انجام گرفت، شاید نیاز به مقدماتی توضیح داشته باشد. چون  $D(j) = d_{j-1}^*$ ، از رابطهٔ (۳۷.۶) نتیجه می‌گیریم

$$D(j) = \sum_{n=1}^{NPOINT} \frac{f_n P_{j-1}(x_n) w(x_n)}{S_{j-1}} \quad (37.6)$$

در حالی که در برنامه، به صورت زیر محاسبه شده است

$$D(j) = \sum_{n=1}^{NPOINT} \frac{ERROR(n) P_{j-1}(x_n) w(x_n)}{S_{j-1}} \quad (38.6)$$

که در آن

به ازای تمام مقادیر  $n$

$$ERROR(n) = f_n - D(1)P_0(x_n) - \dots - D(j-1)P_{j-2}(x_n) \quad (39.6)$$

اگر (۳۹.۶) را در (۳۸.۶) قرار دهیم، چون  $P_{j-1}$  بر  $P_0(x), \dots, P_{j-2}(x)$  متعامد است خواهیم داشت

$D(j)$

$$\begin{aligned}
 &= \sum_{n=1}^{NPOINT} \frac{[f_n - D(1)P_0(x_n) - \dots - D(j-1)P_{j-2}(x_n)] P_{j-1}(x_n) w(x_n)}{S_{j-1}} \\
 &= \frac{\sum_n f_n P_{j-1}(x_n) w(x_n) - D(1)\langle P_0, P_{j-1} \rangle - \dots - D(j-1)\langle P_{j-2}, P_{j-1} \rangle}{S_{j-1}} \\
 &= \sum_n \frac{f_n P_{j-1}(x_n) w(x_n)}{S_{j-1}}
 \end{aligned}$$

از این رو در حساب دقیق یا حساب با دقت نامتناهی، هر دو رابطهٔ (۳۷.۶) و (۳۸.۶) مقادیر واحدی برای  $D(j)$  به دست می‌دهند. اما در حساب با دقت متناهی، می‌توان انتظار داشت که رابطهٔ (۳۸.۶) به دلیل زیر دقیقتر باشد؛ چون

## 1. infinite-precision

$$p_r^*(x) = D(1)P_0(x) + \dots + D(r+1)P_r(x)$$

تقریب با روش کوچکترین توانهای دوم (وزین) برای  $f(x)$  با بسجمله‌ایهای از درجه نایبتر از  $r$  است، نتیجه آن که می‌توان انتظار داشت اعداد

$$\text{ERROR}(n) = f_n - P_{j-1}^*(x_n) \quad n = 1, \dots, \text{NPOINT}$$

از اعداد  $f_n$ ،  $n = 1, \dots, \text{NPOINT}$ ، کوچکتر باشند. از این رو، به نظر می‌رسد که محاسبه رابطه (۳۸.۶) در مقایسه با محاسبه رابطه (۳۷.۶)، به علت تفریق مقادیر تقریباً مساوی از یکدیگر، احتمالاً موجب از دست رفتن ارقام با معنا می‌شود (تمرین ۴-۱ را ببینید)

□ مثال ۱۳.۶: مقادیر  $f_n$  از تابع  $f(x) = e^x$  در

$$x_n = (n-1)/10 - 1, (n = 1, \dots, 21)$$

و گرد شده تا دو رقم بعد از ممیز، داده شده‌اند. سعی کنید که اطلاعات مربوط به  $f(x)$  را که در این داده‌ها گنجانیده شده‌اند بازیابی کنید:

ما سعی می‌کنیم این مسئله را با محاسبه آن بسجمله‌ای  $p_4^*(x)$  از مجموعه تمام بسجمله‌ایهای  $p_4(x)$  از درجه نایبتر از ۳ که عبارت

$$\sum_{n=1}^{21} [f_n - p_4(x_n)]^2$$

را مینیمم می‌سازد، حل کنیم. برنامه فورترن زیر به کمک زیر برنامه ORTPOL، که قبلاً اشاره شد،  $p_4^*(x)$  را بر آورد می‌کند و سپس با استفاده از FUNCTION ORTVL، که اساس آن الگوریتم ۱.۶ است،  $p_4^*(x)$  را در  $x_n$  محاسبه می‌کند.

```

C PROGRAM FOR EXAMPLE 6.13 .
  PARAMETER NPMAX=100
  INTEGER NTERMS, I,J,NPOINT
  REAL B,C,D,ERROR(NPMAX),F(NPMAX),PJ(NPMAX),PJMI(NPMAX),W(NPMAX)
  * ,X(NPMAX)
  COMMON /POLY/ NTERMS,B(20),C(20),D(20)
  NPOINT = 21
  DO 1 I=1,NPOINT
    W(I) = 1.
    X(I) = -1. + FLOAT(I-1)/10.
  1 F(I) = FLOAT(IFIX(EXP(X(I))*100. + .5))/100.
  NTERMS = 4
  CALL ORTPOL( X, F, W, NPOINT, PJMI, PJ, ERROR )
  PRINT 601, (J,B(J),C(J),D(J),J=1,NTERMS)
601 FORMAT(12,3E16.8)
  DO 60 I=1,NPOINT
    PJMI(I) = EXP(X(I))
  60 PJ(I) = ORTVL(X(I))
  PRINT 660, (X(I),F(I),PJ(I),ERROR(I),PJMI(I),I=1,NPOINT)
660 FORMAT(F5.1,F8.3,F10.5,E13.3,F10.5)
  STOP
END

```

REAL FUNCTION ORTVL (X)

C RETURNS THE VALUE AT X OF THE POLYNOMIAL OF DEGREE .LT. NTERMS

C GIVEN BY

C D(1)\*P0(X) + D(2)\*P1(X) + ... + D(NTERMS)\*PNTERMS-1(X),

C WITH THE SEQUENCE P0, P1, ... OF ORTHOGONAL POLYNOMIALS GENERATED

C BY THE THREE-TERM RECURRENCE

C PJPI(X) = (X - B(J+1))\*PJ(X) - C(J+1)\*PJMI(X), ALL J .

C

```

COMMON /POLY/ NTERMS,B(20),C(20),D(20)
PREV = 0.
QRTVAL = D(INTERMS)
IF (INTERMS.EQ. 1) RETURN
DO 10 K=NTERMS-1,1,-1
  PREV2 = PREV
  PREV = QRTVAL
  ORTVAL = D(K) + (X - B(K))*PREV - C(K+1)*PREV2
10 CONTINUE
RETURN
END

```

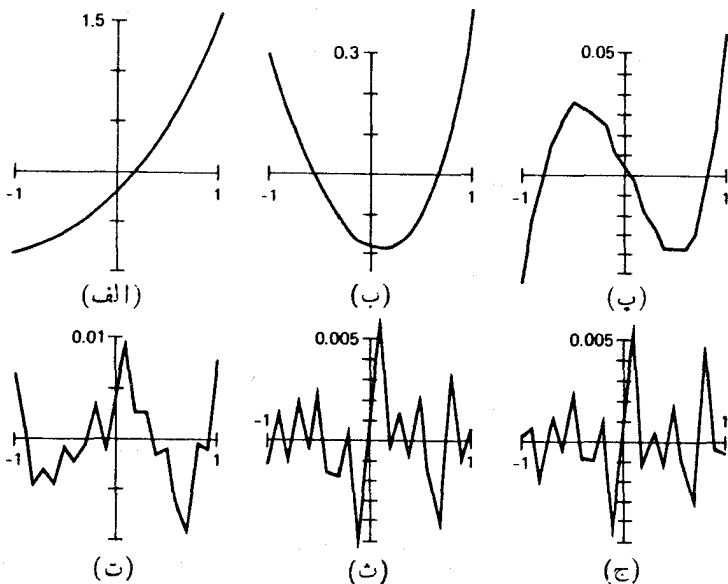
## جدول ۴.۶ نتایج کامپیوتری برای مثال ۱۳.۶

$x_n$	$f_n$	$p_3^*(x_n)$	$f_n - p_3^*(x_n)$	$p_4^*(x_n)$	$f_n - p_4^*(x_n)$	$e_n$
-1.0	0.370	0.36387	6.130E-03	0.37115	-1.154E-03	0.36788
-0.9	0.410	0.40874	1.263E-03	0.40874	1.263E-03	0.40657
-0.8	0.450	0.45481	-4.806E-03	0.45097	-9.719E-04	0.44933
-0.7	0.500	0.50315	-3.148E-03	0.49804	1.964E-03	0.49659
-0.6	0.550	0.55484	-4.836E-03	0.55021	-2.134E-04	0.54881
-0.5	0.610	0.61094	-9.436E-04	0.60789	2.108E-03	0.60653
-0.4	0.670	0.67524	-2.542E-03	0.67156	-1.565E-03	0.67032
-0.3	0.740	0.74070	-7.045E-04	0.74183	-1.832E-03	0.74082
-0.2	0.820	0.81650	3.497E-03	0.81940	6.029E-04	0.81873
-0.1	0.900	0.90101	-1.010E-03	0.90507	-5.070E-03	0.90484
0.0	1.000	0.99530	4.710E-03	0.99976	2.358E-04	1.00000
0.1	1.110	1.10044	9.558E-03	1.10450	5.499E-03	1.10517
0.2	1.220	1.21751	2.490E-03	1.22040	-4.045E-04	1.22140
0.3	1.350	1.34758	2.422E-03	1.34871	1.294E-03	1.34986
0.4	1.490	1.49172	-1.717E-03	1.49074	-7.399E-04	1.49182
0.5	1.650	1.65100	-1.000E-03	1.64795	2.052E-03	1.64872
0.6	1.820	1.82650	-6.499E-03	1.82188	-1.876E-03	1.82212
0.7	2.010	2.01929	-9.287E-03	2.01418	-4.176E-03	2.01375
0.8	2.230	2.23044	-4.368E-04	2.22660	0.397E-03	2.22554
0.9	2.460	2.46102	-1.020E-03	2.46102	-1.020E-03	2.45960
1.0	2.720	2.71211	7.890E-03	2.71939	6.061E-04	2.71828

جدول ۴.۶، نتایج محاسباتی را که روی کامپیوتر ۶۵۰۰ CDC انجام گرفته به دست می‌دهد. خطای  $f_n - p_4^*(x_n)$  در شکل ۴.۶، رسم شده است و نشان می‌دهد خطا تاحدی به صورت منظم عمل می‌کند که این خودگویای آن است که  $p_4^*(x)$  تمام اطلاعات گنجانده شده در داده‌های مفروض را نشان نمی‌دهد. از این رو تقریب با روش کوچکترین مربعات  $p_4^*(x)$  برای داده‌های مفروض، به وسیلهٔ بسجمله‌ایهایی از درجهٔ نایبتر از ۴ نیز محاسبه شده‌اند. این نتایج نیز در جدول ۴.۶ داده شده‌اند. خطای  $f_n - p_4^*(x)$  در شکل ۴.۶، رسم شده است و دیده می‌شود که خطا کاملاً نامنظم عمل می‌کند. از این رو می‌توان فرض کرد که  $p_4^*(x)$  تمام اطلاعات گنجانده شده در داده‌های مفروض  $f_n$  را نشان می‌دهد. افزایش هرچه بیشتر درجهٔ بسجمله‌ای تقریب کننده تنها موجب می‌شود که به تابع تقریب کننده، آزادی بیشتری برای تقریب زدن نوفه در داده‌ها نیز بدهد. □

## تمرین

۴-۱۰ اگر  $f(x) = 6000 + x$ ، آنگاه هر تقریب برای  $f(x)$  از راه کوچکترین مربعات به وسیلهٔ خطهای راست، خود  $f(x)$  است. بسجمله‌ای



شکل ۶.۶ خطا در تقریب زدن با روش کوچکترین مربعات برای داده‌های مثال ۳.۶ به وسیلهٔ بسجمله‌ایهای درجهٔ (الف) صفر، (ب) یک، (پ) دو، (ت) سه، (ث) چهار، (ج) پنج.

$$p_n^*(x) = d_n^* + d_n^* x$$

که عبارت

$$\sum_{n=-2}^2 [f(x) - p_n(n)]^2$$

را مینیمم می‌سازد حساب کنید. توجه کنید که ۱ و  $x$  متعامدند، لذا تنها باید  $d_0^*$  و  $d_1^*$  را محاسبه کنید. تفاضل بین (۳۷.۶) و (۳۸.۶) را با محاسبهٔ  $d_1^*$  به دو طریق و با استفاده از حساب با ممیز شناور با چهار رقم اعشار نشان دهید.

۴-۴.۶ از مجموعهٔ تمام بسجمله‌ایهای  $p(x)$  از درجهٔ نایبتر از ۲، آن یک را که

$$\int_{-1}^1 [\sin \pi x - p(x)]^2 dx$$

را مینیمم می‌سازد محاسبه کنید. از بسجمله‌ی لژاندر استفاده کنید و تمام محاسبات را با پنج رقم اعشار انجام دهید. (توجه:  $\pi = ۳.۱۴۱۵۹۳$ )

۴-۴.۶ زیر برنامهٔ ORTPOL را روی کامپیوتر خود پیاده کنید و سپس از آن برای حل مسئلهٔ زیر استفاده کنید. از جدولی از مقادیر  $f(x) = \sin \pi x$  در  $x_n = (n-1)/۱۰ - ۱$



$(n = 1, \dots, 21)$ ، مقادیر  $f_n(x) = \sin \pi x_n$  را گرد شده تا سه رقم اعشاری به دست آورید. سپس از مجموعهٔ تمام بسجمله‌ایهای نایبتر از درجهٔ ۴،  $p_4(x)$ ، آن بسجمله‌ای  $p_4^*(x)$  را که

$$\sum_{n=1}^{21} [f_n - p_4(x_n)]^2$$

را مینیمم می‌سازد به دست آورید.

### ۵.۶\* تقریب به وسیلهٔ بسجمله‌ایهای مثلثاتی

بسیاری از پدیده‌های فیزیکی، مانند نور و صوت، مشخصهٔ دوره‌ای<sup>۱</sup> دارند. این پدیده‌ها با توابع  $f(x)$  که دوره‌ای هستند، یعنی توابعی که به ازای کلیهٔ مقادیر  $x$  و یک عدد ثابت  $\tau$ ، به نام دورهٔ تناوب<sup>۲</sup>، در رابطهٔ

$$f(x + \tau) = f(x)$$

صدق می‌کنند، بیان می‌شوند. از آنجا که فقط بسجمله‌ایهای دوره‌ای توابع ثابت هستند، لذا باید دسته‌ای دیگر از توابع را برای تقریب‌زدن مؤثر توابع دوره‌ای به کار گیریم و برای انجام این کار بسجمله‌ایهای مثلثاتی از همه مناسبترند. بنا بر تعریف، یک بسجمله‌ای مثلثاتی از مرتبهٔ  $n$  هر تابعی است به شکل

$$p(x) = a_0/2 + \sum_{j=1}^n [a_j \cos jx + b_j \sin jx] \quad (40.6)$$

که در آن  $a_0, \dots, a_n$  و  $b_1, \dots, b_n$  ثابتهای حقیقی با هم‌تافت باشند. این گونه بسجمله‌ای مثلثاتی، یک بسجمله‌ای  $2\pi$ -دوره‌ای<sup>۳</sup> است. بنا بر این وقتی که می‌خواهیم برای تابع  $2\pi$ -دوره‌ای  $f(x)$ ، به ازای  $\tau \neq 2\pi$  تقریب پیدا کنیم، لازم است که اصلاحاتی انجام دهیم. در چنین مواردی قرار می‌گذاریم که تابع  $2\pi$ -دوره‌ای  $g(x) = f(\tau x / (2\pi))$  را مورد توجه قرار دهیم. آنگاه پس از ساختن یک بسجمله‌ای مثلثاتی تقریبی  $p(x)$  برای  $g(x)$ ، از آن یک تابع تقریبی  $2\pi$ -دوره‌ای برای  $f(x)$  به شکل  $p(2\pi x / \tau)$  به دست می‌آوریم. با در نظر گرفتن مطلب فوق، از این به بعد فرض می‌کنیم که تابع  $f(x)$  که باید تقریب زده شود قبلاً  $2\pi$ -دوره‌ای باشد.

چنانچه معلوم خواهد شد، غالباً مناسبتر است که بسجمله‌ایهای مثلثاتی از مرتبهٔ  $n$  را به شکل هم‌تافت هم‌ارز آنها، یعنی به صورت زیر بنویسیم

$$p(x) = \sum_{j=-n}^n c_j e^{ijx} \quad (41.6)$$

1. periodic

2. period

3.  $2\pi$ -periodic

در اینجا و در بقیه این بخش و بخشهای بعدی علامت  $i$  معرف واحد انگاری<sup>۱</sup>

$$i = \sqrt{-1}$$

است و ارتباط بین روابط (۴۰.۶) و (۴۱.۶) به وسیله فرمول اویلر

$$e^{ix} = \cos x + i \sin x \quad (۴۲.۶)$$

برقرار می شود (که اثبات آن را می توان در تمرین ۷.۱-۹ پیدا کرد). با توجه به فرمول اویلر خواهیم داشت [بارعایت  $\cos(-jx) = \cos jx$  و  $\sin(-jx) = -\sin jx$ ] که

$$\begin{aligned} \sum_{j=-n}^n c_j e^{ijx} &= \sum_{j=-n}^n c_j [\cos jx + i \sin jx] \\ &= c_0 + \sum_{j=1}^n [(c_j + c_{-j}) \cos jx + i(c_j - c_{-j}) \sin jx] \end{aligned}$$

این رابطه نشان می دهد که رابطه (۴۱.۶) به همان شکل رابطه (۴۰.۶) است، با فرض

$$a_j = c_j + c_{-j} \quad b_j = i(c_j - c_{-j}) \quad j = 0, \dots, n \quad (\text{الف } ۴۳.۶)$$

این رابطه را به آسانی می توان عکس کرد و رابطه (۴۰.۶)، را که به شکل رابطه (۴۱.۶) است به دست آورد، با شرایط

$$c_j = (a_j - ib_j)/2 \quad c_{-j} = (a_j + ib_j)/2 \quad j = 0, \dots, n \quad (\text{ب } ۴۳.۶)$$

باید توجه کرد که (۴۱.۶) معرف يك تابع حقیقی است اگر، فقط اگر خود مزدوج همثافت خود باشد. اما از آنجا که داریم

$$\overline{\sum_{j=-n}^n c_j e^{ijx}} = \sum_{j=-n}^n \bar{c}_j e^{-ijx} = \sum_{j=-n}^n \bar{c}_{-j} e^{ijx}$$

این رابطه گویای این است که (۴۱.۶) يك تابع حقیقی است اگر، فقط اگر، تساوی

$$c_j = \bar{c}_{-j} \quad \text{به ازای همه مقادیر } j \quad (۴۴.۶)$$

برقرار باشد. بنابراین اگر (۴۰.۶) یا (۴۱.۶) تابع حقیقی باشد، آنگاه (الف ۴۳.۶) به شکل زیر ساده می شود

$$a_j = 2 \operatorname{Re} c_j \quad b_j = -2 \operatorname{Im} c_j \quad (۴۵.۶)$$

تقریب زدن به وسیله بسجمله ایهای مثلثاتی اساساً توسط سری فوریه<sup>۲</sup>

$$f(x) \approx \sum_{j=-\infty}^{\infty} \hat{f}(j)e^{ijx} \quad (۴۶.۶)$$

انجام می‌پذیرد که در آن ضرایب فوریه  $\hat{f}(j)$  با انتگرال زیر محاسبه می‌شوند

$$\hat{f}(j) = \frac{1}{2\pi} \int_0^{2\pi} f(x)e^{-ijx} dx \quad (۴۷.۶)$$

در شرایط نسبتاً معتدل، سری فوق به سمت  $f(x)$  همگرا می‌شود [اما نه برای هر  $f(x)$ ]. برای مثال، اگر  $f(x)$  تابعی پیوسته و دارای اولین مشتق پیوسته-تکه‌ای<sup>۱</sup> باشد، سری به‌طور یکنواخت همگرا می‌شود.

سری فوریه با توجه به دستور مسلم زیر به دست می‌آید

$$\int_0^{2\pi} \overline{e^{ijx}} e^{ikx} dx = \int_0^{2\pi} e^{i(k-j)x} dx = \begin{cases} \int_0^{2\pi} 1 dx = 2\pi & k=j \\ \frac{1}{i(k-j)} e^{i(k-j)x} \Big|_0^{2\pi} = 0 & k \neq j \end{cases}$$

این امر نشان می‌دهد که توابع  $1$  و  $e^{\pm i\pi x}$  و  $e^{\pm i2\pi x}$  نسبت به ضرب اسکالر یا ضرب داخلی متعامند. به عبارت دیگر

$$\langle g, h \rangle = \frac{1}{2\pi} \int_0^{2\pi} \overline{h(x)} g(x) dx$$

متعامند. به عبارت دیگر

$$\langle e^{ikx}, e^{ijx} \rangle = \begin{cases} 1 & k=j \\ 0 & k \neq j \end{cases}$$

که این مطلب اثبات قضیهٔ زیر است.

**قضیه ۴.۶** حاصلجمع جزئی<sup>۲</sup>

$$\sum_{j=-n}^n \hat{f}(j)e^{ijx}$$

سری فوریه برای  $f(x)$  نسبت به نرم

$$\|g\| = \|g\|_2 = \left[ \frac{1}{2\pi} \int_0^{2\pi} |g(x)|^2 dx \right]^{1/2}$$

بهترین تقریب برای  $f(x)$  به وسیلهٔ بسجمله‌ایهای مثلثاتی از مرتبهٔ  $n$  است. بعلاوه می‌توان نشان داد که رابطهٔ پارسوال<sup>۱</sup>

$$\sum_j |\hat{f}(x)|^2 = \frac{1}{2\pi} \int_0^{2\pi} |f(x)|^2 dx \quad (۴۸.۶)$$

برقرار است.

ضرایب فوریه  $\hat{f}(j)$  مربوط به تابع  $f(x)$  برای «درک» این تابع به شرح زیر به کار گرفته می‌شوند. فرض کنید که تابع  $f(x)$  حقیقی و  $2\pi$ -دوره‌ای باشد. اگر  $f(x)$  به صورت موقعیت جسم در حال حرکتی بزرگ خط در زمان  $x$  در نظر گرفته شود، آنگاه تابع  $2\pi$ -دوره‌ای  $f(x)$  یک حرکت دوره‌ای را بیان می‌کند. حال اگر

$$\hat{f}(j) = |\hat{f}(j)| e^{i\theta_j}$$

درست باشد [شکل قطبی عدد مختلط  $\hat{f}(j)$ ]، آنگاه می‌توان سری فوریه برای  $f(x)$  را به صورت زیر نوشت

$$f(x) \approx \sum_{j=0}^{\infty} |\hat{f}(j)| \cos(\theta_j + jx)$$

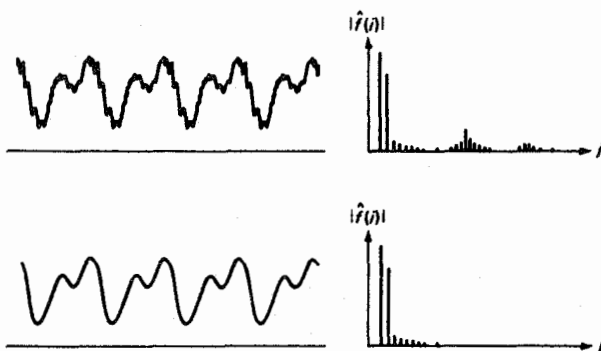
(به تمرین ۵.۶-۷ نگاه کنید). بدین طریق حرکت دوره‌ای مساکسه با  $f(x)$ ، بیان شده به صورت حاصلجمع یا برهمتهی<sup>۲</sup> نوسانهای همساز<sup>۳</sup> ساده نمایش داده شده است. زامین حرکت از این نوع

$$2|\hat{f}(j)| \cos(\theta_j + jx)$$

دارای دامنهٔ  $2|\hat{f}(j)|$ ، فراوانی  $(2\pi)/z$ ، فراوانی زاویه‌ای  $j$ ، دوره یا طول موج  $2\pi/z$  و زاویهٔ فاز  $\theta_j$  می‌باشد. عدد  $|\hat{f}(j)|$  معرف میزان حرکت همساز ساده با فراوانی زاویه‌ای  $j$  است که در کل حرکت موجود است. تمامی دنبالهٔ  $|\hat{f}(0)|$  و  $|\hat{f}(1)|$ ،... (یا احتمالاً، دنبالهٔ مربعات آنها) را طیف توان<sup>۴</sup> یا به طور ساده طیف  $f(x)$  گویند. باید توجه داشت که به موجب رابطهٔ پارسوال (۴۸.۶) طیف  $f(x)$  با کران  $\|f\|_2$  محدود می‌شود، اما بسته به اینکه چگونه «کل انرژی»  $\|f\|_2$  روی طیف  $|\hat{f}(0)|$  و  $|\hat{f}(1)|$  و... توزیع شده باشد،  $f(x)$  ممکن است رفتار بسیار متفاوتی داشته باشد. یک تابع «نوفه‌ای» به ازای مقادیر بزرگ  $j$ ، دارای  $|\hat{f}(j)|$  قابل ملاحظه‌ای است، در صورتی که در یک تابع «هموار» طیف همراه با افزایش  $j$  به سرعت کاهش می‌یابد، به شکل ۷.۶ نگاه کنید.

یک روش مورد توجه ما برای هموار کردن، عبارت است از تولید ضرایب فوریهٔ

1. Parseval's relation
2. superposition
3. harmonic oscillation
4. power spectrum



**شکل ۷.۶** دو تابع  $2\pi$ -دوره‌ای حقیقی و طیف توان آنها. تابع دومی از کوچک کردن اثر فراوانیهای بالا در تابع اولی به دست آمده است.

تابع داده شده از روی داده‌ها، و سپس پالایش این ضرایب یعنی نادیده گرفتن برخی فراوانیها، معمولاً فراوانیهای زیاد، و سپس بازسازی تابع مزبور به صورت سری فوریه با استفاده از این ضرایب «خالص شده» یا «پالایش شده». برای مثال به شکل ۷.۶ نگاه کنید. می‌توان نشان داد که در صورتی که  $f(x)$  دارای  $k-1$  مشتق پیوسته (مانند یک تابع دوره‌ای) و مشتق  $k$ ام به طور تکه‌ای پیوسته (یا حتی فقط تغییرات کرانی محدود) باشد، رابطه

$$|\hat{f}(j)| = O(|j|^{-k-1}) \quad (۷.۹۶)$$

برقرار است. برای مثال، «موج مربعی»

$$f(x) = \text{signum}(\sin x) = \begin{cases} 1 & 0 < x < \pi \\ -1 & \pi < x < 2\pi \end{cases}$$

فقط به طور تکه‌ای پیوسته است. بنابراین انتظار می‌رود که وقتی  $|z| \rightarrow \infty$ ،  $|\hat{f}(j)|$  با سرعتی نه بیشتر از  $1/|z|$  به سمت صفر میل نماید. این امر با محاسبات مستقیم زیر تأیید می‌شود:

$$\begin{aligned} \hat{f}(j) &= \frac{1}{2\pi} \int_0^{2\pi} \text{signum}(\sin x) e^{-ijx} dx \\ &= \frac{1}{2\pi} \left\{ \int_0^{\pi} e^{-ijx} dx - \int_{\pi}^{2\pi} e^{-ijx} dx \right\} \\ &= \frac{1}{2\pi} \frac{1}{-ij} \left\{ \gamma e^{-ij\pi} + \gamma \right\} = i \frac{\gamma}{\pi} \begin{cases} 0 & \text{زوج} \\ 1/j & \text{فرد} \end{cases} \end{aligned}$$

باید توجه داشت که طیف تابع

$$f(x) = x$$

سویعتر از  $1/j$  تحلیل نمی‌رود، ولو اینکه تابع مزبور غالباً بی‌نهایت بار مشتقپذیر باشد. علت آن این است که تحلیل فوریه (آن گونه که در اینجا تشریح شده است) بسا این تابع به صورت یک تابع  $2\pi$ -دوره‌ای، که مقدارش به ازای  $0 < x < 2\pi$  برابر  $x$  است، عمل می‌نماید. اما این تابع اخیر در تمام مضارب  $2\pi$  دارای یک جهش ناپایسته است!

معمولاً بعید است که بتوان ضرایب فوریه (۴۷.۶) را به طور دقیق محاسبه کرد، بدین علت که یا انتگرال را نمی‌توان به شکل بسته محاسبه کرد و یا تابع  $f(x)$  دقیقاً مشخص نیست. در هر حال از انتگرالگیری عددی استفاده می‌شود. در فصل ۷، ضمن یک متن کلی، مدخلی برای این مبحث قدیمی و دامنه‌دار آورده شده است. به منظور کنونی ما، قاعده تقریب زدن بسیار ساده‌ای زیر کفایت می‌کند

$$\int_0^{2\pi} g(x) dx \approx \frac{2\pi}{N} \sum_{n=0}^{N-1} g\left(\frac{2\pi n}{N}\right) \quad (50.6)$$

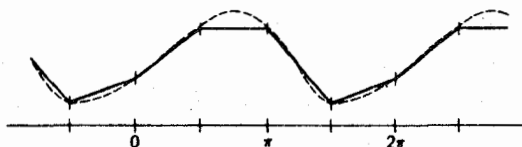
این قانون ذوزنقه‌ای مرکب (۴۹.۷) است که برای انتگرال فوق اعمال شده است، البته بسا در نظر گرفتن اینسکه تابع  $g(x)$  یک تابع  $2\pi$ -دوره‌ای است، و بنا بر این، به ویژه،  $g(2\pi) = g(0)$ . قاعده فوق را می‌توان از گذاردن یک تابع درونیاب به طور تکه‌ای خطی که در نقاط انفعال و متساوی الفاصله  $0, \pm 2\pi/N, \pm 4\pi/N, \dots$  با  $g(x)$  تطابق دارد به جای تابع  $2\pi$ -دوره‌ای  $g(x)$  زیر علامت انتگرال، به دست آورد. به شکل ۸.۶ نگاه کنید.

تقریب متناظر با  $\hat{f}(j)$  را به وسیله  $\hat{f}_N(j)$  در زیر مشخص می‌کنیم:

$$\hat{f}_N(j) = \frac{1}{N} \sum_{n=0}^{N-1} f(x_n) e^{-ix_n} \quad (51.6)$$

که در آن

$$x_n = 2\pi n/N \quad n = 0, \dots, N-1$$



شکل ۸.۶ یک تابع  $2\pi$ -دوره‌ای (نقطه‌چین) و یک تابع درونیاب به طور تکه‌ای خطی (تمام خط) در  $N=4$  نقطه در هر دوره

این نقاط  $x_n$  را نقاط نمونه‌گیری<sup>۱</sup> و اعداد  $f(x_n)$  متناظر را مقادیر نمونه<sup>۲</sup> گویند. عدد  $2\pi/N$  را بازهٔ نمونه‌گیری و عدد  $N/(2\pi)$  را فراوانی نمونه‌گیری<sup>۳</sup> گویند. سؤال این است که  $\hat{f}_N(j)$  با چه دقتی  $\hat{f}(j)$  را تقریب می‌زند؟ برای پاسخ به این سؤال، این امر را متذکر می‌شویم که توابع  $1$  و  $e^{\pm ix}$  و  $e^{\pm i2x}$ ، برخی ویژگیهای تعاملی نسبت به نوع دیگری از ضرب اسکالر یا داخلی دارند که آن را ضرب داخلی گسسته می‌نامند، یعنی

$$\langle g, h \rangle_N = \frac{1}{N} \sum_{n=0}^{N-1} \overline{h(x_n)} g(x_n) \quad (52.6)$$

روشنتر بگوییم، داریم

$$\langle e^{ikx}, e^{ijx} \rangle_N = \begin{cases} 1 & \text{اگر } k = j \pmod{N} \\ 0 & \text{اگر } k \neq j \pmod{N} \end{cases} \quad (53.6)$$

و اثبات این مطلب در حد اثبات محاسبهٔ حاصلجمع یک سری هندسی متناهی است (به تمرین ۵.۶-۸ نگاه کنید).

با توجه به مطلب فوق، ملاحظه می‌کنیم که

$$\hat{f}_N(j) = \langle f, e^{ijx} \rangle_N$$

بنابراین، با فرض اینکه سری فوریه  $\sum \hat{f}(j) e^{ijx}$  به سمت  $f(x)$  مطلقاً همگرا می‌شود (که این خود مستلزم چیزی جز وجود  $\lim_{\infty} \sum_{-\infty}^{\infty} |\hat{f}(j)|$  نیست) از رابطهٔ (۵۳.۶) نتیجه می‌شود که

$$\hat{f}_N(j) = \left\langle \sum_{k=-\infty}^{\infty} \hat{f}(k) e^{ikx}, e^{ijx} \right\rangle_N = \sum_k \hat{f}(k) \langle e^{ikx}, e^{ijx} \rangle_N$$

یا

$$\hat{f}_N(j) = \sum_{k=j \pmod{N}} \hat{f}(k) \quad (54.6)$$

یعنی: ضریب تقریبی فوریهٔ  $\hat{f}_N(j)$ ، مشکل از همهٔ ضرایب درست فوریهٔ  $\hat{f}(k)$  است که تابع مربوطه‌شان،  $e^{ikx}$ ، را نمی‌توان به وسیلهٔ ضرب داخلی (۵۲.۶) از تابع  $e^{ijx}$  تمیز داد.

این پدیده دگرناهی<sup>۴</sup> نامیده شده است. اگر  $k = j \pmod{N}$ ، آنگاه تساوی  $k = j + mN$  به ازای برخی از اعداد صحیح  $m$  برقرار است. اما در این صورت، به ازای

- 
- |                       |                  |
|-----------------------|------------------|
| 1. sampling points    | 2. sample values |
| 3. sampling frequency | 4. aliasing      |

هر  $n$  داریم

$$e^{ikx_n} = e^{i(j+mN)x_n} = e^{ijx_n} e^{imNx_n}$$

و

$$e^{imNx_n} = e^{(i\gamma\pi)mn} = 1$$

در این حال رابطه فوق بیانگر تساوی زیر است

$$e^{ikx} = e^{ijx} \quad \text{برای } x = x_n = 2\pi n/N \text{ و به ازای همه مقادیر } n$$

یعنی، دو تابع  $e^{ikx}$  و  $e^{ijx}$  در کلیه نقاط نمونه گیری که در محاسبه  $\hat{f}_N(j)$ ، یعنی در ضرب داخلی گسسته  $\langle \cdot, \cdot \rangle_N$  به کار برده شده اند، تطابق دارند. اگر تنها مقادیر تابع را در نقاط نمونه گیری  $x_n$ ، به ازای جمیع مقادیر  $n$ ، در نظر گیریم، آنگاه نمی توان بین توابع  $e^{ikx}$  و  $e^{ijx}$  تفاوتی قائل شد.

یک مثال جالب از نتیجه فوق، مثال چرخهای واگن در فیلمهاست، که بدون حرکت و یا حتی در حال چرخش در جهت مخالف با حرکت واگن به نظر می رسند. در اینجا هر  $1/20$  ثانیه از یک حرکت دوره ای نمونه برداری می شود. و در دید تماشاگر به صورت کندترین حرکت سازگار با شواهد جلوه می کند.

به همین ترتیب معمولاً (هنگام نمونه برداری در  $N$  نقطه با فاصله یکنواخت در  $[0, 2\pi]$ )، تابع  $e^{jx}$  را با  $e^{j'x}$  یکی می گیرند، البته در اینجا  $j' = j \pmod{N}$  و  $j'$  فراوانی (زاویه ای)  $|j'|$  کوچکترین مقدار ممکن است. باید توجه کرد که بدین طریق  $j'$  به طور منحصر به فردی به وسیله  $j$  و  $N$  معین می شود، جز در حالت زیر: اگر  $N$  زوج و  $j$  مضرب فردی از  $N/2$  باشد، آنگاه هم  $N/2$  و  $-N/2$  هر دو می توانند  $j'$  گرفته شوند. در مورد اخیر چنین متداول شده است که میانگین دو تابع  $e^{i(N/2)x}$  و  $e^{-i(N/2)x}$ ، یعنی تابع  $\cos(N/2)x$  را به عنوان نماینده رده خود انتخاب می کنند.

به گونه ای مشابه، گرچه رابطه (۵۱.۶)، به ازای کلیه مقادیر  $j$ ، تابع

$$\hat{f}_N(j) = \langle f, e^{ijx} \rangle_N$$

را به عنوان تقریبی برای  $\hat{f}(j)$  در نظر می گیرند ولی معمولاً تابع فوق تنها به ازای  $|j| \leq N/2$  تقریبی برای  $\hat{f}(j)$  محسوب می شود. به ویژه وقتی  $f(x)$  تابعی هموار و  $|j|$  از  $N/2$  خیلی کوچکتر باشد تعبیر خوبی دارد. در این صورت از ترکیب روابط (۴۹.۶) و (۵۴.۶)، درحالتی که  $f(x)$  دارای  $k-1$  مشتق پیوسته، و  $k$  امین مشتق آن به طور تکه ای پیوسته باشد، خواهیم داشت

$$\hat{f}_N(j) = \hat{f}(j) + O(N^{-k}) \quad (55.6)$$

در واقع، وقتی تابعی در  $N$  نقطه متساوی الفاصله در بازه  $[0, 2\pi]$  نمونه برداری می شود



اثر دگرنامی در فراوانیهای بالاتر از  $(2\pi)/(N/2)$ ، مانع مشاهده شدن پدیده‌های دوره‌ای در  $f(x)$  می‌شود. به‌طور قطعی بگوییم، اگر بخواهیم یک پدیدهٔ دوره‌ای با فراوانی  $\nu$  را مشاهده کنیم، باید حداقل با یک فراوانی به بزرگی  $2\nu$  نمونه برداری کنیم اکنون به اجمال تقریب ساز مربوط به بسجمله‌ای مثلثاتی

$$p(x) = \sum_{|j| < N/2} \hat{f}_N(j) e^{ijx} + \text{Re}[\hat{f}_N(N/2) e^{i(N/2)x}]$$

را مورد بحث قرار می‌دهیم. در اینجا تنها زمانی آخرین جمله وجود دارد که  $N/2$  یک عدد صحیح، یعنی  $N$  یک عدد زوج باشد. اما با توجه به اینکه جملهٔ اخیر را برای تکمیل بحث آوردیم (به تمرین ۵.۶-۱۱ نگاه کنید) اکنون تنها به حالتی که  $N$  فرد باشد،

$$N = 2n + 1$$

می‌پردازیم. در این حال به موجب (۵.۶)، تعداد  $N = 2n + 1$  تابع  $1$  و  $e^{\pm ix}, \dots, e^{\pm inx}$  وجود دارد که نسبت به ضرب داخلی گسسته  $\langle \cdot, \cdot \rangle_N$  متعامند، یعنی

$$\langle e^{ikx}, e^{ijx} \rangle_N = \begin{cases} 1 & k = j \\ 0 & k \neq j \end{cases} \quad k, j = -n, \dots, n \quad (5.6.6)$$

بنا به دلایل مذکور در بخش ۲.۶، قضیهٔ زیر حاصل می‌شود.

قضیهٔ ۵.۶ به ازای هر  $m \leq n$ ، بسجمله‌ای مثلثاتی مرتبهٔ  $m$

$$p_m(x) = \sum_{j=-m}^m \hat{f}_N(j) e^{ijx}$$

بهترین تقریب به وسیلهٔ بسجمله‌ایهای مثلثاتی مرتبهٔ  $m$  نسبت به نرم گسستهٔ میانگین مربعات زیر برای  $f(x)$  است.

$$\|g\|_2 = (\langle g, g \rangle_N)^{1/2} = \frac{1}{N} \left( \sum_{j=0}^{N-1} |g(j2\pi/N)|^2 \right)^{1/2}$$

به ازای  $m = n$ ، این بدان معنی است که بسجمله‌ایهای مثلثاتی مرتبهٔ  $m$

$$p_n(x) = \sum_{j=-n}^n \hat{f}_N(j) e^{ijx}$$

تابع  $f(x)$  را در نقاط نمونه‌گیری  $N$   $x_j = 2\pi j/N$ ، به ازای کلیهٔ مقادیر  $j$ ، درونیابی می‌کند.

اگر  $f(x)$  یک تابع حقیقی باشد، آنگاه می‌توانیم به موجب رابطهٔ (۲.۵.۶) بسجمله‌ایهای درونیاب را به شکل حقیقی زیر بنویسیم

$$P_n(x) = a_0/2 + \sum_{m=1}^n [a_m \cos mx + b_m \sin mx] \quad (۵۷.۶)$$

که در آن

$$a_m = 2 \operatorname{Re} \hat{f}_N(m) = \frac{2}{N} \sum_{k=0}^{N-1} f(x_k) \cos mx_k \quad (۵۸.۶ \text{ الف})$$

$$b_m = -2 \operatorname{Im} \hat{f}_N(m) = \frac{2}{N} \sum_{k=0}^{N-1} f(x_k) \sin mx_k \quad (۵۸.۶ \text{ ب})$$

□ مثال ۱۴.۶: برای تابع  $f(x) = \sin x$  يك درونیاب مثلثاتی مرتبه ۱ به دست می آوریم. در این صورت  $N = 3$ ، و کمیت‌های مربوطه عبارتند از:

$$\omega = e^{i2\pi/3} = \cos 2\pi/3 + i \sin 2\pi/3 = \frac{1}{2} + i \sqrt{\frac{3}{4}}$$

این مقادیر حائز اهمیت اند زیرا  $(\omega^m)^j = (e^{imx_j}) = \omega^{mj}$ . بعلاوه

$j$	$x_j$	$f(x_j)$	$\mathbf{w}^{(0)}$	$\mathbf{w}^{(1)}$	$\mathbf{w}^{(-1)}$
۰	۰	۰	۱	۱	۱
۱	$2\pi/3$	$\sqrt{\frac{3}{4}}$	۱	$\omega$	$\omega^{-1} = \bar{\omega}$
۲	$4\pi/3$	$-\sqrt{\frac{3}{4}}$	۱	$\omega^2$	$\omega^{-2} = \bar{\omega}^2$

اما داریم  $c_j = \hat{f}_N(j) = (\mathbf{w}^{(j)})^H f/M$ . بنا بر این  $\omega^2 = \omega^{-1} = \bar{\omega}$  و داریم

$$c_0 = \frac{1}{3} \left[ 0 + \sqrt{\frac{3}{4}} - \sqrt{\frac{3}{4}} \right] = 0$$

$$\begin{aligned} c_1 &= \frac{1}{3} \left[ 0 + \sqrt{\frac{3}{4}} \omega^{-1} - \sqrt{\frac{3}{4}} \omega^{-2} \right] = -\frac{1}{3} \sqrt{\frac{3}{4}} (\omega - \bar{\omega}) = -\frac{2}{3} i \left( \sqrt{\frac{3}{4}} \right)^2 \\ &= -i/2 \end{aligned}$$

$$c_{-1} = \frac{1}{3} \left[ 0 + \sqrt{\frac{3}{4}} \omega - \sqrt{\frac{3}{4}} \omega^2 \right] = \bar{c}_1 = i/2$$

بنابراین  $a_1 = 2\operatorname{Re}(c_1) = c_1 + c_{-1} = 0$  و  $b_1 = -2\operatorname{Im}(c_1) = i(c_1 - c_{-1}) = 1$  که نشان می‌دهند

$$p_1(x) = 0 + 0 \cdot \cos 1x + 1 \cdot \sin 1x = \sin x$$

□

همان‌گونه که انتظار داشتیم.

در این رهگذر تذکر داده می‌شود که درونیایی به وسیلهٔ بسجمله‌ایهای مثلثاتی مرتبهٔ  $n$  در هر یک از  $2n+1$  نقطهٔ متمایز در  $[0, 2\pi]$  به‌طور منحصر به‌فردی انجام می‌گیرد. در مورد بسجمله‌ای درونیاب حاصلهٔ  $p_n(x)$  برای  $f(x)$  می‌توان نشان داد که

$$\|f - p_n\|_\infty \leq \operatorname{const} \operatorname{dist}_\infty(f, \mathring{\pi}_n) \quad (59.6)$$

در اینجا نرم ماکسیمم در بازهٔ  $[0, 2\pi]$  گرفته شده‌است

$$\|g\|_\infty = \max_{0 \leq x \leq 2\pi} |g(x)|$$

و

$$\operatorname{dist}_\infty(f, \mathring{\pi}_n) = \min_{p \in \mathring{\pi}_n} \|f - p\|_\infty$$

که در آن نماد  $p \in \mathring{\pi}_n$  به‌جای ملخص عبارت « $p(x)$  یک بسجمله‌ای مثلثاتی از مرتبهٔ  $n$  است» به‌کار برده شده‌است. می‌توان نشان داد که رابطهٔ (59.6) به‌نامساوی متناظر (17.6) برای درونیایی به‌وسیلهٔ بسجمله‌ای بسیار شبیه است. به‌ویژه اینکه عدد ثابت  $(\operatorname{const})$  به‌نقاط درونیایی بستگی دارد. در این عبارات نقاط درونیایی با فواصل یکنواخت که در اینجا منحصرأ استفاده شده‌اند، بهینه هستند، بدین تعبیر که عدد ثابت  $(\operatorname{const})$  رابطهٔ (59.6) را تا حد امکان کوچک می‌کند، به‌کتاب دبور و پینکوس<sup>۱</sup> [۳۹] نگاه کنید. بهترین مقدار این عدد ثابت توسط الیش و تسلر<sup>۲</sup> [۳۸] به‌صورت زیر محاسبه شده‌است

$(\operatorname{const})_{\text{uniform}}$  (ثابت یکنواخت)

$$= 1 + \frac{1}{N} \left\{ 1 + 2 \sum_{k=1}^n \frac{1}{\sin \left( \frac{(2k-1)\pi}{(2n+1)2} \right)} \right\} \approx \frac{2}{\pi} \ln N + c \quad (60.6)$$

بنابراین برای مقادیری از  $n$  که جنبهٔ عملی داشته باشند، درونیایی در نقاط با فواصل یکنواخت، تقریبی به‌دست می‌دهد که از بهترین تقریب یکنواخت احتمالی از  $\mathring{\pi}_n$  خیلی بدتر نیست. بنا بر این در صورتی که بسجمله‌ای درونیاب به‌آسانی به‌دست آید، معمولاً نیازی به انجام روند پیچیده برای به‌دست آوردن بهترین تقریب یکنواخت نیست. این مسئلهٔ اخیر در بخش

بعد مورد بحث قرار خواهد گرفت.

در این مورد، اشاره می‌شود که در صورتی که  $f(x)$  دارای  $k$  مشتق بسا مشتق  $k$ امی به طور تکه‌ای پیوسته باشد، رابطه (۴۹.۶) ایجاب می‌کند که

$$\text{dist}_{\infty}(f, \hat{r}_n) = O(n^{-k}) \quad (۶۱.۶)$$

### تمرین

۱-۵.۶ سری فوریه برای تابع  $2\pi$ -دوره‌ای  $f(x)$  را که به وسیله  $x = f(x)$  داده شده است در بازه  $[0, 2\pi)$  محاسبه کنید.

۲-۵.۶ تحقیق کنید که تابع  $2\pi$ -دوره‌ای  $f(x)$  که مقادیرش در بازه  $[0, 2\pi)$  به شرح زیر داده شده است

$$f(x) = \begin{cases} (x/\pi)^2 - x/\pi & 0 \leq x \leq \pi \\ (x-\pi)/\pi - ((x-\pi)/\pi)^2 & \pi \leq x \leq 2\pi \end{cases}$$

پیوسته بوده و مشتق اول پیوسته نیز دارد (مانند یک تابع  $2\pi$ -دوره‌ای)، اما در مشتق دوم جهشهایی\* دارد. سپس طیف  $f(x)$  را تشکیل و نشان دهید که وقتی  $j \rightarrow \infty$ ، مانند  $j^{-3}$ ، (ونه سریعتر از آن) به سمت صفر میل می‌کند.

۳-۵.۶ سری فوریه را که در تمرین ۲-۵.۶ به دست آمد بر حسب  $\sin$  و  $\cos$  بنویسید. چرا انتظار دارید که کلیه  $a_j$  برابر با صفر باشند؟

۴-۵.۶ اگر  $f(x)$  یک تابع  $2\pi$ -دوره‌ای باشد، آنگاه به ازای هر عدد صحیح  $m$ ، تابع  $gm(x) = f(mx)$  نیز  $2\pi$ -دوره‌ای است. چه رابطه‌ای بین  $\hat{f}(j)$  و  $\hat{g}_m(j)$  وجود دارد؟

۵-۵.۶ اگر تابع  $f(x)$  یک تابع  $2\pi$ -دوره‌ای باشد، آنگاه به ازای هر عدد  $\alpha$ ، تابع  $g_{\alpha}(x) = f(x - \alpha)$  نیز  $2\pi$ -دوره‌ای است. چه رابطه‌ای بین  $\hat{f}(j)$  و  $\hat{g}_{\alpha}(j)$  وجود دارد؟

۶-۵.۶ فرض کنید که تابع  $f(x)$  تابعی بسیار هموار با دوره (تناوب)  $\pi$  باشد. اما در تبدیل آن به یک تابع  $2\pi$ -دوره‌ای  $g(x) = f(\pi x / (2\pi))$  اشتباهاً  $\pi'$  به جای  $\pi$  به ازای  $\pi' \neq \pi$  به کار گرفته شده است. اثرات احتمالی این اشتباه روی ضرایب فوریه  $\hat{g}(j)$  چه هستند؟

۷-۵.۶ ثابت کنید که اگر  $f(x)$  یک تابع حقیقی باشد، آنگاه رابطه

$$\hat{f}(j)e^{ijx} + \hat{f}(-j)e^{-ijx} = 2|\hat{f}(j)|\cos(\theta_j + jx)$$

\* توضیح مترجم؛ جهش به معنی ناپیوستگی مقطعی است.

به ازای يك اختلاف فاز مناسب  $\theta$  برقرار است. (داهنمایی: از این واقعیت استفاده کنید که هر عدد همتافت  $z$  را می توان به ازای يك  $\theta$  مناسب به شکل قطبی  $|z|e^{i\theta}$  نوشت.)

۸-۵۰۶ قضیهٔ (۵۳.۶) را اثبات کنید (داهنمایی: بیاد آورید که چگونه مجموع يك سری هندسی به دست می آید.)

۹-۵۰۶ رابطهٔ (۶.۵) را اثبات کنید.

۱۰-۵۰۶ فرمولهای جمع برای  $\sin(\alpha+\beta)$  و  $\cos(\alpha+\beta)$  را از فرمول اولیبر (۴۲.۶) و قانون جمع نماها،  $e^{A+B} = e^A e^B$ ، به دست آورید.

۱۱-۵۰۶ ثابت کنید که اگر  $N$  زوج و  $f(x)$  حقیقی باشد، آنگاه تابع

$$p(x) = \sum_{|j| < N/2} \hat{f}_N(j) e^{ijx} + \text{Re}[\hat{f}(N/2) e^{i(N/2)x}]$$

$f(x)$  را در نقاط نمونه گیری  $x_k$  به ازای کلیهٔ مقادیر  $k$  درونیایی می کند.

۱۲-۵۰۶ چگونه می توان بسجمله ای درونیایی مثلثاتی برای  $f(x)$  را در نقاط

$$\alpha + k2\pi/N, k = 0, \dots, N-1,$$

$\alpha$  عدد مثبتی کمتر از  $2\pi/N$ ، به دست آورد؟

### ۶.۶۰ تبدیلهای سریع فوریه<sup>۱</sup>

در آنالیز همساز گسسته، یا آنالیز فوریه، برای تجزیهٔ حرکت  $2\pi$ -دوره ای (که با  $f(x)$  بیان شده است) به همسازهای ساده، اعداد

$$c_j = \hat{f}_N(j) = \frac{1}{N} \sum_{k=0}^{N-1} f(x_k) e^{-ijx_k} \quad (۶۲.۶)$$

$$\text{با شرط } x_k = 2\pi k/N, \text{ به ازای جميع مقادیر } k, \quad (۶۳.۶)$$

را محاسبه می کنند. همان طوری که در بخش ۵.۶ دیدیم اگر  $f(x)$  دارای مشتق  $k$ ام پیوستهٔ تکه ای باشد، آنگاه

$$f(x) \approx \sum_{j=-\infty}^{\infty} \hat{f}(j) e^{ijx}$$

با فرض

$$\hat{f}(j) = \hat{f}_N(j) + \mathcal{O}(N^{-k})$$

ما به آن حالتی که فراوانیهای در  $f(x)$  وجود داشته باشند، و به توان این فراوانیها علاقه مندیم.

## 1. Fast fourier transform

اما به علت اثر دگر نامی، به ازای  $|z| > N/2$ ،  $\hat{f}_N(j)$  به عنوان تقریبی برای  $\hat{f}(j)$  مفید فایده‌ای نیست و معمولاً تقریب خوبی تنها برای موقعی است که مقادیر  $|z|$  خیلی کمتر از  $N/2$  باشد. با توجه بدین امر باید محاسبه  $\hat{f}_N(j)$  به ازای مقادیر «بزرگ»  $N$ ، و لذا مسئله مهم چگونگی محاسبه  $\hat{f}_N(j)$  به نحو کارا، مورد توجه قرار گیرد.

واضح است که محاسبه هر تابع خاص  $\hat{f}_N(j)$  مستلزم تعداد  $\mathcal{O}(N)$  عمل ضرب و جمع است. بنا بر این محاسبه مستقیم  $N$  عدد از این نوع اعداد (مثلاً، اعداد  $\hat{f}_N(j)$  به ازای  $|z| \leq N/2$ ) نیاز به  $\mathcal{O}(N^2)$  عمل دارد. از این رو برای ۱۰۰۰ نقطه نمونه، به میلیونها عمل نیاز خواهد بود و تا این اواخر، این امر مانع اصلی بر سر راه آنالیز گسسته فوریه بود.

وقتی کاملاً معلوم شد که محاسبه همزمان  $N$  شماره متوالی از  $\hat{f}_N(j)$ ها به دلیل روابط قوی بین این اعداد تنها به  $\mathcal{O}(N \log N)$  عمل حساب نیاز دارد، وضعیت به نحو مؤثری دگرگونی پیدا کرد. کلید حل این مشکل تبدیل سریع فوریه (ت س ف) یا FFT بود که این محاسبات را برای  $N \leq 1000$  به یک امر عادی مبدل کرد و حتی این امکان را به وجود آورد که  $N$ های برابر با دهها هزار نیز به کار گرفته شوند.

در اینجا ما می‌توانیم تنها به آن طرح اساسی که به چنین افزایش مؤثری در کارایی منجر شده اشاره کنیم. تازه‌ترین مطالب مذکور در این باب را می‌توان در مقاله ۱۹۷۸ س. وینگراد<sup>۲</sup> [۳۶] به دست آورد. به ویژه کارهایی که قبل و بعد از انتشار مقاله نهادی و بی‌سابقه کولی و تیوکی<sup>۳</sup> [۳۷] انجام گرفته، از مدت‌ها قبل روشن کرده‌اند که FFT (ت س ف)های زیادی وجود دارند و برای به دست آوردن کارایی بیشتر، لازم (و مقرون به صرفه) است که به ازای مقادیر مختلفی از  $N$  که می‌خواهند به کار برند، برنامه‌های مختلفی بنویسند.

برای تحلیل محاسبات مربوط به اعداد  $\hat{f}_N(j)$  به ازای  $|z| \leq N/2$  از روی اعداد  $f(x_0), \dots, f(x_{N-1})$ ، آسانتر است که تبدیل گسسته فوریه  $F_N^z$  معرفی شود. این تبدیل  $N$ -بردار

$$\mathbf{z} = [z_1, z_2, \dots, z_N]^T$$

را به  $N$ -بردار

$$F_N \mathbf{z} = \hat{\mathbf{z}}$$

منتقل می‌کند که با

$$\hat{z}_j = \sum_{n=0}^{N-1} z_n \omega_N^{(j-1)(n-1)} \quad j = 1, \dots, N \quad (64.6)$$

1. fast Fourier transform.

2. S. Winograd

3. Cooley and Tukey

4. discrete Fourier transform

معین می‌شود و در آن  $\omega_N$ ،  $n$  امین ریشهٔ واحد است

$$\omega_N = e^{-i2\pi/N}$$

ارتباط بین محاسبهٔ  $\hat{f}_N(j)$  و این تبدیل گسستهٔ فوریه به شرح زیر است. اگر  $N$ -بردار خاص زیر را در نظر بگیریم

$$\mathbf{z} = [f(x_0)f(x_1) \dots f(x_{N-1})]^T$$

آنگاه

$$\hat{f}_N(j) = \frac{1}{N} \begin{cases} \hat{z}_{j+1} & j = 0, 1, 2, \dots \\ \hat{z}_{j+1+N} & j = -1, -2, \dots \end{cases} \quad |j| \leq N/2 \quad (۶۵.۶)$$

بنا بر این کافی است که سعی خود را بر محاسبهٔ کارای تبدیل گسستهٔ فوریه متمرکز سازیم. این بحث را با ملاحظهٔ این امر آغاز می‌کنیم که  $z_j$  به صورتی که با رابطهٔ (۶۴.۶) داده شده، یک بسجمله‌ای است از درجهٔ کمتر از  $N$  بر حسب کمیت  $\omega_N^{-1}$ ، بنا بر این می‌توان آن را طی  $N$  عمل با ضرب تودرتو محاسبه کرد. در اینجا یک عمل ضرب بعلاوهٔ یک عمل جمع را به عنوان یک عمل<sup>۱</sup> تلقی می‌کنیم. بنا بر این ارزیابی سری است رابطهٔ (۶۴.۶)، به ازای جمیع مقادیر  $j$ ، به تعداد  $N^2$  عمل نیازمند است.

گسترده‌ترین طرحی که برای یک (ت س ف) داریم، طرحی است که کولی و تیوکی آن را به زبان ساده در آورده‌اند و هنگامی قابل استفاده است که  $N$  مساوی با حاصلضرب اعداد صحیح باشد. اکنون این طرح را اول برای حالت زیر مورد بحث قرار می‌دهیم

$$N = P \cdot Q$$

$N$ -بردار  $\mathbf{z}$  را که به شیوهٔ فورتون در یک آرایهٔ یک بعدی ذخیره شده در نظر می‌گیریم. در این صورت می‌توانیم این آرایه را به شیوهٔ فورتون به صورت یک آرایهٔ دو بعدی  $Z$  با ابعاد  $(P, Q)$  نشان دهیم. این بدان معناست که داریم

$$Z(p, q) = z_{p+P(q-1)}$$

همچنین حاصلجمع  $\hat{z}_v$  را به حاصلجمع مضاعف:  $\sum z_n \omega_N^{(p-1)(n-1)}$  را به حاصلجمع مضاعف:

$$\begin{aligned} \hat{z}_v &= \sum_{p=1}^P \sum_{q=1}^Q Z(p, q) \omega_N^{(p-1)(p-1+P(q-1))} \\ &= \sum_{p=1}^P \left[ \sum_{q=1}^Q Z(p, q) \omega_Q^{(p-1)(q-1)} \right] \omega_N^{(p-1)(p-1)} \end{aligned}$$

تجزیه می‌کنیم. در اینجا از تساوی  $\omega_N^P = \omega_Q$  استفاده کرده‌ایم، این تساوی این واقعیت بحرانی را روشن می‌سازد که جمع داخلی در آخرین قسمت سمت راست در  $v$  باید  $Q$ -دوره‌ای باشد، یعنی از قراردادن  $Q+v$  به جای  $Q$ ، مقدار آن به علت واقعیت  $\omega_Q^Q = 1$  تغییر نمی‌کند. معنی این عبارت این است که تنها به محاسبه این حاصلجمع به ازای  $Q, v = 1, \dots, Q$  (و هر  $p$ ) نیاز داریم. بنا بر این به ازای هر  $p = 1, \dots, P$  از  $Q$ -بردار  $Z(p, \cdot)$  برداری را محاسبه می‌کنیم که مؤلفه‌های آن اعداد

$$\sum_{q=1}^Q Z(p, q) \omega_Q^{(v-1)(q-1)} \quad v = 1, \dots, Q$$

باشند. یعنی، تبدیل گسسته فوریه مربوط به  $Q$ -بردارهای  $Z(p, \cdot)$  را به ازای  $p = 1, \dots, P$ ، کلاً با صرف تعداد  $Q \cdot P = N$  عمل، محاسبه می‌کنیم. اما، می‌توانستیم تبدیل  $Z(p, \cdot)$  را بر روی  $Z(p, \cdot)$  ذخیره کنیم. اما برای جلوگیری از توسعه‌های بیشتر، تبدیل  $Z(p, \cdot)$  را در  $Z_1(\cdot, p)$  ذخیره می‌کنیم، که  $Z_1$  یک آرایه دوبعدی است که اندازه آن به جای  $(P, Q)$  برابر  $(Q, P)$  است. با توجه به این امر، محاسبه  $\hat{z}_v$  به محاسبه حاصلجمع

$$\hat{z}_v = \sum_{p=1}^P Z_1(v_Q, p) \omega_N^{(v-1)(p-1)} \quad v = 1, \dots, N$$

تبدیل می‌شود. در اینجا نماد  $v_Q$  معرف عدد صحیحی است بین  $1$  و  $Q$  که به ازای آن  $v - v_Q$  بر  $Q$  قابل قسمت است. بنابراین به ازای برخی از اعداد صحیح  $v'$  بین  $1$  و  $P$  داریم

$$v = v_Q + Q(v' - 1)$$

زیرا اگر بردار  $\hat{z}$  را به شیوه فورتین به یک آرایه دوبعدی  $Z_0$  با اندازه  $(Q, P)$  نشان دهیم، داریم

$$\hat{z}_v = Z_0(v_Q, v') \quad (66.6)$$

با توجه به این رابطه باید مقدار

$$Z_0(v_Q, v') = \sum_{p=1}^P Z_1(v_Q, p) \omega_N^{(p_Q-1+Q(v'-1))(p-1)} \quad v' \text{ و } v_Q \text{ مقادیر}$$

را محاسبه کنیم. در اینجا قسمت سمت راست، یک بسجمله‌ای است از درجه کمتر از  $P$  بر حسب کمیت  $\omega_N^{-1} = \omega_N^{p_Q-1+Q(v'-1)}$ . این کمیت را می‌توان مرحله به مرحله به ترتیب ساده محاسباتی زیر تولید کرد.



$$\begin{array}{l}
 x := 1 \\
 \text{for } v' = 1, \dots, P, \text{ do:} \\
 \quad \text{for } v_Q = 1, \dots, Q, \text{ do:} \\
 \quad \quad Z_0(v_Q, v') := \sum_{p=1}^P Z_1(v_Q, p) x^{p-1} \\
 x := x \cdot \omega_N
 \end{array} \quad (۶۷.۶)$$

البته در داخلی ترین حلقه، حاصلجمع را باید به وسیله ضرب تودرتو محاسبه نمود. در این صورت تعداد کل عمل مصروف در این مرحله برابر است با  $P^2 = NP$ .  $Q$  (در صورتی که از تعداد  $N$  عمل ضرب لازم برای تولید  $P$ های مختلف لازم صرف نظر کنیم). بدین ترتیب در  $Z_0$ ، تبدیل گسسته فوریه  $\hat{Z}$  مربوط به  $Z$  تنها با صرف تعداد  $N(P+Q)$  عمل، در مقایسه با  $N^2$  عمل که با روش ساده لازم بود، به دست آمده است. حال اگر  $N$ ، مساوی با حاصلضرب سه یا چند عدد صحیح بزرگتر از یک، مثلا

$$N = P_1 \dots P_m$$

باشد، آنگاه با استفاده از مرحله دوم (۶۷.۶) به طریق اندکی پیچیده تر، می توانیم تبدیل گسسته فوریه  $Z$  راحتی با تعداد عمل کمتری محاسبه کنیم.

برای توضیح مطلب فوق، به اندکی علامتگذاری نیاز داریم تا بتوانیم نشان دهیم که چگونه یک آرایه یک بعدی را که به شیوه فورترن نشان داده شده، می توان به صورت یک آرایه دو یا سه بعدی نمایش داد. اگر  $Z$  آرایه ای یک بعدی به طول  $N$  باشد، آنگاه آرایه دو بعدی هم ارز آن به ابعاد  $(A, N/A)$  را با  $Z^A$  و آرایه سه بعدی هم ارز آن به ابعاد  $(A, B, N/(AB))$  را با  $Z^{A,B}$  نشان می دهیم. بدین ترتیب داریم

$$\begin{aligned}
 Z^{A,B}(a, b, c) &= Z^A(a, b + B(c-1)) = Z^{AB}(a + A(b-1), c) \\
 &= Z(a + A(b-1 + B(c-1)))
 \end{aligned}$$

حال گیریم، مانند قبل،  $Z$  یک آرایه یک بعدی مشتمل بر  $z$  و به ازای  $k=0, \dots, m$   $Z_k$  یک آرایه یک بعدی مشتمل بر تبدیل گسسته فوریه قسمتهای  $Z$  به شرح زیر باشد:

$$Z_k^A(\dots, c) = F_A Z^{BP}(c, \dots), \quad c = 1, \dots, BP \quad (۶۸.۶)$$

که در آن

$$B := B_k := P_1 \dots P_{k-1}$$

$$P := P_k \quad (الف \ ۶۸.۶)$$

$$A := A_k = P_m \dots P_{k+1}$$

باید توجه داشت که  $Z$  نقش  $Z_m$  را به خوبی ایفا می کند و  $Z_0$  شامل  $F_N Z$  است. برای

رسیدن از  $Z_k$  به  $Z_{k-1}$ ، از شکل اندکی بسط داده شده (۶۷.۶) که در آن  $A$  و  $P$ ،  $B$  و  $P$  با رابطه (۶۸.۶) مشخص می شوند، استفاده می کنیم

$$\begin{array}{l}
 x := 1 \\
 \text{for } p = 1, \dots, P, \text{ do:} \\
 \quad \text{for } a = 1, \dots, A, \text{ do:} \\
 \quad \quad \text{for } b = 1, \dots, B, \text{ do:} \\
 \quad \quad \quad Z_{k-1}^{A,P}(a, p, b) := \sum_{\pi=1}^P Z_k^{A,B}(a, b, \pi) \cdot x^{\pi-1} \\
 \quad \quad \quad x := x \cdot \omega_{AP}
 \end{array} \quad (۶۹.۶)$$

در حقیقت الگوریتم فوق مقادیر زیر را تولید می کند

$$Z_{k-1}^{A,P}(a, p, b) = \sum_{\pi=1}^P Z_k^{A,B}(a, b, \pi) \omega_{AP}^{[a-1+A(p-1)](\pi-1)}$$

از طرف دیگر رابطه (۶۸.۶) مستلزم رابطه

$$Z_k^{A,B}(\cdot, b, \pi) = F_A Z^{B,P}(b, \pi, \cdot) = \sum_{\alpha=1}^A Z^{B,P}(b, \pi, \alpha) \omega_A^{(\cdot-1)(\alpha-1)}$$

است. بنابراین

$$Z_{k-1}^{A,P}(a, p, b) = \sum_{\pi=1}^P \sum_{\alpha=1}^A Z^{B,P}(b, \pi, \alpha) \omega_{AP}^{P(a-1)(\alpha-1) + [a-1+A(p-1)](\pi-1)}$$

اما اکنون به موجب تساوی  $\omega_{AP}^{AP} = 1$  می توان به نمای طرف راست هر ضرب صحیحی از  $AP$  را افزود و این امر به نتیجه زیر منجر می شود

$$Z_{k-1}^{A,P}(a, p, b) = \sum_{\pi=1}^P \sum_{\alpha=1}^A Z^{B,P}(b, \pi, \alpha) \omega_{AP}^{[a-1+A(p-1)][\pi-1+P(\alpha-1)]}$$

و بدین ترتیب ثابت می شود که  $Z_{k-1}$  همان طوری که توسط (۶۹.۶) تولید شده، در رابطه (۶۸.۶) (که در آن  $k-1$  جایگزین  $k$  شده است) صدق می کند.

به ویژه،  $Z_0$  شامل تبدیل گسسته فوری  $Z$  است. با  $m$  بار استفاده از الگوریتم (۶۹.۶) با شروع از  $Z_m = Z$  به  $Z_0$  خواهد رسید.

زیر برنامه فورترن زیر الگوریتم مذکور را پیاده می کند.

```

SUBROUTINE FFT ( Z1, Z2, N, INZEE )
CONSTRUCTS THE DISCRETE FOURIER TRANSFORM OF Z1 (OR Z2) IN THE COOLEY-
C TUKEY WAY, BUT WITH A TWIST.
      INTEGER INZEE,N, AFTER,BEFORE,NEXT,NEXTMX,NOW,PPIME(12)
      COMPLEX Z1(N),Z2(N)
C***** I N P U T *****
C Z1, Z2 COMPLEX N-VECTORS
C N LENGTH OF Z1 AND Z2
C INZEE INTEGER INDICATING WHETHER Z1 OR Z2 IS TO BE TRANSFORMED
C = 1 , TRANSFORM Z1
C = 2 , TRANSFORM Z2

```

```

C***** W O R K A R E A S *****
C 21, 22 ARE BOTH USED AS WORKARRAYS

C***** O U T P U T *****
C 21 OR 22 CONTAINS THE DESIRED TRANSFORM (IN THE CORRECT ORDER)
C INZEE INTEGER INDICATING WHETHER Z1 OR Z2 CONTAINS THE TRANSFORM,
C   = 1, TRANSFORM IS IN Z1
C   = 2, TRANSFORM IS IN Z2
C***** M E T H O D *****
C THE INTEGER N IS DIVIDED INTO ITS PRIME FACTORS (UP TO A POINT).
C FOR EACH SUCH FACTOR P, THE P-TRANSFORM OF APPROPRIATE P-SUBVECTORS
C OF Z1 (OR Z2) IS CALCULATED IN F F T S T P AND STORED IN A SUIT-
C ABLE WAY IN Z2 (OR Z1). SEE TEXT FOR DETAILS.
C
DATA NEXTMX,PRIME / 12, 2,3,5,7,11,13,17,19,23,29,31,37 /
AFTER = 1
BEFORE = N
NEXT = 1

C
10 IF ((BEFORE/PRIME(NEXT))*PRIME(NEXT) .LT. BEFORE) THEN
    NEXT = NEXT + 1
    IF (NEXT .LE. NEXTMX) THEN
        GO TO 10
    ELSE
        NOW = BEFORE
        BEFORE = 1
    END IF
ELSE
    NOW = PRIME(NEXT)
    BEFORE = BEFORE/PRIME(NEXT)
END IF

C
IF (INZEE .EQ. 1) THEN
    CALL FFTSTP( Z1, AFTER, NOW, BEFORE, Z2 )
ELSE
    CALL FFTSTP( Z2, AFTER, NOW, BEFORE, Z1 )
END IF
INZEE = 3 - INZEE
IF (BEFORE .EQ. 1)
    RETURN
AFTER = AFTER*NOW
GO TO 10
END

SUBROUTINE FFTSTP ( ZIN, AFTER, NOW, BEFORE, ZOUT )
CALLED IN F F T .
CARRIES OUT ONE STEP OF THE DISCRETE FAST FOURIER TRANSFORM.
INTEGER AFTER,BEFORE,NOW, IA,IB,IN,J
REAL ANGLE,RATIO,TWOPI
COMPLEX ZIN(AFTER,BEFORE,NOW),ZOUT(AFTER,NOW,BEFORE), ARG,OMEGA,
* VALUE
DATA TWOPI / 6.2831 85307 17958 64769 /
ANGLE = TWOPI/FLOAT(NOW*AFTER)
OMEGA = CMPLX(COS(ANGLE),-SIN(ANGLE))
ARG = CMPLX(1.,0.)
DO 100 J=1,NOW
    DO 90 IA=1,AFTER
        DO 80 IB=1,BEFORE
            VALUE = ZIN(IA,IB,NOW)
            DO 70 IN=NOW-1,1,-1
                VALUE = VALUE*ARG + ZIN(IA,IB,IN)
            70 ZOUT(IA,J,IB) = VALUE
            80 ARG = ARG*OMEGA
        90
    100 CONTINUE
RETURN
END

```

اگر  $N$  مساوی با حاصلضرب  $m$  عدد صحیح باشد،

$$N = P_1 P_2 \dots P_m$$

آنگاه برنامه‌ای مانند برنامه فوق امکان محاسبه تبدیل  $F_N Z = \hat{Z}$  را با تعداد عملیاتی برابر با

$$W = N(P_1 + P_2 + \dots + P_m)$$

(به جای  $N^2$ ) فراهم می‌سازد. از آنجا که به ازای اعداد صحیح  $Q$  و  $R$  بزرگتر از ۱، بجز  $Q = R = 2$ ، تا مساوی  $Q + R < QR$  برقرار است، این عدد  $W$  زمانی به کمترین مقدار می‌رسد که از هر عامل  $N$  بجز آن عواملی از ۲ که از ترکیب آنها ۴ حاصل می‌شود، عملاً استفاده شود. بعلاوه

$$W/N = P_1 + \dots + P_m \quad \text{و} \quad \log N = \log P_1 + \dots + \log P_m$$

بنابراین

$$\frac{W}{N \log N} = \frac{P_1 + \dots + P_m}{\log P_1 + \dots + \log P_m} = \left( \sum_{j=1}^m \frac{P_j}{\log P_j} \log P_j \right) / \sum_{j=1}^m \log P_j$$

این رابطه نشان می‌دهد که عدد  $W/(N \log N)$  یک میانگین وزین اعداد  $p_j / \log p_j$ ،  $j = 1, \dots, m$  است. به آسانی دیده می‌شود که  $P / \log P$  به عنوان تابعی از عدد صحیح  $P$  در  $P = 3$  دارای مقدار حداقل  $1.389 \dots$  در  $P = 2$  و  $P = 4$  مقدار ۲ و در  $P = 10$  تنها برابر با  $3.01 \dots$  است. بنابراین نامساوی

$$1.389 N \log_2 N \leq W$$

برقرار است در حالی که حتی برای عوامل  $P_j$  به بزرگی ۱۰، مقدار  $W$  بزرگتر از  $N \log_2 N$  ۳۰۲ نیست.

در مواردی که بردار داده شده  $Z$  حقیقی باشد، صرفه‌جویی بیشتری امکانپذیر است.

زیرا تساوی

$$\hat{z}_{N+1-j} = \hat{z}_{j+1} \quad (70.6)$$

برقرار است. به روابط (۶۴.۶) و (۶۵.۶) نگاه کنید.

زمانی که  $N$  عددی اول و یا  $N$  مساوی حاصلضرب اعداد صحیحی باشد که دو به دو نسبت به هم اول هستند، (ت س ف) های دیگری هم وجود دارند. به مقاله وینگراد [۳۶] نگاه کنید.

## نهمین

۶۰۶-۱ مستقیماً با توجه به تعریف (۶۴.۶) ثابت کنید که (۷۰.۶) برای حالتی که  $Z$  حقیقی باشد صادق است.

۶۰۶-۲ با استفاده از برنامه (ت س ف) (مثلاً به ازای  $N = 81$ ) جوابهای خود را برای تمرینهای ۱-۵.۰۶ و ۲-۵.۰۶ امتحان کنید. [این کار شما را وادار خواهد ساخت که به تمام جزئیات موجود در رابطه (۶۵.۶) دقیقاً توجه کنید!]

۳-۶.۶ با استفاده از برنامه (ت س ف) ضرایب فوریه  $f(j)$  را برای توابع زیر (به طور تقریبی) محاسبه کنید

$$\text{الف. } f(x) = \sin 3x \quad \text{ب. } f(x) = \sin(\pi x)$$

برای مثال  $N$  را برابر با ۸۱ یا ۳۲۴ یا هر عدد دیگری که می‌خواهید بگیرید. چرا با افزایش  $|j|$  ضرایب فوریه برای  $f(x) = \sin(\pi x)$  به سرعت کاهش نمی‌یابد؟

۴-۶.۶ برنامه (ت س ف) تیلر<sup>۱</sup> برای حالت خاص  $N = ۳۰۴$  را درز بگیرید و تا آنجا که می‌توانید در محاسبات و مقدار حافظه صرفه‌جویی کنید.

۵-۶.۶ با افزودن عبارت‌های مخصوص به‌جای دامنهٔ حلقهٔ DO روی IB، مثلاً در موردی که  $N = ۲, ۳$  یا  $N = ۲$ ، زیر برنامهٔ FFTSTP را بهتر کنید.

۶-۶.۶ در باب استفاده از (ت س ف)، برای محاسبهٔ حاصلجمع مثلثاتی

$$i(x) = \sum_{j=-n}^n a(j)e^{ijx}$$

در نقاط  $x_j = \alpha + 2\pi j / (2n + 1)$ ،  $j = 0, \dots, n$ ، به‌ازای یک مقدار ثابت  $\alpha$  در بازهٔ  $[0, 2\pi / (2n + 1)]$  بحث کنید.

۷-۶.۶ با استفاده از (ت س ف) بسجمله‌ای درونیاب مثلثاتی را در  $N = 2n + 1$  نقطهٔ  $x_j = 2\pi j / N$ ،  $j = 0, \dots, N - 1$  برای موج‌مربعی  $f(x) = \text{signum}(\sin x)$  و با استفاده از  $N = ۳۵$  بسازید. سپس با ردیگر از (ت س ف) برای محاسبهٔ بسجمله‌ای درونیاب در ۱۰۵ نقطهٔ  $j = 0, \dots, ۱۰۴$  و  $y_j = 2\pi j / ۱۰۵$  استفاده کنید. (دانه‌یابی: از تمرین ۶-۶.۶ استفاده کنید).

۸-۶.۶ با استفاده از (ت س ف) تقریبی بسازید برای طیف یک تابع  $f(x)$  که مقادیرش در نقاط  $x_j = 2\pi j / N$ ،  $j = 0, \dots, N - 1$  و با مثلاً  $N = ۱۲۸$ ، از یک مولد عددی (شبه) تصادفی<sup>۲</sup> که اعدادی با توزیع یکنواخت بین ۱۰ و ۱۰۰ تولید می‌کند، به‌دست آمده باشند. نتیجه را با طیف تابعی که در تمرین ۵-۶.۶ بررسی کردیم مقایسه کنید.

۹-۶.۶ با استفاده از تمرین ۸-۶.۶ در چگونگی استفاده از (ت س ف) برای بازیابی مقادیر  $f(x_j)$ ،  $j = 0, \dots, N - 1$ ، از یک تابع «هموار»  $2\pi$ -دوره‌ای  $f(x)$  از داده‌های مفروض  $f(x_j) + \varepsilon_j$ ، به‌ازای کلیهٔ مقادیر  $j$ ، که در آن  $\varepsilon_j$  نوفه‌ای است با توزیع یکنواخت، بحث کنید.

۱۰-۶.۶ نشان دهید که  $F_N^{-1}Z = 1/N \overline{F_N Z}$  (یعنی  $N$ -بردار  $Z$  را از تبدیل مجرد فوریه

آن  $\hat{z} = FNz$  به دست آورید، بدین ترتیب که (الف) با تبدیل کلیه درایه‌های  $\hat{z}$  به فرد و زوجهای همتافتشان و سپس (ب) با تشکیل تبدیل گسسته فوریه بردار حاصله  $\hat{z}$  و آنگاه (پ) با تقسیم هر درایه بردار حاصل بر  $N$ .

۶-۱۱ چگونگی استفاده از (ت س ف) را برای تشکیل بسجمله‌ای درونیایی از درجه  $n$  تا بیشتر از  $n$  در نقاط چبیشف (۱۸.۶) برای داده‌های مفروض شرح دهید (دانهمایی: بسجمله‌ای درونیاب را از یک ترکیب خطی از  $n+1$  بسجمله‌ای چبیشف،  $T_0, \dots, T_n$ ، با استفاده از (۶.۹) بسازید. ارزیابی بعدی، البته، از راه تابع چبیشف صورت خواهد گرفت).

### ۷.۶ تقریب به وسیله بسجمله‌ای - تکه‌ای

یک مثال ساده و آشنا از تقریب زدن به وسیله بسجمله‌ای - تکه‌ای، درونیایی خطی درجدولی است از مقادیر  $f(x_i)$ ،  $i = 1, \dots, N+1$ ، که در آن

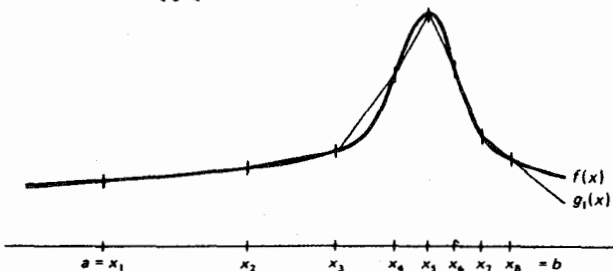
$$a = x_1 < x_2 < \dots < x_{N+1} = b$$

در اینجا  $f(x)$  در یک نقطه  $\bar{x}$ ، با تعیین بازه  $[x_k, x_{k+1}]$  که شامل  $\bar{x}$  است و گرفتن

$$p_1(\bar{x}) = f(x_k) + f[x_k, x_{k+1}](\bar{x} - x_k)$$

به عنوان تقریبی برای  $f(\bar{x})$ ، تقریب زده می‌شود. در نتیجه،  $f(x)$  توسط تابع «خط شکسته» یا تابع خطی - تکه‌ای  $g_1(x)$  روی  $[a, b]$  (به شکل ۹.۶ نگاه کنید) با نقاط انفصال  $x_1, \dots, x_N$ ، که  $f(x)$  را در نقاط  $x_1, \dots, x_{N+1}$  درونیایی می‌کند، تقریب زده شده است. از مثال ۶.۲ که برای هر یک از زیر بازه‌های  $[x_k, x_{k+1}]$ ،  $k = 1, \dots, N$ ، به کار برده شده نتیجه می‌شود که به ازای جمیع مقادیر  $x \in [a, b]$  داریم

$$|f(x) - g_1(x)| \leq \frac{1}{\lambda} \max_k \left\{ (\Delta x_k)^2 \max_{x_k \leq \xi \leq x_{k+1}} |f'''(\xi)| \right\} \leq \max_{a \leq \xi \leq b} |f'''(\xi)| \frac{1}{\lambda} (\max_k \Delta x_k)^2 \quad (71.6)$$



شکل ۹.۶ درونیایی به وسیله خط شکسته

به شرطی که  $f(x)$  در بازه  $[a, b]$  دوبار مشتق پذیر باشد. باید توجه داشت که با کوچک نمودن  $\Delta x_k$ ، به ازای جمیع مقادیر  $k$ ، می توان خطای درونبایی را تا حد ممکن تقلیل داد. بعد هم توجه می کنید که این گونه افزایش در نقاط درونبایی، کار کردن بعدی با  $g_1(x)$  را پیچیده تر نمی کند، زیرا  $g_1(x)$  «به طور محلی» تابعی است بسیار ساده.

با استفاده از یک تابع بسجمله‌ای-تکه‌ای  $g_r(x)$  از درجه بزرگتر از یک، به جای تابع خطی-تکه‌ای  $g_1(x)$ ، می توان برای تابع  $f(x)$  تقریبی تولید کرد که جمله خطای مربوط به  $\max_k \Delta x_k$  در آن شامل تسوان  $(r+1)$ ام باشد؛ از این رو با کوچک شدن  $\max \Delta x_k$ ، این جمله خطا نسبت به خطای  $(71.6)$  برای درونبایی خطی-تکه‌ای، سریعتر به سمت صفر میل می کند. مخصوصاً تقریب درجه سوم-تکه‌ای<sup>۱</sup> متداولتر است. اکنون چندین نمونه از درونبایی درجه سوم-تکه‌ای را مورد بحث قرار می دهیم.

گیریم  $f(x)$  تابعی حقیقی مقدار و معین در یک بازه  $[a, b]$  باشد. می خواهیم یک تابع (بسجمله‌ای) درجه سوم-تکه‌ای  $g_3(x)$  بسازیم که  $f(x)$  در نقاط  $x_1, \dots, x_{N+1}$  با

$$a = x_1 < x_2 < \dots < x_{N+1} = b \quad (72.6)$$

درونبایی کند؛ همان گونه که با درونبایی تکه‌ای - خطی عمل کردیم، نقاط درونبایی داخل  $x_2, \dots, x_N$  را در نقاط انفصالی  $g_3(x)$  انتخاب می کنیم، یعنی در هر بازه  $[x_i, x_{i+1}]$ ،  $g_3(x)$  را مانند یک بسجمله‌ای درجه سوم  $P_i(x)$ ،  $i = 1, \dots, N$ ، می سازیم. برای تسهیل استفاده از  $g_3(x)$  در محاسبات بعدی هر قطعه درجه سوم  $P_i(x)$  از  $g_3(x)$  را به صورت زیر می نویسیم

$$P_i(x) = c_{1,i} + c_{2,i}(x-x_i) + c_{3,i}(x-x_i)^2 + c_{4,i}(x-x_i)^3 \quad (73.6)$$

زمانی که ضرایب  $c_{j,i}$ ،  $i = 1, \dots, N$ ، و  $j = 1, \dots, 4$  معلوم باشند، آنگاه برنامه تابع فورترن موسوم به PCUBIC را به ازای هر نقطه خاص  $x = \bar{x}$  به نحو چشمگیری محاسبه می نماید.

```

REAL FUNCTION PCUEIC ( XBAR, XI, C, N )
C RETURNS THE VALUE AT XBAR OF THE PIECEWISE CUBIC FUNCTION ON N
C INTERVALS WITH BREAKPOINT SEQUENCE XI AND COEFFICIENTS C .
  INTEGER N, I, J
  REAL C(4,N), XBAR, XI(N+1), DX
  DATA I /1/
  IF (XBAR .GE. XI(I)) THEN
    DO 10 J=1,N
      IF (XBAR .LT. XI(J+1)) GO TO 30
10  CONTINUE
    J = N
  ELSE
    DO 20 J=I-1,1,-1
      IF (XBAR .GE. XI(J)) GO TO 30
20  CONTINUE
    J = 1
  END IF
30  I = J
    DX = XBAR - XI(I)
    PCUBIC = C(1,I) + DX*(C(2,I) + DX*(C(3,I) + DX*(C(4,I)))
    RETURN
END

```

## 1. piecewise\_cubic approximation

اکنون به تعیین تابع درونیابی درجه سوم-تکه‌ای  $g_3(x)$  می‌پردازیم. از آنجایی که می‌خواهیم

$$g_3(x_i) = f(x_i) \quad i = 1, \dots, N+1$$

باید داشته باشیم

$$P_i(x_i) = f(x_i) \quad P_i(x_{i+1}) = f(x_{i+1}) \quad i = 1, \dots, N \quad (74.6)$$

باید توجه کرد که روابط (74.6) ایجاب می‌کنند که داشته باشیم

$$P_{i-1}(x_i) = P_i(x_i) \quad i = 2, \dots, N$$

لذا پیوسته بودن  $g_3(x)$  در بازه  $[a, b]$  تضمین می‌شود.

با توجه به قضیه ۱۰۲ یا ۴۰۲ خاطر نشان می‌سازیم که همیشه می‌توانیم تابع مفروضی را در چهار نقطه توسط یک بسجمله‌ای درجه سوم درونیابی کنیم. تا اینجا لازم بوده که هر یک از قطعه‌های درجه سوم  $P_i(x)$ ، تابع  $f(x)$  را تنها در دو نقطه درونیابی کند و از این رو در انتخاب  $P_i(x)$  مقداری آزادی عمل داشته‌ایم و تفاوت روشهای مختلف درونیابی تنها در چگونگی استفاده از این آزادی عمل است.

در درونیابی درجه سوم-تکه‌ای هر میت،  $P_i(x)$  چنان تعیین می‌شود که  $f(x)$  در نقاط  $x_i, x_{i+1}, x_i, x_{i+1}$  و  $x_{i+1}$  درونیابی شود، یعنی چنان که روابط زیر نیز برقرار باشند

$$P'_i(x_i) = f'(x_i) \quad P'_i(x_{i+1}) = f'(x_{i+1}) \quad i = 1, \dots, N \quad (75.6)$$

در این حال از فرمول نیوتن (۳۲۰۲) نتیجه می‌شود که به ازای  $i = 1, \dots, N$

$$P_i(x) = f(x_i) + f[x_i, x_i](x - x_i) + f[x_i, x_i, x_{i+1}](x - x_i)^2 + f[x_i, x_i, x_{i+1}, x_{i+1}](x - x_i)^3(x - x_{i+1})$$

چون داریم  $(x - x_{i+1}) = (x - x_i) + (x_i - x_{i+1})$  نتیجه می‌شود که

$$P_i(x) = f(x_i) + f'(x_i)(x - x_i) + (f[x_i, x_i, x_{i+1}] - f[x_i, x_i, x_{i+1}, x_{i+1}]\Delta x_i) \times (x - x_i)^2 + f[x_i, x_i, x_{i+1}, x_{i+1}](x - x_i)^3$$

که در آن  $\Delta x_i = x_{i+1} - x_i$ ، و از روی آن ضرایب  $c_{1,i}, c_{2,i}, c_{3,i}, c_{4,i}$  را می‌توان به‌طور مستقیم به‌دست آورد. با استفاده از فرمولهای اختصاری

$$f_i = f(x_i) \quad s_i = f'(x_i) \quad i = 1, \dots, N+1 \quad (76.6)$$

داریم

$$c_{1,i} = f_i \quad c_{2,i} = s_i$$



$$c_{\varphi, i} = f[x_i, x_i, x_{i+1}] - f[x_i, x_i, x_{i+1}, x_{i+1}] \Delta x_i$$

$$= \frac{f[x_i, x_{i+1}] - s_i}{\Delta x_i} - c_{\varphi, i} \Delta x_i \quad (77.6)$$

$$c_{\varphi, i} = \frac{f[x_i, x_{i+1}, x_{i+1}] - f[x_i, x_i, x_{i+1}]}{\Delta x_i}$$

$$= \frac{s_{i+1} + s_i - 2f[x_i, x_{i+1}]}{(\Delta x_i)^2}$$

با ذخیره شدن  $f_i$  در  $c_{1,i}$  و  $s_i$  در  $c_{\varphi, i}$  به ازای  $i = 1, \dots, N+1$  در زیر برنامه فورترن زیر برای محاسبه  $c_{\varphi, i}$  و  $c_{1,i}$  در  $i = 1, \dots, N$  از رابطه (77.6) استفاده شده.

```

SUBROUTINE CALCCF ( XI, C, N )
  INTEGER N, I
  REAL C(4,N+1), XI(N+1), DIVDF1, DIVDF3, DX
  C***** I N P U T *****
  C XI(1), ..., XI(N+1) STRICTLY INCREASING SEQUENCE OF BREAKPOINTS.
  C C(1,I), C(2,I), VALUE AND FIRST DERIVATIVE AT XI(I), I=1,...,N+1.
  C OF THE PIECEWISE CUBIC FUNCTION.
  C***** O U T P U T *****
  C C(1,I), C(2,I), C(3,I), C(4,I) POLYNOMIAL COEFFICIENTS OF THE FUNC-
  C TION ON THE INTERVAL (XI(I), XI(I+1)) I=1,...,N .
  C
  DO 10 I=1,N
    DX = XI(I+1) - XI(I)
    DIVDF1 = (C(1,I+1) - C(1,I))/DX
    DIVDF3 = C(2,I) + C(2,I+1) - 2.*DIVDF1
    C(3,I) = (DIVDF1 - C(2,I) - DIVDF3)/DX
  10 C(4,I) = DIVDF3/(DX*DX)
  RETURN
END

```

□ مثال ۱۵۰۶: مسئله درونیایی مثال ۴.۲ را با استفاده از درونیایی درجه سوم-تکه‌ای هر میت حل کنید، یعنی به ازای  $N = 2, 4, \dots, 16$  بگیرید

$$x_i = \frac{(i-1)10}{N} - 5 \quad i = 1, \dots, N+1$$

$$f(x) = (1+x^2)^{-1}$$

را در این نقاط درونیایی کنید و مانند قبل ماکسیم خطای درونیایی را در بازه  $[-5, 5]$  برآورد کنید.

برنامه فورترن زیر این مسئله را حل می‌کند:

## تقریب ۳۷۳

```

C PROGRAM FOR EXAMPLE 6.15 .
  INTEGER I,J,K,N
  REAL C(4,17),ERRMAX,H,X(17),Y
C   PIECEWISE CUBIC HERMITE INTERPOLATION AT EQUALLY SPACED POINTS
C   TO THE FUNCTION
      F(Y) = 1./(1. + Y*Y)
C
  PRINT 600
600 FORMAT('1  N',5X,'MAXIMUM ERROR')
  DO 40 N=2,16,2
    H = 10./FLOAT(N)
    DO 10 I=1,N+1
      X(I) = FLOAT(I-1)*H - 5.
      C(1,I) = F(X(I))
C   10  C(2,I) = F'(X(I))
      C(2,I) = -2.*X(I)*C(1,I)**2
    CALL CALCCF ( X, C, N )
C   ESTIMATE MAXIMUM INTERPOLATION ERROR ON (-5,5).
    ERRMAX = 0.
    DO 30 I=1,101
      Y = .1*I - 5.
      ERRMAX = MAX(ERRMAX, ABS(F(Y)-PCUBIC(Y,X,C,N)))
    30  CONTINUE
    40  PRINT 640, N,ERRMAX
640 FORMAT(15,E18.7)
                                STOP
END

```

## نتایج کامپیوتری برای مثال ۱۵.۶

N	MAXIMUM ERROR
2	4.9188219E - 01
4	2.1947326E - 01
6	9.1281965E - 02
8	3.5128250E - 02
10	1.2705882E - 02
12	4.0849234E - 03
14	1.6011164E - 03
16	1.6953134E - 03

در مقایسه با درونیایی به وسیلهٔ بسجمله‌ای (مثال ۴.۲ را ببینید)، در این حال با افزایش

□

$N$ ، ما کمترین خطا به نحو محسوسی کاهش می‌یابیم.

بر آورد خطا در درونیایی درجهٔ سوم-تکه‌ای هر میت به آسانی صورت می‌گیرد.

چون به‌ازای  $x \in [x_i, x_{i+1}]$  داریم  $g_p(x) = P_i(x)$ ، که در اینجا  $P_i(x)$  تابع  $f(x)$

را در نقاط  $x_i, x_{i+1}, x_i, x_{i+1}$  درونیایی می‌کند، از رابطهٔ (۳۷.۲) نتیجه می‌شود که

به‌ازای  $x \in [x_i, x_{i+1}]$  داریم

$$f(x) - g_p(x) = f[x_i, x_i, x_{i+1}, x_{i+1}, x](x - x_i)^2(x - x_{i+1})^2$$

$$= \frac{1}{4!} f^{(4)}(\xi_x)(x - x_i)^2(x - x_{i+1})^2 \quad \xi_x \in (x_i, x_{i+1})$$

به شرط آنکه  $f(x)$  چهار بار به‌طور پیوسته مشتق‌پذیر باشد، بعلاوه

$$\max_{x \in [x_i, x_{i+1}]} |(x - x_i)^2(x - x_{i+1})^2| = \left(\frac{1}{2} \Delta x_i\right)^4 \leq \frac{(\max_j \Delta x_j)^4}{16}$$

پس برای  $a \leq x \leq b$  به‌ازای

$$|f(x) - g_r(x)| \leq \frac{1}{384} \max_k \left\{ (\Delta x_k)^4 \max_{x_k \leq \xi \leq x_{k+1}} |f^{(4)}(\xi)| \right\}$$

$$\leq \max_{\xi \in [a, b]} |f^{(4)}(\xi)| \frac{(\max_i \Delta x_i)^4}{384} \quad (78.6)$$

درونیابی درجه سوم-تکه ای هر میت، به اطلاعاتی در باب  $f'(x)$  نیاز دارد. در عمل، اغلب به دست آوردن اعداد مورد نیاز  $f'(x_i)$ ،  $i = 1, \dots, N+1$ ، اگر غیرممکن نباشد دشوار است. در چنین حالتی برای محاسبه  $s_i$ ، تقریب مناسبی برای  $f'(x_i)$ ،  $i = 1, \dots, N+1$ ، در نظر می گیرند. بنابراین در درونیابی درجه سوم-تکه ای پس، به جای  $s_i = f'(x_i)$  از رابطه

$$s_i = \frac{\Delta x_{i-1} f[x_i, x_{i+1}] + \Delta x_i f[x_{i-1}, x_i]}{\Delta x_{i-1} + \Delta x_i} \quad (79.6)$$

استفاده می کنند، اما مانند قبل، با تعیین ضرایب  $c_{j,i}$  برای تکه های درجه سوم طبق رابطه (77.6) به گونه دیگر عمل می کنند. باید توجه کرد که در به دست آوردن اعدادی برای مشتقاتی کرانه ای  $g_r(x)$ ،  $s_1$  و  $s_{N+1}$ ، تساوی (79.6) به دو نقطه اضافی  $x_0$  و  $x_{N+2}$  نیاز دارد. این نقاط به گونه ای، مثلا به صورت

$$x_0 = x_2 \quad x_{N+2} = x_{N-1}$$

انتخاب می شوند. یا متناظر با انتخاب  $x_0 = a$  و  $x_{N+2} = b$ ، اعداد

$$s_1 = f'(a) \quad s_{N+1} = f'(b) \quad (80.6)$$

را به کار می برند به شرط آنکه این اعداد در دسترس باشند. بازیک امکان دیگر آن است که  $s_1$  و  $s_{N+1}$  به طریقی انتخاب شوند که شرایط «پایانی آزاد»<sup>۱</sup>

$$g_r''(a) = g_r''(b) = 0 \quad (81.6)$$

برقرار باشند.

اگر استفاده از  $f_i = f(x_i)$ ،  $i = 1, \dots, N+1$ ، را در رابطه (77.6) ادامه دهیم، آنگاه اعداد خاص انتخابی  $s_i$ ،  $i = 1, \dots, N+1$ ، هر چه باشند تابع درجه سوم-تکه ای  $g_r(x)$ ، تابع  $f(x)$  را در نقاط  $x_1, \dots, x_{N+1}$ ، علاوه بر  $g_r(x)$  گذشته از پیوسته بودن در بازه  $[a, b]$ ، به طور پیوسته مشتق پذیر نیز هست، زیرا (77.6) ایجاب می کند که داشته باشیم

$$P'_{i-1}(x_i) = s_i = P'_i(x_i) \quad i = 2, \dots, N$$

همان گونه که اکنون نشان خواهیم داد، همواره ممکن است اعداد  $s_1, \dots, s_{N+1}$  را چنان تعیین کرد که  $g_3(x)$  حاصل، حتی دو بار به طور پیوسته مشتقپذیر باشد. این روش تعیین  $g_3(x)$  به درونیابی درجه سوم قلمی<sup>۱</sup> معروف است. در این حالت اطلاق نام «قلمی» به درونیاب  $g_3(x)$  از این جهت است که نمودار آن، تقریب زدن وضعیتی است که «قلم ظریف و انعطاف پذیر» یک نقشه کش، در صورت اجبار به گذشتن از نقاط  $\{x_i, f_i\}$ ،  $i = 1, \dots, N+1$ ، اختیار می کند.

شرط اینکه  $g_3(x)$  دو بار به طور پیوسته مشتقپذیر باشد، با شرط

$$P''_{i-1}(x_i) = P''_i(x_i) \quad i = 2, \dots, N$$

یا با توجه به رابطه (۷۳.۶) با شرط

$$2c_{3,i-1} + 6c_{4,i-1}\Delta x_{i-1} = 2c_{3,i} \quad i = 2, \dots, N$$

هم ارز است. از این رو با داشتن (۷۷.۶) می خواهیم داشته باشیم

$$\frac{2(f[x_{i-1}, x_i] - s_{i-1})}{\Delta x_{i-1}} + 2c_{4,i-1}\Delta x_{i-1} = \frac{2(f[x_i, x_{i+1}] - s_i)}{\Delta x_i} - 2c_{4,i}\Delta x_i$$

$$i = 2, \dots, N$$

اگر برای بیان  $c_{4,i-1}$  و  $c_{4,i}$  بر حسب  $f_j$ ها و  $s_j$ ها، از رابطه (۷۷.۶) استفاده، و آن را ساده کنیم خواهیم داشت

$$\begin{aligned} (\Delta x_i)s_{i-1} + 2(\Delta x_{i-1} + \Delta x_i)s_i + (\Delta x_{i-1})s_{i+1} \\ = 3(f[x_{i-1}, x_i]\Delta x_i + f[x_i, x_{i+1}]\Delta x_{i-1}) \quad i = 2, \dots, N \end{aligned} \quad (۸۲.۶)$$

این، یک دستگاه  $N-1$  معادله خطی است با  $N+1$  مجهول  $s_1, \dots, s_{N+1}$ . اگر  $s_1$  و  $s_{N+1}$  را به طریقی، مثلاً به کمک رابطه (۷۹.۶) یا رابطه (۸۰.۶)، انتخاب کنیم می توانیم دستگاه (۸۲.۶) را نسبت به  $s_2, \dots, s_N$  یا روش حذف گاوس (فصل ۴ را ببینید) حل کنیم. در این صورت ماتریس ضرایب دستگاه (۸۲.۶) نافذ قطری (سطراً مؤکداً) و بنابراین (تمرین ۳-۶.۴ را ببینید) وارونپذیر خواهد بود، و لذا دستگاه (۸۲.۶) یک جواب منحصر-به فرد خواهد داشت. وقتی  $s_2, \dots, s_N$ ، یعنی جوابهای دستگاه خطی (۸۲.۶) را به دست آوریم، برای ساختن ضرایب بسجمله ای محلی تابع درجه سوم قلمی درونیاب، از آنها و شیبهای کرانی  $s_1$  و  $s_{N+1}$  در برنامه CALCCF استفاده می کنیم.

می توان نشان داد (مثلاً دبور [V(۶); ۴۰] را ببینید) که خطا در درونیاب درجه

سوم قلمی، به ازای  $a \leq x \leq b$ ، در رابطه

$$|f(x) - g_2(x)| \leq \max_{\xi \in [a, b]} |f^{(4)}(\xi)| \frac{\Delta(\max_i \Delta x_i)^4}{384} \quad (۸۳.۶)$$

صدق می‌کند. این خطای کرانی تنها پنج برابر خطای کرانی (۷۸.۶) برای درونبایی درجه سوم هرमित است، اگرچه در درونبایی درجه سوم هرमित از دو برابر اطلاعات در باره تابع  $f(x)$ ، یعنی مقادیر  $f'(x_i)$ ،  $i = 2, \dots, N$ ، و نیز مقادیر تابع استفاده می‌شود. این بدان معنی است که شیبهای  $g'_2(x_i)$  مربوط به تابع درونبای قلمی باید برای شیبهای متناظر  $f'(x_i)$  مربوط به  $f(x)$  تقریبهای خوبی باشند. می‌توان نشان داد (برای مثال کتاب دیور [۴۰؛ ۱۱۱ - V(۱۲)] را ببینید) که به ازای  $a \leq x \leq b$ :

$$|f'(x) - g'_2(x)| \leq \max_{\xi \in [a, b]} |f^{(4)}(\xi)| \frac{(\max_i \Delta x_i)^3}{24} \quad (۸۴.۶)$$

در مورد يك دنبالهٔ یکنواخت از نقاط یعنی  $x_i = x_0 + ih$ ، به ازای جمیع مقادیر  $i$ ، حتی داریم

$$i: 2, \dots, N \quad \text{به ازای } N$$

$$|f'(x_i) - g'_2(x_i)| \leq \max_{\xi \in [a, b]} |f^{(4)}(\xi)| h^3 / 60 \quad (۸۵.۶)$$

که این مسئله، موجب تعمیم درونبایی تابع درجه سوم قلمی به عنوان وسیله‌ای برای مشتقگیری عددی (فصل ۷ را ببینید) شده است.

در زیر برنامهٔ فورترن SPLINE (قلمی) زیر روش حذف گاوس پذیرفته شده است. تسا از امتیاز مشخصه سه قطری بودن ماتریس ضرایب (۸۲.۶) الگوریتم ۳.۴ را ببینید) برای محاسبهٔ  $S_i = c_{\nu, i}$ ،  $i = 2, \dots, N$ ، به عنوان جواب دستگاه (۸۲.۶) استفاده کنند، بسا فرض اینکه اعداد،  $c_{\nu, i} = f_i$ ،  $i = 1, \dots, N+1$ ،  $c_{\nu, 1} = S_1$  و  $c_{\nu, N+1} = S_{N+1}$  داده شده باشند.

```

SUBROUTINE SPLINE ( XI, C, N )
PARAMETER NP1MAX=50
INTEGER N, M
REAL C(4,N+1), XI(N+1), D(NP1MAX), DIAG(NP1MAX), G
C***** I N P U T *****
C XI(1), ..., XI(N+1) STRICTLY INCREASING SEQUENCE OF BREAKPOINTS
C C(1,I), C(2,I), VALUE AND FIRST DERIVATIVE AT XI(I), I=1,...,N+1,
C OF THE CUBIC SPLINE.
C***** O U T P U T *****
C C(1,I), C(2,I), C(3,I), C(4,I) POLYNOMIAL COEFFICIENTS OF THE SPLINE
C ON THE INTERVAL (XI(I), XI(I+1)), I=1,...,N.
DATA DIAG(1), D(1) /1., 0./
DO 10 M=2, N+1
  D(M) = XI(M) - XI(M-1)
10 DIAG(M) = (C(1,M) - C(1,M-1))/D(M)
DO 20 M=2, N
  C(2,M) = 3.*(D(M)*DIAG(M+1) + D(M+1)*DIAG(M))
20 DIAG(M) = 2.*(D(M) + D(M+1))
DO 30 M=2, N
  G = -D(M+1)/DIAG(M-1)

```

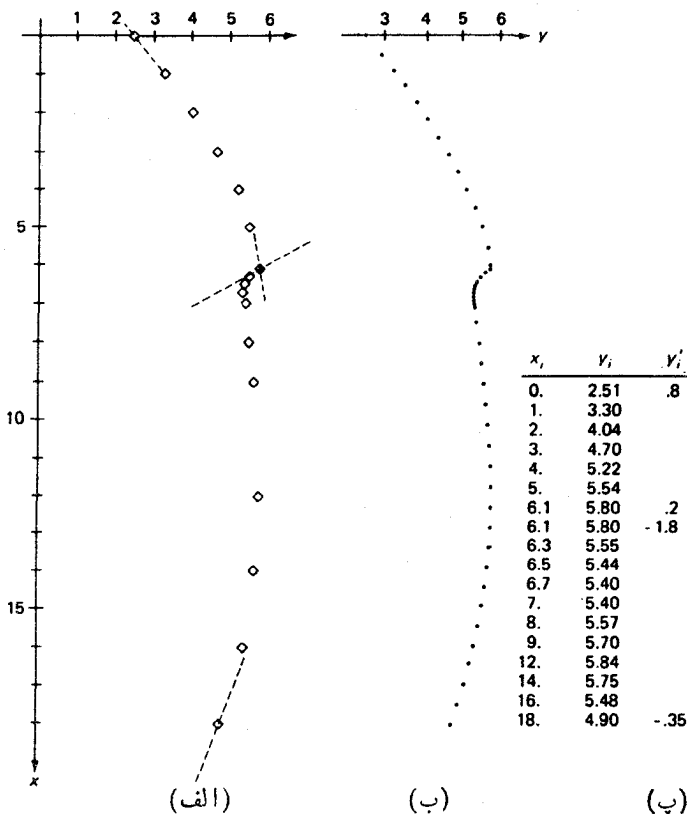
```

DIAG(M) = DIAG(M) + G*D(M-1)
30 C(2,M) = C(2,M) + G*C(2,M-1)
DO 40 M=N,2,-1
40 C(2,M) = (C(2,M) - D(M)*C(2,M+1))/DIAG(M)
RETURN
END

```

□ مثال ۱۶۰۶: تقریب‌زدن یک منحنی طرح توسط یک تابع درجه سوم قلمی یک منحنی طرح، مثلاً قسمتی از مقطع عرضی یک درب ماشین، به گونه‌ای که در شکل ۱۰۰۶ (الف) آمده، داده شده است. این منحنی در  $x = 61$  یک ناپیوستگی شیبی دارد. همان‌گونه که در شکل‌های ۱۰۰۶ (الف) و (ج) نشان داده شده است، اندازه‌گیری‌هایی انجام گرفته و شیب‌های انتهایی به‌طور نموداری برآورد شده‌اند. مسئله ما پیدا کردن یک تابع  $s(x)$  است که به داده‌های مفروض بپردازد و «هموار جلوه کند».

با درونیایی درجه سوم قلمی برای داده‌های مفروض استفاده از دو تابع درجه سوم قلمی که به‌طور پیوسته بهم وصل شده‌اند ولی در  $x = 61$  شیب‌های متفاوتی دارند، جواب این مسئله به آسانی به دست می‌آید. این کار در برنامه فورترن زیر با استفاده از زیر برنامه‌های SPLINE و CALCCF که قبلاً مورد بحث قرار گرفت، انجام می‌گیرد.



شکل ۱۰۰۶ تقریب درجه سوم قلمی برای یک منحنی طرح.

این برنامه داده‌های تا  $x=۶۱$ ، به انضمام دوشیب انتهایی مفروض، را می‌دهد و ضرایب بسجمله‌ای محاسبه شده مربوط به شش تکه نخست بسجمله‌ای را در

$$C(J, I), J=۱, \dots, ۴ \quad I=۱, \dots, ۶$$

ذخیره می‌کند. سپس داده‌ها از  $x=۶۱$  تا  $x=۱۸$  همراه با دوشیب انتهایی خوانده می‌شوند و با استفاده مجدد از برنامه‌های SPLINE و CALCCF، ضرایب

$$C(J, I), J=۱, \dots, ۴ \quad I=۷, \dots, ۱۶$$

متعلق به ۱۰ تکه بسجمله‌ای باقیمانده به دست می‌آیند. سرانجام تابع درجه سوم-تکه‌ای محاسبه شده  $s(x)$  با به کارگیری برنامه PCUBIC، به ازای مقادیر مختلف  $x$  ارزیابی می‌شود. برخی از این مقادیر در شکل ۱۰.۶ (ب)، مشخص شده‌اند. حتی بدون ناپیوستگی شیبی، درونیابی بسجمله‌ای برای این داده‌ها تقریبی «ناهموار» یعنی تقریبی نوسانی تولید می‌کند، زیرا این ناحیه با انحنای نسبتاً زیاد در نزدیک نقطه  $x=۶۱$ ، توسط یک قسمت مبهم و نسبتاً صاف دنبال می‌شود (تمرین ۷.۰۶-۲ را ببینید).

برنامه فورترن برای درونیابی قلمی درجه سوم مثال (۱۶.۶)

```

C PROGRAM FOR EXAMPLE 6.16
  PARAMETER NP1MAX = 50
  INTEGER I, IEND, N, N1, N2
  REAL C(4, NP1MAX), FX, X, XI(NP1MAX)
  READ 500, N1
500 FORMAT(I2)
  READ 501, (XI(I), C(1, I), I=1, N1), C(2, 1), C(2, N1)
501 FORMAT(2E10.3)
  N = N1 - 1
  CALL SPLINE(XI, C, N)
  CALL CALCCF(XI, C, N)
C
  READ 500, N2
  IEND = N + N2
  READ 501, (XI(I), C(1, I), I=N1, IEND), C(2, N1), C(2, IEND)
  N = N2 - 1
  CALL SPLINE(XI(N1), C(1, N1), N)
  CALL CALCCF(XI(N1), C(1, N1), N)
C
  N = IEND - 1
  X = XI(1)
  DO 10 I=1, 40
    FX = PCUBIC(X, XI, C, N)
    PRINT 600, I, X, FX
600 FORMAT(I5, F10.1, E20.9)
  10 X = X + .5
                                STOP
END

```



در اینجا ما تنها مقدمه کوتاهی برای تقریب به وسیله بسجمله‌ای-تکه‌ای داده‌ایم برای توضیحات بیشتر، مثلاً، به کتاب دبور [۴۰] مراجعه کنید. تقریب به وسیله بسجمله‌ای و تقریب درجه سوم-تکه‌ای از جهات مختلف و مهمی که هنگام درونیابی آشکار می‌شوند، با هم متفاوت‌اند. اگر داده‌ها در نقاط متساوی الفاصله

داده شده باشند، آنگاه همان گونه که در مثال ۴.۲ نشان داده شده، با افزایش تعداد این نقاط، درونیایی به وسیلهٔ بسجمله‌ای، قدرت خود را از دست می‌دهد. این گونه مشکلات حتی در درونیایی درجهٔ سوم قلمی اصلاً وجود ندارد. (توجه می‌کنید که در درونیایی به وسیلهٔ بسجمله‌ای مثلثاتی نیز هیچ مشکلی وجود ندارد.) و نیز با افزایش تعداد نقاط، بسجمله‌ای (و بسجمله‌ای مثلثاتی) پیچیده‌تر، یعنی ارزیابی آن پرهزینه‌تر می‌شود. همچنین وقتی درجهٔ بسجمله‌ای بیشتر از ۱۵ باشد، بدعلت بد شرطی شکل توانی، یا باید از دقت مضاعف استفاده کرد، یا بسجمله‌ای را به شکل‌های دیگری، مثلاً برحسب بسجمله‌ایهای چیشف، نوشت. چنین مشکلاتی در درونیایی درجهٔ سوم-تکه‌ای وجود ندارد. زیرا تعداد نقاط درونیایی هر قدر زیاد باشد، درونیاب همیشه به‌طور محلی یک تابع خیلی ساده یعنی یک بسجمله‌ای درجهٔ سوم است. و سرانجام اگر تابعی که می‌باید تقریب زده شود در بعضی جاها رفتار نامناسب داشته باشد، آنگاه بهترین بسجمله‌ای تقریب، در همه جا تقریب ناتوانی خواهد بود (به‌تمرین ۱۰۶-۱۰۵ و ۱۱-۱۰۶ نگاه کنید). در تقریب بسجمله‌ای-تکه‌ای این امکان وجود دارد که با انتخاب مناسب نقاط ناپیوستگی، این اثرات به یک بازه در محدودهٔ نقاطی که رفتار تابع نامطلوب است محدود شوند، و در بقیهٔ جاها تقریب خوبی به دست آید.

## تمرین

۱-۷۰۶ به موجب نمادگذاری در این قسمت، معادله‌ای بیابید که  $f_1, f_2, f_3, f_4, f_5$  باید در آن صدق کنند تا شرط «پایانی آزاد»

$$g_7''(a) = 0$$

برقرار شود.

۲-۷۰۶ یک بسجمله‌ای از درجهٔ مناسب پیدا کنید که منحنی طرح مثال ۱۶.۶ را در تمام نقاط داده‌های مفروض از ۱ تا ۱۸ (به انضمام شیبه‌ها) درونیایی کند، و آن را با تقریب قلمی که در مثال ۱۶.۶ حساب شده، مقایسه کنید.

۳-۷۰۶ داده‌های مثال ۱۶.۶ را توسط درونیایی درجهٔ سوم بس، درونیایی کنید و نتایج حاصله را مقایسه کنید.

۴-۷۰۶ درونیایی درجهٔ سوم بس محلی است، بدین معنی که مقدار تابع درونیاب  $g_3(x)$  در هر نقطهٔ  $x$ ، تنها به چهار مقدار تابع مفروض در نزدیکی‌ترین نقاط به نقطهٔ  $x$  بستگی دارد. برعکس، درونیایی درجهٔ سوم قلمی کلی<sup>۱</sup> است، یعنی مقدار  $g_3(x)$  در هر نقطهٔ مفروض به کلیهٔ اطلاعات داده شده در بازهٔ  $f(x)$  بستگی دارد. این دو حکم را ثابت کنید.

۵-۷۰۶ سعی کنید از درونیایی یک تابع مفروض توسط یک تابع سهمی-تکه‌ای  $g_2(x)$  الگوی مناسبی بسازید. آیا می‌توانید تابع  $g_2(x)$  را به‌طور پیوسته مشتق‌پذیر سازید؟





## مشتقگیری و انتگرالگیری

در فصل ۲، تکنیکهایی را برای تقریب زدن يك تابع مفروض به وسیله يك بسجمله‌ای، نوعاً به روش درونیایی، بسط دادیم. در این فصل، يك استفاده مهم از این بسجمله‌ایهای تقریبی یعنی جایگزینی تحلیلی را مورد مطالعه قرار می‌دهیم. در اینجا، سروکار، با گذاردن يك بسجمله‌ای تقریبی به جای يك تابع پیچیده یا فقط يك تابع جدول‌بندی شده است به طوری که اعمال بنیادی دیفرانسیل و انتگرال بتوانند آسانتر انجام گیرند، یا اساساً بتوانند انجام گیرند. این اعمال شامل عملهای

$$I(f) = \int_a^b f(x) dx \quad D(f) = f'(a)$$

$$S_n(f) = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) \sin nx \, dx$$

و حتی

$$Z(f) = \lim_{h \rightarrow 0} f(h)$$

هستند. به گونه‌ای مجرد، اگر  $L$  معرف یکی از اعمال فوق (یا عملی مشابه) روی توابع باشد، و به ازای يك  $f(x)$  مفروض،  $p(x)$  تقریبی برای  $f(x)$  باشد، آنگاه عدد  $L(f)$  را با عدد  $L(p)$  تقریب می‌زنیم. امید ما این است که عمل  $L$  بتواند به آسانی روی  $p(x)$

انجام گیرد و اگر  $p(x)$  يك بسجمله‌ای و  $L$  یکی از اعمال فوق‌الذکر باشد، این امید قابل توجیه باشد.

خطی بودن عمل  $L$  (همان گونه که اعمال مذکور در بالا هستند) در برآورد خطای  $L(f) - L(p)$ ، تا حدی مفید است. معنی خطی بودن این است که

$$L(f(x) + g(x)) = L(f) + L(g)$$

$$L(af(x)) = aL(f)$$

که در اینجا  $a$  يك عدد است و  $f(x)$  و  $g(x)$  دو تابع. خطی بودن ایجاب می‌کند که داشته باشیم

$$L(f) - L(p) = L(e)$$

که  $e(x)$  عبارت است از خطای موجود در تابع تقریب  $p(x)$  از تابع  $f(x)$ ، یعنی

$$f(x) = p(x) + e(x)$$

ما معمولاً  $p(x)$  را به عنوان يك بسجمله‌ای درونیاب انتخاب می‌کنیم. مثلاً  $p(x)$  يك بسجمله‌ای از درجه نایبتر از  $k$  است که  $f(x)$  را در نقاط  $x_0, \dots, x_k$  درونیابی می‌کند. اگر این نقاط متمایز باشند، آنگاه بنا بر (۷.۲) داریم

$$p(x) = \sum_{i=0}^k f(x_i) l_i(x)$$

که در آن  $l_i(x)$ ها بسجمله‌ایهای لاگرانژ به ازای نقاط  $x_0, \dots, x_k$  هستند. اما اگر عمل  $L$  خطی باشد، نتیجه می‌شود که

$$L(p) = \sum_{i=0}^k f(x_i) w_i$$

که در آن اعداد  $w_i$  با روابط

$$w_i = L(l_i) \quad i = 0, \dots, k$$

مشخص می‌شوند و به  $f(x)$  بستگی ندادند؛ و بنا بر این می‌توانند (به ازای هر مجموعه از نقاط خاص  $x_0, \dots, x_k$ ) یکبار برای همیشه محاسبه شوند. در این شکل نمایش، معمولاً تقریب  $L(p)$  يك قاعده [برای تقریب  $L(f)$ ] نامیده می‌شود، و اعداد  $w_i$  اوزان یا ضرایب آن و نقاط  $x_0, \dots, x_k$  رؤس آن نام دارند. در این قاعده، با اعمال عمل  $L$  بر تابع خطای بسجمله‌ای درونیابی، همان گونه که در (۱۸.۲) و (۳۷.۲) آمده است، و با

استفاده از این حقیقت که تفاضل منقسم تابعی است خوش تعریف از شناسه‌هایش، عبارت زیر را برای خطا به دست می‌آوریم

$$E(f) = L(f) - L(p)$$

## ۱.۷ مشتگیری عددی

در وهله اول بعضی از تکنیک‌های عددی برای تقریب زدن مشتق يك تابع داده شده، یعنی  $f'(x)$  را بررسی می‌کنیم. قواعد حاصل در حل عددی معادلات دیفرانسیل از اهمیت درجه اولی برخوردارند، و علت اساسی شرح آنها در اینجا همین امر است. این قواعد برای به دست آوردن تقریب‌های عددی يك مشتق از روی مقادیر تابع نیز به کار می‌روند. اما باید متذکر شویم که مشتگیری عددی بر مبنای بسجمله‌ای درونیاب، اساساً فرایندی است ناپایدار و نمی‌توانیم دقت خوبی را انتظار داشته باشیم حتی اگر داده‌های اولیه دقیق باشند. به طوری که خواهیم دید، خطای  $f'(x) - p'(x)$  ممکن است بسیار بزرگ باشد، بخصوص زمانی که مقادیر  $f(x)$  در نقاط درونیابی «نوفه‌ای» باشند. این مطالب در زیر دقیقتر توضیح داده خواهد شد.

گیریم  $f(x)$  تابعی باشد در بازه  $[c, d]$  به طور پیوسته مشتق پذیر. اگر در بازه  $[c, d]$ ، نقاط  $x_0, \dots, x_k$ ، نقاطی متمایز باشند، می‌توانیم  $f(x)$  را به موجب (۳۷.۲) به صورت زیر بنویسیم

$$f(x) = p_k(x) + f[x_0, \dots, x_k, x] \psi_k(x) \quad (۱۰۷)$$

که در آن  $p_k(x)$  بسجمله‌ای است از درجه نایبتر از  $k$  که  $f(x)$  در نقاط  $x_0, \dots, x_k$  درونیابی می‌کند و

$$\psi_k(x) = \prod_{j=0}^k (x - x_j)$$

بنابر (۳۸.۲)، اگر  $f(x)$  به اندازه کافی هموار باشد، داریم

$$\frac{d}{dx} f[x_0, \dots, x_k, x] = f[x_0, \dots, x_k, x, x]$$

لذا، در يك چنین حالتی، می‌توانیم از (۱۰۷) مشتق بگیریم و چنین به دست آوریم

$$f'(x) = p'_k(x) + f[x_0, \dots, x_k, x, x] \psi_k(x) + f[x_0, \dots, x_k, x] \psi'_k(x) \quad (۲۰۷)$$

عملگر  $D$  را در نقطه غیر مشخص  $a$  روی  $[c, d]$  به شرح زیر تعریف می‌کنیم

$$D(f) = f'(a)$$

اگر  $D(f)$  را با  $D(p_k)$  تقریب زنیم، آنگاه به موجب (۳.۷)، خطا در این تقریب به شرح زیر است

$$E(f) = D(f) - D(p_k)$$

$$= f[x_0, \dots, x_k, a, a] \psi_k(a) + f[x_0, \dots, x_k, a] \psi'_k(a)$$

یا به ازای مقادیری مانند  $\xi, \eta \in (c, d)$

$$E(f) = \frac{f^{(k+2)}(\xi) \psi_k(a)}{(k+2)!} + \frac{f^{(k+1)}(\eta) \psi'_k(a)}{(k+1)!} \quad (3.7)$$

در مشتقگیری عددی، عبارت (۳.۷) در حالت کلی برای خطای  $E(f)$  اطلاعات بسیار کمی از خطای واقعی به دست می‌دهد، زیرا به ندرت مشتقهای  $f^{(k+1)}$  و  $f^{(k+2)}$  موجود در  $E(f)$  را در دست داریم و تقریباً هیچ وقت شناسه‌های  $\xi$  و  $\eta$  بر ما معلوم نیستند. در برخی موارد این عبارت خطا یا با انتخاب نقطه  $a$  که در آن مشتق باید محاسبه شود و یا با گزینش مناسب نقاط درونیایی  $x_0, \dots, x_k$  می‌تواند بسیار ساده شود.

در ابتدا حالتی را که  $a$  یکی از نقاط درونیایی باشد بررسی می‌کنیم. گزینش به ازای مقداری مانند  $i$ ، داشته باشیم  $a = x_i$ . در این صورت چون  $\psi_k(x)$  متضمن عامل  $(x - x_i)$  است، نتیجه می‌شود که  $\psi_k(a) = 0$ ، اولین جمله در رابطه خطای (۳.۷) حذف می‌شود. بعلاوه  $\psi'_k(a) = q(a)$  که در آن

$$q(x) = \frac{\psi_k(x)}{x - x_i} = (x - x_0) \dots (x - x_{i-1})(x - x_{i+1}) \dots (x - x_k)$$

بنابراین، اگر به ازای مقداری مانند  $i$ ،  $a = x_i$  را انتخاب کنیم، آنگاه معادله (۳.۷) به معادله زیر بدل می‌شود

$$E(f) = \frac{1}{(k+1)!} f^{(k+1)}(\eta) \prod_{\substack{j=0 \\ j \neq i}}^k (x_i - x_j) \quad \eta \in (c, d) \text{ مانند ای نقطه‌ای} \quad (4.7)$$

راه دیگر برای ساده کردن عبارت خطای (۳.۷) انتخاب  $a$  است به طوری که  $\psi'_k(a) = 0$ ، زیرا در این صورت جمله دوم رابطه (۳.۷) صفر می‌شود. اگر  $k$  عدد فردی باشد، می‌توانیم این عمل را با قراردادن  $x_j$ ‌ها به طوری متقارن پیرامون  $a$ ، یعنی به طوری که داشته باشیم

$$x_{k-j} - a = a - x_j \quad j = 0, \dots, \frac{k-1}{2} \quad (5.7)$$

انجام دهیم. زیرا

$$\begin{aligned}(x-x_j)(x-x_{k-j}) &= (x-a+a-x_j)(x-a+a-x_{k-j}) \\ &= (x-a)^2 - (a-x_j)^2 \quad j=0, \dots, \frac{k-1}{2}\end{aligned}$$

بنابراین خواهیم داشت

$$\psi_k(x) = \prod_{j=0}^{(k-1)/2} [(x-a)^2 - (a-x_j)^2]$$

زیرا به ازای کلیه مقادیر  $j$

$$\frac{d}{dx} [(x-a)^2 - (a-x_j)^2] \Big|_{x=a} = 2(x-a) \Big|_{x=a} = 0$$

بنابراین نتیجه می شود که  $\psi'_k(a) = 0$ . خلاصه، اگر رابطه (۵.۷) برقرار باشد، آنگاه (۳.۷) به معادله زیر بدل می شود

$$E(f) = \frac{1}{(k+2)!} f^{(k+2)}(\xi) \prod_{j=0}^{(k-1)/2} [-(a-x_j)^2] \quad (۶.۷)$$

باید توجه کرد که درجه مشتق  $f(x)$  در رابطه (۶.۷) یک مرتبه بالاتر از درجه مشتق در رابطه (۴.۷) است.

اکنون مثالهای خاصی را مورد بررسی قرار می دهیم. اگر  $k=0$ ، آنگاه داریم  $D(p_k) = 0$  که نکته اطمینان بخشی است، اما (معمولا) تقریب خیلی خوبی برای  $D(f) = f'(a)$  نیست. بنابراین  $k \geq 1$  را انتخاب می کنیم. به ازای  $k=1$ ، داریم

$$p_k(x) = f(x_0) + f[x_0, x_1](x-x_0)$$

از این رو بدون در نظر گرفتن  $a$ ، خواهیم داشت

$$D(p_k) = f[x_0, x_1]$$

اگر  $a = x_0$ ، آنگاه به ازای  $h = x_1 - x_0$  روابط (۲.۷) و (۴.۷) فرمول تفاضل-پیشرو<sup>۱</sup> زیر را به دست می دهد

$$f'(a) \approx f[a, a+h] = \frac{f(a+h) - f(a)}{h}$$

$$E(f) = -\frac{1}{2} h f''(\eta)$$

(۷.۷)

از سوی دیگر، اگر بگیریم  $a = (1/2)(x_0 + x_1)$ ، آنگاه  $x_0$  و  $x_1$  پیرامون  $a$  متقارن بوده و به ازای  $x_0 = a - h$  و  $x_1 = a + h$  و  $h = (1/2)(x_1 - x_0)$  رابطه (۶.۷) خود فرمول معمولی تفاضل مرکزی<sup>۱</sup> زیر را به دست می‌دهد

$$f'(a) \approx f[a-h, a+h] = \frac{f(a+h) - f(a-h)}{2h} \quad (۸.۷)$$

$$E(f) = -\frac{h^2}{6} f'''(\eta)$$

بنابراین، اگر  $x_0$  و  $x_1$  «نزدیک به هم باشند» آنگاه تقریب برای  $f'(a)$  در نقطه میانی  $a = (1/2)(x_0 + x_1)$  بهتر از تقریب در یکی از دوسر  $a = x_0$  یا  $a = x_1$  است. این امر چندان شگفت آور نیست زیرا به موجب قضیه مقدار میانگین برای مشتقات (به قسمت ۷.۱ نگاه کنید) داریم

$$f[x_0, x_1] = f'(a) \quad \text{به ازای نقطه‌ای مانند } a \text{ بین } x_0 \text{ و } x_1$$

این نکته در شکل ۱۰.۷ نیز نشان داده شده است.

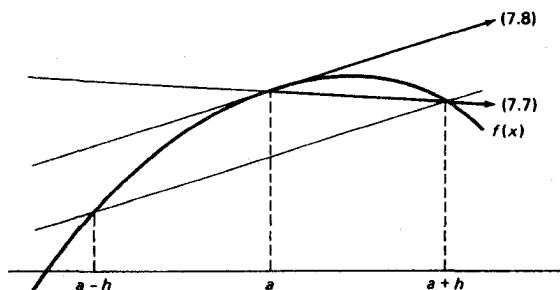
سپس استفاده از سه نقطه درونیابی را مورد بررسی قرار می‌دهیم به طوری که

$k = 2$ . بنابراین

$$p_k(x) = f(x_0) + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1)$$

لذا

$$p_k'(x) = f[x_0, x_1] + f[x_0, x_1, x_2](2x - x_0 - x_1)$$



شکل ۱۰.۷ مشتقگیری عددی

بنابراین اگر  $a = x_0$ ، آنگاه از روابط (۲.۷) و (۲.۷) به دست می‌آوریم

$$f'(a) = f[a, x_1] + f[a, x_1, x_2](a - x_1) + \frac{1}{6}(a - x_1)(a - x_2)f'''(\eta) \quad (۹.۷)$$

اکنون در حالت خاص، گیرسیم  $x_1 = a + h$  و  $x_2 = a + 2h$ ، بنا بر این (۹.۷) به صورت زیر بدل می‌شود

$$f'(a) \approx \frac{-3f(a) + 2f(a+h) - f(a+2h)}{2h} \quad (۱۰.۷)$$

$$E(f) = \frac{h^3}{3} f'''(\xi) \quad \text{به‌ازای نقطه‌ای مانند } \xi \text{ بین } a \text{ و } a+2h$$

از طرف دیگر، اگر  $x_1 = a - h$  و  $x_2 = a + h$ ، آنگاه داریم

$$f'(a) \approx \frac{f(a+h) - f(a-h)}{2h}$$

$$E(f) = -\frac{h^3}{6} f'''(\xi) \quad |\xi - a| < |h| \quad (۱۱.۷)$$

که عیناً همان معادلهٔ (۸.۷) است.

فرمولهایی برای تقریب‌زدن مشتقهای مراتب بالاتر  $f(x)$  را می‌توان به‌طریقی مشابه به دست آورد. بنا بر این با دوبار مشتق گرفتن از (۱۰.۷)، خواهیم داشت

$$\begin{aligned} f''(x) &= p_k''(x) + 2f[x_0, \dots, x_k, x, x, x]\psi_k'(x) \\ &\quad + 2f[x_0, \dots, x_k, x, x]\psi_k''(x) \\ &\quad + f[x_0, \dots, x_k, x]\psi_k'''(x) \end{aligned} \quad (۱۲.۷)$$

به‌ازای  $k=2$  و  $a=x_0$  از این رابطه نتیجه می‌شود

$$\begin{aligned} f''(a) &= 2f[a, x_1, x_2] + 2f[a, x_1, x_2, a, a](a - x_1)(a - x_2) \\ &\quad + f[a, x_1, x_2, a]2(a - x_1 + a - x_2) \end{aligned}$$

از این دو به‌ازای  $x_1 = a + h$  و  $x_2 = a + 2h$ ، داریم

$$f''(a) \approx \frac{f(a) - 2f(a+h) + f(a+2h)}{h^2}$$

$$E(f) = \frac{h^4}{6} f^{iv}(\xi) - h f'''(\eta) \quad (13.7)$$

در عوض، با انتخاب  $x_2 = a+h$  و  $x_1 = a-h$ ، به گونه‌ای که نقاط درونیابی پیرامون  $a$  متقارن باشند، خواهیم داشت

$$f''(a) \approx \frac{f(a-h) - 2f(a) + f(a+h)}{h^2}$$

$$E(f) = -\frac{h^4}{12} f^{iv}(\xi) \quad |\xi - a| < |h| \quad \text{به ازای} \quad (14.7)$$

ملاحظه می‌کنیم که قراردادن نقاط درونیابی به‌طور متقارن پیرامون  $a$ ، مجدداً فرمولی از مرتبه بالاتر را نتیجه داده‌است.

و سرانجام از (۱۷.۲) نتیجه می‌گیریم که

$$k! f[x_0, \dots, x_k]$$

تقریب «خوبی» برای  $f^{(k)}(a)$  است، به شرط آنکه کلیه  $p_k$  «به اندازه کافی نزدیک» به  $a$  باشند. فرمولهای (۷۰.۷)، (۸۰.۷)، (۱۰۰.۷) همگی به‌صورت کلی زیر هستند.

$$D(f) = D(p_k) + \text{const} h^r f^{r+1}(\xi) \quad (15.7)$$

که در آن  $h^r f^{(r+1)}(\xi)$  مقداری است ثابت،  $D(f) = f'(a)$ ، و  $h$  فواصل نقاطی است که برای درونیابی به‌کار می‌روند. علاوه، عدد  $D(p_k)$  تنها در تعدادی متناهی<sup>۱</sup> از نقاط مجزا<sup>۲</sup> شامل مقادیر  $f(x)$  است. روند قراردادن  $D(f)$  به‌جای  $D(p_k)$  را عمل جداسازی<sup>۳</sup> نامند و ثابت جمله خطای یعنی  $h^r f^{(r+1)}(\xi)$  را خطای جداسازی گویند.

از رابطه (۱۵.۷) نتیجه می‌شود که تنها با محاسبه  $D(p_k)$  به‌ازای  $h$  به‌اندازه کافی کوچک، باید بتوانیم بسا هردقت مطلوبی،  $D(f)$  را محاسبه کنیم. ولی محدودیت طول کلمه‌ها در کامپیوترها و از دست رفتن ارقام با معنی هنگام تفریق کمیت‌های تقریباً مساوی، توأمأ حصول دقت از مرتبه بالا را دشوار می‌سازند. زیرا برای کامپیوتری با طول کلمه



ثابت و برای يك تابع مفروض، مقدار بهینه‌ای از  $h$  وجود دارد که کمتر از آن، تقریب را بدتر می‌سازد. برای مثال، مقادیر جدول ۱۰۷ را در نظر می‌گیریم. این مقادیر به صورت محاسبات با دقت ساده در حساب با ممیز شناور و بسا استفاده از IBM مدل ۷۰۹۴ انجام گرفته‌اند. در این جدول، ستونی که با عنوان  $D_h$  مشخص شده است مقادیر  $f'(a)$  را که توسط (۸.۷) برآورد می‌شود به دست می‌دهد و ستون  $D_h^2$  مقادیر  $f''(a)$  را، که بسا (۱۴.۷) برآورد شده است. تابع  $f(x)$  برابر  $e^x$  اختیار شده، و واضح است که به ازای  $a=0$  مقادیر دقیق  $f'(a)$  و  $f''(a)$  برابر یک است. از جدول فوق چنین برمی‌آید که با کوچک شدن  $h$  تا مقدار  $h=0.001$ ، مقادیر  $D_h$  و  $D_h^2$  مرتباً رويه بهتر شدن می‌روند. پس از این، نتایج بدتر می‌شوند. به ازای  $h=0.00001$  در  $D_h$  چهار رقم بسا معنی و در  $D_h^2$  هفت رقم بسا معنی از دست می‌روند. تنها راه چاره برای از دست نرفتن این ارقام با معنی، وقتی  $h$  کوچکتر می‌شود، افزودن تعدادی ارقام با معنی به ارقام حساب شدهٔ  $f(x)$  است. معمولاً این عمل روی بیشتر کامپیوترها غیرممکن است. بعلاوه  $f(x)$  خود، نتیجهٔ محاسبات دیگری است که خطاهای دیگری را وارد می‌کنند.

برای تحلیل این مسئله، مطالعهٔ فرمول (۱۱.۷) رابطهٔ زیر را به دست می‌دهد

$$f'(a) = \frac{f(a+h) - f(a-h)}{2h} - \frac{h^2 f'''(a)}{6}$$

در حقیقت به دلیل گرد کردن، اعداد  $f(a-h) + E_-$  و  $f(a+h) + E_+$  به جای اعداد  $f(a-h)$  و  $f(a+h)$  در محاسبات به کار می‌روند، بنابراین محاسبات زیر را انجام می‌دهیم

$$\begin{aligned} f'_{\text{comp}} &= \frac{f(a+h) + E_+ - f(a-h) - E_-}{2h} \\ &= \frac{f(a+h) - f(a-h)}{2h} + \frac{E_+ - E_-}{2h} \end{aligned}$$

### جدول ۱۰۷

$h$	EXP( $h$ )	EXP( $-h$ )	$D_h$	$D_h^2$
1.0	0.27182817E 01	0.36787944E 00	0.11752012E 01	0.10861612E 01
0.1	0.11051708E 01	0.90483743E 00	0.10016673E 01	0.10008334E 01
0.01	0.10100501E 01	0.99004984E 00	0.10000161E 01	0.10000169E 01
0.001	0.10010005E 01	0.99990050E 00	0.99999458E 00	0.99837783E 00
0.0001	0.10000999E 01	0.99990001E 00	0.99994244E 00	0.14901161E 01

لذا، با توجه به رابطه (۱۱.۷) خواهیم داشت

$$f'(a) = f'_{\text{comp}} - \frac{E_+ - E_-}{2h} - \frac{h^2 f'''(\xi)}{6} \quad (16.7)$$

از این رو مشاهده می‌شود که خطا در تقریب محاسبه  $f'_{\text{comp}}$  برای  $f'(a)$  از دو قسمت تشکیل شده است، یک قسمت بر اثر گرد کردن و قسمت دیگر بر اثر جداسازی. اگر  $f'''(x)$  کراندار باشد، آنگاه به ازای  $h \rightarrow 0$  خطای جداسازی به سمت صفر میل می‌کند، اما اگر فرض کنیم (همان گونه که در عمل باید انجام داد) که  $E_+ - E_-$  کم نمی‌شود، آنگاه خطای گرد کردن زیاد می‌شود (اما به تدریج  $10^{-7}$  -  $5$  نگاه کنید).

مقدار بهینه  $h$  به موجب تعریف، مقداری است که به ازای آن، حاصل جمع مقادیر خطای گرد کردن و خطای جداسازی حداقل می‌شود. برای آنکه روش پیدا کردن مقدار بهینه  $h$  را نشان دهیم، مسئله فوق را در رابطه با محاسبه  $f'(0)$  در حالت  $f(x) = e^x$  بررسی می‌کنیم. فرض کنید که در محاسبه  $e^x$  خطای حاصل برابر با  $\pm 1 \times 10^{-8}$  و مقدار  $E_+ - E_-$  محدود و تقریباً برابر با  $\pm 2 \times 10^{-8}$  باشد. در این صورت با توجه به رابطه (۱۶.۷) خطای گرد کردن  $R$  به طور تقریبی برابر است با

$$R = \pm \frac{2 \times 10^{-8}}{2h}$$

چون  $f'''(\xi)$  تقریباً برابر یک است، خطای جداسازی  $T$  تقریباً برابر با

$$T = -\frac{1}{6} h^2$$

خواهد بود. برای پیدا کردن مقدار بهینه  $h$  می‌باید رابطه زیر را مینیمم سازیم

$$|R| + |T| = \frac{10^{-8}}{h} + \frac{1}{6} h^2 = g(h)$$

برای یافتن مقداری از  $h$  که به ازای آن مقدار  $g(h)$  مینیمم باشد، از رابطه  $g(h)$  نسبت به  $h$  مشتق می‌گیریم و ریشه آن را به دست می‌آوریم. بنابراین خواهیم داشت:

$$g'(h) = \frac{-10^{-8}}{h^2} + \frac{h}{3} = 0$$

که جواب مثبت آن برابر است با

$$h^3 = 3 \times 10^{-8}$$

$$h = 10^{-3} \sqrt[3]{30} \approx 0.003$$

که این مقدار بهینهٔ  $h$  است. دانشجویان می‌توانند با بررسی جدول ۱۰۷ تحقیق کنند که بهترین مقدار برای  $h$  بین ۰.۰۰۱ و ۰.۰۰۱ قرار دارد.

فرمولهایی که در این قسمت در مشتقگیری عددی به دست آمد در مطالعهٔ روشهای حل عددی معادلات دیفرانسیل بسیار مفیدند (به فصل ۸ و ۹ نگاه کنید). اما تحلیلهای فوق‌نشان می‌دهند که برای محاسبهٔ تقریبی مشتقها، این فرمولها از مطلوبیت محدودی برخوردارند. این تحلیلهای نشان می‌دهند که می‌توانیم با به کار بستن محاسبات با دقت عمل «به اندازهٔ کافی» بالا، با اثر خطای گرد کردن مقابله کنیم. اما وقتی  $f(x)$  را تنها به طور تقریبی در چندین نقطه داشته باشیم، این امر غیرممکن است.

اگر محاسبهٔ عددی مشتقها غیر قابل اجتناب باشد، معمولاً ترجیح می‌دهند که  $D(f)$  را از روی  $D(p_k)$  بر او رد کنند، که در آن  $p_k(x)$  تقریب  $f(x)$  با کوچکترین توانهای دوم به کمک بسجمله‌هایی از درجهٔ پایین (به بخش ۴.۶ نگاه کنید) است. روش مناسب دیگر برای تقریب زدن  $D(f)$ ، استفاده از  $D(g_p)$  است، که در آن  $g_p(x)$  یک درونیاب درجهٔ سوم قلمی است که  $f(x)$  را در تعدادی از نقاط درونیابی می‌کند، یا به بهترین نحو،  $f(x)$  را به تعبیر کوچکترین توانهای دوم تقریب می‌زند.

## تمرین

۱-۱۰۷ با استفاده از جدول زیر و روابط  $(7.7)$ ،  $(8.7)$  و  $(10.7)$ ،  $f'(1.4)$  را تعیین کنید. همچنین با استفاده از رابطهٔ  $(14.7)$ ،  $f''(1.4)$  را به دست آورید. نتایج به دست آمده را با  $2.1509$  و  $f'(1.4) = \cosh 1.4 = 1.90433$  و  $f''(1.4) = \sinh 1.4 = 1.90433$  مقایسه نمایید. که تا تعداد ارقام داده شده دقیق هستند، مقایسه نمایید.

$x$	$f(x)$
۱.۲	۱.۵۰۹۵
۱.۳	۱.۶۹۸۴
۱.۴	۱.۹۰۴۳
۱.۵	۲.۱۲۹۳
۱.۶	۲.۳۷۵۶

۲-۱۰۷ با استفاده از جدول مقادیر زیر که به تابع  $f(x) = \sinh x$  مربوط است و استفاده

از رابطه (۸.۷) به ازای  $h = 0.001$  و  $h = 0.002$ ،  $f'(0.4)$  را تعیین کنید، کدامیک از مقادیر به دست آمده دقیق هستند؟ جواب درست برابر است با

$$f'(0.4) = \cosh 0.4 = 1.0810720.$$

$x$	$f(x)$
0.398	0.408591
0.399	0.409671
0.400	0.410752
0.401	0.411834
0.402	0.412915

۳-۱۰۷ در معادله (۱۶.۷)،  $f(x) = \sinh x$ ، گیریم و فرض می‌کنیم که در طی محاسبه  $\sinh x$ ، میزان خطای گرد کردن ثابت باقی بماند به طوری که داشته باشیم

$$E_+ - E_- = 0.5 \times 10^{-7}.$$

اگر فرمول (۸.۷) برای محاسبه  $f'(0)$  به کار رود، مقدار بهینه  $h$  را که می‌باید در فرمول فوق به کار رود معین کنید.

۴-۱۰۷ با سه بار مشتقگیری از رابطه (۱۰.۷) و با انتخاب  $k = 3$ ، به ازای  $a = x_0$ ،  $x_1 = a - h$ ،  $x_2 = a + h$ ،  $x_3 = a + 2h$  فرمولی برای  $f'''(a)$  به دست آورید. همچنین جمله خطا را برای این فرمول به دست آورید.

۵-۱۰۷ به کمک کامپیوتر دنباله اعداد

$$a_n = f[2 - 2^{-n}, 2 + 2^{-n}] \quad n = 1, 2, 3, \dots$$

را که در آن  $f(x) = \ln x$  است، محاسبه کنید. بدون در نظر گرفتن اثرات خطای گرد کردن داریم

$$\lim_{n \rightarrow \infty} a_n = f'(2) = 0.5$$

و با بحثی که در این قسمت آمد، به دلیل خطای گرد کردن، داریم

$$\lim_{n \rightarrow \infty} |a_n| = \infty$$

آیا این امر واقعاً اتفاق می‌افتد؟ اگر نمی‌افتد، چرا؟ این امر بحثی را که در متن انجام گرفته از اعتبار می‌اندازد.

۶-۱۰۷ با بسط  $f(a+h)$  و  $f(a-h)$  به سری تیلر پیرامون نقطهٔ  $a$ ، فرمول (۸.۷) را بررسی کنید.

۷-۱۰۷ با استفاده از بسط سری تیلر، فرمول (۱۲.۷) را برای  $f''(a)$  به دست آورید.

## ۲.۷ انتگرالگیری عددی: برخی قواعد اساسی

مسئلهٔ انتگرالگیری عددی یا ترییع، عبارت است از برآورد کردن عدد

$$I(f) = \int_a^b f(x) dx \quad (17.7)$$

این مسئله زمانی پیش می‌آید که نتوان انتگرالگیری را به‌طور دقیق انجام داد، یا هنگامی که  $f(x)$  تنها در تعداد متناهی از نقاط، داده شده باشد.

بدین دلیل، خط مشیبی را که در آغاز این فصل بیان کردیم پیش می‌گیریم.  $I(f)$  را به وسیلهٔ  $I(p_k)$  تقریب می‌زنیم، که در اینجا  $p_k(x)$  یک بسجمله‌ای است از درجهٔ  $n$  بیشتر از  $k$  که در نقاط  $x_0, \dots, x_k$  بر  $f(x)$  منطبق است. این تقریب معمولاً به شکل یک قاعده یعنی به صورت حاصلجمع وزین

$$I(p_k) = A_0 f(x_0) + A_1 f(x_1) + \dots + A_k f(x_k)$$

از مقادیر تابع  $f(x_0), \dots, f(x_k)$  نوشته می‌شود. وزنه‌های توانند به صورت  $A_i = I(l_i)$  که در آن  $l_i(x)$   $i$ امین بسجمله‌ای لاگرانژ است محاسبه شوند.

اکنون فرض کنید که عامل زیر علامت انتگرال  $f(x)$  در یک بازهٔ  $[c, d]$  که متضمن نقاط  $a$  و  $b$  است، به اندازهٔ کافی هموار باشد به طوری که بتوانیم آن را مانند (۳۷.۲) به شکل

$$f(x) = p_k(x) + f[x_0, \dots, x_k, x] \psi_k(x)$$

بنویسیم، که  $\psi_k(x)$  در تساوی زیر صدق می‌کند

$$\psi_k(x) = \prod_{j=0}^k (x - x_j)$$

در این صورت مقدار خطا در برآورد  $I(p_k)$  برای  $I(f)$  به شرح زیر است

$$E(f) = I(f) - I(p_k) = \int_a^b f[x_0, \dots, x_k, x] \psi_k(x) dx \quad (18.7)$$

که در آن  $f[x_0, \dots, x_k, x]$  پیوسته و از این رو بنا بر قضیه ۵.۲ به صورت تابعی از  $x$  انتگرالپذیر است.

گاهی این عبارت خطا می‌تواند ساده شود. برای مثال اگر  $\psi_k(x)$  روی  $(a, b)$  دارای یک علامت باشد، آنگاه بنا بر قضیه مقدار میانگین برای انتگرالها (به قسمت ۷.۱ نگاه کنید) داریم،

به ازای مقداری مانند  $\xi \in (a, b)$

$$\int_a^b f[x_0, \dots, x_k, x] \psi_k(x) dx = f[x_0, \dots, x_k, \xi] \int_a^b \psi_k(x) dx \quad (19.7)$$

بعلاوه اگر  $f(x)$  به طور پیوسته  $(k+1)$  مرتبه بر  $(c, d)$  مشتقپذیر باشد، با توجه به (۱۸.۷) و (۱۹.۷) خواهیم داشت

$$E(f) = \frac{1}{(k+1)!} f^{(k+1)}(\eta) \int_a^b \psi_k(x) dx \quad \eta \in (c, d) \quad \text{مانند } (20.7)$$

حتی اگر  $\psi_k(x)$  دارای یک علامت نباشد، تا حدودی ساده کردن عبارت خطای (۱۸.۷) امکان پذیر است. یک مورد خاص مطلوب از این نوع وقتی است که داشته باشیم

$$\int_a^b \psi_k(x) dx = 0 \quad (21.7)$$

در چنین حالتی می‌توانیم از همانی<sup>۱</sup>

$$f[x_0, \dots, x_k, x] = f[x_0, \dots, x_k, x_{k+1}] + f[x_0, \dots, x_{k+1}, x](x - x_{k+1})$$

که به ازای مقدار دلخواهی مانند  $x_{k+1}$  معتبر است، برای یافتن معادله

$$\begin{aligned} E(f) &= \int_a^b f[x_0, \dots, x_{k+1}] \psi_k(x) dx \\ &\quad + \int_a^b f[x_0, \dots, x_{k+1}, x](x - x_{k+1}) \psi_k(x) dx \\ &= \int_a^b f[x_0, \dots, x_{k+1}, x] \psi_{k+1}(x) dx \end{aligned}$$

استفاده کنیم، زیرا طرف دوم معادلهٔ اخیر به صورت زیر تبدیل می‌شود

$$\int_a^b f[x_0, \dots, x_{k+1}] \psi_k(x) dx = f[x_0, \dots, x_{k+1}] \int_a^b \psi_k(x) dx = 0$$

اکنون اگر بتوانیم  $x_{k+1}$  را به طریقی انتخاب کنیم که تابع  $\psi_{k+1}(x) = (x - x_{k+1})\psi_k(x)$  در بازهٔ  $(a, b)$  دارای یک علامت و  $f(x)$  تابعی  $(k+2)$  مرتبه مشتق‌پذیر باشد، آنگاه (مانند قبل) نتیجه می‌شود که

$$E(f) = \frac{1}{(k+2)!} f^{(k+2)}(\eta) \int_a^b \psi_{k+1}(x) dx \quad \eta \in (c, b) \text{ مانند } \quad (22.7)$$

باید توجه کرد که درجهٔ مشتق  $f(x)$  که در (۲۲.۷) ظاهر می‌شود، یک درجه بالاتر از درجهٔ مشتقی است که در (۲۰.۷) آمده بود. این امر مانند مشتق‌گیری عددی معرف آن است که مرتبهٔ (۲۲.۷) از مرتبهٔ (۲۰.۷) بالاتر است.

اکنون مثالهای خاصی را مورد بررسی قرار می‌دهیم. گیریم  $k=0$ ، آنگاه

$$f(x) = f(x_0) + f[x_0, x](x - x_0)$$

بنابراین

$$I(p_k) = (b-a)f(x_0)$$

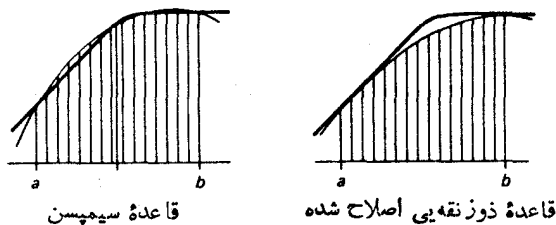
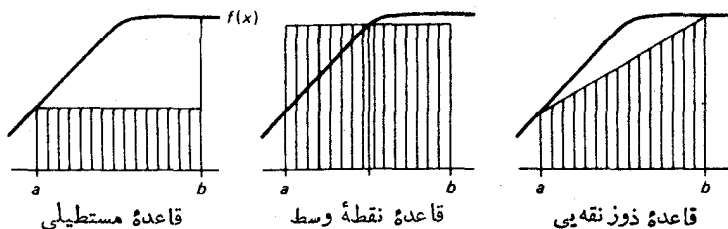
اگر  $x_0 = a$ ، آنگاه این تقریب به صورت زیر درمی‌آید

$$I(f) \approx R = (b-a)f(a) \quad (23.7)$$

که به قاعدهٔ مستطیلی<sup>۱</sup> معروف است (به شکل ۲.۷ نگاه کنید). چون در این حالت تابع  $\psi_0(x) = x - a$  روی  $(a, b)$  دارای یک علامت است، خطای  $E^R$  در قاعدهٔ مستطیلی می‌تواند از رابطهٔ (۲۰.۷) محاسبه شود. لذا داریم

$$E^R = f'(\eta) \int_a^b (x-a) dx = \frac{f'(\eta)(b-a)^2}{2} \quad (24.7)$$

اگر  $x_0 = (a+b)/2$ ، آنگاه  $\psi_0(x)$  دارای یک علامت نخواهد بود. ولی در این صورت داریم



شکل ۲.۷ انتگرال عددی

$$\int_a^b (x - x_0) dx = 0$$

در حالی که  $(x - x_0)^2$  دارای یک علامت است. بنابراین، در این حالت مقدار خطا در  $I(p_k)$  را می‌توان از (۲۲.۷) به‌ازای  $x_1 = x_0$  محاسبه کرد. که خواهیم داشت

$$I(f) \approx M = (b-a)f\left(\frac{a+b}{2}\right)$$

$$E^M = \frac{f'''(\eta)(b-a)^3}{24}$$

به‌ازای مقداری مانند  $\eta \in (a, b)$

(۲۵.۷)

که به‌قاعدۀ نقطۀ وسط موسوم است. سپس گیریم  $k=1$ ، آنگاه

$$f(x) = f(x_0) + f[x_0, x_1](x - x_0) + f[x_0, x_1, x]\psi_1(x)$$

برای بدست آوردن  $\psi_1(x) = (x - x_0)(x - x_1)$  که روی  $(a, b)$  دارای یک علامت

## 1. midpoint rule



باشد، مقادیر  $x_0 = a$  و  $x_1 = b$  را انتخاب می‌کنیم. در این صورت، با توجه به رابطه (۲۰.۷) خواهیم داشت

$$I(f) = \int_a^b \{f(a) + f[a, b](x-a)\} dx + \frac{1}{4} f''(\eta) \int_a^b (x-a)(x-b) dx$$

یا

$$I(f) \approx T = \frac{1}{2}(b-a)[f(a) + f(b)]$$

$$E^T = -\frac{f''(\eta)(b-a)^3}{12}$$

به ازای مقداری مانند  $\eta \in (a, b)$

(۲۶.۷)

که به قاعدهٔ ذوزنقه‌ای موسوم است (به شکل ۲.۷ نگاه کنید). اکنون، گیریم  $k=2$ ، در این صورت

$$f(x) = p_2(x) + f[x_0, x_1, x_2, x]\psi_2(x)$$

باید توجه داشت که برای نقاط متمایز  $x_0, x_1, x_2$  در بازهٔ  $(a, b)$ ، تابع

$$\psi_2(x) = (x-x_0)(x-x_1)(x-x_2)$$

روی  $(a, b)$  دارای يك علامت نیست. اما اگر  $x_0 = a$  و  $x_1 = (a+b)/2$  و  $x_2 = b$  انتخاب کنیم، آنگاه با انتگرالگیری مستقیم یا استدلالهای مقارن، می‌توان نشان داد که

$$\int_a^b \psi_2(x) dx = \int_a^b (x-a)(x-(a+b)/2)(x-b) dx = 0$$

که در این حالت خطا به شکل رابطه (۲۲.۷) خواهد بود. اگر  $x_2 = x_1 = (a+b)/2$  انتخاب کنیم، آنگاه داریم

$$\psi_2(x) = (x-a)\left(x - \frac{a+b}{2}\right)^2(x-b)$$

که روی  $(a, b)$  دارای يك علامت است. بنابراین از روابط (۱۸.۷) و (۲۲.۷) نتیجه می‌شود که

$$I(f) = I(p_2) + \frac{1}{4!} f^{iv}(\eta) \int_a^b \psi_2(x) dx$$

## 1. trapezoid (al) rule

رابطه زیر مستقیماً محاسبه می‌شود

$$\int_a^b \psi_2(x) dx = \int_a^b (x-a) \left(x - \frac{a+b}{2}\right)^2 (x-b) dx = -\frac{2}{15} \left(\frac{b-a}{2}\right)^5$$

لذا مقدار خطا در این فرمول به صورت زیر درخواهد آمد

$$E^s(f) = \frac{-1}{90} \left(\frac{b-a}{2}\right)^5 f^{iv}(\eta) \quad \eta \in [a, b]$$

اکنون برای به دست آوردن فرمولی متناظر با حالت  $k=2$ ، با انتخاب نقاط درونیاب  $x_0 = a$ ،  $x_1 = (a+b)/2$ ،  $x_2 = b$ ، تابع  $I(p_2)$  را مستقیماً محاسبه می‌کنیم. بهتر است بسجمله‌ای درونیاب را به شکل زیر بنویسیم

$$p_2(x) = f(a) + f[a, b](x-a) + f\left[a, b, \frac{a+b}{2}\right](x-a)(x+b)$$

بنابراین داریم

$$\begin{aligned} \int_a^b p_2(x) dx &= f(a)(b-a) + f[a, b](b-a)^2/2 \\ &+ f\left[a, b, \frac{a+b}{2}\right] \int_a^b (x-a)(x-b) dx \end{aligned}$$

اما همان گونه که هنگام به دست آوردن رابطه (۲۶.۷)، پیدا کردیم،

$$\int_a^b (x-a)(x-b) dx = -(b-a)^3/6$$

بنابراین

$$\begin{aligned} \int_a^b p_2(x) dx &= f(a)(b-a) + f[a, b](b-a)^2/2 \\ &- f\left[a, \frac{a+b}{2}, b\right](b-a)^3/6 \end{aligned} \quad (27.7)$$

که در رابطه بالا از تساوی

$$f\left[a, b, \frac{a+b}{2}\right] = f\left[a, \frac{a+b}{2}, b\right]$$

که به علت وجود تقارن در تفاضل منقسم برقرار است، استفاده شده است. اما اکنون داریم

$$f[a, b](b-a) = f(b) - f(a)$$

در حالی که

$$\begin{aligned} f\left[a, \frac{a+b}{2}, b\right](b-a)^2 &= \left(f\left[\frac{a+b}{2}, b\right] - f\left[a, \frac{a+b}{2}\right]\right)(b-a) \\ &= 2\left(f(b) - f\left(\frac{a+b}{2}\right) - \left(f\left(\frac{a+b}{2}\right) - f(a)\right)\right) \end{aligned}$$

از گذاردن این عبارتها در (۲۷.۷) خواهیم داشت

$$\begin{aligned} \int_a^b p_2(x) dx &= (b-a) \left\{ f(a) + (f(b) - f(a))/2 \right. \\ &\quad \left. - 2\left(f(b) - 2f\left(\frac{a+b}{2}\right) + f(a)\right) / 6 \right\} \\ &= \frac{b-a}{6} \left\{ f(a) + 2f\left(\frac{a+b}{2}\right) + f(b) \right\} \end{aligned}$$

بنابراین درست به قاعدهٔ سیمپسن رسیدیم که مقدار خطای متناظر آن چنین است

$$I(f) \approx S = \frac{b-a}{6} \left\{ f(a) + 2f\left(\frac{a+b}{2}\right) + f(b) \right\}$$

$$E^S = -\frac{f^{iv}(\eta)[b-a]^5/2^5}{90} \quad (28.7)$$

و سرانجام می‌گیریم  $k=3$ . در این صورت

$$f(x) = p_3(x) + f[x_0, x_1, x_2, x_3, x]\psi_4(x)$$

یا انتخاب  $x_2 = x_3 = b$ ,  $x_0 = x_1 = a$  می‌توان اطمینان داشت که تابع

$$\psi_4(x) = (x-a)^2(x-b)^2$$

در بازهٔ  $(a, b)$  يك علامت دارد و از این رو از رابطهٔ (۲۵.۷)، نتیجه می‌گیریم که خطا را می‌توان به صورت زیر بیان کرد

$$E(f) = \frac{1}{4!} f^{iv}(\eta) \int_a^b (x-a)^2(x-b)^2 dx = \frac{f^{iv}(\eta)(b-a)^5}{720}$$

برای به دست آوردن يك فرمول انتگرالگیری متناظر با انتخاب نقاط  $x_0 = x_1 = a$ ،  $x_2 = x_3 = b$  در ابتدا مشاهده می‌کنیم که

$$p_3(x) = f[a] + f[a, a](x-a) + f[a, a, b](x-a)^2 + f[a, a, b, b](x-a)^2(x-b)$$

لذا

$$\int_a^b p_3(x) dx = f(a)(b-a) + f[a, a] \frac{(b-a)^2}{2} + f[a, a, b] \frac{(b-a)^3}{3} + f[a, a, b, b] \left\{ \frac{(b-a)^4}{2} - \frac{(b-a)^4}{3} \right\} \quad (29.7 \text{ الف})$$

بر اساس درونیابی بوسانی بخش ۷.۲، داریم

$$f[a, a] = f'(a)$$

$$f[a, a, b] = \{f[a, b] - f'(a)\} / (b-a)$$

$$f[a, a, b, b] = \{f'(b) - 2f[a, b] + f'(a)\} / (b-a)^2$$

با گذاردن روابط بالا در (29.7 الف) و ساده کردن آن خواهیم داشت

$$\int_a^b p_3(x) dx = f(a)(b-a) + f'(a) \frac{(b-a)^2}{2} + \{f[a, b] - f'(a)\} \frac{(b-a)^3}{3} - \{f'(b) - 2f[a, b] + f'(a)\} \frac{(b-a)^4}{12}$$

و بالاخره از قرار دادن  $(f(b) - f(a)) / (b-a)$  به جای  $f[a, b]$  و مرتب کردن رابطه حاصله بر حسب توانهای  $(b-a)$  به فرمول زیر می‌رسیم

$$I(f) \approx CT = \frac{(b-a)}{2} [f(a) + f(b)] + \frac{(b-a)^2}{12} [f'(a) - f'(b)]$$

(29.7 ب)

که به دلایلی که محتاج به توضیح نیست، به قاعدهٔ دوزنقه‌ی اصلاح شده معروف است. خطای قاعدهٔ دوزنقه‌ی تصحیح شده چنین است

$$E^{CT} = \frac{f^{iv}(\eta)(b-a)^5}{720}$$

البته اگر قواعد فوق‌الذکر در مورد انتگرالگیری عددی، تقریب رضایتبخشی برای  $I(f)$  به دست ندهند، می‌توان درجهٔ  $k$ ی بسجمله‌ای درونیاب به کار رفته را افزایش داده خطرات انجام چنین عملی در بخش ۷.۶ مورد بحث قرار گرفت و در آنجا به کارگیری درونیابی بسجمله‌ای تکه‌ای<sup>۱</sup> پیشنهاد شد که روش بسیار منطقی و مناسبی برای دستیابی به دقت بالاست. بنا بر این  $I(f)$  را به کمک  $I(g_k)$  تقریب می‌زنیم که در آن  $g_k(x)$  يك تابع بسجمله‌ای تکه‌ای از درجهٔ «پایین»  $k$  است که  $f(x)$  را درونیابی می‌کند. قواعد انتگرالگیری حاصل، که معمولاً قواعد مرکب<sup>۲</sup> نامیده می‌شوند، در قسمت ۴.۷ مورد بحث قرار خواهند گرفت.

در این بخش پنج قاعدهٔ اصلی انتگرالگیری را به دست آورده‌ایم. این قاعده‌ها عبارت‌اند از قاعدهٔ مستطیلی (۲۴.۷)، قاعدهٔ نقطهٔ وسط (۲۵.۷)، قاعدهٔ دوزنقه‌یی (۲۶.۷)، قاعدهٔ سیمپسن (۲۸.۷) و قاعدهٔ دوزنقه‌یی اصلاح شده (۲۹.۷). از میان قواعد مذکور تنها، قاعدهٔ دوزنقه‌یی اصلاح شده است که به مشتق  $f(x)$  نیاز دارد و این امر يك عیب آشکار بر این روش خاص است. عبارتهای خطای مربوط به این قواعد نشان می‌دهند که هر گاه تابع  $f(x)$  به اندازهٔ کافی هموار باشد، استفاده از قاعدهٔ سیمپسن و یا قاعدهٔ دوزنقه‌یی اصلاح شده، می‌بایست ارجح باشد. با وجود این، توابعی وجود دارند که فرمولهای پایین مرتبه نسبت به فرمولهای بالا مرتبه، نتیجهٔ بهتری را به دست می‌دهند. [به تمرین ۲-۲.۷ نگاه کنید].

□ مثال ۱۰۷: برای پیدا کردن برآوردی برای انتگرال زیر، از هر يك از پنج قاعدهٔ مذکور در بالا استفاده کنید

$$I = \int_0^1 e^{-x^2} dx$$

مقادیر  $a=0$ ،  $b=1$ ،  $(a+b)/2 = 1/2$  را برمی‌گزینیم و با توجه به جدول مقادیر خواهیم داشت

$$f(0) = 1 \quad f(1) = e^{-1} = 0.36788 \quad f\left(\frac{1}{2}\right) = e^{-1/4} = 0.77880$$

در محاسبات به مقادیر زیر نیز احتیاج داریم

$$f'(0) = 0 \quad f'(1) = -2e^{-1} = -0.73576$$

در این صورت می‌توانیم مقادیر زیر را از فرمولهای مناسب، محاسبه کنیم

$$R = 1 \times e^0 = 1$$

$$M = 1 \times e^{-1/4} = 0.77880$$

$$T = \frac{1}{4} [e^0 + e^{-1}] = 0.68292$$

$$S = \frac{1}{6} [e^0 + 2e^{-1/4} + e^{-1}] = 0.74718$$

$$CT = \frac{1}{4} [e^0 + e^{-1}] + \frac{1}{12} [0 + 2e^{-1}] = 0.74525$$

مقدار انتگرال تسا پنج رقم صحیح اعشاری برابر با  $I = 0.74682$  است. چنانچه از بررسی عبارتهای خطا برمی‌آید، و این حقیقت که چند مشتق اول تابع  $e^{-x^2}$  از نظر مقدار خیلی تغییر نمی‌کنند، بدیهی است که قاعدهٔ دوزنقه‌ی اصلاح شده (CT) و قاعدهٔ سیمپسن (S) بهترین نتایج را به‌دست می‌دهند. □

## تمرین

۱-۲۰۷ اگر تابع

$$\psi_2(x) = (x-a)(x-(a+b)/2)(x-b)$$

مفروض باشد، با روش انتگرالگیری مستقیم تحقیق کنید  $\int_0^1 \psi_2(x) dx = 0$ .

۲-۲۰۷ با استفاده از هریک از پنج قاعده‌ای که در این بخش بیان کردیم، تقریبی برای  $I = \int_0^1 x \sin x dx$  به‌دست آورید. نتایج حاصله را با مقدار صحیح

$$I = \sin 1 - \cos 1 = 0.301169$$

مقایسه کنید.

۳-۲۰۷ تابع  $f(x)$  بر بازهٔ  $[0, 1]$  به‌شرح زیر تعریف شده‌است

$$f(x) = \begin{cases} x & 0 \leq x \leq 1/2 \\ 1-x & 1/2 \leq x \leq 1 \end{cases}$$

با استفاده از قواعد زیر، مقدار  $\int_0^1 f(x) dx$  را محاسبه کنید:  
الف) قاعدهٔ دوزنقه‌ی بر بازهٔ  $[0, 1]$ .

(ب) قاعدهٔ ذوزنقه‌یی ابتدا بر بازهٔ  $[0, 1/2]$  و سپس بر بازهٔ  $[1/2, 1]$ .

(پ) قاعدهٔ سیمپسن بر بازهٔ  $[0, 1]$ .

(ت) قاعدهٔ ذوزنقه‌یی اصلاح‌شده بر بازهٔ  $[0, 1]$ .

علت وجود تفاوت در نتایج حاصله را بیان کنید.

۴-۲۰۷ با ملاحظهٔ اینکه  $p_2(x)$  یک بسجمله‌ای از درجهٔ ۳ است و  $p_2''(x) = 0$  و نیز اینکه قاعدهٔ سیمپسن (۲۸-۷) می‌تواند دقیقاً برای محاسبهٔ  $I(p_2)$  به‌کار رود، در نتیجه قاعدهٔ ذوزنقه‌یی اصلاح‌شده را می‌توان بسیار ساده‌تر به‌دست آورد. بنابراین

$$I(p_2) = \frac{b-a}{6} \left\{ p_2(a) + 4p_2\left(\frac{a+b}{2}\right) + p_2(b) \right\}$$

چون  $p_2(x)$  درونیاب  $f(x)$  در نقاط  $a, b, a, b$  است، باید داشته باشیم

$$p_2(a) = f(a), \quad p_2(b) = f(b)$$

با استفاده از نتایج حاصله از بخش ۷-۲ در مورد درونیابی بوسانی، نشان دهید که

$$p_2\left(\frac{a+b}{2}\right) = \frac{1}{4} [f(a) + f(b)] + \frac{b-a}{8} [f'(a) - f'(b)]$$

سپس رابطهٔ بالا را در رابطهٔ  $I(p_2)$  قرار دهید تا قاعدهٔ ذوزنقه‌یی اصلاح‌شده (۲۹-۷) به‌دست آید.

۵-۲۰۷ برای برآورد مقدار انتگرال

$$I = \int_0^1 (1-x^2)^{3/2} dx$$

از قاعدهٔ سیمپسن استفاده کنید.

۶-۲۰۷ برای برآورد مقدار انتگرال  $I = \int_0^1 xe^{-x^2} dx$ ، قاعدهٔ ذوزنقه‌یی را به‌کار برید. کرانی را در خطای حاصل از قاعدهٔ ذوزنقه‌یی (۲۶-۷) به‌دست آورید و آن را با خطای واقعی مقایسه کنید.

### ۳-۷ انتگرالگیری عددی: قواعد گاوسی

همهٔ قواعدی را که در بخش ۲-۷ استخراج کردیم، به‌استثنای قانون ذوزنقه‌یی اصلاح‌شده، می‌توان به‌شکل زیر نوشت

$$I(g) \approx A_0 g(x_0) + A_1 g(x_1) + \dots + A_k g(x_k) \quad (30.7)$$

که در آن وزنه‌های  $A_0, \dots, A_k$  به تابع خاص  $g(x)$  بستگی ندارند. تا به حال به نحوی رأسهای  $x_0, \dots, x_k$  را انتخاب کرده‌ایم، مثلاً آنها را در جدولی با فواصل مساوی در نظر گرفتیم و سپس به ازای جمیع مقادیر  $z$ ، وزنه‌های  $A_i$  را به صورت  $I(I_i)$  محاسبه کردیم. این عمل تضمین می‌کند که قاعده مزبور برای بسجمله‌هایی ناپیشتتر از درجه  $k$  دقیق باشد. اما این امکان وجود دارد که با انتخاب رأسهای مناسب، این قاعده را برای بسجمله‌هایی ناپیشتتر از درجه  $1 - 2k$  نیز دقیق درآورد. این مسئله زیربنای فکری قواعد گاوس است.

قواعد حاصله پیچیده‌تر از قواعد مستخرجه در بخش ۲.۷ به نظر می‌رسند. در این قواعد، هم رأسها و هم وزنها، در حالت کلی، اعداد گنگ  $2$  هستند. این امر موجب می‌شود که وقتی بنا باشد افراد محاسبات را به وسیله دست انجام دهند، عموماً از به کار بردن این قواعد خودداری کنند. اما با یک کامپیوتر، معمولاً تفاوتی ندارد که محاسبه  $y_k$  تابع در  $x = 3$  انجام گیرد یا  $0.269057735058 \approx 1/\sqrt{3}$ . هنگامی که رأسها و وزنها در چنین قواعدی به شکل خاصی (مثلاً به صورت، زیریرنامه LGNDRE که در زیر آمده است) ذخیره شده باشند، این قواعد به همان آسانی قاعده دوزنقه‌ی یا قاعده سیمپسن به کار برده می‌شوند. در عین حال، بر اساس تعداد مقادیری که برای تابع به کار برده شده، معمولاً قواعد گاوسی در مقایسه با قواعد بخش ۲.۷ بسیار دقیق‌ترند.

ما این قواعد گاوسی را در قالب کلیتر یک انتگرال  $\int_a^b f(x) dx$  که در آن عامل زیر علامت انتگرال،  $f(x)$ ، به اندازه کافی برای توجیه کاربرد قواعد بخش ۲.۷ مشتقبذیر نیست، مورد بحث قرار می‌دهیم. مثلاً  $f(x)$  به ازای مقداری مانند  $-1 < \alpha$  و در نزدیکی نقطه  $a$ ، ممکن است به صورت  $(x-a)^\alpha$  عمل کند یا آنکه  $a$  و  $b$  ممکن است نامتناهی باشند. غالباً در این گونه موارد می‌توان انتگرال را به صورت

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

نوشت، که در آن  $w(x)$  تابعی است نامنفی انتگرالپذیر و

$$g(x) = \frac{f(x)}{w(x)}$$

تابعی است هموار. در مثال بالا،  $w(x) = (x-a)^\alpha$ ، همین وضع را دارد. انتخابهای دیگری که برای  $w(x)$  صورت می‌گیرد در زیر مورد بحث قرار می‌گیرد. حالتی که تابع زیر علامت انتگرال تابع بی‌دردسری است نیز در اینجا با انتخاب ساده  $1 \equiv w(x)$ ، مورد بحث قرار می‌گیرد.

حال محاسبه تقریبی انتگرال وزین

$$I(g) = \int_a^b g(x)w(x) dx \quad (31.7)$$



را با استفاده از قاعده‌ای به شکل (۳۰.۷) در نظر می‌گیریم. گوئیم که قاعده (۳۰.۷) برای تابع خاص  $p(x)$  صحیح است، اگر بنا بر قرار دادن  $p(x)$  به جای  $g(x)$  در (۳۰.۷)، این رابطه به یک تساوی بدل شود. برای مثال، قاعدهٔ دوزنقه‌یی

$$\int_a^b g(x) dx \approx \frac{b-a}{2} g(a) + \frac{b-a}{2} g(b)$$

برای همهٔ بسجمله‌های از درجهٔ نایبتر از یک، صحیح است. برای بررسی این امر فقط کافی است به عبارت خطای این قاعده، یعنی

$$E^T = \frac{g''(\eta)(b-a)^3}{12}$$

نگاه کنیم. از آنجا که این عبارت خطا متضمن مشتق دوم  $g(x)$  است، و چون مشتق دوم هر بسجمله‌ای از درجهٔ نایبتر از یک برابر با صفر است، لذا نتیجه می‌شود که هر گاه  $g(x)$  یک بسجمله‌ای از درجهٔ نایبتر از یک باشد، خطا برابر با صفر است. به صورت کلیتر، اگر عبارت خطای (۳۰.۷) به شکل

$$E = g^{(r+1)}(\eta) \times (\text{تابعی از } x_0, \dots, x_k) \quad (32.7)$$

باشد، آنگاه قاعده (۳۰.۷) باید برای تمامی بسجمله‌ایهای نایبتر از درجهٔ  $r$ ، صحیح باشد.

بنا بر این، اگر بخواهیم قاعده‌ای به شکل (۳۰.۷) بسازیم که، با  $k$  ثابت برای بسجمله‌ایهای از درجات بالاتر صحیح باشد، می‌باید قاعده‌ای بسازیم که عبارت خطای آن به شکل (۳۲.۷)، یعنی به صورت یک عدد صحیح تاحد امکان بزرگ، باشد. این کار را می‌توانیم با ترفندی که در بخش ۲.۷ به کار بردیم انجام دهیم.

مانند بخش (۲.۷)، از جایگذاری تحلیلی استفاده می‌کنیم یعنی نقاط  $x_0, \dots, x_k$  را در بازهٔ  $(a, b)$  انتخاب می‌کنیم و می‌نویسیم

$$g(x) = p_k(x) + g[x_0, \dots, x_k, x] \psi_k(x)$$

که در آن  $p_k(x)$  یک بسجمله‌ای است از درجهٔ نایبتر از  $k$  که درونیاب  $g(x)$  در نقاط  $x_0, \dots, x_k$  است و داریم

$$\psi_k(x) = (x - x_0) \dots (x - x_k)$$

که نتیجه می‌دهد

$$I(g) = I(p_k) + \int_a^b g[x_0, \dots, x_k, x] \psi_k(x) w(x) dx$$

واضح است که تقریب  $I(p_k)$  برای  $I(g)$  به شکل (۳۰.۷) است. زیرا اگر  $p_k(x)$  را به شکل لاگرانژ بنویسیم (به بخش ۲۰.۲ نگاه کنید)

$$p_k(x) = g(x_0)l_0(x) + g(x_1)l_1(x) + \dots + g(x_k)l_k(x)$$

که در آن داریم

$$l_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^k \frac{x - x_j}{x_i - x_j} \quad i = 0, \dots, k$$

در این صورت

$$\begin{aligned} I(p_k) &= \int_a^b p_k(x)w(x) dx \\ &= g(x_0) \int_a^b l_0(x)w(x) dx + \dots + g(x_k) \int_a^b l_k(x)w(x) dx \end{aligned}$$

بنابراین

$$I(p_k) = A_0 g(x_0) + A_1 g(x_1) + \dots + A_k g(x_k) \quad (33.7)$$

که در آن

$$A_i = \int_a^b l_i(x)w(x) dx \quad i = 0, \dots, k \quad (34.7)$$

سپس خطای

$$I(g) - I(p_k) = \int_a^b g[x_0, \dots, x_k, x] \psi_k(x) w(x) dx$$

را در نظر می‌گیریم. فرض کنید که

$$\int_a^b \psi_k(x) w(x) dx = 0$$

سپس همان گونه که در بخش ۲.۷ بحث شد، به ازای هر انتخابی از  $x_{k+1}$  داریم

$$I(g) - I(p_k) = \int_a^b g[x_0, \dots, x_k, x_{k+1}, x] \psi_{k+1}(x) w(x) dx$$

اکنون اگر تساوی

$$\int_a^b \psi_{k+1}(x) w(x) dx = 0$$

نیز برقرار باشد، آنگاه بهمان دلیل داریم

$$I(g) - I(p_k) = \int_a^b g[x_0, \dots, x_{k+2}, x] \psi_{k+2}(x) w(x) dx$$

بنابراین، در حالت کلی، اگر برای نقاطی مانند  $x_0, \dots, x_{k+m}$  داشته باشیم

$$\int_a^b \psi_k(x) (x - x_{k+1}) \dots (x - x_{k+1+i}) w(x) dx = 0$$

$$i = 0, \dots, m-1 \quad (35.7)$$

آنگاه برای هر انتخابی از  $x_{k+m+1}$  داریم

$$I(g) - I(p_k) = \int_a^b g[x_0, \dots, x_{k+m+1}, x] \psi_{k+m+1}(x) w(x) dx \quad (36.7)$$

اکنون با توجه به بخش ۶.۶ یادآوری می‌کنیم که برای بسیاری از  $w(x)$ ‌ها می‌توانیم یک بسجمله‌ای  $p_{k+1}(x)$  چنان پیدا کنیم که به‌ازای تمام بسجمله‌ایهای  $q(x)$  نا بیشتر از درجهٔ  $k$  (به‌ویژگی ۳ی بسجمله‌ایهای متعامد در بخش ۶.۶ مراجعه کنید) داشته باشیم

$$\int_a^b p_{k+1}(x) q(x) w(x) dx = 0 \quad (37.7)$$

بعلاوه، بنا بر ویژگی ۲ی بسجمله‌ایهای متعامد، داریم

$$p_{k+1}(x) = \alpha_{k+1} (x - \xi_0)(x - \xi_1) \dots (x - \xi_k)$$

که در آن  $\xi_0, \dots, \xi_k$   $k+1$  نقطهٔ متمایز در بازهٔ  $(a, b)$  هستند که در آنها بسجمله‌ای  $p_{k+1}$  صفر می‌شود. لذا اگر بگیریم

$$x_j = \xi_j \quad j = 0, \dots, k \quad (38.7)$$

و به‌ازای  $z = 1, \dots, k+1$   $x_{k+z}$  را نقاط انتخابی در بازهٔ  $(a, b)$  برگزینیم، آنگاه به‌ازای  $m = k$ ، رابطهٔ (۳۵.۷) و در نتیجه رابطهٔ (۳۶.۷) برقرار خواهد بود. به‌ازای تمام بسجمله‌ایهای نا بیشتر از درجهٔ  $k$ ، رابطهٔ (۳۵.۷) به‌شکل (۳۷.۷) خواهد بود که در آن به‌ازای  $m \leq k$  داریم

$$q(x) = \frac{(x - x_{k+1}) \dots (x - x_{k+1+i})}{\alpha_{k+1}} \quad i = 0, \dots, m-1$$

بنابراین

$$I(g) - I(p_k) = \int_a^b g[x_0, \dots, x_{k+1}, x] \psi_{k+1}(x) w(x) dx \quad (39.7)$$

برای به دست آوردن این خطا به شکل (۳۲.۷)،  $x_{k+j}$  ها را به صورت زیر انتخاب می کنیم

$$x_{k+j} = \xi_{j-1} \quad j = 1, \dots, k+1$$

در این صورت خواهیم داشت

$$\begin{aligned} \psi_{\nu_{k+1}}(x) &= (x - x_0) \dots (x - x_{\nu_{k+1}}) \\ &= (x - \xi_0) \dots (x - \xi_k)(x - \xi_0) \dots (x - \xi_k) \\ &= \left[ \frac{P_{k+1}(x)}{\alpha_{k+1}} \right]^2 \end{aligned}$$

به طوری که  $\psi_{\nu_{k+1}}(x)w(x)$  دارای یک علامت، یعنی علامت نامنفی در بازه  $(a, b)$  خواهد بود. بنابراین می توان قضیه مقدار میانگین برای انتگرال را به کار برد (به بخش ۷۰۱ نگاه کنید) و رابطه

$$I(g) - I(p_k) = g[x_0, \dots, x_{\nu_{k+1}}, \eta] \int_a^b \left[ \frac{1}{\alpha_{k+1}} P_{k+1}(x) \right]^2 w(x) dx$$

را به دست آورد. و سرانجام اگر  $g(x)$  به طور پیوسته،  $\nu_{k+2}$  مرتبه مشتقپذیر باشد، می توانیم قضیه ۵.۲ را به کار گیریم و خطا را به شکل زیر بیان کنیم

$$I(g) - I(p_k) = \frac{1}{(\nu_{k+2})!} g^{(\nu_{k+2})}(\xi) \frac{S_{k+1}}{\alpha_{k+1}^{\nu_{k+2}}} \quad (40.7)$$

که در آن داریم

$$S_{k+1} = \int_a^b [P_{k+1}(x)]^2 w(x) dx$$

به طور خلاصه، نشان دادیم که اگر نقاط  $x_0, \dots, x_k$  را در رابطه (۳۳.۷) به عنوان ریشه های بسجمله ای  $P_{k+1}(x)$ ، از درجه  $(k+1)$  ام، برگزینیم و این بسجمله ای نسبت به تابع وزن  $w(x)$  در بازه  $(a, b)$  بر هر بسجمله ای نایبتر از درجه  $k$  متعامد باشد و اگر ضرایب  $A_i$  ( $i = 0, \dots, k$ ) در رابطه (۳۳.۷) بر اساس رابطه (۳۴.۷) انتخاب شوند، آنگاه فرمول گاوس حاصل، (۳۳.۷)، برای کلیه بسجمله ایهای نایبتر از درجه  $\nu_{k+1}$  صحیح است. قواعد انتگرالگیری از این نوع را، به تعبیر ما و در شرایط مذکور در بالا، «بهترین قاعده ممکن» گویند.

اکنون مثالهایی را بررسی می کنیم. در ابتدا می گیریم  $w(x) = 1$ . اگر  $(a, b)$  بازه متناهی باشد، آنگاه برای تغییر حدود انتگرالگیری از  $(a, b)$  تا  $(-1, 1)$ ، از تغییر خطی متغیرها به صورت  $x = [(b-a)t + (b+a)]/2$  استفاده می کنیم. با انجام این کار

خواهیم داشت

$$\int_a^b f(x) dx = \int_{-1}^1 f(x(t))x'(t) dt = \int_{-1}^1 f(x(t)) \frac{b-a}{2} dt \quad (۴۱.۷)$$

با فرض اینکه این تغییر قبلا انجام گرفته باشد، انتگرال (۴۱.۷) را که به شکل زیر است در نظر می‌گیریم

$$I(g) = \int_{-1}^1 g(x)w(x) dx$$

از آنجایی که داریم  $w(x) = 1$ ، بسجمله‌ایهایی متعادل مناسب بسجمله‌ایهای لژاندر هستند (به مثال ۶.۶ نگاه کنید). در این حالت

$$P_1(x) = x \quad \xi_0 = 0$$

$$P_2(x) = \frac{3}{2} \left( x^2 - \frac{1}{3} \right) \quad \xi_0 = -\frac{1}{\sqrt{3}}, \xi_1 = \frac{1}{\sqrt{3}}$$

$$P_3(x) = \frac{5}{2} \left( x^3 - \frac{3}{5}x \right) \quad \xi_0 = -\sqrt{\frac{3}{5}}, \xi_1 = 0, \xi_2 = \sqrt{\frac{3}{5}}$$

و غیره. اگر  $k = 1$  انتخاب شود، آنگاه  $x_0 = \xi_0 = -1/\sqrt{3}$  و  $x_1 = \xi_1 = 1/\sqrt{3}$  و با قراردادن این مقادیر در روابط (۳۳.۷) و (۴۰.۷) خواهیم داشت

$$\int_{-1}^1 g(x) dx \approx A_0 g\left(\frac{-1}{\sqrt{3}}\right) + A_1 g\left(\frac{1}{\sqrt{3}}\right) \quad (۴۲.۷)$$

$$E = Cg^{iv}(\eta)$$

که در آن

$$A_1 = \int_{-1}^1 \frac{x - (-1/\sqrt{3})}{1/\sqrt{3} - (-1/\sqrt{3})} dx = 1$$

$$A_0 = \int_{-1}^1 \frac{x - 1/\sqrt{3}}{(-1/\sqrt{3}) - 1/\sqrt{3}} dx = 1$$

$$C = \frac{1}{4!} \frac{S_2}{\alpha_2^2} = \frac{1}{24} \times \frac{2}{5} \times \frac{4}{9} = \frac{1}{135}$$

چون

$$S_r = \int_{-1}^1 [P_r(x)]^2 dx = \frac{2}{5} \quad \alpha_r = \frac{3}{2}$$

با گذاردن این مقادیر ثابت در (۴۲.۷)، فرمول انتگرال دو نقطه‌ای گاوس

$$\int_{-1}^1 g(x) dx \approx g\left(\frac{-1}{\sqrt{3}}\right) + g\left(\frac{1}{\sqrt{3}}\right) \quad (۴۳.۷)$$

را با خطای آن

$$E = \frac{1}{135} g^{iv}(\eta) \quad (۴۴.۷)$$

به دست خواهیم آورد.

به ازای  $k > 1$ ، هم نقاط  $\xi_i$  و هم وزنهای  $A_i$ ، هر دو، اعداد گنگ می‌شوند. ولی محاسبه آنها ساده است. این رأسها و وزنهای  $k$  از  $0, \dots, 5$  در زیر برنامه فرعی فورترن به نام LGNDRE ذخیره می‌شوند. توجه کنید که محدودیتهای (قبلی) فورترن ما را وادار نموده که رأسها و وزنهای  $k$  از  $0$  تا  $k-1$  تا  $k+1$  NP شماره گذاری کنیم. بنابراین پارامتر ورودی NP به جای مشخص نمودن درجه بسجمله‌ای مورد بررسی، به مشخص کردن تعداد نقاط می‌پردازد.

```

SUBROUTINE LGNDRE ( NP , POINT , WEIGHT )
C SUPPLIES POINTS AND WEIGHTS FOR GAUSS-LEGENDRE QUADRATURE
C INTEGRAL(F(X), -1 .LE. X .LE. 1) IS APPROXIMATELY EQUAL TO
C SUM(F(POINT(I))*WEIGHT(I), I=1,...,NP) .
      INTEGER NP, I
      REAL POINT(NP), WEIGHT(NP)
      IF (NP .GT. 6) THEN
        PRINT 600, NP
600  FORMAT(' THE GIVEN NUMBER NP = ', I2, ' IS GREATER THAN 6.')
      *  /' EXECUTION STOPPED IN SUBROUTINE LGNDRE .'
      STOP
      END IF
      GO TO (1,2,3,4,5,6), NP
1  POINT(1) = 0.
   WEIGHT(1) = 2.
      GO TO 99
2  POINT(2) = .57735 02691 89626 D0
   WEIGHT(2) = 1.
      GO TO 95
3  POINT(2) = 0.
   POINT(3) = .7459 66692 41483 D0
   WEIGHT(2) = .88888 88888 88888 9 D0
   WEIGHT(3) = .55555 55555 55555 6 D0
      GO TO 95
4  POINT(3) = .33998 10435 84856 D0
   POINT(4) = .86113 63115 94053 D0
   WEIGHT(3) = .65214 51548 62546 D0
   WEIGHT(4) = .34785 48451 37454 D0
      GO TO 95
5  POINT(3) = 0.
   POINT(4) = .53846 93101 05683 D0
   POINT(5) = .90617 98459 38664 D0
   WEIGHT(3) = .56888 88888 88888 9 D0
   WEIGHT(4) = .47862 86704 99366 D0
   WEIGHT(5) = .23692 68850 56189 D0
      GO TO 95

```

```

6 POINT(4) = .23861 91860 83197 D0
POINT(5) = .66120 93864 66265 D0
POINT(6) = .93246 95142 03152 D0
WEIGHT(4) = .46791 39345 72691 D0
WEIGHT(5) = .36076 15730 48139 D0
WEIGHT(6) = .17132 44923 79170 D0

```

```

C
95 DO 96 I=1,NP/2
POINT(I) = -POINT(NP+1-I)
96 WEIGHT(I) = WEIGHT(NP+1-I)
99 RETURN
END

```

□ مثال ۲۰۷: به منظور مقایسه، می‌خواهیم دوباره مقدار  $I = \int_0^1 \exp(-x^2) dx$  را، ولی این بار با استفاده از فرمول پنج نقطه‌ای گاوس، محاسبه کنیم ( $k=4$ ). در اینجا تغییر متغیرهای لازم ( $41.7$ )، به صورت  $x = (t+1)/2$  است، بنابراین

$$I = \int_0^1 e^{-x^2} dx = \int_{-1}^1 \frac{1}{2} e^{-(t+1)^2/4} dt = \sum_{i=1}^4 A_i \exp(-(\xi_i+1)^2/4)/2$$

طبیعتاً برای انجام محاسبات از برنامه‌ای به صورت زیر استفاده خواهیم کرد.

```

EXAMPLE 7.2 GAUSSIAN INTEGRATION
REAL INTGRL,P(5),WEIGHT(5)
F(T) = EXP(-(1.+T)**2/4.)/2.
CALL LGNDRE ( 5, P, WEIGHT )
INTGRL = WEIGHT(1)*(F(P(1))+F(P(5))) + WEIGHT(2)*(F(P(2))+F(P(4)))
* + WEIGHT(3)*F(P(3))
PRINT 600,INTGRL
600 FORMAT(' EXAMPLE 7.2, GAUSS QUADRATURE',/ ' INTEGRAL = ',1PE14.7)
STOP
END

```

خروجی برنامه فوق به گونهٔ زیر است

$$\text{INTEGRAL} = 774682413.001$$

برای دسترسی به دقتی که نسبت به قانون وزن‌نمایی قابل قیاس باشد به ۲۸۰۰ قسمت فرعی  $\square$  احتیاج است، در حالی که قانون سیمپسن به تقریباً ۲۰ قسمت فرعی نیاز دارد.

□ مثال ۳۰۷: با استفاده از قاعدهٔ انتگرال‌گیری گاوس به ازای  $k=3$ ، برای انتگرال زیر تقریبی بیابید

$$I = \int_1^2 \frac{(\sin x)^2}{x} dx$$

(مقدار درست  $I = 0.79482518 \dots$ )

دوباره محدودهٔ انتگرال را به بازهٔ  $[-1, 1]$  تغییر می‌دهیم و این کار با تغییر متغیر  $x = t+2$  انجام می‌گیرد. این امر موجب می‌شود که داشته باشیم

$$I = \int_{-1}^1 \frac{(\sin(t+2))^2}{t+2} dt$$

بعد از آنکه برنامه اصلی<sup>۱</sup> در مثال ۲.۷ را به صورت زیر تغییر دادیم

$$F(T) = \text{SIN}(T + 2.) * 2 / (T + 2.)$$

$$\text{CALL LGNDRE}(4, P, \text{WEIGHT})$$

$$\text{INTGRL} = \text{WEIGHT}(1) * (F(P(1)) + F(P(4))) + \text{WEIGHT}(2) * (F(P(2)) + F(P(3)))$$

خروجی به گونه زیر خواهد شد

$$\square \quad \text{INTEGRAL} = ۷۹۴۸۲۸۳۳۰۰۱$$

هنگامی که با انتگرالهای تکین<sup>۲</sup> سروکار داریم، به ویژه فرمولهایی از نوع گاوس بسیار مفیدند. اگر برای مثال، محاسبه  $\int_a^b f(x) dx$ ، که در آن  $f(x)$  در  $a$  و  $b$  دارای تکینگی<sup>۳</sup> جبری می باشد مورد نظر باشد، آنگاه انتگرال را به انتگرال زیر تبدیل می کنند

$$\int_{-1}^1 g(x)w(x) dx$$

که در آن  $w(x)$  به ازای توانهای مناسب  $\alpha$  و  $\beta$  چنین است:

$$w(x) = (1-x)^\alpha (1+x)^\beta$$

در این حالت  $\xi_i$ ها، ریشه های بسجمله ای مناسب ژاکوبی هستند. در حالت خاص  $\alpha = \beta = -(1/2)$ ، بسجمله ایهای حاصل، درست همان بسجمله ایهای چیبیشف می شوند که در مثال ۷.۶ معرفی شدند و در قسمت ۱۰.۶ مورد بحث قرار گرفتند. برای این حالت خاص، قاعده جالب زیر به دست می آید

$$\int_{-1}^1 \frac{g(x)}{(1-x^2)^{1/2}} dx \approx \frac{\pi}{k+1} \sum_{i=0}^k g(\xi_i) \quad (۴۵.۷)$$

که در آن همه وزنهای  $A_i$  بر هم منطبق و  $\xi_i$ ها نقاط چیبیشف هستند [رابطه (۱۸.۶) را ببینید].

$$\xi_i = \cos\left(\frac{2i+1}{k+1} \times \frac{\pi}{2}\right) \quad i = 0, \dots, k \quad (۴۶.۷)$$

اگر بازه انتگرالگیری نیمه نامتناهی<sup>۴</sup> باشد، گاهی بهتر است که انتگرال را به شکل زیر تبدیل کنیم

$$\int_0^\infty g(x)w(x) dx$$

- 
1. body
  2. singular
  3. singularity
  4. semi-infinite



که در آن

$$w(x) = x^a e^{-x}$$

در این حالت،  $\xi$ ها ریشه‌های بسجمله‌ای خاص لاگرانژ هستند؛ به مثال ۸.۶ نگاه کنید. و سرانجام انتگرالهایی به شکل

$$\int_{-\infty}^{\infty} g(x) e^{-x^2} dx$$

را می‌توان اغلب با استفاده از ریشه‌های بسجمله‌ای مناسب هرمیت به خوبی برآورد نمود (به مثال ۸.۶ نگاه کنید).

برای کلیهٔ این مثالها (و مثالهای دیگر) جدولهایی هم برای  $\xi$ ها و هم برای  $A_i$ ها وجود دارند و احتمالاً جدیدترین و وسیعترین آنها «فرمولهای انتگرالگیری گاوسی» استرود و سیکرست<sup>۲</sup> هستند، [۲۰] و نیز [۲۷] را ببینید.

## تمرین

۱-۳۰۷ قاعدهٔ سیمپسن برای چه بسجمله‌هایی صحیح است؟

۲-۳۰۷ قاعده‌ای به شکل

$$I(f) = \int_{-1}^1 f(x) dx \approx A_0 f\left(-\frac{1}{4}\right) + A_1 f(0) + A_2 f\left(\frac{1}{4}\right)$$

سازید که برای تمام بسجمله‌های نایبتر از درجهٔ ۲ صحیح باشد.

۳-۳۰۷ مقدار انتگرال زیر را تا چهار رقم اعشاری دقیق محاسبه کنید

$$\int_0^1 \frac{\sin \pi x}{[x(1-x)]^{3/4}} dx$$

[دانهمایی: انتگرال را به نحو مناسبی تغییر دهید و از رابطهٔ (۴۵.۷) و (۴۶.۷) استفاده کنید].

۴-۳۰۷ برآوردی برای انتگرال  $\int_0^{\infty} e^{-x^2} dx$  پیدا کنید.

۵-۳۰۷ به ازای  $k=3$  و با استفاده از ریشه‌های  $\xi$  که در زیر برنامهٔ LGNDRE داده شده‌اند و زنها  $A_i$  در فرمول گاوسی را بیابید.

۳-۴.۷ با استفاده از فرمول پنج نقطه‌ای گاوسی، برای انتگرالهایی که در تمرینهای ۴.۷-۳ و ۴.۷-۴ داده شده‌اند برآوردی تعیین کنید.

۴-۳.۷ از تمرین ۳.۶-۷ استفاده کنید تا نشان دهید که می‌توان رابطه (۳.۷) را به ازای  $i = 0, \dots, k$  به صورت  $A_i = \int_a^b [l_i(x)]^2 w(x) dx$  نیز نوشت. از اینجا نتیجه‌گیری کنید که وزنهای فرمول گاوسی همواره مثبت‌اند.

۸-۳.۷ قاعده لو باتو<sup>۱</sup> برای انتگرالگیری  $I = \int_{-1}^1 f(x) dx$  يك فرمول گاوسی است، به استثنای آنکه این قاعده شامل  $\pm 1$  به عنوان دو طول<sup>۲</sup> ثابت است. این قاعده به شکل زیر است [رابطه (۳.۷) را ببینید]

$$I(f) \approx A_0 f(-1) + A_1 f(x_1) + \dots + A_{k-1} f(x_{k-1}) + A_k f(1)$$

قاعده لو باتو را برای حالت  $k = 2$ ، به دست آورید و نشان دهید که این قاعده برای تمام بسجمله‌ایهای نایبتر از درجه ۳، صحیح است.

۹-۳.۷ برنامه فرعی LGNDRE را با به کار بردن آن در محاسبه

$$\int_{-1}^1 x^n dx = \frac{1}{(n+1)} (1 - (-1)^{n+1}) \quad n = 0, 1, 2, \dots$$

امتحان کنید. به ازای چه مقادیری از  $n$ ، قاعده گاوس-لژاندر روی نقاطی به تعداد NP، مقدار انتگرال را دقیقاً به دست می‌دهد؟

## ۴.۷ انتگرالگیری عددی: قاعده‌های مرکب<sup>۳</sup>

قواعد ساده محاسبه انتگرال که در بخشهای پیشین برای برآورد

$$I = \int_a^b f(x) dx$$

پرورانیده شده بود، بخصوص وقتی بازه  $[a, b]$  نسبتاً بزرگ باشد، معمولاً برآوردی به قدر کافی دقیق به دست نمی‌دهد. در عمل مرسوم است که بازه مفروض  $[a, b]$  را به  $N$  بازه کوچکتر تقسیم می‌کنند و قواعد ساده محاسبه انتگرال را برای هر يك از این زیر بازه‌ها به کار می‌برند. بنا بر این بازه  $[a, b]$  را به بازه‌های کوچکتری چنان تقسیم می‌کنیم که

$$a = x_0 < x_1 < x_2 < \dots < x_N = b$$

و يك تابع بسجمله‌ای-تکه‌ای (به بخش ۴.۶ مراجعه کنید) در نقاط انفصال  $\{x_i\}$  را،

$(i = 1, \dots, N) P_{i,k}(x)$  بگیریم. بعلاوه،  $g_k(x)$  نشان می‌دهیم. بعلاوه،  $(i = 1, \dots, N-1)$  معرف یک بسجمله‌ای نایبتر از درجه  $k$  باشد که با  $g_k(x)$  در بازه  $(x_{i-1}, x_i)$  تطابق دارد. بنا بر قواعد انتگرالگیری می‌دانیم که

$$I(f) = \int_a^b f(x) dx = \sum_{i=1}^N \int_{x_{i-1}}^{x_i} f(x) dx$$

و

$$I(g_k) = \sum_{i=1}^N \int_{x_{i-1}}^{x_i} g_k(x) dx = \sum_{i=1}^N \int_{x_{i-1}}^{x_i} P_{i,k}(x) dx$$

بنابراین تقریب زدن  $I(f)$  با  $I(g_k)$ ، مانند تقریب زدن

$$\int_{x_{i-1}}^{x_i} P_{i,k}(x) dx \quad \text{به وسیلهٔ} \quad \int_{x_{i-1}}^{x_i} f(x) dx \quad i = 1, \dots, N$$

و جمع کردن نتایج است. بدیهی است که در هر زیر بازه  $(x_{i-1}, x_i)$  دقیقاً مانند بخشهای ۲.۷ و ۳.۷ عمل می‌کنیم. بخصوص اینکه می‌توانیم در هر زیر بازه با گذاردن یک بسجمله‌ای به جای تابع زیر علامت انتگرال، هر یک از قواعد حاصله از بخشهای ۲.۷ و ۳.۷ را به کار بریم و سپس نتایج حاصله را جمع کنیم.  $x$ ها را با فواصل مساوی به صورت زیر انتخاب می‌کنیم مگر آنکه برای انتخاب دیگر دلیل موجهی وجود داشته باشد

$$x_i = a + ih \quad i = 0, \dots, N$$

$$h = \frac{b-a}{N}$$

همچنین مانند بخش ۶.۲ علامت اختصاری زیر را به کار می‌بریم

$$f_s = f(a + sh)$$

لذا  $f_i = f(x_i)$  به ازای  $i = 0, \dots, N$ .

اکنون مثالهای خاصی را در نظر می‌گیریم. اگر قاعدهٔ مستطیلی (۲۳.۷) را در هر زیر بازه به کار بریم، آنگاه برای زیر بازه  $(x_{i-1}, x_i)$  خواهیم داشت

$$\begin{aligned} \int_{x_{i-1}}^{x_i} f(x) dx &= (x_i - x_{i-1})f(x_{i-1}) + \frac{f'(\eta_i)(x_i - x_{i-1})^2}{2} \\ &= hf(x_{i-1}) + \frac{f'(\eta_i)h^2}{2} \end{aligned}$$

از محاسبهٔ حاصلجمع این مقادیر خواهیم داشت

$$I(f) \approx R_N = h \sum_{i=1}^N f_{i-1} \quad (۴۷.۷ \text{ الف})$$

که رابطه فوق، قاعده مستطیلی مرکب (بر  $N$  زیر بازه است). خطای آن نیز درست برابر با مجموع خطاهایی است که در هر زیر بازه ایجاد می شود

$$E_N^R = \sum_{i=1}^N \frac{f'(\eta_i)h^2}{2}$$

که در آن  $\eta_i \in (x_{i-1}, x_i)$  می باشد. اگر  $f'(x)$  پیوسته باشد (همان گونه که فرض کرده ایم)، رابطه فوق را می توان با استفاده از قضیه ۲.۱ از بخش ۷.۱، به صورت زیر ساده کرد:

$$\sum_{i=1}^N \frac{f'(\eta_i)h^2}{2} = f'(\eta) \sum_{i=1}^N \frac{h^2}{2} = \frac{f'(\eta)Nh^2}{2}$$

لذا به ازای  $Nh = b - a$ ، خواهیم داشت

$$E_N^R = \frac{f'(\eta)(b-a)h}{2}$$

(ب)  $\eta \in (a, b)$  به ازای مقداری مانند

سهس قاعده مرکب سیمپسن را به دست می آوریم. اگر در رابطه (۲۸.۷) بگیریم  $a = x_{i-1}$  و  $b = x_i$ ، آنگاه برای یک زیر بازه تنها خواهیم داشت

$$\int_{x_{i-1}}^{x_i} f(x) dx = \frac{h}{6} [f_{i-1} + 4f_{i-1/2} + f_i] - \frac{f^{iv}(\eta_i)(h/2)^5}{90}$$

$$x_{i-1} < \eta_i < x_i$$

با محاسبه حاصل جمع به ازای  $i = 1, \dots, N$ ، رابطه زیر حاصل می شود

$$\begin{aligned} I(f) &= \sum_{i=1}^N \int_{x_{i-1}}^{x_i} f(x) dx \\ &= \frac{h}{6} \sum_{i=1}^N [f_{i-1} + 4f_{i-1/2} + f_i] - \sum_{i=1}^N \frac{f^{iv}(\eta_i)(h/2)^5}{90} \\ &= S_N + E_N^S \end{aligned}$$

تقریب مرکب سیمپسن  $S_N$  را می توان ساده کرد و نتیجه زیر را به دست آورد

$$S_N = \frac{h}{6} \left[ f_0 + f_N + 2 \sum_{i=1}^{N-1} f_i + 2 \sum_{i=1}^N f_{i-1/2} \right] \quad (\text{الف } ۲۸.۷)$$

درعین حال عبارت خطا را نیز می توان با استفاده از قضیه ۲.۱ از بخش ۷.۱، به صورت زیر ساده کرد:

$$E_N^S = -\frac{f^{iv}(\xi)(h/2)^4(b-a)}{180} \quad a < \xi < b \quad (\text{ب } ۲۸.۷)$$

باید توجه داشت که در قاعده سیمپسن باید بتوانیم تابع را در نقاط میانی  $x_{i-1/2}$ ،  $(i = 1, \dots, N)$ ، و نیز در نقاط انفصال  $x_i$ ،  $(i = 0, 1, \dots, N)$ ، محاسبه کنیم. این امر مخصوصاً ایجاب می کند که همواره به تعداد فردی از نقاط متساوی الفاصله، که در آنها مقدار تابع زیر علامت انتگرال را می دانیم، نیاز داشته باشیم. به همین ترتیب برای قاعده نقطه وسط مرکب چنین به دست می آوریم:

$$I(f) \approx M_N = h \sum_{i=1}^N f_{i-1/2} \quad E_N^M = \frac{f''(\xi)h^2(b-a)}{24} \quad (\text{الف } ۲۹.۷)$$

با توجه به قاعده نقطه وسط (۲۵.۷)، و قاعده دوزنقه‌ی مرکب، و (۲۶.۷) داریم

$$I(f) \approx T_N = h \sum_{i=1}^{N-1} f_i + \frac{h}{2} (f_0 + f_N)$$

$$E_N^T = -\frac{f''(\eta)h^2(b-a)}{12} \quad (\text{ب } ۲۹.۷)$$

از قاعده دوزنقه‌ی اصلاح شده (۲۹.۷)، رابطه زیر به دست می آید

$$I(f) \approx CT_N = h \sum_{i=1}^{N-1} f_i = \frac{h}{2} (f_0 + f_N) + \frac{h^2}{12} [f'(a) - f'(b)]$$

(۵۰.۷)

$$E_N^{CT} = \frac{f^{iv}(\eta)h^4(b-a)}{720}$$

توجه داشته باشید که کلیه مشتق‌های درونی  $f'(x_i)$ ،  $i = 1, 2, \dots, N-1$ ، هنگام جمع نتایج حاصل از قاعدهٔ ذوزنقه‌یی اصلاح شده بر هر زیر بازه، از یکدیگر حذف می‌شوند. بنا بر این قاعدهٔ ذوزنقه‌یی اصلاح شدهٔ مرکب، در حقیقت یک قاعدهٔ ذوزنقه‌یی مرکب اصلاح شده است، یعنی

$$CT_N = T_N + \frac{h^2[f'(a) - f'(b)]}{12} \quad (51.7)$$

البته قاعدهٔ ذوزنقه‌یی اصلاح شده دارای این اشکال است که مشتق  $f(x)$  باید در دست و یا قابل محاسبه باشد [به استثنای حالتی که در آن  $f(x)$  یک تابع  $(b-a)$ -دوره‌ای است]. اگر قرار باشد که یکی از این قواعد مرکب به کار روند، در آغاز باید مقدار مناسبی مانند  $N$ ، و یا معادل با آن، مقداری مانند  $h = (b-a)/N$  تعیین شود. اگر اطلاعی از مقدار مشتقی که در عبارت خطا ظاهر می‌شود در دست باشد، اصلاً می‌توان  $h$  یا  $N$  را طوری تعیین کرد که مقدار خطای حاصل از تحمل (تولرانس) مقرر کمتر باشد.

□ مثال ۴۰۷:  $N$  را طوری تعیین کنید که قاعدهٔ ذوزنقه‌یی مرکب (۳۳.۵ ب)، مقدار انتگرال

$$\int_0^1 e^{-x^2} dx$$

را (با فرض اینکه بتوان  $e^{-x^2}$  را دقیقاً محاسبه کرد) تا شش رقم اعشار صحیح به دست دهد و سپس، تقریب این انتگرال را نیز محاسبه کنید.

در این مثال داریم  $f(x) = e^{-x^2}$ ،  $a = 0$ ،  $b = 1$  و  $h = 1/N$ . بنا بر این به ازای مقداری مانند  $\eta \in (a, b)$  مقدار خطا در قاعدهٔ ذوزنقه‌یی مرکب برابر با  $1/12 N^{-2} f''(\eta)$  است. از آنجا که  $\eta$  را نداریم، بهترین فرضی که می‌توانیم بکنیم این است که بگوئیم مقدار خطا از لحاظ قدر مطلق نباید بزرگتر از مقدار

$$\max_{0 \leq \eta \leq 1} \frac{|f''(\eta)| N^{-2}}{12}$$

باشد. محاسبهٔ زیر را انجام می‌دهیم

$$f''(x) = e^{-x^2}(4x^2 - 2)$$

بعلاوه  $f'''(x) = e^{-x^2}4x(3 - 2x^2)$  در نقطه‌های  $x = 0$  و  $x = \pm\sqrt{3/2}$  برابر با صفر است. لذا  $\max |f''(x)|$  در بازهٔ  $[0, 1]$  باید در نقطهٔ  $x = 0$  یا در نقاط انتهایی  $x = 0, 1$  موجود باشد. از این رو

$$\max_{0 \leq \eta \leq 1} |f'''(\eta)| = \max \{|f'''(0)|, |f'''(1)|\} = \max \{2, 2e^{-1}\} = 2$$

بنابراین اگر  $N$  را چنان انتخاب کنیم که

$$\frac{2N^{-2}}{12} < 5 \times 10^{-7}$$

یا

$$N^2 > \frac{10^6}{3} = \frac{10^7}{6 \times 5}$$

یا

$$N > \frac{10^3}{\sqrt{3}} \approx 577$$

دقت تا شش رقم اعشاری (بعد از ممیز) برای ما تضمین شده است. نتایج کامپیوتری زیر نشان می‌دهد که  $N$  اندکی بیش از اندازه برآورد شده است.

محاسبات روی IBM مدل ۷۰۹۴ با دقت مضاعف (DP) ودقت ساده (SP) انجام گرفته و برنامه حاصل برای مقادیر مختلف  $N$  به شرح زیر است.

$N$	$I(SP)$	$I(DP)$	ERROR(SP)	ERROR(DP)
50	7.467994E-01	7.4670061D-01	2.466E-05	2.452D-05
100	7.4681776E-01	7.4681800D-01	6.37 E-06	6.13 D-06
200	7.4682212E-01	7.4682260D-01	2.01 E-06	1.53 D-06
400	7.4682275E-01	7.4682375D-01	1.56 E-06	3.8 D-07
800	7.4682207E-01	7.4682404D-01	2.06 E-06	9. D-08

مقدار صحیح  $I$  تا هشت رقم با معنی برابر است با  $I = 0.74682275$ . بنابراین، چنین به نظر می‌رسد که تعداد تقسیمها هر چه اختیار شوند ظاهراً در محاسبه با دقت ساده نمی‌توانیم به شش رقم اعشاری دقیق دست یابیم. درحقیقت نتایج برای  $N = 800$  نسبت به  $N = 400$  بدتر است. این امر نشان می‌دهد که خطای ناشی از گرد کردن روی سه رقم آخر اثر گذاشته است. چنانچه از قبل تا حدی پیشبینی شده بود نتایج با دقت مضاعف نشان می‌دهند که به ازای  $N = 400$  دقت تا شش رقم اعشاری به دست می‌آید. برنامه فورترن محاسبات فوق به صورت زیر است:

برنامه فورترن برای مثال ۴.۷ (دقت ساده)

```
C EXAMPLE 7.4 . TRAPEZOID RULE .
INTEGER I,N
REAL A,B,H,T
F(X) = EXP(-X*X)
1 PRINT 601
```

```

601 FORMAT(' EXAMPLE 7.4 TRAPEZOIDAL INTEGRATION')
    READ 501, A,B,N
501 FORMAT(2E20.0,I5)
    IF (N .LT. 2) STOP
    T = F(A)/2.
    H = (B - A)/FLOAT(N)
    DO 2 I=1,N-1
      2 T = F(A + FLOAT(I)*H) + T
    T = (F(B)/2. + T)*H
    PRINT 602, A,E,N,T
602 FORMAT(' INTEGRAL FROM A = ',1PE14.7,' TO B = ',E14.7,
          * ' FOR N = ',I5,' IS ',E14.7)
          GO TO 1
END

```

اگر از قاعدهٔ ذوزنقه‌یی اصلاح شده (۵۰.۷) استفاده کنیم،  $N$  مورد نیاز به‌طور قابل ملاحظه‌ای کاهش می‌یابد. اما خطا به‌وسیلهٔ

$$\max_{0 \leq \eta \leq 1} \frac{|f^{iv}(\eta)| N^{-4}}{720}$$

کراندار شده است. مشتق مرتبهٔ چهارم  $f^{iv}(x) = 4e^{-x^2}(3 - 12x^2 + 4x^4)$  را محاسبه می‌کنیم. بنا براین

$$\max_{0 \leq \eta \leq 1} |f^{iv}(\eta)| = |f^{iv}(0)| = 12$$

لذا برای دستیابی به‌دقت با شش رقم، کافی است که داشته باشیم

$$\frac{12N^{-4}}{720} < 5 \times 10^{-7}$$

یا

$$N^4 > (10/3)10^4 = \frac{10^7}{300}$$

یا

$$N > 1325 \approx \sqrt[4]{(10/3)(10^7)}$$

به‌طوری که درمقایسه با ۵۷۸ زیربازهٔ مورد نیاز برای قانون ذوزنقه‌یی مرکب بدون تصحیح  $\square$  دیگر انسیل نهایی، تنها به ۱۴ زیربازه نیاز است.

همان‌طوری که مثال فوق نشان می‌دهد اگر مشتقهای مرتبه بالاتر تابع زیر علامت انتگرال یا مشتقهای مرتبهٔ پایینتر آن تقریباً هم‌اندازه باشند، آنگاه فرمولهای مرتبه بالاتر می‌توانند تعداد تسایع لازم برای محاسبه را نسبت به فرمولهای مرتبهٔ پایینتر به‌طور قابل ملاحظه‌ای کاهش دهند. بالاخص قواعد گاوسی در این مورد می‌توانند بسیار مؤثر باشند. اگر دربارهٔ اندازهٔ مشتق مناسب  $f(x)$  اطلاعاتی در دست نباشد، قواعد مرکب را تنها



می‌توان برای مقادیر مختلف  $N$  به‌کار برد، و بدین ترتیب یک دنباله از تقریبات  $I_N$  برای  $I(f)$ ، به‌دست می‌آیند که اگر  $f(x)$  به‌اندازهٔ کافی هموار باشد این مقادیر از لحاظ نظری وقتی  $\infty \rightarrow N$ ، به‌سمت  $I(f)$  همگرا می‌شوند. این فرایند هنگامی به‌پایان می‌رسد که تفاوت بین برآوردهای متوالی «به‌اندازهٔ کافی کوچک» شود. خطرات این شیوهٔ عمل در بخش ۶.۱ مورد بحث قرار گرفت. اشکال دیگری که در این حالت پیش می‌آید آثار ناشی از گرد کردن است که با زیاد شدن  $N$  افزایش می‌یابد. نتایج کامپیوتری در مثال ۴.۷ این امر را به‌خوبی نشان می‌دهند.

□ مثال ۵.۷: برنامه‌ای برای قاعدهٔ دوزنقه‌یی اصلاح شده بنویسید و مثال ۴.۷ را بسا استفاده از این برنامه حل کنید.

### برنامهٔ فورترن

```
C      EXAMPLE 7.5 . CORRECTED TRAPEZOID RULE
      INTEGER I,N
      REAL A,B,CORTRP,H,TRAP
      F(X) = EXP(-X*X)
      FPRIME(X) = -2.*X*F(X)
      DATA A,B /0., 1. /
      PRINT 600
600  FORMAT(9X,'N',7X,'TRAPEZOID SUM',7X,'CORR.TRAP.SUM')
      DO 10 N = 10,15
         H = (B - A)/FLOAT(N)
         TRAP = (F(A) + F(B))/2.
         DO 1 I=1,N-1
            TRAP = TRAP + F(A + FLOAT(I)*H)
            TRAP = H*TRAP
            CORTRP = TRAP + H*H*(FPRIME(A) - FPRIME(B))/12.
10     PRINT 610, N,TRAP,CORTRP
610  FORMAT(110,2E20.7)
      STOP
END
```

نتایج با دقت ساده

N	TRAPEZOID SUM	CORR.TRAP.SUM
10	0.7462108E 00	0.7468239E 00
11	0.7463173E 00	0.7468240E 00
12	0.7463983E 00	0.7468240E 00
13	0.7464612E 00	0.7468240E 00
14	0.7465112E 00	0.7468240E 00
15	0.7465516E 00	0.7468241E 00

نتایج با دقت مضاعف

N	TRAPEZOID SUM	CORR.TRAP.SUM
10	7.4621080E-01	7.4682393E-01
11	7.4631727E-01	7.4682399E-01
12	7.4639825E-01	7.4682403E-01
13	7.4646126E-01	7.4682406E-01
14	7.4651126E-01	7.4682408E-01
15	7.4655159E-01	7.4682409E-01

□ مثال ۶.۷: برنامه‌ای برای قاعدهٔ سیمپسن بنویسید و مثال ۴.۷ را با این برنامه، هم با دقت ساده و هم با دقت مضاعف حل کنید.

برنامهٔ فورتون و نتایج به دست آمده از IBM ۷۰۹۴ در زیر به‌زای تعداد زیر تقسیمه‌هایی برابر با ۱۰۰ و ۵۰ و  $N = 25$  داده شده‌اند. توجه کنید که نتایج با دقت سادهٔ ۱۰۰ و  $N = 50$  نسبت به  $N = 25$ ، بدتر است که این امر حاکی از اثرات خطای ناشی از گرد کردن است. نتایج با دقت مضاعف، تماماً با تعداد ارقام داده شده صحیح‌اند. از مقایسهٔ این نتایج با نتایجی که از مثال ۴.۷ و ۵.۷ به دست آمد، مشاهده می‌شود که هر دو قاعدهٔ سیمپسن و قاعدهٔ دوزنقه‌یی اصلاح شده به‌مراتب کاراتر از قاعدهٔ دوزنقه‌یی هستند.

```
C PROGRAM FOR EXAMPLE 7.6 . SIMPSON'S RULE .
  INTEGER I,N
  REAL A,B,H,HALF,HOVER2,S,X
  F(X) = EXP(-X*X)
  PRINT 600
600 FORMAT(' EXAMPLE 7.6 SIMPSON'S RULE')
  1 READ 501, A,B,N
  501 FORMAT(2E20.0,I5)
  IF (N .LT. 2) STOP
  H = (B - A)/FLOAT(N)
  HOVER2 = H/2.
  S = 0.
  HALF = F(A + HOVER2)
  DO 2 I=1,N-1
    X = A + FLOAT(I)*H
    S = S + F(X)
  2 HALF = HALF + F(X+HOVER2)
  S = (H/6.)*(F(A) + 4.*HALF + 2.*S + F(B))
  PRINT 602, A,B,N,S
602 FORMAT(' INTEGRAL FROM A = ',1PE14.7,' TO B = ',E14.7,
  * ' FOR N = ',I5,' IS ',E14.7)
  GO TO 1
  4 FORNAT(2E20.0,I5)
  END
```

نتایج کامپیوتری برای مثال ۶.۷

$N$	$I(SP)$	ERROR (SP)	$I(DP)$	ERROR (DP)
25	7.4682406E-01	7. E-07	7.4682413D-01	0.
50	7.4682400E-01	1.3E-06	7.4682413D-01	0.
100	7.4682392E-01	2.1E-06	7.4682413D-01	0.

□

سرانجام، قواعد مرکبی را که بر فرمولهای گاوس استوارند نیز می‌توان به دست آورد. برای آنکه این قواعد با قواعد مرکبی که تاکنون بحث شد سازگار باشند، مسا خود را به انتگرالهای معین به شکل زیر محدود خواهیم ساخت

$$I = \int_a^b f(x) dx$$

بار دیگر بازهٔ  $(a, b)$  را به  $N$  زیر بازهٔ مساوی تقسیم می‌کنیم به طوری که

به ازای  $h = (b - a) / N$

$$x_i = a + ih \quad i = 0, 1, \dots, N$$

می‌خواهیم انتگرال‌گیری گاوسی را در بازهٔ  $I_i$  برای این انتگرال یعنی برای انتگرال

$$I_i = \int_{x_{i-1}}^{x_i} f(x) dx \quad (52.7)$$

به‌کار بریم. وزنها و نقاط گاوسی که بر اساس بسجمله‌ریهای لژاندر استوار و در قسمت ۳.۷ داده شده‌اند. حدود انتگرال‌گیری را از  $-1$  تا  $+1$  می‌پذیرند. از این رو نخست تبدیل متغیرهای خطی به‌صورت زیر انجام می‌دهیم

$$x = \frac{h}{2}t + x_{i-1/2}$$

که در آن

$$x_{i-1/2} = (x_i + x_{i-1}) / 2$$

از قرارداد این مقادیر در (52.7) رابطهٔ زیر به‌دست می‌آید

$$I_i = \frac{h}{2} \int_{-1}^1 f\left(\frac{ht}{2} + x_{i-1/2}\right) dt = \int_{-1}^1 g_i(t) dt$$

که در آن

$$g_i(t) = \frac{h}{2} f\left(\frac{ht}{2} + x_{i-1/2}\right)$$

حال در  $k+1$  نقطه، انتگرال  $I_i$  را با فرمول گاوسی تقریب می‌زنیم و رابطهٔ زیر را به‌دست می‌آوریم

$$I_i \approx A_0 g_i(\xi_0) + A_1 g_i(\xi_1) + \dots + A_k g_i(\xi_k) \quad (53.7)$$

که در آن وزنها و طولها<sup>۱</sup> از زیر برنامهٔ LGNDRE در قسمت ۳.۷ به‌دست می‌آیند. سرانجام با به‌دست آوردن حاصلجمع این مقادیر تقریبی در  $N$  زیربازه، خواهیم داشت

$$I \approx \sum_{i=1}^N \int_{x_{i-1}}^{x_i} f(x) dx = \sum_{i=1}^N \int_{-1}^1 g_i(t) dt = \sum_{i=1}^N I_i$$

که با توجه به (53.7)، این مقدار تقریبی برابر است با

$$I \approx \sum_{i=1}^N \{A_0 g_i(\xi_0) + A_1 g_i(\xi_1) + \dots + A_k g_i(\xi_k)\}$$

$$= \frac{h}{2} \sum_{i=1}^N \left\{ A_0 f\left(\frac{h}{2} \xi_0 + x_{i-1/2}\right) + \dots + A_k f\left(\frac{h}{2} \xi_k + x_{i-1/2}\right) \right\}$$

(الف ۵۴.۷)

توجه دارید که وزنها، مستقل از  $i$  هستند.

بهموجب معادله خطای (۴۰.۷)، در یک بازه  $(x_{i-1}, x_i)$ ، خطا به فرم زیر قابل بیان

است.

$$E_i = C_k g_i^{(k+2)}(\eta_i) \quad [-1, 1] \text{ در } \eta_i \text{ مانند } \eta_i \text{ مقدار می مانند}$$

اما این بدان معنی است که

$$E_i = C_k \left(\frac{h}{2}\right)^{k+2} f^{(k+2)}(\eta'_i) \quad x_{i-1} < \eta'_i < x_i$$

بنابراین خطا در بازه  $(a, b)$  می تواند به صورت زیر بیان شود

$$E_N^G = \frac{1}{2} C_k \left(\frac{h}{2}\right)^{k+2} f^{(k+2)}(\eta) \quad \text{(ب ۵۴.۷)}$$

□ مثال ۷.۷: با استفاده از انتگرالگیری گاوسی به ازای  $k=3$  و به تعداد  $N=2$  زیر بازه از بازه  $[1, 3]$ ، انتگرال  $I = \int_1^3 [(\sin x)^2/x] dx$  را برآورد کنید. مثال ۳.۷ را ببینید.

C PROGRAM FOR EXAMPLE 7.7. COMPOSITE FOUR-POINT GAUSS-LEGENDRE .

```

INTEGER I,N
REAL A,B,H,HOVER2,P1,P2,POINT(2),S,S1,S2,WEIGHT(2),X
DATA POINT,WEIGHT / .33998 10436, .86113 63116,
.65214 51549, .34785 48451 /
F(X) = SIN(X)**2/X
PRINT 600
600 FORMAT(' EXAMPLE 7.7 FOUR-POINT GAUSS-LEGENDRE'/)
1 READ 501, A,B,N
501 FORMAT(2E20.0,I5)
IF (N.LT. 1) STOP
H = (B - A)/FLOAT(N)
HOVER2 = H/2.
P1 = POINT(1)*HOVER2
P2 = POINT(2)*HOVER2
S1 = 0.
S2 = 0.
DO 2 I=1,N
X = A + FLOAT(I)*H - HOVER2
S1 = S1 + F(-P1+X) + F(P1+X)
S2 = S2 + F(-P2+X) + F(P2+X)
S = HOVER2*(WEIGHT(1)*S1 + WEIGHT(2)*S2)
PRINT 602, A,B,N,S
602 FORMAT(' INTEGRAL FROM A = ',1PE14.7,' TO B = ',E14.7,
' FOR N = ',I3,' IS ',E14.7)
GO TO 1

```

END

جواب حاصل از محاسبات با دقت ساده که روی کامپیوتر ۱۱۱۰ UNIVAC انجام گرفته برابری با ۰۷۹۴۸۲۵۱۷ که میزان خطا در آخرین رقم آن کمتر از ۳ واحد\* است. □

### تمرین

۱-۴.۷ قاعدهٔ وزنقهبی مرکب  $T_N$  (۲۹.۷) و قاعدهٔ نقطهٔ میانی مرکب  $M_N$  (۴۸.۷) را به دست آورید.

۲-۴.۷ قاعدهٔ وزنقهبی اصلاح شدهٔ مرکب  $CT_N$  (۵۰.۷) را به دست آورده و تحقیق کنید که مشتقهای درونی  $f'(x_i)$ ،  $(i = 1, \dots, N-1)$ ، در حاصلجمع حذف می شوند.

۳-۴.۷ برنامه‌ای برای قاعدهٔ مرکب سیمپسن بنویسید. ورودیهای این برنامه می باید  $f(x)$  باشد، بازهٔ آن  $[a, b]$ ، و تعداد زیرفاصله‌ها  $N$ . از این برنامه برای محاسبهٔ انتگرال زیر با  $N = 10$  و  $N = 20$  زیرفاصله استفاده کنید

$$I = \int_1^2 x \ln x \, dx$$

۴-۴.۷ برنامهٔ مربوط به قاعدهٔ سیمپسن را در محاسبهٔ تقریبی انتگرالهای زیر که تا شش رقم اعشاری صحیح باشند به کار برید

$$I = \int_0^1 x e^{-x} \, dx \quad I = \int_0^1 x \cos x \, dx \quad I = \int_0^1 (1+x^2)^{3/2} \, dx$$

این عمل را با  $N = 10$  شروع و هر دفعه  $N$  را دو برابر کنید تا به دقتی که می خواهید دست یابید.

۵-۴.۷ برنامه‌ای برای قاعدهٔ وزنقهبی اصلاح شده بنویسید. در این حالت، ورودی شامل  $f(x)$ ،  $f'(x)$ ،  $[a, b]$  و  $N$  خواهد بود. این برنامه را برای انتگرال تمرین ۳-۴.۷ به کار برید و نتایج حاصله را با نتایج به دست آمده از قاعدهٔ سیمپسن مقایسه کنید.

۶-۴.۷ به ازای  $k = 3$ ، برنامه‌ای برای قاعدهٔ مرکب گاوسی (۵۴.۷ الف) بنویسید. برنامهٔ فوق را برای محاسبهٔ انتگرال تمرین ۳-۴.۷ به کار گیرید. این محاسبات را در ابتدا به ازای  $N = 2$  زیربازه و سپس به ازای  $N = 4$  زیربازه انجام دهید. مقدار محاسبات و دقت به دست آمده را با موارد مشابه در قاعدهٔ سیمپسن مقایسه کنید.

۷-۴.۷ تابع خطای  $\text{erf}(x)$  به صورت زیر تعریف می شود

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$

به ازای  $N = 2$  و  $N = 3$  زیر بازه و  $k = 3$  برای محاسبه  $\operatorname{erf}(0.5)$  از قاعده مرکب گاوسی استفاده کنید. دقت نتیجه حاصله را بر آورد و با مقدار دقیق

$$\operatorname{erf}(0.5) = 0.520499876$$

مقایسه کنید.

۸-۴۰۷ برای تعیین میزان چگالش یک بخار خالص در خارج از یک لوله افقی سرد شده، به محاسبه میانگین ضریب انتقال دمای  $Q$  نیاز داریم. برای محاسبه این ضریب، همراه با سایر پارامترها، محاسبه انتگرال زیر لازم می آید

$$I = \int_0^{\pi} (\sin x)^{1/3} dx$$

به ازای  $N = 5, 10, 15, 20$  زیر فاصله و بسا به کارگیری قانون سیمپسن، مقدار این انتگرال را محاسبه کنید.

$$I \approx 2.5286949, N = 5 \text{ جواب:}$$

## ۵.۲ انتگرال تطبیقی

قواعد مرکبی که تا به حال مورد بحث قرار گرفتند، همه بر پایه  $N$  زیر بازه مساوی نهاده شده اند. انتخاب این گونه زیر بازه های مساوی، هر گاه تابع زیر علامت انتگرال در یک رشته از نقاط با فواصل مساوی معلوم باشد، مثلاً اگر  $f(x)$  تنها به صورت جدولی از مقادیر تابع داده شده باشد، کاملاً طبیعی و حتی اغلب لازم است. اما اگر  $f(x)$  در هر نقطه از بازه انتگرالگیری به راحتی محاسبه شود، معمولاً به صرفه نزدیکتر است که از زیر بازه هایی استفاده کنیم که طول آنها با توجه به رفتار محلی تابع زیر علامت انتگرال، معین شود. به عبارت دیگر، اگر به جای اصرار در مساوی الفاصله بودن زیر بازه ها، آنها نامساوی ولی مناسب انتخاب شوند، آنگاه می توان  $I(f)$  را بسا دقتی از پیش تعیین شده و محاسباتی کمتر به دست آورد.

برای مثال، قاعده کلی ذوزنقه ای مرکب

$$I(f) = \sum_{i=1}^N \frac{x_i - x_{i-1}}{2} [f(x_{i-1}) + f(x_i)] - \sum_{i=1}^N \frac{f''(\eta_i)(x_i - x_{i-1})^3}{12}$$

را در نظر بگیرید، که در آن لازم نیست، نقاط انفصال  $a = x_0 < x_1 < \dots < x_N = b$

به فواصل مساوی قرار گیرند، در بازهٔ  $(x_{i-1}, x_i)$  به ازای نقطه‌ای مانند  $\eta_i \in (x_{i-1}, x_i)$  سهم مقدار

$$\frac{f''(\eta_i)(x_i - x_{i-1})^3}{12}$$

در رابطه با کل خطا، هم به مقدار  $f''(x)$  در بازهٔ  $(x_{i-1}, x_i)$  بستگی دارد و هم به اندازهٔ زیر بازهٔ  $|x_i - x_{i-1}|$ . بنا بر این اگر بخواهیم سهمی که هر زیر بازه در کل خطا دارد مساوی باشد، می‌توانیم در قسمتهایی از بازهٔ انتگرالگیری  $(a, b)$  که مقدار  $|f''(x)|$  «کوچک» است، زیر بازه‌ها را «بزرگ» بگیریم و در ناحیه‌هایی که  $|f''(x)|$  «بزرگ» است، زیر بازه‌ها را «کوچک» می‌توان نشان داد که اگر منظور به حداقل رساندن تعداد زیر بازه‌ها و در پی آن به حداقل رساندن تعداد محاسباتی از تابع باشد که برای به دست آوردن  $I(f)$  تسادقت معینی لازم است، آنگاه روش مزبور بهترین روش است.

آن دسته از روشهای انتگرالگیری که طول زیر بازه‌ها را با رفتار محلی تابع زیر علامت انتگرال وفق می‌دهند، روشهای انتگرالگیری تطبیقی نام دارند. اشکال مهمی که در این گونه روشها وجود دارد بی‌اطلاعی از مشتقی است که در عبارت خطا ظاهر می‌شود و این بدان معنی است که در چنین روشهایی، رفتار محلی تابع زیر علامت انتگرال را باید از روی مقادیر آن در چند نقطه حدس زد.

در اینجا یک روش الگوی انتگرالگیری تطبیقی را که بر استفاده از قاعدهٔ سیمپسن به عنوان یک فرمول انتگرال پایه متکی است، به اجمال شرح می‌دهیم. فرض آن است که تابع  $f(x)$ ، بازهٔ  $[a, b]$  و یک معیار خطای  $\varepsilon$  داده شده‌اند. هدف، محاسبهٔ تقریبی  $P$  انتگرال  $I = \int_a^b f(x) dx$  با حداقل تعداد محاسبات برای تابع است به طوری که

$$|P - I| \leq \varepsilon \quad (55.7)$$

برای این کار، عمل را تا آنجا که ممکن است، با استفاده از عدد کمی از مقادیر تابع انجام می‌دهیم.

عملیات را با تقسیم بازهٔ  $[a, b]$  به  $N$  زیر بازه که معمولاً، ولی نه الزاماً، برابرند آغاز می‌کنیم. گیریم  $x_i$  و  $x_{i+1}$  دوسریک چنین زیر بازه‌ای باشند و  $h = x_{i+1} - x_i$ . اکنون برای انتگرال

$$I_i = \int_{x_i}^{x_{i+1}} f(x) dx$$

دو تقریب با قاعدهٔ سیمپسن به دست می‌آوریم. یکی از این دو تقریب را که با  $I$  نشان می‌دهیم بر اساس استفاده از دو قطعهٔ مبتنی شده است، و دیگری را که با  $I_1$  نشان می‌دهیم بر اساس

استفاده از ۴ قطعه پی‌ریزی شده است. به موجب فرمول (۲۸.۷) این تقریبات با روابط زیر معین می‌شوند

$$S_i = \frac{h}{6} \left\{ f(x_i) + 4f\left(x_i + \frac{h}{4}\right) + f(x_{i+1}) \right\} \quad (\text{الف } ۵۶.۷)$$

$$\bar{S}_i = \frac{h}{12} \left\{ f(x_i) + 4f\left(x_i + \frac{h}{4}\right) + 2f\left(x_i + \frac{h}{2}\right) + 4f\left(x_i + \frac{3h}{4}\right) + f(x_{i+1}) \right\} \quad (\text{ب } ۵۶.۷)$$

از این دو تقریب، می‌توانیم میزان خطا را در تقریب دقیقتر  $\bar{S}_i$  به گونه‌ای زیر محاسبه کنیم. به موجب عبارت خطا در قاعدهٔ سیمپسن (۲۸.۷)، داریم

$$I_i - S_i = -\frac{f^{iv}(\eta)}{90} \left(\frac{h}{4}\right)^5 \quad (\text{الف } ۵۷.۷)$$

$$I_i - \bar{S}_i = -\frac{2f^{iv}(\eta)}{90} \left(\frac{h}{4}\right)^5 \quad (\text{ب } ۵۷.۷)$$

ظاهر شدن ضریب ۲ در رابطه (ب ۵۷.۷) ناشی از آن حقیقت است که عمل انتگرالگیری روی دو زیربازه انجام می‌گیرد که هر یک پهنایی برابر با  $h/2$  دارد. با فرض اینکه مشتق  $f^{iv}(x)$  در بازه  $[x_i, x_{i+1}]$  تقریباً ثابت است، می‌توان رابطه (ب ۵۷.۷) را از (الف ۵۷.۷) کم کرد و بعد از ساده کردن نتیجه خواهیم داشت

$$\bar{S}_i - S_i = \frac{f^{iv} \cdot h^5}{25 \cdot 90} \left(\frac{1-2^4}{2^4}\right)$$

که از این رابطه، نتیجه می‌گیریم

$$\frac{f^{iv} \cdot h^5}{25 \cdot 90} = \frac{2^4(\bar{S}_i - S_i)}{1 - 2^4} \quad (\text{۵۸.۷})$$

با گذاردن (۵۸.۷) در سمت راست (ب ۵۷.۷)، مقدار برآورد خطا برابر می‌شود با

$$I_i - \bar{S}_i = \frac{\bar{S}_i - S_i}{2^4 - 1} = \frac{1}{15} (\bar{S}_i - S_i) \quad (\text{۵۹.۷})$$

یا به صورت گفتاری، میزان خطا در تقریب دقیقتر  $\bar{S}_i$  تقریباً برابر است با  $1/15$  تفاضل بین دو تقریب  $\bar{S}_i$  و  $S_i$ ، که کمیتی است به آسانی قابل محاسبه.

اگر بازه  $[a, b]$  با  $N$  زیربازه پوشانده شده باشد، و اگر در هر یک از این زیربازه‌ها،



بر آورد خطا در رابطه

$$E_i = \frac{1}{15} |\bar{S}_i - S_i| \leq \frac{h}{b-a} \varepsilon \quad (60.7)$$

صدق کند، آنگاه می‌توان نشان داد که تقریب حاصل برای انتگرال  $I$  از حاصلجمع

$$P = \sum_{i=1}^N \bar{S}_i$$

در کل بازهٔ  $[a, b]$  در ملاء خطای مورد نیاز (۵۵.۷) صدق خواهد کرد. در رابطه (۶۰.۷) توجه به این امر که مقدار  $h = x_{i+1} - x_i$  با تغییر پهنای زیر بازه تغییر خواهد کرد، ضروری است.

انتگرالگیری تطبیقی اساساً، عبارت است از به کار بردن فرمولهای (۵۶.۷ الف) و (۵۶.۷ ب) برای هر زیر بازه‌ای که  $[a, b]$  را طی می‌کند تا نامساوی (۶۰.۷) برقرار شود. اگر نامساوی (۶۰.۷) در یک یا چند زیر بازه صدق نکند، آنگاه زیر بازه‌های مزبور به بازه‌های کوچکتر تقسیم و کل فرایند تکرار می‌شود.

هر زیر بر نامه‌ای که بر اساس انتگرالگیری تطبیقی نوشته می‌شود باید باریکهٔ همهٔ زیر بازه‌ها را حفظ کند تا اطمینان حاصل شود که بازهٔ  $[a, b]$  طی شده است، و باید پهنای زیر بازه‌ها، یعنی  $h$  مورد نیاز در فرمولهای (۵۶.۷ الف)، (۵۶.۷ ب)، (۶۰.۷) به طور مناسب برگزیده شود. پیچیدگی زیر بر نامه‌های انتگرالگیری تطبیقی ناشی از ساماندهی وسیع لازم برای حفظ باریکهٔ زیر بازه‌های تودرتو، و نیاز به جریانه‌های متفاوت عمل هنگام مواجهه با مشکلات است. با توجه به اینکه در فرمولهای (۵۶.۷ الف) و (۵۶.۷ ب) نقاطی که  $f(x)$  در آنها در (۵۶.۷ الف) محاسبه شده‌اند در همان نقاطی هستند که در (۵۶.۷ ب) وجود دارند، زیر بر نامه‌های تطبیقی که بر قاعدهٔ سیمپسن متکی هستند نیز کارا تر می‌شوند. بنابراین مقادیر  $f(x)$  در این نقاط را می‌توان کم کرد. مثال زیر روشی را که در اینجا بیان شد، روشن خواهد ساخت.

□ مثال ۸۰۷: با به کارگیری روش انتگرال تطبیقی بر اساس قاعدهٔ سیمپسن، تقریبی برای انتگرال

$$I = \int_0^1 \sqrt{x} dx$$

دقیق تا خطایی برابر با  $\varepsilon = 0.00005$  بیابید.

جواب دقیق برابر با  $I = 2/3$  به آسانی محاسبه می‌شود. ولی روشن است که برای محاسبهٔ انتگرال فوق حتی الامکان نباید از روش قاعدهٔ تطبیقی سیمپسن استفاده کرد. بسا

رسم نمودار تابع  $f(x) = \sqrt{x}$  مشاهده می‌شود که در نزدیکی مبدأ [در واقع  $f'(0) = \infty$ ] منحنی دارای شیب خیلی تندی است، درحالی که وقتی  $x \rightarrow 1$ ، منحنی نسبتاً هموار می‌شود. از این رو در یک بازه نزدیک به مبدأ در مقایسه با یک بازه نزدیک به ۱،  $x$ ، اشکالات زیادتری در امر انتگرالگیری وجود دارد.

عملیات را با تقسیم بازه  $[0, 1]$  به دو زیر بازه  $[0, 1/2]$  و  $[1/2, 1]$  آغاز می‌کنیم. در ابتدا فرمولهای (۵۶.۷ الف) و (۵۶.۷ ب) را در بازه  $[1/2, 1]$  به کار می‌بریم. در اینجا  $h = 1/2$ ، و از این رو خواهیم داشت

$$S\left[\frac{1}{2}, 1\right] = \frac{1}{12} \{ \sqrt{1/2} + 2\sqrt{3/4} + \sqrt{1} \} = 0.43093403$$

$$\bar{S}\left[\frac{1}{2}, 1\right] = \frac{1}{24} \{ \sqrt{1/2} + 2\sqrt{5/8} + 2\sqrt{3/4} + 2\sqrt{7/8} + \sqrt{1} \}$$

$$= 0.43096219$$

در اینجا برای مشخص ساختن زیر بازه مورد نظر، تقریباً نشانه گذاری متفاوتی به کار رفته است. به موجب فرمول خطای (۶۰.۷)، داریم

$$E\left[\frac{1}{2}, 1\right] = \frac{1}{15} (\bar{S} - S) = 0.0000018775 < \frac{1/2}{1} (0.00005)$$

$$= 0.000025$$

از آنجایی که ملاک خطا برقرار شده است، مقدار  $\bar{S}[1/2, 1]$  را می‌پذیریم و آن را در یک ثبات  $\sum$  کنار می‌گذاریم. سپس فرمولهای (۵۶.۷ الف) و (۵۶.۷ ب) را برای بازه  $[0, 1/2]$  به کار می‌گیریم. به ازای  $h = 1/2$ ، خواهیم داشت

$$S\left[0, \frac{1}{2}\right] = \frac{1}{12} \{ 0 + 2\sqrt{1/4} + \sqrt{1/2} \} = 0.22559223$$

$$\bar{S}\left[0, \frac{1}{2}\right] = \frac{1}{24} \{ 0 + 2\sqrt{1/8} + 2\sqrt{1/4} + 2\sqrt{3/8} + \sqrt{1/2} \} = 0.223211709$$

و

$$E\left[0, \frac{1}{2}\right] = 0.000043499 < 0.000025$$

در اینجا آزمون  $\sum$  خطا برقرار نیست و بنا بر این بازه  $[0, 1/2]$  تقسیم می‌شود. بانصف کردن این بازه، دو زیر بازه  $[0, 1/4]$  و  $[1/4, 1/2]$  به دست می‌آید. به ازای  $h = 1/4$  به کارگیری فرمولهای (۵۶.۷ الف) و (۵۶.۷ ب)، خواهیم داشت

$$S\left[\frac{1}{4}, \frac{1}{2}\right] = 0.15235819$$

$$\bar{S}\left[\frac{1}{4}, \frac{1}{2}\right] = 0.15236814$$

$$E\left[\frac{1}{4}, \frac{1}{2}\right] = 0.664 \times 10^{-6} < \frac{1/4}{1} (0.0005) = 0.000125$$

آشکار است که ملاک خطا برقرار است، از این رو مقدار  $\bar{S}[1/4, 1/2]$  را به محتوای ثبات SUM می‌افزاییم تا تقریب جزئی زیر را به دست آوریم

$$\text{SUM}\left[\frac{1}{4}, 1\right] = 0.23096219 + 0.15236814 = 0.38333033$$

با به کارگیری مجدد فرمول پایهٔ  $(56.7)$  در بازهٔ  $[0, 1/4]$  و به ازای  $h = 1/4$  خواهیم داشت

$$S\left[0, \frac{1}{4}\right] = 0.07975890$$

$$\bar{S}\left[0, \frac{1}{4}\right] = 0.08206578$$

$$D\left[0, \frac{1}{4}\right] = (0.0001537922) < 0.000125$$

آزمون خطا برقرار نبوده از این رو بازهٔ  $[0, 1/4]$  را به دو زیر بازهٔ  $[0, 1/8]$  و  $[1/8, 1/4]$  تقسیم می‌کنیم. با انجام عملیاتی مشابه و به ازای  $h = 1/8$  خواهیم داشت

$$S\left[\frac{1}{8}, \frac{1}{4}\right] = 0.05386675$$

$$\bar{S}\left[\frac{1}{8}, \frac{1}{4}\right] = 0.05387027$$

$$E\left[\frac{1}{8}, \frac{1}{4}\right] = 0.0000002346 < 1/8 (0.0005) = 0.0000625$$

$$S \left[ 0, \frac{1}{8} \right] = 0.002819903$$

$$\bar{S} \left[ 0, \frac{1}{8} \right] = 0.002901264$$

$$E \left[ 0, \frac{1}{8} \right] = 0.000005437 < 0.00000625$$

چون آزمون خطا در هر دو بازه برقرار است، می‌توانیم این مقادیر را به‌محتوای ثبات SUM اضافه کنیم، که با انجام این کار خواهیم داشت

$$\begin{aligned} P = \text{SUM} [0, 1] &= 0.58331033 + 0.05387027 + 0.002901264 \\ &= 0.639081864 \end{aligned}$$

از آنجایی که مقدار دقیق  $I$  برابر  $0.6666666666666666$  است، دیده می‌شود که در تمام بازه  $[0, 1]$  تقریب  $P$  برای  $I$  درملاک خطای مطلوب صدق می‌کند، یعنی

$$\square \quad |P - I| = 0.000045142 < 0.00005$$

همان گونه که این مثال نشان می‌دهد، در جاهایی که منحنی  $f(x)$  به آرامی تغییر می‌کند، روشهای انتگرالگیری تطبیقی بازه‌های بزرگتری را می‌طلبد، اما در جاهایی نظیر قله‌های تیز یا نزدیک به نقاط تکین، که منحنی  $f(x)$  سریعاً تغییر می‌کند، باید بازه‌های انتگرالگیری خیلی کوچکتر را به کار گرفت تا دقت عمل لازم به دست آید.

ما در اینجا زیر برنامه‌ای بر اساس روش انتگرالگیری تطبیقی عرضه نمی‌کنیم. همان گونه که قبلاً بیان شد، اگر قرار باشد که رده‌های بسیاری از توابع با این برنامه اجرا شوند، چنین زیر برنامه‌ای مطمئناً بسیار پیچیده خواهد بود. برای بسیاری از کامپیوترهای جدید، روشهای انتگرالگیری تطبیقی عالی تهیه کرده‌اند.

## تمرین

۱-۵۰۷ با استفاده از یک ماشین حساب جیبی نتایج به‌آمده از مثال ۸۰۷ را برای  $\bar{S}[0, 1/8]$ ،  $\bar{S}[1/8, 1/4]$ ،  $\bar{S}[1/4, 1/2]$  تحقیق کنید.

۲-۵۰۷ در مثال ۸۰۷، ملاک خطا را به  $\epsilon = 0.00001$  بدل کنید. کدامیک از بازه‌هایی که تا به حال تعیین شده است ملاک خطای مطلوب را برقرار می‌کند و کدامیک نمی‌کند؟ مادامی که ملاک خطای جدید برقرار نشده است، بازه  $[0, 1/8]$  را تقسیم و مانند مثال مذکور، انتگرال را محاسبه کنید.

۳-۵۰۷ با استفاده از روش انتگرالگیری تطبیقی براساس قاعده سیمپسن، تقریبی برای انتگرال

$$I = \int_0^1 (1-x^2)^{3/2} dx$$

با دقتی سه رقم اعشار بیابید. نخست منحنی  $f(x)$  را رسم و سپس سعی کنید که محلی را که انتظار می رود با اشکالاتی مواجه شوید، معین نمایید.

۴-۵۰۷ با استفاده از روش انتگرالگیری تطبیقی، تقریبی برای انتگرال زیر، با دقتی تا شش رقم اعشاری بیابید

$$I = \int_0^1 \frac{\sin x}{x^{3/2}} dx$$

۵-۵۰۷ برنامه‌ای برای انتگرالگیری تطبیقی براساس قاعده سیمپسن بنویسید و در این برنامه محدودیتهای زیر را در نظر بگیرید.

۱. ورودی شامل تابع  $f(x)$ ، یک بازه متناهی  $[a, b]$ ، و ملاک خطای مطلق  $\epsilon$  باشد.  
 ۲. برنامه می‌باید بازه  $[a, b]$  را به دو قسمت مساوی تقسیم کند و فرمولهای (۵۶.۷ الف)، (۵۶.۷ ب) و (۶۰.۷) را برای محاسبه  $S$  و  $\bar{S}$  و  $E$  به‌ازای هر قسمت به‌کار ببرد.

۳. اگر  $E$  در شرایط خطای مطلوب در یک زیر بازه صدق می‌کند،  $\bar{S}$  ذخیره شود؛ در غیر این صورت بازه مزبور نصف و مرحله ۲ تکرار شود.

۴. اگر تقسیم بازه لازم باشد این کار را حداکثر ۴ مرحله در یک زیر بازه انجام دهید.  
 ۵. خروجی این برنامه به‌صورت زیر باشد.

(i) اگر آزمون خطا در مجموعه‌ای از بازه‌هایی که  $[a, b]$  را طی می‌کند، برقرار شود یک متغیر به‌صورت عدد صحیح  $IFLAG = 1$  در خروجی داشته باشیم، و اگر آزمون خطا به‌ازای یک یا چند زیر بازه برقرار نشود، متغیر  $IFLAG = 2$  در خروجی ظاهر شود.

(ii) اگر  $IFLAG = 1$ ،  $P = \sum \bar{S}_i$  چاپ شود.

اگر  $IFLAG = 2$ ، در آن فواصلی که آزمون خطا برقرار می‌شود حاصلجمع جزئی  $PP = \sum \bar{S}_i$  چاپ شود، و فهرستی از بازه‌های  $[x_i, x_{i+1}]$  که در آنها آزمون خطا برقرار نیست نیز چاپ شود.

۶-۵۰۷ این حکم موجود در متن را تحقیق کنید: اگر رابطه خطای (۶۰.۷) در هر یک از  $N$  زیر بازه‌ای که بازه  $[a, b]$  را طی می‌کند برقرار باشد، آنگاه  $P = \sum_{i=1}^N \bar{S}_i$  در شرط خطای مطلوب (۵۵.۷) در کل بازه  $[a, b]$  صادق خواهد بود.

## ۶.۷\* برونیایی به سمت حد

در قسمتهای قبل، تلاش زیادی برای به دست آوردن عبارتهایی برای خطای قاعده‌های مختلف انتگرالگیری و مشتقگیری تقریبی به عمل آوردیم. خلاصه بگوییم: با  $L(f)$  که انتگرال  $f(x)$  در بازه  $[a, b]$  بود، یا با مقدار مشتق  $f(x)$  در نقطه  $a$ ، تقریب  $L_h(f)$  را برای  $L(f)$  که بستگی به پارامتر  $h$  دارد و در رابطه

$$\lim_{h \rightarrow 0} L_h(f) = L(f)$$

صدق می‌کند، به دست می‌آوریم. روشنتر بگوییم، مطابق معمول، رابطه

$$L(f) = L_h(f) + ch^r f^{(s)}(\xi)$$

را که در آن  $c$  یک عدد ثابت و  $s$  و  $r$  اعداد صحیح مثبت و  $\xi = \xi(h)$  یک نقطه نامعلوم از بازه مورد بحث بودند، ثابت و اشاره کردیم که حد مستقیم برای مقدار عبارت خطا مستلزم دانستن مقدار  $|f^{(s)}(\xi)|$  است، که این خود غالباً نمی‌تواند به اندازه کافی دقیق به دست آید (اگر نخواهیم بگوییم که اصلاً دقیق به دست نمی‌آید).

با وجود این، یک چنین عبارت خطایی میزان نزدیکی  $L_h(f)$  به  $L(f)$  را (وقتی  $h \rightarrow 0$ ) بیان می‌کند. این اطلاعات را می‌توان گاهی برای برآورد خطا از روی مقادیر متوالی  $L_h(f)$  به کار برد. امکان یک چنین برآوردی در قسمت ۶.۱ به طور مختصر ذکر شد و در بخش ۴.۳ مثال خاصی، فرایند  $\Delta^2$  ای تکن، مورد بحث قرار گرفت و مثال دیگری در بخش قبلی ۵.۷ آورده شد.

به عنوان یک مثال ساده، تقریب

$$D_h(f) = \frac{f(a+h) - f(a-h)}{2h}$$

برای مقدار

$$D(f) = f'(a)$$

را که مربوط به مشتق اول  $f(x)$  در نقطه  $x = a$  است در نظر می‌گیریم. اگر  $f(x)$  دارای سه مشتق پیوسته باشد، آنگاه به موجب (۸.۷) یا (۱۱.۷)، داریم

$$D(f) = D_h(f) - \frac{1}{6} h^2 f'''(\xi) \quad | \xi - a | < |h| \text{ با شرط}$$

چون وقتی  $h \rightarrow 0$ ، خواهیم داشت  $\xi(h) \rightarrow a$ ، و چون  $f'''(x)$  پیوسته است، وقتی  $h \rightarrow 0$ ، داریم

$$f'''(\xi) \rightarrow f'''(a)$$

از این رو

$$\frac{[f'''(\xi) - f'''(a)]h^2}{6}$$

سرریز از  $h^2$  به سمت صفر میل می‌کند. بنا بر این با استفاده از علامت گذاریهای مرتبه‌ای که در بخش ۶.۱ بیان شد، خواهیم داشت

$$D(f) = D_h(f) + C_1 h^2 + o(h^2) \quad (۶.۷)$$

که در آن ثابت  $C_1 = -f'''(a)/6$  به  $h$  بستگی ندارد. يك مثال عددی می‌تواند به روشن کردن اهمیت معادلهٔ (۶.۷) كمك كند. به ازای  $f(x) = \sin x$  و  $a = 1$  خواهیم داشت

$$D(f) = 0.5840402$$

$$C_1 = 0.0090050$$

در جدول ۲.۷ مقادیر  $D_h(f)$  خطای  $E_h(f) = -h^2 f'''(\xi)/6$ ، و دو مؤلفهٔ آن، یعنی  $C_1 h^2$  و  $o(h^2)$  به ازای مقادیر مختلف  $h$  داده شده‌اند. برای جلوگیری از اشکالهای جزئی ناشی از خطای گرد کردن، کلیهٔ عناصر این جدول با دقت مضاعف محاسبه و سپس گرد شده‌اند. چنانچه این جدول نشان می‌دهد،  $C_1 h^2$  سریعاً به مؤلفهٔ نافذاً خطا تبدیل می‌شود زیرا اگرچه  $C_1 h^2$  (همراه با  $h$ ) به سمت صفر میل می‌کند، ولی مؤلفهٔ  $o(h^2)$  سرریز از آن به سمت صفر نزدیک می‌شود. اما این امر بدین معناست که می‌توانیم بر آورد خوبی در مورد خطا برای مؤلفهٔ اصلی یعنی  $C_1 h^2$ ، به صورت زیر داشته باشیم: در رابطهٔ (۶.۷) به جای  $h$  مقدار  $2h$  را می‌گذاریم و داریم

جدول ۲.۷

$h$	$D_h(f)$	$E_h(f)$	$C_1 h^2$	$o(h^2)$	$(D_h - D_{2h})/3$	$R_h$
6.4	0.009839	0.530463	3.688464	-3.158001		
3.2	-0.009856	0.550158	0.922116	-0.371957	-0.065652	-0.57
1.6	0.337545	0.202757	0.230529	-0.027772	0.115800	2.37
0.8	0.484486	0.055816	0.057632	-0.001816	0.048980	3.54
0.4	0.526009	0.014293	0.014408	-0.000115	0.013841	3.88
0.2	0.536707	0.003594	0.003602	-0.000007	0.003566	3.97
0.1	0.539402	0.000900	0.000901	-0.0000005	0.000898	

۱. dominant (مؤلفهٔ نافذ = آن مؤلفه‌ای که تأثیر بیشتری روی خطا می‌گذارد). -م.

$$D(f) = D_{\gamma h}(f) + {}^4C_1 h^2 + o(h^2)$$

از کم کردن این معادله از (۶۱.۷) خواهیم داشت

$$0 = D_h(f) - D_{\gamma h}(f) - {}^3C_1 h^2 + o(h^2)$$

یا

$$C_1 h^2 = \frac{D_h(f) - D_{\gamma h}(f)}{3} + o(h^2) \quad (62.7)$$

معادلهٔ اخیر بیان می‌دارد که، به ازای  $h$  به قدر کافی کوچک، عدد قابل محاسبه

$$\frac{D_h(f) - D_{\gamma h}(f)}{3} \quad (63.7)$$

بر آورد خوبی است برای مؤلفهٔ خطای نافذ  $C_1 h^2$  که معمولاً مجهول است. این امر در جدول ۲.۷ که در آن اعداد به دست آمده از رابطهٔ (۶۳.۷) نیز ثبت شده اند، به خوبی مشاهده می‌شود.

البته نکته‌ای که در بررسیهای فوق وجود دارد، عبارت «به ازای  $h$  به قدر کافی کوچک» است. در حقیقت از جدول ۲.۷ مشاهده می‌شود که در مثال عددی فوق وقتی  $h = 1.6$ ،  $(D_h - D_{\gamma h})/3$  تنها بر آوردی برای حدود دامنهٔ  $C_1 h^2$  است. و هنگامی که  $h = 3.2$ ،  $(D_h - D_{\gamma h})/3$  حتی در حدود مقدار دامنه هم نیست. بنابراین عدد (۶۳.۷) را نباید بدون بررسی به عنوان بر آوردی از خطا قبول کرد. بلکه، می‌باید عملیات خود را در برابر اشتباه بزرگتر حفظ کرد که این امر بایک بررسی ساده بر اساس بحث زیر امکان پذیر است: اگر  $C_1 h^2$  مؤلفهٔ خطای نافذ باشد، یعنی اگر  $o(h^2)$  در مقایسه با  $C_1 h^2$  «کوچک» باشد، آنگاه به موجب (۶۲.۷) داریم

$$C_1 h^2 \approx \frac{D_h(f) - D_{\gamma h}(f)}{3}$$

و همچنین

$$C_1 \left(\frac{h}{\gamma}\right)^2 \approx \frac{D_{h/\gamma}(f) - D_h(f)}{3}$$

بنابراین خواهیم داشت

$$\frac{D_h(f) - D_{\gamma h}(f)}{D_{h/\gamma}(f) - D_h(f)} \approx \frac{C_1 h^2}{C_1 h^2 / \gamma^2} = \gamma^2$$

به گونهٔ توضیحی، اگر  $C_1 h^2$  مؤلفهٔ خطا باشد، آنگاه نسبت تفاضلهای قابل محاسبه یعنی



$$R_h = \frac{D_h(f) - D_{2h}(f)}{D_{h/2}(f) - D_h(f)} \quad (۶۴.۷)$$

می‌باید تقریباً برابر با ۴ باشد. این امر برای مثال عددی ما در جدول ۲.۷، که نسبتهای  $R_h$  را نیز داده‌ایم کاملاً آشکار است.

به محض رسیدن به این باور که (۶۳.۷) برآورد خوبی برای خطا در  $D_h(f)$  است و اطمینان مجدد به اینکه ۴  $R_h \approx$ ، آنگاه می‌توان انتظار داشت که

$$D_h^*(f) = D_h(f) + \frac{D_h(f) - D_{2h}(f)}{3} \quad (۶۵.۷)$$

در مقایسه با  $D_h(f)$ ، تقریب خیلی بهتری برای  $D(f)$  باشد. بخصوص در این صورت، این باور حاصل می‌شود که

$$|D(f) - D_h^*(f)| < \frac{|D_h(f) - D_{2h}(f)|}{3} \quad (۶۶.۷)$$

برای اینکه ببینیم تا چه حد تقریب  $D_h^*(f)$  بهتر است، اکنون بیان مشروحتری از عبارت خطای

$$E_h(f) = -\frac{1}{6}h^2 f'''(\xi)$$

برای  $D_h(f)$  به دست می‌آوریم. به منظور تنوع، به جای تفاضلهای منقسم، سری تیلر را به کار می‌گیریم. اگر  $f(x)$  دارای پنج مشتق پیوسته باشد، آنگاه با بسط  $f(a+h)$  و  $f(a-h)$  به سری تیلر در حول  $x=a$  خواهیم داشت

$$f(a+h) = f(a) + f'(a)h + \frac{f''(a)h^2}{2} + \frac{f'''(a)h^3}{6} + \frac{f^{iv}(a)h^4}{24} \\ + \frac{f^v(a)h^5}{120} + o(h^5)$$

$$f(a-h) = f(a) - f'(a)h + \frac{f''(a)h^2}{2} - \frac{f'''(a)h^3}{6} + \frac{f^{iv}(a)h^4}{24} \\ - \frac{f^v(a)h^5}{120} + o(h^5)$$

معادلهٔ دوم را از معادلهٔ اول کم، و سپس نتیجه را بر  $2h$  تقسیم می‌کنیم، که در نتیجه خواهیم داشت

$$D_h(f) = f'(a) + \frac{f'''(a)h^2}{6} + \frac{f^{(4)}(a)h^3}{120} + o(h^3)$$

بنابراین

$$D(f) = D_h(f) + C_1 h^2 + C_2 h^3 + o(h^3) \quad (67.7)$$

که در آن ثابتهای

$$C_1 = \frac{-f'''(a)}{6} \quad C_2 = \frac{-f^{(4)}(a)}{120}$$

به  $h$  بستگی ندارند. بنابراین با گذاردن  $2h$  به جای  $h$  در (67.7) داریم

$$D(f) = D_{2h}(f) + 4C_1 h^2 + 16C_2 h^3 + o(h^3) \quad (68.7)$$

از کم کردن  $\frac{1}{3}$  معادله (68.7) از  $\frac{2}{3}$  معادله (67.7)، خواهیم داشت

$$D(f) = D_h^*(f) + C_2^* h^3 + o(h^3) \quad (69.7)$$

که در آن داریم

$$C_2^* = -4C_2 = \frac{f^{(4)}(a)}{30}$$

و چون به موجب (65.7) داریم

$$D_h^*(f) = \frac{4D_h(f) - D_{2h}(f)}{3}$$

یک مقایسه (69.7) با (67.7) نشان می‌دهد که  $D_h^*(f)$  در مقایسه با  $D_h(f)$  تقریبی از مرتبه بالاتر برای  $D(f)$  است: اگر  $C_1 \neq 0$ ، آنگاه  $D(f) - D_h(f)$  همراه با  $h$  تنها با همان سرعت  $h^2$  به سمت صفر میل می‌کند، در حالی که  $D(f) - D_h^*(f)$  حداقل با سرعت  $h^3$  به سمت صفر میل خواهد کرد.

این روند به دست آوردن یک تقریب مرتبه بالاتر را از دو تقریب مرتبه پایینتر، معمولاً برون‌یابی به سمت حد یا برون‌یابی به سمت صفر اندازه‌ها نامند. (برای توضیح این اصطلاح به تمرین 6.7-3 نگاه کنید).

برون‌یابی به سمت حد به هیچ وجه به تقریبهای خطای  $O(h^2)$  محدود نمی‌شود، برای مثال، از رابطه (69.7) با قراردادن  $2h$  به جای  $h$ ، خواهیم داشت

$$D(f) = D_{2h}^*(f) + 16C_2^* h^3 + o(h^3)$$

بنابراین، از کم کردن این رابطه از (69.7) و مرتب کردن نتیجه داریم

$$C_{\psi}^{\setminus} h^{\psi} = \frac{D_h^{\setminus}(f) - D_{\psi h}^{\setminus}(f)}{15} + o(h^{\psi})$$

بنا بر این، اگر قرار دهیم

$$D_h^{\setminus}(f) = D_h(f) + \frac{D_h^{\setminus}(f) - D_{\psi h}^{\setminus}(f)}{15}$$

خواهیم داشت

$$D(f) = D_h^{\setminus}(f) + o(h^{\psi})$$

که نشان می‌دهد،  $D_h^{\setminus}(f)$  در مقایسه با  $D_h^{\setminus}(f)$  برای  $D(f)$  حتی تقریبی از مرتبه بالاتر است. روشنتر بگوییم، می‌توان نشان داد که اگر  $f(x)$  به قدر کافی هموار باشد، آنگاه

$$D(f) = D_h^{\setminus}(f) + C_{\psi}^{\setminus} h^{\psi} + o(h^{\psi}) \quad (70.7)$$

اما باید توجه کرد که به ازای هر مقدار خاص  $h$ ، نمی‌توان انتظار داشت که  $D_h^{\setminus}(f)$  در مقایسه با  $D_h^{\setminus}(f)$  تقریب بهتری برای  $D(f)$  باشد، مگر اینکه

$$\frac{D_h^{\setminus}(f) - D_{\psi h}^{\setminus}(f)}{15}$$

بر آورد خوبی برای خطا در  $D_h^{\setminus}(f)$  باشد، یعنی مگر اینکه  $C_{\psi}^{\setminus} h^{\psi}$  قسمت نافذ خطا در  $D_h^{\setminus}(f)$  باشد. این مطلب فقط زمانی صحیح است که تساوی زیر برقرار باشد

$$R_h^{\setminus} = \frac{D_h^{\setminus}(f) - D_{\psi h}^{\setminus}(f)}{D_{h/\psi}^{\setminus}(f) - D_h^{\setminus}(f)} \approx \frac{C_{\psi}^{\setminus} h^{\psi}}{C_{\psi}^{\setminus} (h/\psi)^{\psi}} = 16$$

بنا بر این، پیش از رسیدن به این باور که رابطه

$$|D(f) - D_h^{\setminus}(f)| < \frac{|D_h^{\setminus}(f) - D_{\psi h}^{\setminus}(f)|}{15}$$

برقرار است، این شرط باید بررسی شود. در جدول ۳.۷ نتایج دوبار به کار بردن برونیاپی به سمت حد را برای دنباله  $D_h^{\setminus}(f)$ ، که در جدول ۲.۷ محاسبه شده بود، فهرست کرده‌ایم. همچنین مقادیر مختلف  $R_h^{\setminus}$  را نیز داده‌ایم. کلیه محاسبات با گرد کردن تا شش رقم اعشاری انجام گرفته‌اند.

سرانجام، هیچ اشکالی برای به کار گرفتن عدد ۲، که در بالا به کار بردیم، در کلیه برونیاپیها وجود ندارد. برای مثال، اگر  $q$  در واقع عددی ثابت باشد، آنگاه از رابطه (۶۷.۷) نتیجه می‌گیریم

$$D(f) = D_{qh}(f) + q^{\psi} C_{\psi}^{\setminus} h^{\psi} + q^{\psi} C_{\psi}^{\setminus} h^{\psi} + o(h^{\psi})$$

از کم کردن این معادله از (۶۷.۷) و مرتب کردن نتیجه، خواهیم داشت

## جدول ۳.۷

$h$	$D_h(f)$	$R_h$	$D'_h(f)$	$R'_h$	$D''_h(f)$
6.4	0.009839				
3.2	-0.009856	-0.57	-0.075508		
1.6	0.337545	2.37	0.453345	6.1	0.488602
0.8	0.484486	3.54	0.533466	12.5	0.538807
0.4	0.526009	3.88	0.539850	15.1	0.540276
0.2	0.536707	3.97	0.540273	15.7	0.540301
0.1	0.539402		0.540300		0.540302

$$C_\gamma h^\gamma = \frac{D_h(f) - D_{qh}(f)}{q^\gamma - 1} - (1 + q^\gamma) C_\gamma h^\gamma + o(h^\gamma)$$

بنابراین، به ازای

$$D_{h,q}(f) = D_h(f) + \frac{D_h(f) - D_{qh}(f)}{q^\gamma - 1}$$

خواهیم داشت

$$D(f) = D_{h,q}(f) - q^\gamma C_\gamma h^\gamma + o(h^\gamma)$$

که نشان می‌دهد  $D_{h,q}(f)$  یک تقریب  $\Theta(h^\gamma)$  برای  $D(f)$  است. برای مثال از جدول ۳.۷ محاسبه می‌کنیم که

$$D_{0.1, 4}(f) = 0.539402 + \frac{0.539402 - 0.526009}{16 - 1} = 0.540295$$

که خطای آن تنها برابری با هفت واحد در آخرین رقم است. نکات اصلی بحثهای فوق در الگوریتم زیر گنجانیده شده است.

**الگوریتم ۱۰۷: پرونیاسایی به سمت حد.** برای محاسبه تقریب  $L_h(f)$  برای عدد  $L(f)$  به ازای هر  $h > 0$  راه‌حلهایی داده شده‌اند، که در آنها  $L_h(f)$  در رابطه زیر صدق می‌کند

$$L(f) = L_h(f) + Ch^r + o(h^r) \quad h > 0$$

در این رابطه  $C$  عددی است ثابت مستقل از  $h$  و  $r$  عددی است مثبت.

یک  $h$  و یک عدد  $q > 1$  (برای مثال  $q = 2$ ) را انتخاب می‌کنیم و محاسبه زیر را با دو عدد  $L_h(f)$  و  $L_{qh}(f)$  انجام می‌دهیم

$$L_{h,q}(f) = L_h(f) + \frac{L_h(f) - L_{qh}(f)}{q^r - 1}$$

آنگاه

$$L(f) = L_{h,q}(f) + o(h^r)$$

لذا برای  $h$  به قدر کافی کوچک داریم

$$|L(f) - L_{h,q}(f)| < |L_{h,q}(f) - L_h(f)| \quad (۷۱.۷)$$

پیش از اینکه اطمینان کامل از رابطه (۷۱.۷) به دست آید، می باید تحقیق کنیم که رابطه

$$\frac{L_h(f) - L_{qh}(f)}{L_{hp}(f) - L_h(f)} \approx \frac{Ch^r(q^r - 1)}{Ch^r(1 - p^{-r})} = p^r \frac{q^r - 1}{p^r - 1}$$

به ازای مقداری مانند  $p > 1$  (مثلا،  $p = q$ ) برقرار است.

## تمرین

۱-۶۰۷ با  $f(x) = x + x^2 + x^5$  و  $a = 0$ ، مقادیر  $D_h(f)$  و  $D'_h(f)$  را به ازای مقادیر مختلف  $h$  محاسبه کنید. چرا  $D'_h(f)$  در مقایسه با  $D_h(f)$  همواره تقریب بسدی برای  $D(f) = f'(0)$  است؟ (از حساب با دقت بسیار زیاد استفاده کنید تا خطای گرد کردن را به علت نارسا بودن رد کنید یا یک عبارت روشنی برای  $D'_h$  و  $D_h$  بر حسب  $h$  به دست آورید.)

۲-۶۰۷ با استفاده از روش برونمایی به سمت حد، مقدار  $f'(0.۴)$  را برای داده‌های زیر پیدا کنید

$x$	$\sinh x = f(x)$
۰.۳۹۸	۰.۴۰۸۵۹۱
۰.۳۹۹	۰.۴۰۹۶۷۱
۰.۴۰۰	۰.۴۱۰۷۵۲
۰.۴۰۱	۰.۴۱۱۸۳۴
۰.۴۰۲	۰.۴۱۲۹۱۵

در این حالت مقدار برونمایی شده تقریب خوبی نیست. شرح دهید که چرا چنین است. [توجه: مقدار صحیح  $f'(0.۴)$  برابر است با ۱.۱۰۵۸۱۰۷۲].

۳-۶۰۷ نشان دهید که برونمایی به سمت حد می تواند اساس جایگزینی تحلیلی گرفته شود

بخصوص با قراردادهای الگوریتم ۱.۷ نشان دهید که

$$L_{h,q}(f) = \lim_{x \rightarrow 0} p(x) \approx \lim_{x \rightarrow 0} L_x(f) = L(f)$$

کسه در آن تقریب  $p(x)$  برای  $g(x) = L_x(f)$  از راه پیدا کردن  $A$  و  $B$ ، به طوری که رابطه

$$p(x) = A + Bx^r$$

در  $x = h$  و  $x = qh$  با  $g(x)$  تطابق داشته باشد، به دست آمده است. چگونه این عمل نام «برونیایی به سمت حد» را توجیه می کند؟

### ۷.۷\* انتگرالگیری رامبرگ

برونیایی به سمت حد احتمالاً بیشتر به علت کار بردش با قاعده ذوزنقه ای مرکب، که نام انتگرالگیری رامبرگ به خود گرفته، بیشتر معروف است. با تقریب قاعده ذوزنقه ای مرکب (به قسمت ۴.۷ نگاه کنید)

$$T_N = T_N(f) = h \sum_{i=1}^{N-1} f_i + \frac{h(f_0 + f_N)}{2} \quad (72.7)$$

برای عدد

$$I = I(f) = \int_a^b f(x) dx$$

شروع می کنیم. در اینجا  $N$  عددی است صحیح و مثبت که با رابطه زیر به  $h$  بستگی دارد

$$h = \frac{b-a}{N}$$

و

$$f_i = f_{i,N} = f(a+ih) \quad i = 0, \dots, N$$

اگر  $f(x)$  چهار مرتبه به طور پیوسته مشتق پذیر باشد، از روابط (۵۰.۷) و (۵۱.۷) نتیجه می گیریم

$$I(f) = T_N(f) + C_1 h^2 + \theta(h^4) \quad (73.7)$$

که در آن ثابت  $C_1 = [f'(a) - f'(b)]/12$  مستقل از  $h$  است. بنابراین برونیایی به سمت حد، را می توان به کار برد و خواهیم داشت

$$T_{N,q}(f) = T_N + \frac{T_N(f) - T_{N/q}(f)}{q^2 - 1}$$

که يك تقريب  $\mathcal{O}(h^4)$  برای  $I(f)$  است، در حالی که در حالت کلی  $T_N(f)$  تنها دارای خطایی از مرتبه  $\mathcal{O}(h^2)$  است.

باید توجه کرد که انتخاب  $q$  یا  $N$  مقید به این شده است که  $N/q$  عدد صحیح باشد. معمولاً  $q=2$  انتخاب می شود (لذا  $N$  باید زوج باشد). این انتخاب  $q$  از نظر محاسباتی این مزیت را دارد که کلیه مقادیر تابع، که برای محاسبه  $T_{N/q}$  به کار گرفته می شود، برای محاسبه  $T_N$  نیز می تواند به کار رود. بخصوص ثابت می کنیم که هر گاه  $N$  عددی زوج باشد، داریم

$$T_N(f) = \frac{T_{N/2}(f)}{2} + h \sum_{i=1}^{N/2} f(a + (2i-1)h) \quad (74.7)$$

چون به موجب (74.7) داریم

$$\begin{aligned} T_N(f) &= h \sum_{i=1}^{N-1} f(a+ih) + \frac{h(f(a)+f(b))}{2} \\ &= h \sum_{i=1}^{N/2} f(a+(2i-1)h) + h \sum_{i=1}^{N/2} f(a+2ih) + \frac{h(f(a)+f(b))}{2} \end{aligned}$$

در اینجا، مجموعیابی اول بر نقاط «فرد» و مجموعیابی دوم بر نقاط «زوج» تعمیم داده شده است. دو جمله آخر را می توان به صورت زیر نوشت

$$\left[ 2h \sum_{j=1}^{N/2-1} f(a+j(2h)) + \frac{2h(f(a)+f(b))}{2} \right] / 2$$

بنابراین، از آنجایی که داریم

$$2h = \frac{2(b-a)}{N} = \frac{b-a}{N/2}$$

این دو جمله آخر به  $T_{N/2}(f)/2$  اضافه می شود. و همین امر رابطه (74.7) را ثابت می کند. باید توجه داشت که (74.7) را می توان به شکل ساده تری به صورت زیر نوشت

$$T_N(f) = \frac{T_{N/2}(f) + M_{N/2}(f)}{2}$$

که در آن  $M$  معرف قاعده مرکب نقطه میانی (۴۹.۷ الف) است. اگر تابع زیر علامت انتگرال دارای  $2k+2$  بار مشتق پیوسته باشد، می توان به طریقی

روشنتر از (۷۳.۷) نشان داد

$$I(f) = T_N(f) + C_1 h^2 + C_2 h^4 + C_3 h^6 + \dots + C_k h^{2k} + \mathcal{O}(h^{2k+2})$$

که در آن ثابتهای  $C_1, \dots, C_k$  مستقل از  $h$  اند. از این رو با توجه به

$$T_N^1(f) = T_N(f) + \frac{T_N(f) - T_{N/2}(f)}{3}$$

خواهیم داشت

$$I(f) = T_N^1(f) + C_1^1 h^2 + C_2^1 h^4 + \dots + C_k^1 h^{2k} + \mathcal{O}(h^{2k+2})$$

که در آن ثابتهای  $C_1^1, \dots, C_k^1$  مستقل از  $h$  هستند. بنا بر این بر رویایی بعدی معنی پیدا می کند. با قراردادن

$$T_N^2(f) = T_N^1(f) + \frac{T_N^1(f) - T_{N/2}^1(f)}{15}$$

داریم

$$I(f) = T_N^2(f) + C_1^2 h^2 + \dots + C_k^2 h^{2k} + \mathcal{O}(h^{2k+2})$$

به شکل کلیتر، دیده می شود که به ازای  $k, \dots, 1, m$

$$T_N^m(f) = T_N^{m-1}(f) + \frac{T_N^{m-1}(f) - T_{N/2}^{m-1}(f)}{4^m - 1}$$

یک  $\mathcal{O}(h^{2m+2})$  تقریبی برای  $I(f)$  است.

توجه داریم که محاسبه  $T_N^m$  مستلزم محاسبه  $T_{N/2}^{m-1}$  و  $T_N^{m-1}$  و بنا بر این مستلزم محاسبه  $T_N^{m-2}, T_{N/2}^{m-3}, T_{N/4}^{m-3}, \dots$  و سرانجام  $T_{N/2^m}$  است. بنا بر این  $N/2^m$  باید عدد صحیحی مثلا به صورت

$$\frac{N}{2^m} = M$$

برای  $T_N^m$  تعریف شده باشد. ساده تر این است که این تقریبهای مختلف برای  $I(f)$  را به صورت عناصر یک آرایه مثلثی، موسوم به جدول  $T$  مجسم کنیم

$T_M^0$				
$T_{2M}^1$	$T_M^1$			
$T_{4M}^2$	$T_{2M}^2$	$T_M^2$		
.....				
$T_{2^m M}^m$	$T_{2^{m-1} M}^m$	$T_{2^{m-2} M}^m$	...	$T_M^m$



در اینجا به جای  $T_N^0$ ،  $T_N$  نوشته شده است.

**الگوریتم ۲۰۷:** انتگرالگیری رامبرگ تابع  $f(x)$  که در بازه  $[a, b]$  تعریف شده و عدد صحیح مثبت  $M$  (معمولا  $M = 1$ ) داده شده اند.

$$h := (b - a) / M$$

$$\text{Calculate } T_M^0 = h \sum_{i=1}^{M-1} f(a + ih) + h(f(a) + f(b)) / 2$$

For  $k = 1, 2, 3, \dots$ , do:

$$h := h / 2$$

$$\text{Calculate } T_{2^k M}^0 = \frac{1}{2} T_{2^{k-1} M}^0 + h \sum_{i=1}^{2^k M - 1} f(a + (2i - 1)h)$$

For  $m = 1, \dots, k$ , do:

$$\text{Calculate } T_{2^m M}^m = T_{2^{m-1} M}^{m-1} + (T_{2^m M}^{m-1} - T_{2^{m-1} M}^{m-1}) / (4^m - 1)$$

اگر  $f(x)$  دارای  $2m + 2$  بار مشتق پیوسته باشد، آنگاه

$$I(f) = \int_a^b f(x) dx = T_{\sqrt[k]{k}M}^m + O\left(\left[\frac{b-a}{\sqrt[k]{k}M}\right]^{2m+2}\right) \quad k = m, m+1, \dots$$

همچنین، اگر  $k$  به اندازه کافی بزرگ باشد، آنگاه

$$|I(f) - T_{\sqrt[k]{k}M}^m| < |T_{\sqrt[k]{k}M}^m - T_{\sqrt[k]{k}M}^{m-1}|$$

اما پیش از اینکه اطمینان کاملی از این نامساوی داشته باشیم، بررسی می کنیم که دست کم رابطه زیر برقرار باشد

$$R_k^{m-1} = \frac{T_{\sqrt[k]{k-1}M}^{m-1} - T_{\sqrt[k]{k}M}^{m-1}}{T_{\sqrt[k]{k}M}^{m-1} - T_{\sqrt[k]{k-1}M}^{m-1}} \approx \frac{1}{k}$$

□ **مثال ۹۰۷:** انتگرالگیری رامبرگ را برای مثال ۱۰۷ به کار گیرید. انتگرال مورد نظر به صورت زیر است

$$I(f) = \int_0^1 e^{-x^2} dx$$

برنامه فورترنی که در زیر داده شده، شش ردیف اول جدول  $T$  و جدول نسبتهای  $R_k^m$  متناظر با آن را به صورت زیر پدید می آورد:

جدول  $T$  رامبرگ

---

0.7313700E 00					
0.7429838E 00	0.7468551E 00				
0.7458653E 00	0.7468258E 00	0.7468238E 00			
0.7465842E 00	0.7468238E 00	0.7468237E 00	0.7468237E 00		
0.7467639E 00	0.7468237E 00	0.7468237E 00	0.7468237E 00	0.7468237E 00	
0.7468069E 00	0.7468212E 00	0.7468210E 00	0.7468210E 00	0.7468209E 00	

---

## جدول نسبتها

4.03				
4.01	14.88			
4.00	16.50	0.0		
4.17	0.05	0.0	0.0	

مقدار  $M$  برابر با ۲ انتخاب شده بود، لذا اولین عنصر جدول  $T$ ،  $T_4(f)$  است. توجه کنید که اولین ستون نسبتها به نحو مطلوبی به سمت ۴ همگرا می شود، ولی از آن به بعد شروع می کند به دور شدن از ۴. این نتیجه حتی در ستون دوم نسبتها که به سمت ۱۶ میل می کند (چنانکه باید باشد) مشخصتر است و سپس، آن گونه که آخرین عنصر نشان می دهد غیر قابل پیشبینی می شود. نتیجه اینکه: خطا در درایه های آخرین سطر جدول  $T$ ، عمدتاً بر اثر گرد کردن حاصل شده نه بر اثر تجزیه، بنا بر این به نظر می رسد که

$$0.77468237$$

بهترین بروردی برای  $I(f)$  است که از این محاسبات می توانسته به دست آید. چون

$$R_4 = 16.5 \approx 16 = 4^2$$

و

$$T_{4,M}^4 = T_{4,M}^4 = 0.77468237$$

که این تساویها تا تعداد ارقام نشان داده شده برقرارند، و می توان نتیجه گرفت که این برورد تا تعداد ارقام نشان داده شده صحیح است. ولی به طور دقیق داریم

$$\int_0^1 e^{-x^2} dx = 0.774682413279 \dots$$

که اختلاف بین این دو عدد و برورد «صحیح» انجام شده ناشی از این حقیقت است که در محاسبات با تابع زیر علامت انتگرال

$$f(x) = e^{-x^2}$$

سروکار نداریم بلکه با  $f(x)$  گرد شده، یعنی با تابع

$$F(X) = \text{EXP}(-X * X)$$

مواجه هستیم. کلیه محاسبات با دقت ساده روی IBM ۳۶۰ که مخصوصاً مشخصه نامطلوبی در گرد کردن دارد، انجام گرفته است.

## برنامه فورترن برای مثال ۹.۷

```

REAL T(100)
EXTERNAL FERR
CALL RMBERG( FERR, 0., 1., 2, T, 6 )
                                STOP
END
SUBROUTINE RMBERG ( F, A, B, MSTART, T, NROW )
C  CONSTRUCTS AND PRINTS OUT THE FIRST NROW ROWS OF THE ROMBERG T-
C  TABLE FOR THE INTEGRAL OF F(X) FROM A TO B , STARTING WITH THE
C  TRAPEZOIDAL SUM ON MSTART INTERVALS.
  INTEGER MSTART,NROW, I,K,M
  REAL A,B,T(NROW,NROW), H,SUM
  M = MSTART
  H = (B-A)/M
  SUM = (F(A) + F(B))/2.
  IF (M .GT. 1) THEN
    DO 10 I=1,M-1
      10  SUM = SUM + F(A+FLOAT(I)*H)
    END IF
    T(1,1) = SUM*H
    PRINT 610
610  FORMAT('1',10X,'ROMBERG T-TABLE'//)
    PRINT 611, T(1,1)
611  FORMAT('E15.7')
    IF (NROW .LT. 2) RETURN
C
    DO 20 K=2,NROW
      H = H/2.
      M = M*2
      SUM = 2.
      DO 11 I=1,M,2
        11  SUM = SUM + F(A+FLOAT(I)*H)
      T(K,1) = T(K-1,1)/2. + SUM*H
      DO 12 J=1,K-1
        C  SAVE DIFFERENCES FOR LATER CALC. OF RATIOS
        12  T(K-1,J) = T(K,J) - T(K-1,J)
        20  PRINT 611, (T(K,J),J=1,K)
      IF (NROW .LT. 3) RETURN
C
      CALCULATE RATIOS
      PRINT 620
620  FORMAT('///11X','TABLE OF RATIOS'//)
      DO 30 K=1,NROW-2
        DO 25 J=1,K
          IF (T(K+1,J) .EQ. 0.) THEN
            RATIO= 0.
          ELSE
            RATIO = T(K,J)/T(K+1,J)
          END IF
          25  T(K,J) = RATIO
        30  PRINT 630, (T(K,J),J=1,K)
630  FORMAT(8F10.2)
                                RETURN
      END
      REAL FUNCTION FERR(X)
      REAL X
      FERR = EXP(-X*X)
                                RETURN
    END
  END

```

□

## تمرین

۱-۷۰۷ ثابت کنید که در انتگرالگیری دامبرگ تساوی  $T_{M}^1 = S_M$  برقرار است که در آن  $S_M$ ، قاعدهٔ مرکب سیمپسن است، ر.ک. (۴۸۰۷).

۲-۷۰۷ سعی کنید مقدار  $I(f) = \int_a^b f(x) dx$  را تا تقریب  $10^{-6}$  برآورد کنید. برای

انجام این کار، انتگرالگیری رامبرگ را برای هر یک از موارد زیر به کار گیرید:

$$M \text{ اختیاری}, b=1, a=0 \quad f(x)=x^2 \quad (\text{الف})$$

$$M=1, b=1, a=0 \quad f(x)=\sin 101\pi x \quad (\text{ب})$$

$$M=1, b=1, a=0 \quad f(x)=1+\sin 10\pi x \quad (\text{پ})$$

$$M=3 \text{ و } M=1, b=1, a=0 \quad f(x)=\left|x-\frac{1}{3}\right| \quad (\text{ت})$$

$$M \text{ اختیاری}, b=1, a=0 \quad f(x)=\sqrt{x} \quad (\text{ث})$$

$$b=1, a=0 \quad f(x)=\begin{cases} \frac{\sin x}{x} & x \neq 0 \\ 1 & x = 0 \end{cases} \quad (\text{ج})$$

۳-۷.۷ با استفاده از انتگرالگیری رامبرگ و به ازای داده‌های زیر، انتگرال  $\int_1^{1.8} f(x) dx$  را تا آنجا که ممکن است با دقت محاسبه کنید. با شروع از  $M=1$ ، جدول  $T$  را پر کنید.

$x$	$f(x)$
۱.۰	۰.۳۶۷۸۷۹۴۴
۱.۱	۰.۳۶۶۱۵۸۱۹
۱.۲	۰.۳۶۱۴۳۳۰۵
۱.۳	۰.۳۵۴۲۹۱۳۳
۱.۴	۰.۳۴۵۲۳۵۷۴
۱.۵	۰.۳۳۴۶۹۵۲۴
۱.۶	۰.۳۲۳۰۳۴۴۳
۱.۷	۰.۳۱۰۵۶۱۹۹
۱.۸	۰.۲۹۷۵۳۸۰۰

۴-۷.۷ با برونیابی از قاعده نقطه میانی و قاعده ذوزنقه‌ای، برای

$$I_h(f) = \int_{-h}^h f(x) dx$$

قاعدهٔ سیمپسن را به دست آورید. (دانهمایی: ترکیب خطی مناسبی از دو معادلهٔ

$$I_h(f) = T(f) + C_T h^2 + \mathcal{O}(h^4) \quad I_h(f) = M(f) + C_M h^2 + \mathcal{O}(h^4)$$

برای حذف جملات  $h^2$  تشکیل دهید، برای انجام این کار، باید مقادیر ثابتهای  $C_M$  و  $C_T$  را پیدا کنید).



## حل معادلات دیفرانسیل

بسیاری از مسائل مهندسی و علوم را می‌توان برحسب معادلات دیفرانسیل بیان کرد. قسمت اعظم انگیزه ساختن کامپیوترهای ابتدایی ناشی از نیاز به محاسبات دقیق و سریع مسیر پرتابه‌ها بوده است. امروزه کامپیوترها در حل معادلات مربوط به پرتاب موشکها و نظریهٔ اعمار مصنوعی و نیز در مورد نظریهٔ شبکه‌های الکتریکی، خمش پرتوهای نوری، پایداری هواپیما، نظریهٔ ارتعاش و نظریه‌های دیگر به‌طور وسیعی به‌کار گرفته می‌شوند.

فرض بر آن است که دانشجویان با نظریهٔ مقدماتی معادلات دیفرانسیل آشنایی دارند. در ریاضیات اول، تکنیکهای مختلفی برای حل رستهٔ مشخصی از معادلات دیفرانسیل به‌یک شکل معین آموزش داده می‌شود. ولی اکثر معادلاتی را که در عمل با آنها مواجه می‌شویم نمی‌توانیم با روشهای تحلیلی حل کنیم و بناچار لازم می‌شود از روشهای عددی استفاده کنیم. خوشبختانه برای حل معادلات دیفرانسیل بساکامپیوتر، روشهای بسیار خوبی وجود دارند. در این فصل، چند رده از این روشها را استخراج و آنها را از جنبهٔ کارایی محاسباتی ارزیابی خواهیم کرد.

### ۱۰۸ مقدمات ریاضی

در اینجا مرور برخی تعاریف و مفاهیم مقدماتی نظریهٔ معادلات دیفرانسیل را بسیار مفید می‌دانیم. معادله‌ای را که شامل رابطه‌ای بین یک تابع مجهول و یک یا چند مشتق آن باشد

معادلهٔ دیفرانسیل نامند. ما همواره فرض خواهیم کرد که می توان این معادله را نسبت به مشتق بالاترین مرتبه اش حل کرد. يك معادلهٔ دیفرانسیل معمولی از مرتبهٔ  $n$  به شکل زیر خواهد بود

$$y^{(n)}(x) = f(x, y(x), y'(x), \dots, y^{(n-1)}(x)) \quad (۱.۸)$$

منظور ما از جواب معادلهٔ (۱.۸) تابعی است به صورت  $\phi(x)$  که به طور پیوسته در يك بازهٔ معلوم،  $n$  بار مشتقپذیر باشد و در معادلهٔ (۱.۸) صدق کند، یعنی  $\phi(x)$  باید در رابطه‌های به صورت

$$\phi^{(n)}(x) = f(x, \phi(x), \phi'(x), \phi''(x), \dots, \phi^{(n-1)}(x))$$

صدق نماید. معمولاً جواب عمومی معادلهٔ (۱.۸) دارای  $n$  ثابت اختیاری است و بنابراین يك خانوادهٔ  $n$  پارامتری از جوابها وجود دارد. اگر  $y(x_0)$ ،  $y'(x_0)$ ،  $\dots$ ، و  $y^{(n-1)}(x_0)$  در نقطهٔ  $x = x_0$  تعریف شده باشند، آنگاه يك مسئله با مقدار اولیه خواهیم داشت. همواره فرض ما بر این است که تابع  $f$  در شرایط کافی برای تضمین يك جواب یگانه برای این مسئله با مقدار اولیه، صدق می کند.  $y' = y$  يك مثال ساده از معادلهٔ مرتبهٔ يك است، جواب عمومی آن  $y(x) = Ce^x$  است که در آن  $C$  يك ثابت اختیاری است. در صورتی که شرط اولیهٔ  $y(x_0) = y_0$  معلوم باشد، جواب رami توان به شکل  $y(x) = y_0 e^{(x-x_0)}$  نوشت. علاوه معادلات دیفرانسیل به دو ردهٔ خطی و غیر خطی تقسیم بندی شده اند. يك معادله را زمانی خطی گویند که تابع  $f$  در معادلهٔ (۱.۸) به طور خطی شامل  $y$  و مشتقهای آن باشد. معادلات دیفرانسیل خطی این ویژگی مهم را دارند که اگر  $y_1(x)$ ،  $y_2(x)$ ،  $\dots$ ،  $y_m(x)$  جوابهایی از (۱.۸) باشند، آنگاه  $C_1 y_1(x) + C_2 y_2(x) + \dots + C_m y_m(x)$  به ازای ثابتهای اختیاری  $C_i$  نیز يك جواب خواهد بود.  $y'' = y$  يك معادلهٔ ساده از مرتبهٔ دو است. می توان به آسانی تحقیق کرد که  $e^x$  و  $e^{-x}$  جوابهای این معادله هستند و بنابراین به موجب ویژگی خطی بودن، حاصل جمع زیر نیز يك جواب آن است:

$$y(x) = C_1 e^x + C_2 e^{-x} \quad (۲.۸)$$

دو جواب  $y_1$  و  $y_2$  از معادلهٔ دیفرانسیل مرتبهٔ دو را مستقل خطی گویند، اگر در مینان ورونیسکی<sup>۱</sup> این جوابها مخالف صفر باشد، که بنا بر تعریف، در مینان ورونیسکی به صورت زیر است

$$W(y_1, y_2) = y_1 y_2' - y_2 y_1' = \begin{vmatrix} y_1 & y_1' \\ y_2 & y_2' \end{vmatrix} \quad (۳.۸)$$

مفهوم استقلال خطی را می توان برای جوابهای معادلات از مرتبهٔ بالاتر نیز تعمیم داد. اگر  $y_1(x)$ ،  $y_2(x)$ ،  $\dots$ ،  $y_n(x)$   $n$  جواب به طور خطی مستقل از يك معادلهٔ دیفرانسیل همگن

درجه  $n$  باشند، آنگاه

$$y(x) = C_1 y_1(x) + C_2 y_2(x) + \dots + C_n y_n(x)$$

جواب عمومی نامیده می‌شود.

در بین معادلات خطی، آنهایی که ضرایب ثابت دارند به‌ویژه مفیدند، از این جهت که راه‌حل آنها ساده‌تر است. معادله دیفرانسیل مرتبه  $n$ ام با ضرایب ثابت به‌صورت زیر نوشته می‌شود

$$Ly = y^{(n)} + a_{n-1}y^{(n-1)} + \dots + a_0y^{(0)} = 0 \quad (۴.۸)$$

که در آن  $a_i$ ها حقیقی فرض شده‌اند. اگر به دنبال جوابهایی از (۴.۸) که به‌صورت  $e^{\beta x}$  هستند باشیم، از قرار دادن مستقیم این مقدار در معادله می‌بینیم که  $\beta$  باید در بسجمله‌ای زیر صدق کند

$$\beta^n + a_{n-1}\beta^{n-1} + \dots + a_0 = 0 \quad (۵.۸)$$

این معادله را معادله مشخصه معادله دیفرانسیل مرتبه  $n$ ام (۴.۸) نامند. اگر معادله (۵.۸)  $n$  ریشه متمایز  $\beta_i (i=1, \dots, n)$  داشته باشد، آنگاه می‌توان نشان داد که عبارت

$$y(x) = C_1 e^{\beta_1 x} + C_2 e^{\beta_2 x} + \dots + C_n e^{\beta_n x} \quad (۶.۸)$$

جواب عمومی معادله (۴.۸) است، که در آن  $C_i$ ها مقادیر ثابت اختیاری هستند. اگر  $\beta_1 = \alpha + i\beta$  یک ریشه همگام از معادله (۵.۸) باشد، مزدوج آن یعنی  $\beta_2 = \alpha - i\beta$  نیز یک ریشه است. متناظر با این ریشه‌های همگام مزدوج، دو جواب  $y_1 = e^{\alpha x} \cos \beta x$  و  $y_2 = e^{\alpha x} \sin \beta x$  وجود دارد که مستقل خطی هستند. هنگامی که رابطه (۵.۸) دارای ریشه‌های مکرر باشد، تکنیکهای خاصی برای به‌دست آوردن جوابهای مستقل خطی وجود دارند. در حالت خاص، اگر  $\beta_1$  یک ریشه مکرر رابطه (۵.۸) باشد، آنگاه  $y_1 = e^{\beta_1 x}$  و  $y_2 = x e^{\beta_1 x}$  جوابهای مستقل خطی (۴.۸) می‌شوند. برای معادله خاص  $y'' + a^2 y = 0$ ، معادله مشخصه به‌صورت  $\beta^2 = -a^2$  می‌باشد که ریشه‌های آن عبارت انداز  $\beta_{1,2} = \pm ia$  و جواب عمومی آن به‌صورت  $y(x) = C_1 \cos ax + C_2 \sin ax$  است. و سرانجام، اگر معادله (۱.۸) خطی ولی غیر همگن باشد، یعنی اگر

$$Ly = g(x) \quad (۷.۸)$$

و اگر  $\zeta(x)$  یک جواب خصوصی (۷.۸) باشد، یعنی اگر

$$L\zeta = g(x)$$

آنگاه با فرض آنکه ریشه‌های (۵.۸) متمایزند، جواب عمومی (۷.۸) به‌صورت زیر درمی‌آید



$$y = \zeta(x) + C_1 e^{\beta_1 x} + C_2 e^{\beta_2 x} + \dots + C_n e^{\beta_n x} \quad (۱.۸)$$

□ مثال: جواب معادله

$$y'' - 4y' + 3y = x \quad (\text{الف})$$

را که در شرایط اولیه

$$y'(0) = \frac{7}{3}, y(0) = \frac{4}{9} \quad (\text{ب})$$

صدق کند پیدا کنید.

۱. برای پیدا کردن يك جواب خاص  $\zeta(x)$  برای قسمت (الف) تابع  $ax + b$  را به عنوان جواب بررسی می کنیم، زیرا سمت راست تساوی يك بسجمله ای از درجه نا بیشتر از يك است، و سمت چپ نیز يك چنین بسجمله ای است هر گاه  $y = y(x)$  با گذاردن  $\zeta(x)$  در (الف) مقادیر  $a = 1/3$ ،  $b = 4/9$  را به دست خواهیم آورد. بنا بر این

$$\zeta(x) = \frac{1}{3}x + \frac{4}{9}$$

۲. برای پیدا کردن جوابهای معادله همگن

$$y'' - 4y' + 3y = 0$$

معادله مشخصه زیر را بررسی می کنیم

$$\beta^2 - 4\beta + 3 = 0$$

ریشه های این معادله عبارت اند از  $\beta_1 = 3$  و  $\beta_2 = 1$ . بنا بر این دو جواب مستقل خطی از دستگاه همگن، به گونه زیرند

$$y_1(x) = e^{3x} \quad y_2(x) = e^x$$

۳. جواب عمومی معادله (الف) به صورت زیر است

$$y(x) = \frac{x}{3} + \frac{4}{9} + C_1 e^{3x} + C_2 e^x$$

۴. برای پیدا کردن جوابهایی که در شرایط (ب) صدق کنند، باید داشته باشیم

$$y(0) = \frac{4}{9} + C_1 + C_2 = \frac{4}{9}$$

$$y'(0) = \frac{1}{3} + 3C_1 + C_2 = \frac{7}{3}$$

جواب این دستگاه برابر است با  $C_1 = 1$ ،  $C_2 = -1$ . بنا بر این جواب مورد نظر برابر است با

□

$$y(x) = \frac{x}{3} + \frac{2}{9} + e^{3x} - e^x$$

### تمرین

۱-۱۰۸ جواب عمومی معادلات زیر را پیدا کنید

$$y'' - 4y' + 4y = 0 \quad (\text{ب}) \qquad y' = -2y \quad (\text{الف})$$

$$y' - ay = x \quad (\text{ت}) \qquad y''' - 2y'' - y' + 2y = 0 \quad (\text{پ})$$

$$y'' - 2y' + 2y = 0 \quad (\text{ج}) \qquad y' - xy = e^x \quad (\text{ث})$$

۲-۱۰۸ جوابهای مسائل با مقدار اولیه زیر را پیدا کنید:

$$y(0) = 1 \qquad y' + 2y = 1 \quad (\text{الف})$$

$$y'(0) = 1, y(0) = 0 \qquad y'' - a^2y = 0 \quad (\text{ب})$$

$$y'(0) = 1, y(0) = 0 \qquad y'' - 4y' + 4y = x \quad (\text{پ})$$

### ۲.۸ معادلات تفاضلی ساده

برای تحلیل روشهای عددی در حل معادلات دیفرانسیل، لازم است که با نظریه ساده معادلات تفاضلی مختصری آشنایی پیدا کنیم. یک معادله تفاضلی از مرتبه  $N$  رابطه‌ای است بین تفاضلهای  $\Delta^0 y_n, \Delta^1 y_n, \Delta^2 y_n, \dots, \Delta^N y_n$  از یک دنباله عددی، یعنی

$$\Delta^N y_n = f(n, y_n, \Delta y_n, \dots, \Delta^{N-1} y_n) \quad (9.8)$$

یک جواب این معادله تفاضلی به صورت دنباله‌ای است از اعداد  $y_m, y_{m+1}, y_{m+2}, \dots$  به طوری که به ازای  $n = m, m+1, m+2, \dots$  رابطه (۹.۸) برقرار باشد. بنا بر این در صورتی که یک معادله دیفرانسیل متضمن توابعی باشد که خود و مشتقاتشان بر بازه‌ای از اعداد حقیقی تعریف شده باشند، یک معادله تفاضلی متضمن توابعی است که خود و تفاضلاتشان بر «بازه»‌ای از اعداد صحیح تعریف شده باشند.

اگر (۹.۸) یک معادله تفاضلی خطی باشد که سمت راست آن به طور خطی به



معادله تفاضلی خطی همگن از مرتبه  $N$  با ضرایب ثابت

$$y_{n+N} + a_{N-1}y_{n+N-1} + \dots + a_0y_n = 0 \quad (12.8)$$

را به تفصیل مورد بررسی قرار می‌دهیم. همان گونه که در معادلات دیفرانسیل خطی همگن با ضرایب ثابت عمل می‌کردیم، به دنبال جوابهایی می‌گردیم که به شکل  $y_n = \beta^n$  (به ازای جمیع مقادیر) باشند. از گذاردن این مقدار در (۱۲.۸)، نتیجه زیر به دست می‌آید

$$\beta^{n+N} + a_{N-1}\beta^{n+N-1} + \dots + a_0\beta^n = 0$$

که از تقسیم آن بر  $\beta^n$ ، معادله مشخصه

$$p(\beta) = \beta^N + a_{N-1}\beta^{N-1} + \dots + a_0 = 0 \quad (13.8)$$

به دست می‌آید که بسجمله‌ای مشخصه‌ای است از درجه  $N$ . در ابتدا فرض می‌کنیم که ریشه‌های آن یعنی  $\beta_1, \beta_2, \dots, \beta_N$  متمایز باشند، سپس  $\beta_1^n, \beta_2^n, \dots, \beta_N^n$  همگی جوابهای (۱۲.۸) هستند و بنا بر ویژگی خطی بودن، نتیجه می‌شود که مقدار

$$y_n = c_1\beta_1^n + c_2\beta_2^n + \dots + c_N\beta_N^n \quad n \text{ به ازای جمیع مقادیر } n \quad (14.8)$$

که  $c_i$ ها در آن ثابتهای اختیاری هستند، یک جواب معادله (۱۲.۸) است. بعلاوه، در این مورد می‌توان نشان داد که (۱۴.۸) جواب عمومی معادله (۱۲.۸) است. به عنوان مثال، معادله تفاضلی

$$y_{n+3} - 2y_{n+2} - y_{n+1} + 2y_n = 0 \quad (15.8)$$

معادله‌ای است از مرتبه سوم و دارای معادله مشخصه زیر

$$\beta^3 - 2\beta^2 - \beta + 2 = 0$$

ریشه‌های این معادله بسجمله‌ای برابر است با  $1, -1, 2$  و  $2$  و جواب عمومی معادله (۱۵.۸) چنین است

$$\begin{aligned} y_n &= c_1(1)^n + c_2(-1)^n + c_3(2)^n \\ &= c_1 + (-1)^n c_2 + 2^n c_3 \end{aligned} \quad (16.8)$$

اگر اولین  $N-1$  مقدار  $y_n$  داده شده باشند، معادله تفاضلی با مقدار اولیه حاصل را می‌توان آشکارا به ازای جمیع مقادیری که به  $n$  داده می‌شود حل کرد. بنا بر این اگر در (۱۵.۸)، مقادیر اولیه  $y_0 = 0, y_1 = 1, y_2 = 1$  داده شده باشند، مقدار  $c_3$  که از (۱۵.۸) محاسبه می‌شود برابر است با

$$c_3 = 2(1) + 1 - 0 = 3$$

اگر استفاده از (۱۵.۸) را همچنان ادامه دهیم، خواهیم داشت  $y_4 = 5$ ،  $y_5 = 11$  و قس علیهذا. اما این عمل يك فرمول ثابتی برای  $y_n$  به ما نمی دهد. ولی با استفاده از رابطه (۱۶.۸) و قرار دادن شرایط اولیه  $n = 0, 1, 2$ ، دستگاه معادلات زیر را برای  $c_1, c_2, c_3$  به دست می آوریم

$$0 = c_1 + c_2 + c_3$$

$$1 = c_1 - c_2 + 2c_3$$

$$1 = c_1 + c_2 + 4c_3$$

که جواب آن برابر است با  $c_1 = 1/3$ ،  $c_2 = -1/3$ ،  $c_3 = 1/3$ ، لذا جواب ثابت<sup>۱</sup> برای این مسئله با مقدار اولیه به صورت زیر است

$$y_n = -\frac{1}{3}(-1)^n + \frac{2}{3}$$

اگر بسجمله ای مشخصه (۱۳.۸) دارای يك جفت ریشه مزدوج همناف باشد، باز هم جواب می تواند به شکل حقیقی بیان شود. بنا بر این اگر داشته باشیم  $\beta_1 = \alpha + i\beta$  و  $\beta_2 = \alpha - i\beta$ ، در ابتدا  $\beta_1, \beta_2$  را به شکل قطبی<sup>۲</sup> بیان می کنیم

$$\beta_1 = re^{i\theta}$$

$$\beta_2 = re^{-i\theta}$$

که در آن  $r = \sqrt{\alpha^2 + \beta^2}$  و  $\theta = \arctan(\beta/\alpha)$ . در این حال جواب (۱۲.۸) متناظر با این ریشه های مزدوج به صورت زیر خواهد شد

$$\begin{aligned} c_1 \beta_1^n + c_2 \beta_2^n &= c_1 r^n e^{in\theta} + c_2 r^n e^{-in\theta} \\ &= r^n [c_1 (\cos n\theta + i \sin n\theta) + c_2 (\cos n\theta - i \sin n\theta)] \\ &= r^n (C_1 \cos n\theta + C_2 \sin n\theta) \end{aligned}$$

که در آن  $C_1 = c_1 + c_2$  و  $C_2 = i(c_1 - c_2)$ . به عنوان يك مثال ساده، معادله تفاضلی زیر را در نظر می گیریم

$$y_{n+2} - 2y_{n+1} + 2y_n = 0 \quad (17.8)$$

معادله مشخصه آن  $\beta^2 - 2\beta + 2 = 0$  و ریشه های این معادله برابر است با  $\beta_{1,2} = 1 \pm i$ . بنا بر این داریم  $r = \sqrt{2}$  و  $\theta = \pi/4$ ، لذا جواب عمومی معادله (۱۷.۸)

به صورت زیر است

$$y_n = (\sqrt{r})^n \left( C_1 \cos \frac{n\pi}{\varphi} + C_2 \sin \frac{n\pi}{\varphi} \right)$$

اگر  $\beta_1$  يك ریشه مضاعف معادله مشخصه (۱۳.۸) باشد، آنگاه جواب دوم معادله (۱۳.۸) برابر است با  $n\beta_1^n$ . برای تحقیق این امر، ملاحظه می کنیم که اگر  $\beta_1$  يك ریشه مضاعف  $p(\beta)$  باشد، آنگاه  $p(\beta_1) = 0$  و نیز  $p'(\beta_1) = 0$ . حال با قراردادن  $y_n = n\beta_1^n$  در معادله (۱۲.۸) و مرتب کردن آن و با توجه به تساوی  $p(\beta_1) = p'(\beta_1) = 0$  خواهیم داشت

$$\begin{aligned} (n+N)\beta_1^{n+N} + a_{n-1}(n+N-1)\beta_1^{n+N-1} + \dots + a_0 n\beta_1^n \\ = \beta_1^n \{ n(\beta_1^n + a_{N-1}\beta_1^{N-1} + \dots + a_0) \\ + \beta_1 [N\beta_1^{N-1} + a_{N-1}(N-1)\beta_1^{N-2} + \dots + a_1] \} \\ = \beta_1^n [np(\beta_1) + \beta_1 p'(\beta_1)] = 0 \end{aligned}$$

بعلاوه می توان نشان داد که دو جواب  $\beta_1^n$  و  $n\beta_1^n$  مستقل خطی هستند. به عنوان مثال، معادله تفاضلی زیر را در نظر می گیریم

$$y_{n+3} - 5y_{n+2} + 8y_{n+1} - 4y_n = 0$$

ریشه های معادله مشخصه آن عبارت اند از ۱، ۲، ۲، و جواب عمومی آن به صورت زیر است

$$y_n = 2^n(c_1 + nc_2) + c_3$$

و بالاخره جواب معادله تفاضلی خطی غیر همگن با ضرایب ثابت را مورد بررسی قرار می دهیم. جواب عمومی معادله

$$y_{n+N} + a_{N-1}y_{n+N-1} + \dots + a_0 y_n = b_n \quad (18.8)$$

را می توان به شکل زیر نوشت

$$y_n = y_n^c + y_n^p$$

که در آن  $y_n^c$  جواب عمومی دستگاه همگن (۱۲.۸) و  $y_n^p$  جواب خصوصی (۱۸.۸) است. در حالت خاص، وقتی که  $b_n = b$  يك عدد ثابت باشد، جواب خصوصی را می توان با قراردادن  $y_n^p = A$  (يك عدد ثابت) در معادله (۱۸.۸) به آسانی به دست آورد. گذاردن  $y_n = A$  در معادله (۱۸.۸)، به شرطی که مجموع ضرایب آن صفر نباشد، منجر به تعیین مقدار  $A$  زیر خواهد شد

$$A = \frac{b}{1 + a_{N-1} + \dots + a_0}$$

مثلا جواب عمومی معادلهٔ غیر همگن

$$y_{n+2} - 2y_{n+1} + 2y_n = 1$$

به صورت زیر است

$$y_n = (\sqrt{2})^n \left( C_1 \cos \frac{n\pi}{4} + C_2 \sin \frac{n\pi}{4} \right) + 1$$

خواص سادهٔ معادلات تفاضلی که تا اینجا بررسی شدند، برای استفادهٔ در بقیهٔ این فصل کافی هستند.

□ مثال: نشان دهید که جواب عمومی معادلهٔ تفاضلی

$$y_{n+2} - (2+h^2)y_{n+1} + y_n = h^2 \quad h > 0 \quad (\text{الف})$$

را می‌توان به شکل زیر بیان کرد

$$y_n = c_1 \left[ 1 + h + \frac{h^2}{2} + O(h^3) \right]^n + c_2 \left[ 1 - h + \frac{h^2}{2} + O(h^3) \right]^n - 1 \quad (\text{ب})$$

حل:

۱. جواب خصوصی معادلهٔ (الف) از قرارداد  $y_n^p = C$  در (الف) به دست می‌آید که برابر است با  $y_n^p = -1$ .
۲. معادلهٔ مشخصهٔ مربوط به معادلهٔ غیر همگن (الف) به صورت زیر است

$$\beta^2 - (2+h^2)\beta + 1 = 0$$

به موجب فرمول ریشه‌های معادلهٔ درجهٔ دوم، داریم

$$\begin{aligned} \beta_{1,2} &= \frac{2+h^2 \pm \sqrt{4h^2+h^4}}{2} \\ &= \frac{2+h^2 \pm 2h\sqrt{1+h^2/4}}{2} \\ &= 1+h^2/2 \pm h\left(1+\frac{h^2}{4}\right)^{1/2} \end{aligned}$$

از بسط  $(1+t)^{1/2}$  به سری تیلر بیرامون  $t=0$  و قراردادن  $h^2/4$  به جای  $t$  خواهیم داشت

$$\beta_{1,2} = 1 + \frac{h^2}{4} \pm h \left[ 1 + \frac{h^2}{8} + O(h^4) \right]$$

$$\beta_1 = 1 + h + \frac{h^2}{4} + O(h^3)$$

$$\beta_2 = 1 - h + \frac{h^2}{4} + O(h^3)$$

بنابراین جواب عمومی دستگاه همگن به صورت زیر است

$$y_n^c = c_1 \beta_1^n + c_2 \beta_2^n$$

۳. در نتیجه جواب معادله (الف) عبارت است از

$$y_n = y_n^p + y_n^c$$

□

که بودن جواب به شکل (ب) را ثابت می کند.

### تمرین

۱-۲۰۸ جواب عمومی معادلات تفاضلی زیر را پیدا کنید

$$y_{n+1} - 3y_n = 5 \quad (\text{الف})$$

$$y_{n+2} - 4y_{n+1} + 4y_n = n \quad (\text{ب})$$

(دانهمایی: برای پیدا کردن جواب خصوصی،  $y_n^p = an + b$  را آزمایش کنید).

$$y_{n+2} + 2y_{n+1} + 2y_n = 0 \quad (\text{ب})$$

$$y_{n+2} - y_{n+1} + 2y_n = 0 \quad (\text{ت})$$

$$y_{n+2} - y_{n+1} - y_n = 0 \quad (\text{ث})$$

۲-۲۰۸ جواب معادلات تفاضلی با مقدار اولیه زیر را پیدا کنید

$$y_1 = 1 \quad y_0 = 0 \quad y_{n+2} - 4y_{n+1} + 3y_n = 2^n \quad (\text{الف})$$

$$y_1 = 1 \quad y_0 = 0 \quad y_{n+2} - y_{n+1} - y_n = 0 \quad (\text{ب})$$

(دانهمایی: برای پیدا کردن يك جواب خصوصی معادله (الف)،  $y_n^p = A2^n$  را امتحان کنید).]



۳-۲۰.۸ نشان دهید که جواب عمومی معادله تفاضلی

$$y_{n+2} + 2hy_{n+1} - y_n = 2h$$

را که در آن  $h$  عدد ثابت و مثبتی است، می توان به شکل زیر نوشت

$$y_n = c_1 [1 - 2h + \theta(h^2)]^n + c_2 (-1)^n [1 + 2h + \theta(h^2)]^n + \frac{1}{\theta}$$

۴-۲۰.۸ نشان دهید که اگر  $y_0 = 1$ ،  $y_1 = x$ ، آنگاه جمله  $n$ ام،  $y_n = y_n(x)$ ، از جواب معادله

$$y_{n+2} - 2xy_{n+1} + y_n = 0$$

یک بسجمله ای است از درجه  $n$  بر حسب  $x$  که ضریب پیشرو آن  $2^{n-1}$  است. [توجه:  $y_n(x)$  ها بسجمله ایهای چیشف هستند که در بخش ۱۰.۶ دیدیم.]

### ۳.۸ انتگرالگیری عددی به وسیله سری تیلر

اکنون آماده ایم که روشهای عددی انتگرالگیری معادلات دیفرانسیل را بررسی کنیم. ابتدا یک معادله دیفرانسیل مرتبه اول با مقدار اولیه به شکل

$$y' = f(x, y) \quad y(x_0) = y_0 \quad (19.8)$$

را در نظر می گیریم. تابع  $f$  ممکن است خطی یا غیر خطی باشد، اما فرض بر آن است که  $f$  به اندازه کافی هم نسبت به  $x$  و هم نسبت به  $y$  مشتق پذیر است. می دانیم که اگر  $\partial f / \partial y$  در حوزة مورد نظر پیوسته باشد، معادله (۱۹.۸) دارای جواب یگانه است. اگر  $y(x)$  جواب دقیق (۱۹.۸) باشد، می توانیم  $y(x)$  را پیرامون نقطه  $x = x_0$  به سری تیلر بسط دهیم:

$$y(x) = y_0 + (x - x_0)y'(x_0) + \frac{(x - x_0)^2}{2!} y''(x_0) + \dots \quad (20.8)$$

روشن است که در این بسط مشتقها معلوم نیستند زیرا که جواب نامعلوم است. ولی اگر  $f$  به اندازه کافی مشتق پذیر باشد، می توان مشتقات مذکور را با گرفتن مشتق کلی از معادله (۱۹.۸) نسبت به  $x$ ، با توجه به اینکه  $y$  خود تابعی از  $x$  است، به دست آورد (به بخش ۷.۱ نگاه کنید). بنابراین برای چند مشتق اولیه داریم

$$y' = f(x, y)$$

$$y'' = f' = f_x + f_y y' = f_x + f_y f$$

$$y''' = f'' = f_{xx} + f_{xy} f' + f_{yx} f' + f_{yy} f'^2 + f_y f_{xx} + f_y^2 f'$$

$$= f_{xx} + 2f_{xy} f' + f_{yy} f'^2 + f_x f_y + f_y^2 f' \quad (21.8)$$

به همین ترتیب، می‌توانیم هر مشتقی از  $y$  را بر حسب  $f(x, y)$  و مشتقات جزئی آن بیان کنیم. اما کاملاً آشکار است که مادامی که  $f(x, y)$  تابع خیلی ساده‌ای نباشد، مشتقات کلی از درجات بالاتر بسیار پیچیده خواهند شد. لذا به دلایل عملی، می‌باید تعداد جملات را در بسط (۲۰.۸) به تعداد مناسبی محدود سازیم، که این محدودیت به محدودیت مقدار  $x$  در رابطه (۲۰.۸) منجر می‌شود که می‌تواند تقریب موجهی برای (۲۰.۸) باشد. اگر فرض کنیم که سری ناتمام (۲۰.۸) تقریب خوبی برای نمو  $h$ ، یعنی به ازای  $x - x_0 = h$ ، به دست می‌دهد، آنگاه می‌توانیم  $y$  را در  $x_0 + h$  محاسبه و مجدداً مقادیر  $y'$  و  $y''$  و غیره را در  $x = x_0 + h$  حساب کنیم و سپس (۲۰.۸) را برای مرحله بعدی به کار ببریم. اگر عمل را به همین طریق ادامه دهیم، مجموعه گسسته‌ای از مقادیر  $y_n$  خواهیم داشت که تقریبهایی برای جواب حقیقی در نقاط  $x_n = x_0 + nh$  ( $n = 0, 1, 2, \dots$ ) هستند. در این فصل همه‌جا مقدار جواب دقیق در نقطه  $x_n$  را با  $y(x_n)$  و جواب تقریبی را با  $y_n$  نمایش می‌دهیم.

برای آنکه روش فوق را به صورت دستوری بیان کنیم، ابتدا عملگر  $T$  زیر را وارد می‌نماییم

$$T_k(x, y) = f(x, y) + \frac{h}{1!} f'(x, y) + \dots + \frac{h^{k-1}}{k!} f^{(k-1)}(x, y)$$

$$k = 1, 2, \dots \quad (22.8)$$

که در آن فرض بر آن است که از یک نموناب  $h$  استفاده شده و  $f^{(l)}$  معرف مشتق کلی مرتبه  $l$ ام تابع  $f(x, y(x))$  نسبت به  $x$  است. حال می‌توانیم الگوریتم ۱۰.۸ را بیان کنیم.

**الگوریتم ۱۰.۸:** الگوریتم تیلر از مرتبه  $k$  برای پیدا کردن یک جواب تقریبی معادله دیفرانسیل

$$y' = f(x, y)$$

$$y(a) = y_0$$

بر یک بازه  $[a, b]$ :

۱. نمو  $h = (b - a) / N$  را انتخاب و فرض می‌کنیم

$$x_n = a + nh \quad n = 0, 1, \dots, N$$

۲. تقریبهای  $y_n$  را برای  $y(x_n)$  با استفاده از فرمول بازگشتی

$$y_{n+1} = y_n + hT_k(x_n, y_n) \quad n = 0, 1, \dots, N-1$$

که در آن  $T_k(x, y)$  به وسیله رابطه (۲۲.۸) تعریف شده است، پدید می آورییم. غالباً الگوریتم تیلر و روشهای دیگری را که بر اساس این الگوریتم مبتنی هستند و در آنها برای محاسبه  $y$  در نقطه  $x = x_{n+1}$  تنها از اطلاعات مربوط به  $y$  و  $y'$  در نقطه  $x = x_n$  استفاده می کنند، روشهای یک مرحله ای نامند. قضیه تیلر یا باقیمانده (قضیه ۷.۱)، نشان می دهد که خطای موضعی الگوریتم تیلر از مرتبه  $k$  به شرح زیر است

$$E = \frac{h^{k+1} f^{(k)}(\xi, y(\xi))}{(k+1)!} \quad x_n < \xi < x_n + h$$

$$= \frac{h^{k+1}}{(k+1)!} y^{(k+1)}(\xi)$$

الگوریتم تیلر را از مرتبه  $k$  نامند اگر خطای موضعی  $E$ ، که در بالا تعریف شد، برابر با  $\theta(h^{k+1})$  باشد. با قراردادن  $k=1$  در الگوریتم ۱.۸، روش اولیو و خطای موضعی آن به دست می آید

$$y_{n+1} = y_n + hf(x_n, y_n)$$

$$E = \frac{h^2}{2} y''(\xi) \quad (23.8)$$

برای روشن ساختن روش اولیو، مسئله با مقدار اولیه

$$y' = y \quad y(0) = 1$$

را در نظر می گیریم. با به کار گرفتن فرمول (۲۳.۸) به ازای  $h=0.01$  و حفظ شش رقم اعشار، خواهیم داشت

$$y(0.01) \approx y_1 = 1 + 0.01 = 1.01$$

$$y(0.02) \approx y_2 = 1.01 + 0.01(1.01) = 1.0201$$

$$y(0.03) \approx y_3 = 1.0201 + 0.01(1.0201) = 1.030301$$

$$y(0.04) \approx y_4 = 1.030301 + 0.01(1.030301) = 1.040606$$

## 1. one-step

\* خطای موضعی، خطایی است که در هر مرحله به ازای  $n=0, 1, 2, \dots$  وجود دارد. م.

از آنجایی که  $y = e^*$  جواب دقیق این معادله است، مقدار صحیح آن در  $x = 0.04$  برابر با  $1.0408$  خواهد شد. واضح است که برای به دست آوردن دقت بیشتر در روش اولر، باید مقدار  $h$  را به طور قابل ملاحظه‌ای کوچک بگیریم. اگر بگیریم  $h = 0.005$ ، خواهیم داشت

$$y(0.005) \approx y_1 = 1.0050$$

$$y(0.010) \approx y_2 = 1.0100$$

$$y(0.015) \approx y_3 = 1.0151$$

$$y(0.020) \approx y_4 = 1.0202$$

$$y(0.025) \approx y_5 = 1.0253$$

$$y(0.030) \approx y_6 = 1.0304$$

$$y(0.035) \approx y_7 = 1.0356$$

$$y(0.040) \approx y_8 = 1.0408$$

این نتایج تا چهار رقم اعشار صحیح اند. از آنجا که  $h$  می‌باید نسبتاً کوچک باشد، معمولاً برای انتگرالگیری معادلات دیفرانسیل از روش اولر استفاده نمی‌شود.

البته برای به دست آوردن دقت بیشتر می‌توانستیم الگوریتم تیلر از مرتبه بالاتر را به کار ببریم و در حالت کلی باید انتظار داشته باشیم که برای یک نمو معین، هر قدر مرتبه الگوریتم بالاتر باشد، به همان اندازه دقت بیشتری به دست آید. اگر  $f(x, y)$  تابع نسبتاً ساده‌ای از  $x$  و  $y$  باشد، می‌توان مشتقات لازم را غالباً با هزینه نسبتاً کمتری در کامپیوتر به دست آورد. این عمل با به کار بردن مشتقگیری نمادی\* یا با بهره‌گیری از هر یک از ویژگی‌های تابع  $f(x, y)$  که ممکن باشد، می‌تواند انجام پذیرد (به‌تقریب ۳.۸-۴ نگاه کنید). ولی لزوم محاسبه مشتقات مرتبه بالاتر، الگوریتم تیلر را برای انتگرالگیری کلی در کامپیوترهای سریع کاملاً نامناسب می‌سازد. مع‌هذا، الگوریتم تیلر از لحاظ نظری مورد توجه زیاد است زیرا در بسیاری از روشهای عملی سعی بر این است که دقت عمل مشابهی مانند الگوریتم تیلر به دست آورند، بی‌آنکه با وضع نامساعد محاسبه مشتقات مرتبه بالا مواجه باشند. اگرچه الگوریتم کلی تیلر به ندرت برای مقاصد عملی به کار برده می‌شود، ولی حالت خاص آن یعنی روش اولر، با شرح بیشتری برای نشان دادن اهمیت نظری آن مورد بررسی قرار می‌گیرد.

□ مثال ۱۰۸: با استفاده از سری تیلر جواب معادله دیفرانسیل

\* مشتقگیری نمادی به معنای مشتقگیری با استفاده از روشهای غیر عددی است. م.

$$xy' = x - y \quad y(2) = 2$$

را در  $x = 2$  از  $x = 2$  تا پنج رقم اعشاری دقیق به دست آورید.  
چند مشتق اول و مقادیر آنها در  $x = 2$  و  $y = 2$  به شرح زیرند:

$$y' = 1 - \frac{y}{x} \quad y'_0 = 0$$

$$y'' = \frac{-y'}{x} + \frac{y}{x^2} \quad y''_0 = \frac{1}{2}$$

$$y''' = \frac{-y''}{x} + \frac{2y'}{x^2} - \frac{2y}{x^3} \quad y'''_0 = -\frac{3}{4}$$

$$y^{iv} = \frac{-y'''}{x} + \frac{3y''}{x^2} - \frac{6y'}{x^3} + \frac{6y}{x^4} \quad y^{iv}_0 = \frac{3}{2}$$

بسط سری تیلر پیرامون نقطه  $x_0 = 2$  چنین است:

$$y(x) = y_0 + (x-2)y'_0 + \frac{1}{2}(x-2)^2 y''_0 + \frac{1}{6}(x-2)^3 y'''_0$$

$$+ \frac{1}{24}(x-2)^4 y^{iv}_0 + \dots$$

$$= 2 + (x-2) \cdot 0 + \frac{1}{2}(x-2)^2 - \frac{1}{8}(x-2)^3 + \frac{1}{16}(x-2)^4 + \dots$$

در  $x = 2$  خواهیم داشت

$$y(2) = 2 + 0.00025 - 0.0000125 + 0.000000625 - \dots$$

$$\approx 2.000238$$

چون جمله‌ها در این سری تیلر از نظر کمیت کاهش می‌یابند و از لحاظ علامت متناوباً تغییر می‌کنند (به تمرین ۳۰۸-۴ نگاه کنید)، این نتیجه تا پنج رقم اعشاری صحیح است. اکنون اگر بخواهیم  $y(2.2)$  را با همان دقت به دست آوریم، می‌باید دو جمله دیگر از سری فوق را نیز به دست آوریم. یا اینکه می‌توانیم بسط جدیدی پیرامون  $x = 2$  بدسیم و چهار مشتق اول را مجدداً در  $x = 2$  تعیین و سپس  $y(2.2)$  را محاسبه کنیم. □

$$y' = \frac{1}{x^2} - \frac{y}{x} - y^2$$

$$y(1) = -1$$

را از  $x=1$  تا  $x=2$  حل کنید. از الگوریتم مرتبه ۲ی تیلر استفاده می‌کنیم. مسئله را به‌ازای  $h=1/16, 1/32, 1/64, 1/128$  حل و دقت نتایج را برآورد می‌کنیم  
حل: چون

$$f(x, y) = \frac{1}{x^2} - \frac{y}{x} - y^2$$

$$f'(x, y) = -\frac{2}{x^3} - \frac{y'}{x} + \frac{y}{x^2} - 2yy'$$

لذا

$$T_2(x, y) = f + \frac{h}{2} f'$$

و

$$y_{n+1} = y_n + h \left[ f(x_n, y_n) + \frac{h}{2} f'(x_n, y_n) \right]$$

نتایج به‌دست آمده با کامپیوتر IBM ۷۰۹۰ در زیر داده شده‌اند. نمو  $h$  درستون اول و مقادیر  $y(200), y'(200), y(100), y'(100)$  به ترتیب در چهارستون بعدی داده شده‌اند. جواب دقیق این معادله  $y = -1/x$  است، لذا مقدار دقیق  $y(100)$  برابر  $2/3 -$  و مقدار دقیق  $y(200)$  برابر  $1/2 -$  می‌شود، می‌توانیم خطای کلی ناشی از گسسته‌سازی را به‌طریق زیر برآورد کنیم: خطای موضعی الگوریتم مرتبه ۲ی تیلر برابر  $y'''(h^3/6)$  است. چون داریم  $y''' = 6/x^4$ ، مقدار ماکزیم آن در بازه  $[1, 2]$  برابر ۶ است و بنا براین در هر مرحله، خطای موضعی حداکثر برابر با  $h^3$  است. به‌ازای  $h = 1/128$ ، مرحله انتگرال‌گیری داریم به‌طوری که خطای انباشته شده، حداکثر برابر با  $0.000006 \approx (1/128)^2 = 128 h^2$  است. در توافق نزدیک با این مقدار برآورد، ظاهراً خطای حقیقی در  $x = 200$ ، برابر با  $0.000003$  است. در حالت کلی، جواب واقعی را در دست نداریم تا صحت آن را آزمایش کنیم. حتی با نداشتن جواب، به‌ازای  $h \rightarrow 0$  می‌توانیم دقت جواب را از روی تعداد ارقام موافق برآورد کنیم. از آنجایی که با هر بار نصف کردن  $h$ ، تقریباً یک رقم دیگر به‌دقت افزوده می‌شود، ظاهراً در نبود خطای ناشی از گرد کردن، نمو  $1/1024$  می‌باید حداقل هفت رقم دقیق تولید نماید.

## نتایج کامپیوتری برای مثال ۲۰۸

روش ۱- روش مرتبه دوم بسط تیلر

H	Y(1.5)	YPRM(1.5)	Y(2.)	YPRM(2.)
0.6250000E-01	-0.66787238E 00	0.44363917E-00	-0.50187737E 00	0.24905779E-00
0.3125000E-01	-0.66696430E 00	0.44424593E-00	-0.50046334E 00	0.24976812E-00
0.1562500E-01	-0.66674034E 00	0.44439532E-00	-0.50011456E 00	0.24994271E-00
0.7812500E-02	-0.66668454E 00	0.44443253E-00	-0.50002744E 00	0.24998628E-00

روش ۲- روش ساده شده مرتبه دوم رونگه- کوتا

H	Y(1.5)	YPRM(1.5)	Y(2.)	YPRM(2.)
0.6250000E-01	-0.66552725E 00	0.44520275E-00	-0.49822412E-00	0.25088478E-00
0.3125000E-01	-0.66637699E 00	0.44463748E-00	-0.49954852E-00	0.25022554E-00
0.1562500E-01	-0.66659356E 00	0.44449317E-00	-0.49988601E-00	0.25005698E-00
0.7812500E-02	-0.66664808E 00	0.44445683E-00	-0.49997083E-00	0.25001458E-00

روش ۳- روش کلاسیک مرتبه چهارم رونگه- کوتا

H	Y(1.5)	YPRM(1.5)	Y(2.)	YPRM(2.)
0.6250000E-01	-0.66666625E 00	0.44444472E-00	-0.49999941E-00	0.25000029E-00
0.3125000E-01	-0.66666664E 00	0.44444446E-00	-0.49999997E-00	0.25000001E-00
0.1562500E-01	-0.66666666E 00	0.44444444E-00	-0.50000000E 00	0.25000000E-00
0.7812500E-02	-0.66666667E 00	0.44444444E-00	-0.50000001E 00	0.24999999E-00

همین مسئله بعداً با استفاده از دوروش دیگر حل خواهد شد. به منظور مقایسه، نتایج به دست آمده از هر سه روش در اینجا داده شده‌اند.

□

## تمرین

۳-۳۰۸ برای معادله

$$y' = -xy + \frac{1}{y^2} \quad y(1) = 1$$

معادله تفاضلی متناظر با الگوریتم مرتبه ۴ی تیلر را به دست آورید. یک مرحله از انتگرالگیری را به ازای  $h = 0.1$  با دست انجام دهید. برنامه‌ای برای حل این مسئله بنویسید و انتگرالگیری را از  $x = 1$  تا  $x = 2$  به ازای  $h = 1/64$  و  $h = 1/28$  انجام دهید.

۳-۳۰۸ برای معادله

$$y' = 2y \quad y(0) = 1$$

جواب دقیق معادله تفاضلی حاصل از روش اویلر را تعیین کنید. مقدار «به اندازه کافی کوچک»  $h$  را چنان برآورد کنید که این مقدار، جوابی با دقت چهار رقم دارد بازه  $[0, 1]$  تضمین کند. به ازای مقدار مناسبی از  $h$ ، حل این مسئله را برای ۱۰ مرحله انجام دهید.

۳-۸.۸ اگر  $y(x)$  در رابطه زیر صدق کند، از سری تیلر مربوط به  $y(x)$ ، مقدار  $y(0.1)$  را تا شش رقم اعشاری دقیق به دست آورید

$$y' = xy + 1 \quad y(0) = 1$$

۴-۳.۸ ثابت کنید که برای تابع  $f(x, y) = 1 - y/x$  از مثال ۱-۸،  $y'' = (1 - 2y')/x$ ،  $y^{(k)} = -ky^{(k-1)}/x$ ،  $k = 3, 4, \dots$ ، برای این اساس با استفاده از الگوریتم ۱-۸، یک برنامه فورترن بنویسید که مقدار  $y(3)$  را تا دقت  $10^{-6}$  از روی جواب  $y(x)$  در مثال ۱-۸ به دست دهد.

### ۴.۸ برآوردهای خطا و همگرایی روش اویلر

برای حل معادله دیفرانسیل  $y' = f(x, y)$ ،  $y(x_0) = y_0$  با روش اویلر، نمونابند  $h$  را انتخاب و از فرمول

$$y_{n+1} = y_n + hf(x_n, y_n) \quad n = 0, 1, \dots \quad (23.8)$$

که در آن  $x_n = x_0 + nh$ ، استفاده می‌کنیم. جواب دقیق معادله دیفرانسیل در  $x = x_n$  را با  $y(x_n)$  و جواب تقریبی حاصل از کاربرد فرمول (۲۳.۸) را با  $y_n$  مشخص می‌سازیم. می‌خواهیم مقدار خطای گسسته‌سازی  $e_n$  را که به صورت

$$e_n = y(x_n) - y_n \quad (24.8)$$

تعریف می‌شود برآورد کنیم. یادآوری می‌کنیم که اگر  $y_0$ ، همان گونه که فرض خواهیم کرد، دقیق باشد آنگاه مقدار  $e_0$  مساوی ۰ خواهد شد. با فرض وجود داشتن مشتقاتی مناسب، می‌توانیم با استفاده از «قضیه تیلر با باقیمانده»،  $y(x_{n+1})$  را پیرامون  $x = x_n$  بسط دهیم

$$y(x_{n+1}) = y(x_n) + hy'(x_n) + \frac{h^2}{2} y''(\xi_n) \quad x_n \leq \xi_n \leq x_{n+1} \quad (25.8)$$

کمیت  $(h^2/2)y''(\xi_n)$  را خطای گسسته‌سازی موضعی<sup>۱</sup> نامند، یعنی خطایی که تنها در یک مرحله عبور از  $x_n$  به  $x_{n+1}$ ، با فرض معلوم بودن  $y$  و  $y'$  در  $x = x_n$  به طور دقیق، حاصل می‌شود. در یک کامپیوتر نیز هنگام محاسبه  $y_{n+1}$  با استفاده از فرمول (۲۳.۸)، خطای ناشی از گرد کردن وجود دارد. اما این نوع خطاها در این بخش نادیده گرفته خواهند شد.



از کسر کردن (۲۳.۸) از (۲۵.۸) و استفاده از (۲۴.۸) خواهیم داشت

$$e_{n+1} = e_n + h[f(x_n, y(x_n)) - f(x_n, y_n)] + \frac{h^2}{2} y''(\xi_n) \quad (26.8)$$

بنابر قضیهٔ مقدار میانگین حساب دیفرانسیل، داریم

$$\begin{aligned} f(x_n, y(x_n)) - f(x_n, y_n) &= f_y(x_n, \bar{y}_n)(y(x_n) - y_n) \\ &= f_y(x_n, \bar{y}_n)e_n \end{aligned}$$

که در آن  $\bar{y}_n$  بین  $y_n$  و  $y(x_n)$  واقع است. بنابراین رابطهٔ (۲۶.۸) به صورت زیر در خواهد آمد

$$e_{n+1} = e_n + hf_y(x_n, \bar{y}_n)e_n + \frac{h^2}{2} y''(\xi_n) \quad (27.8)$$

اکنون فرض می‌کنیم که در بازهٔ مورد نظر داریم

$$|f_y(x, y)| < L \quad |y''(x)| < Y$$

در اینجا  $L$  و  $Y$  اعداد ثابت و مثبتی هستند. با قدر مطلق گیری از رابطهٔ (۲۷.۸) خواهیم داشت

$$|e_{n+1}| \leq |e_n| + hL|e_n| + \frac{h^2}{2} Y = (1 + hL)|e_n| + \frac{h^2}{2} Y \quad (28.8)$$

اکنون با روش استقرا نشان خواهیم داد که جواب معادلهٔ تفاضلی

$$\xi_{n+1} = (1 + hL)\xi_n + \frac{h^2}{2} Y \quad (29.8)$$

به ازای  $\xi_0 = 0$  نافذ بر جواب (۲۷.۸) است، یعنی نشان خواهیم داد که

$$\xi_n \geq |e_n| \quad n = 0, 1, \dots \quad (30.8)$$

از آنجایی که  $e_0 = \xi_0 = 0$ ، رابطهٔ (۳۰.۸) مطمئناً به ازای  $n = 0$  برقرار است. به فرض اینکه رابطهٔ (۳۰.۸) به ازای عدد صحیح  $n$  برقرار باشد، چون  $\xi_n \geq |e_n|$  و  $(1 + hL) > 1$ ، آنگاه از رابطهٔ (۲۹.۸) نتیجه می‌شود که

$$\xi_{n+1} \geq |e_{n+1}| \quad n = 0, 1, \dots$$

که استقرای ما را کامل می‌کند.

بنابراین جواب  $\xi_n$  از معادلهٔ تفاضلی غیرهمگن (۲۹.۸)، کران بالایی برای خطای

گسسته‌سازی  $e_n$  به دست می‌دهد. از نظریهٔ مربوط به معادلهٔ تفاضلی کسه در بخش ۲۰۸ بیان شد، جواب معادلهٔ (۲۹.۸) به شرح زیر به دست می‌آید

$$\xi_n = c(1 + hL)^n - B \quad (31.8)$$

که در آن  $c$  يك ثابت اختیاری است و

$$B = \frac{hY}{\gamma L}$$

برای برقراری شرط  $\xi_n = 0$ ، ملاحظه می‌کنیم که باید  $c = +B$  انتخاب شود، که با این انتخاب رابطهٔ (۳۱.۸)، به صورت زیر در خواهد آمد

$$\xi_n = B(1 + hL)^n - B$$

با توجه به بخش ۷۰۱ نتیجه می‌گیریم که  $e^x = 1 + x + e^\xi x^2 / 2$ ؛ از این رو به ازای کلیهٔ مقادیر  $x$  رابطهٔ  $e^x \geq 1 + x$  برقرار است، که از این موضوع نتیجه می‌شود که  $1 + hL \leq e^{hL}$  و بنابراین  $(1 + hL)^n \leq e^{nhL}$ . اگر از این مطلب در رابطهٔ (۳۱.۸) استفاده کنیم می‌توانیم بگوییم که

$$\begin{aligned} \xi_n &\leq B(e^{nhL} - 1) \\ &= \frac{hY}{\gamma L} (e^{nhL} - 1) \\ &= \frac{hY}{\gamma L} (e^{(x_n - x_0)L} - 1) \end{aligned}$$

که در آنها از تساوی  $nh = x_n - x_0$  استفاده کرده‌ایم. چون  $|\xi_n| \leq e_n$ ، لذا قضیهٔ زیر را اثبات نموده‌ایم.

**قضیه ۲۰۸** گیریم  $y_n$  جواب تقریبی معادلهٔ (۱۹.۸) باشد که به وسیلهٔ روش اویلر (۲۳.۸) به دست آمده است. اگر جواب دقیق  $y(x)$  از معادلهٔ (۱۹.۸) دارای مشتق دوم پیوسته در بازهٔ  $[x_0, b]$  باشد، و اگر در این بازه نامساویهای

$$|f_y(x, y)| \leq L \quad |y''(x)| < Y$$

به ازای مقادیر ثابت مثبت  $L$  و  $Y$  برقرار باشند، خطای  $e_n = y(x_n) - y_n$  ناشی از روش اویلر در نقطهٔ  $x_n = x_0 + nh$  به صورت زیر کراندار می‌شود

$$|e_n| \leq \frac{hY}{\gamma L} (e^{(x_n - x_0)L} - 1) \quad (32.8)$$

این قضیه نشان می‌دهد که خطا برابر با  $\theta(h)$  است؛ یعنی به ازای  $h \rightarrow 0$ ، خطا به سمت صفر میل می‌کند، همانند  $ch$  که به ازای مقدار ثابتی مانند  $c$  وقتی  $x = x_n$  ثابت نگه داشته شود به سمت صفر میل می‌کند. می‌باید تأکید شود که برآورد (۳۲.۸)، به جای کران واقعی یک کران بالایی به دست می‌دهد. اهمیت اساسی این رابطه بیشتر در اثبات همگرایی این روش است، نه در فراهم کردن یک برآورد واقعی از خطایی که از قبل تعیین شده است.

□ مثال ۳۰.۸: کران بالایی را برای خطای گسسته‌سازی روش اویلر در حل معادله  $y' = y$ ،  $y(0) = 1$  از  $x = 0$  تا  $x = 1$  تعیین کنید.

حل: در اینجا داریم  $f(x, y) = y$  و  $\partial f / \partial y = 1$ . بنا بر این می‌توانیم  $L$  را برابر با یک بگیریم. همچنین چون  $y = e^x$ ، بنا بر این  $y'' = e^x$ ، و به ازای  $0 \leq x \leq 1$  داریم،  $|y''(x)| \leq e$ . برای پیدا کردن کرانی برای خطا در  $x = 1$  داریم  $x_n - x_0 = 1$  و  $y = e^1$  و از رابطه (۳۲.۸) به دست می‌آوریم

$$|e(1)| \leq \frac{he}{4}(e-1)$$

$$< 2.4h$$

بنا بر این خطای  $e(1)$  در  $x = 1$  با  $2.4h$  کراندار شده است. برای اینکه ببینیم این کران چقدر به مقدار واقعی نزدیک است، از روش اویلر جواب دقیقی برای این مسئله به دست می‌آوریم. لذا

$$\begin{aligned} y_{n+1} &= y_n + hf(x_n, y_n) \\ &= (1+h)y_n \end{aligned}$$

جواب این معادلهٔ تفاضلی که  $y(0) = 1$  در آن صدق می‌کند به صورت زیر است

$$y_n = (1+h)^n$$

حال اگر  $h = 0.1$  و  $n = 10$  اختیار شوند با بسط  $(1.1)^{10}$  درمی‌یابیم که روش اویلر مقدار  $y(1) = 2.5937$  را به دست می‌دهد. از کم کردن این مقدار، از جواب دقیق  $y(1) = e = 2.71828$  درمی‌یابیم که در مقایسه با کران  $2.4$  که با استفاده از رابطه (۳۲.۸) به دست آمده، مقدار خطا برابر با  $0.1246$  است. □

تمرین

$$۱-۴۰۸ \text{ در معادله } y' = -2y, y(0) = 1 \text{ و } 0 \leq x \leq 1$$

(الف) کران بالایی برای خطا بر حسب نمو  $h$  با استفاده از رابطه (۳۲.۸) در  $x = 1$  پیدا کنید.

(ب) معادله تفاضلی را که از روش اویلر نتیجه می شود حل کنید.

(پ) کران به دست آمده از (الف) را با خطای حقیقی به دست آمده از (ب) به ازای  $h = 0.01$  و  $h = 0.05$ ، در  $x = 1$  مقایسه کنید.

(ت) با به کار گرفتن روش اویلر، نمو  $h$  به چه اندازه می باید کوچک باشد تا شش رقم دقت در  $x = 1$  به دست آید (فرض کنید که خطای گرد کردن وجود ندارد).

۴-۴.۸ خطای  $e_n$  از یک روش انتگرالگیری در تانماسوی تفاضلی

$$|e_{n+2}| \leq a_1 |e_{n+1}| + a_2 |e_n| + A$$

که در آن  $a_1$ ،  $a_2$  و  $A$  ثابتهای مثبت هستند و  $e_1 = e_0 = 0$ ، صدق می کند. گیریم  $\xi_n$  یک جواب معادله تفاضلی

$$\xi_{n+2} = a_1 \xi_{n+1} + a_2 \xi_n + A$$

باشد و  $\xi_1 = \xi_0 = 0$ . به استقرا نشان دهید که

$$|e_n| \leq \xi_n \quad \text{به ازای جميع مقادیر } n$$

## ۵.۸ روشهای رونگه-کوتا

چنانکه قبلاً ذکر شد، روش اویلر در مسائل عملی چندان مفید نیست، زیرا برای دستیابی به یک دقت مناسب به نمو بسیار کوچکی نیاز دارد. الگوریتم مرتبه بالاتر تیلر نیز به عنوان شیوه ای با هدف کلی غیر قابل قبول است، زیرا به مشتقهای کلی از درجه بالاتر  $y(x)$  احتیاج دارد. در روش رونگه-کوتا سعی بر این است که به دقت بیشتری دست یافته شود و همزمان با آن از محاسبه مشتقهای مرتبه بالاتر اجتناب شود، و این کار با محاسبه تابع  $f(x, y)$  در نقاط انتخاب شده در هر زیربازه انجام می گیرد. در اینجا ساده ترین روش رونگه-کوتا را مورد بحث قرار می دهیم. ما در اینجا به دنبال فرمولی به شکل زیر هستیم

$$y_{n+1} = y_n + ak_1 + bk_2 \quad (33.8)$$

که در آن

$$k_1 = hf(x_n, y_n)$$

$$k_2 = hf(x_n + ah, y_n + \beta k_1)$$

و  $a, b, \alpha, \beta$  ثابت‌هایی هستند که باید چنان تعیین شوند که فرمول (۳۳.۸) با الگوریتم تیلر از مرتبه هر چه بزرگتر توافق داشته باشد. با بسط  $y(x_{n+1})$  به سری تیلر تا جمله  $h^3$  خواهیم داشت

$$\begin{aligned} y(x_{n+1}) &= y(x_n) + h y'(x_n) + \frac{h^2}{2} y''(x_n) + \frac{h^3}{6} y'''(x_n) + \dots \\ &= y(x_n) + h f(x_n, y_n) + \frac{h^2}{2} (f_x + f f_y)_n \\ &\quad + \frac{h^3}{6} (f_{xx} + 2f f_{xy} + f_{yy} f_x^2 + f_x f_y + f_y^2 f)_n + O(h^4) \end{aligned} \quad (34.1)$$

که در آن از بسط (۲۱.۸) استفاده کرده‌ایم، ولی زیر نمایه  $n$  به معنای آن است که کلیه توابع مورد نظر باید در  $\{x_n, y_n\}$  ارزیابی شوند. از طرف دیگر با استفاده از بسط تیلر برای توابع دو متغیره (به بخش ۷.۱ نگاه کنید)، خواهیم داشت

$$\begin{aligned} \frac{k_\gamma}{h} &= f(x_n + \alpha h, y_n + \beta k_\gamma) = f(x_n, y_n) + \alpha h f_x + \beta k_\gamma f_y \\ &\quad + \frac{\alpha^2 h^2}{2} f_{xx} + \alpha h \beta k_\gamma f_{xy} + \frac{\beta^2 k_\gamma^2}{2} f_{yy} + O(h^3) \end{aligned}$$

که در آن کلیه مشتقها در  $\{x_n, y_n\}$  ارزیابی شده‌اند. اکنون اگر این عبارت را در فرمول (۳۳.۸) به جای  $k_\gamma$  بگذاریم و توجه داشته باشیم که  $k_\gamma = h f(x_n, y_n)$ ، آنگاه با مرتب کردن نتیجه حاصله بر حسب توانهای  $h$  خواهیم داشت

$$\begin{aligned} y_{n+1} &= y_n + (a+b)hf + bh^2(\alpha f_x + \beta f f_y) \\ &\quad + bh^3\left(\frac{\alpha^2}{2} f_{xx} + \alpha\beta f f_{xy} + \frac{\beta^2}{2} f_y^2 f_{yy}\right) + O(h^4) \quad (\text{الف } 34.8) \end{aligned}$$

از مقایسه این رابطه با (۳۴.۸)، ملاحظه می‌کنیم که برای اینکه توانهای  $h$  و  $h^2$  در دو رابطه با هم سازگار باشند، باید داشته باشیم

$$a + b = 1$$

$$b\alpha = b\beta = \frac{1}{2} \quad (35.8)$$

اگرچه چهارمجهول داریم، ولی تنها سه معادله در دست داریم، بنابراین هنوز يك درجه آزادی در جواب (۳۵.۸) داریم. می‌توانیم این درجه آزادی را برای به‌دست آوردن هماهنگی در ضرایب جملات  $h^3$  به‌کار ببریم. ولی آشکار است که این امر برای کلیه توابع  $f(x, y)$  شدنی نیست.

برای (۳۵.۸) جوابهای زیادی وجود دارند، که شاید ساده‌ترین آنها به‌صورت زیر

باشد

$$a = b = \frac{1}{4} \quad \alpha = \beta = 1$$

الگوریتم ۲.۸: روش مرتبه دوم رونکه-کوتا برای معادله

$$y' = f(x, y) \quad y(x_0) = y_0$$

تقریبه‌های  $y_n$  را برای  $y(x_0 + nh)$  به‌ازای  $h$  ثابت و  $n = 0, 1, \dots$  با استفاده از فرمولهای بازگشتی زیر پدید آورید

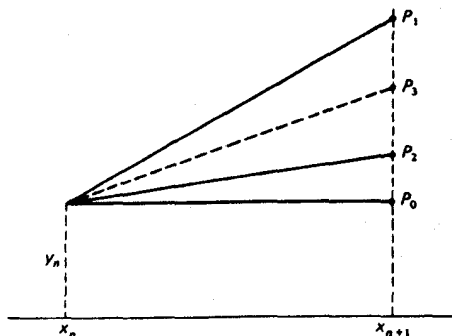
$$\boxed{y_{n+1} = y_n + \frac{1}{4}(k_1 + k_2) \quad k_1 = hf(x_n, y_n) \quad \text{به‌ازای} \quad (۳۶.۸)} \\ k_2 = hf(x_n + h, y_n + k_1)$$

الگوریتم ۲.۸ را می‌توان از لحاظ هندسی به‌صورت شکل ۱.۸ تجسم نمود. روش

اولی برای  $y_n$ ، نمو  $P_1 P_0 = hf(x_n, y_n)$  را تولید می‌کند و

$$P_2 P_0 = hf(x_n + h, y_n + hf(x_n, y_n))$$

نمود دیگری بر پایه‌ی شیب به‌دست آمده در  $x_{n+1}$  است. با به‌دست آوردن مقدار متوسط این دو نمو،



شکل ۱.۸

فرمول (۳۶.۸) به دست خواهد آمد.

خطای موضعی فرمول (۳۶.۸) به شکل زیر است

$$y(x_{n+1}) - y_{n+1} = \frac{h^4}{12} (f_{xx} + 2ff_{xy} + f^2f_{yy} - 2ff_{xy} - 2ff_{yy}^2) + O(h^4)$$

پیچیدگی ضرایب در این عبارت خطا، مشخصه کلیه روشهای رونگه-کوتا است، این امر یکی از جنبه‌های نامطلوب این روشها به شمار می‌آید، زیرا بر آورد خطای موضعی بسیار دشوار است. خطای موضعی فرمول (۳۶.۸) از مرتبه  $h^3$  است، در حالی که این خطا در روش اویلر از مرتبه  $h^2$  است. بنابراین می‌توانیم انتظار داشته باشیم که بتوانیم نمود بزرگتری را با (۳۶.۸) به کار ببریم. زیبایی که ما در این مسئله متحمل می‌شویم این است که تابع  $f(x, y)$  را در هر مرحله انتگرالگیری دوبار محاسبه می‌کنیم. با روش فوق می‌توان فرمولهایی از نوع فرمول رونگه-کوتا از هر مرتبه‌ای را به دست آورد. ولسی، مشتقات حاصله در آنها بسیار پیچیده می‌شود. عادیترین و متداولترین فرمول از این نوع در قالب الگوریتم ۳۰۸، آمده است.

**الگوریتم ۳۰۸:** روش مرتبه ۴ رونگه-کوتا برای معادله  $y' = f(x, y)$ ،  $y(x_0) = y_0$ ، تقریبهایی  $y_n$  برای  $y(x_0 + nh)$  را به ازای مقدار ثابت  $h$ ،  $n = 0, 1, 2, \dots$ ، با استفاده از فرمول بازگشتی زیر تولید کنید

$$y_{n+1} = y_n + \frac{1}{6} (k_1 + 2k_2 + 2k_3 + k_4) \quad (37.8)$$

که در آن

$$k_1 = hf(x_n, y_n)$$

$$k_2 = hf\left(x_n + \frac{h}{4}, y_n + \frac{1}{4}k_1\right)$$

$$k_3 = hf\left(x_n + \frac{h}{4}, y_n + \frac{1}{4}k_2\right)$$

$$k_4 = hf(x_n + h, y_n + k_3)$$

خطای گسسته‌سازی موضعی در الگوریتم ۳۰۸ برابر با  $O(h^5)$  است. و باز ضروری که ما برای این امر متحمل می‌شویم آن است که در هر مرحله به چهار بار ارزیابی تابع نیاز پیدا می‌کنیم. در مسائلی که تابع  $f(x, y)$  پیچیده است، این ضرر ممکن است از نظر وقت کامپیوتر قابل ملاحظه باشد. روشهای رونگه-کوتا معایب دیگری نیز دارند که بعداً مورد بحث قرار خواهند گرفت. در عمل، فرمول (۳۷.۸) به نحوی گسترده با موفقیتی بسیار زیاد

## حل معادلات دیفرانسیل ۴۷۵

به کار برده می‌شود. امتیاز مهم آن خود-آغازگری آن است: یعنی برای یافتن  $y$  و  $y'$  در  $x = x_0 + 1$ ، تنها به مقدار  $y$  در نقطه  $x = x_0$  نیاز دارد.

در زیر یک برنامه فورترن باهدف کلی بر اساس الگوریتم ۲.۸ برای معادله دیفرانسیل داده شده است. برای استفاده از این برنامه، باید زیر برنامه‌ای برای محاسبه تابع  $f(x, y)$  به آن افزوده شود و مقدار اولیه  $y_0 = y(x_0)$ ، نقطه انتهایی  $x_{NSTEP}$  و تعداد کل مراحل NSTEPS مشخص شوند.

## برنامه فورترن برای الگوریتم ۲.۸

```

C FORTRAN PROGRAM TO SOLVE THE FIRST ORDER DIFFERENTIAL EQUATION
C   Y'(X) = F(X,Y)
C WITH INITIAL CONDITION OF
C   Y(XBEGIN) = YBEGIN
C TO THE POINT XEND , USING THE SECOND ORDER RUNGE-KUTTA METHOD.
C A FUNCTION SUBPROGRAM CALLED 'F' MUST BE SUPPLIED.
      INTEGER I,N,NSTEPS
      REAL DERIV,H,K1,K2,XBEGIN,XN,XEND,YBEGIN,YN
      1 READ 501, XBEGIN,YBEGIN,XEND,NSTEPS
501 FORMAT(3F10.5,I3)
      IF (NSTEPS .LT. 1)                STOP
      H = (XEND - XBEGIN)/NSTEPS
      XN = XBEGIN
      YN = YBEGIN
      DERIV = F(XN,YN)
      N = 0
      PRINT 601, N,XN,YN,DERIV
601 FORMAT(1X,I3,3E21.9)
      DO 10 N=1,NSTEPS
          K1 = H*F(XN,YN)
          K2 = H*F(XN+H,YN+K1)
          YN = YN + .5*(K1+K2)
          XN = XBEGIN + N*H
          DERIV = F(XN,YN)
      10 PRINT 601, N,XN,YN,DERIV
                                         GO TO 1
      END
      REAL FUNCTION F(X,Y)
      REAL X,Y
      F = (1./X - Y)/X - Y*Y
                                         RETURN
      END

```

□ مثال ۴.۸: مسئله مثال ۲.۸ را با روش مرتبه دوم رونگه-کوتا (۳.۸) و با روش مرتبه چهارم رونگه-کوتا (۳.۸) حل کنید.

در نتایج کامپیوتری که در بخش ۳.۸ داده شد، فرمول (۳.۸) را روش ۲ و فرمول (۳.۸) را روش ۳ نامیده ایم. مشاهده می‌شود که روش مرتبه دوم رونگه-کوتا نتایجی به دست می‌دهد که با الگوریتم مرتبه دوم تیلر (روش ۱) کاملاً قابل مقایسه است. ولی روش مرتبه چهارم رونگه-کوتا نتایج خیلی بهبود یافته‌ای به دست می‌دهد، که به ازای  $h = 1/16$  تا شش رقم اعشاری و به ازای سایر مقادیر  $h$  تا هفت یا هشت رقم اعشاری دقت دارند.

## 1. self-starting



کارایی محاسباتی روشهای ۲ و ۳ را می توان با بررسی تعداد دفعاتی که تابع می باید برای هر يك محاسبه شود، مقایسه نمود. روش ۲، در هر مرحله نیاز به دو بار محاسبه تابع دارد و به ازای  $h = 1/128$ ، جمعاً ۲۵۶ بار محاسبه تابع لازم است. روش ۳، در هر مرحله نیاز به چهار بار محاسبه تابع دارد و به ازای  $h = 1/16$  جمعاً ۶۴ بار محاسبه تابع لازم است و نتایج دقیقتری به دست می دهد. آشکار است که روش مرتبه چهار رونگه - کوتاه، برای این مسئله روش کارآمدتری است و این مطلب از لحاظ کلی درست است.  $\square$

### تمرین

۱-۵۰۸ برای معادله  $y' = x + y$ ،  $y(0) = 1$ ، خطای موضعی روش (۳۶.۸) را محاسبه کنید. نتایج خود را با خطای الگوریتم مرتبه دوم تیلر مقایسه نمایید. انتظار دارید که کدامیک از این روشها درباره  $[0, 1]$  نتایج بهتری بدهد؟

۲-۵۰۸ با استفاده از فرمول (۳۶.۸) و نمو  $h = 0.01$ ، چند مرحله از انتگرالگیری  $y' = x + y$ ،  $y(0) = 1$  را انجام دهید و سپس برنامه ای بنویسید که این مسئله را از  $x = 1$  تا  $x = 0$  حل کند.

۳-۵۰۸ به معادله (۳۵.۸) شرطی اضافه کنید که ضرایب  $f_{xx}$  در (۳۴.۸) و (الف ۳۴.۸) با هم سازگار شوند. دستگاه معادلات حاصل را نسبت به  $a, b, \alpha, \beta$  حل کنید، عبارت خطا برای روش مرتبه دوم رونگه - کوتاه که با این مقادیر  $a, b, \alpha, \beta$  به دست می آید تعیین کنید.

۴-۵۰۸ می توان نشان داد که خطای روش مرتبه چهارم رونگه - کوتاه به ازای نمو  $h$ ، وقتی  $h$  به سمت صفر میل می کند، رابطه ای به شکل

$$y_n(h) - y(b) = A(b)h^4 + O(h^5)$$

پیدا می کند، که در آن  $b = x_0 + nh$ ، و در نتیجه  $n(h) = (b - x_0)/h$  و ثابت  $A(b)$  به  $h$  بستگی ندارد. مانند انتگرالگیری رامبرگ، با استفاده از روش برونیایی، تقریبی برای  $y(b)$  به دست آورید که خطای آن برابر با  $O(h^5)$  باشد.

### ۶.۸ کنترل اندازه $h$ در روشهای رونگه - کوتاه

در بخش ۵.۸، دور روش رونگه - کوتاه (رنگ) یکی از مرتبه ۲ و دیگری از مرتبه ۴ مورد بررسی قرار گرفت. اگر چه به دست آوردن فرمولهای روشهای رونگه - کوتاه بسیار پیچیده می شود، می توان فرمولهای مذکور را از هر مرتبه ای به دست آورد. نکته مهم در استفاده از روشهای رنگ مرحله ای نوع رونگه - کوتاه آن است که می توان خطای موضعی را بر آورد و در نتیجه مرحله مناسبی را برای حصول دقت لازم انتخاب کرد. هیچ دلیلی وجود ندارد که

ازماً اندازهٔ نمو  $h$ ، به گونه‌ای که در مثال ۴.۸ عمل کردیم، در تمام فاصلهٔ مورد بحث ثابت بماند. به کار بردن نمو‌های ثابت متفاوت برای برآورد دقت، به گونه‌ای که در مثال ۴.۸ انجام دادیم، ممکن است خیلی کارا نباشد. در این بخش روشهایی را برای برآورد خطای موضعی و برای تغییر نمو بر حسب معیار خطایی مورد بررسی قرار می‌دهیم.

روش اول بر پایهٔ نصف کردن بازه استوار است. فرض می‌کنیم که روش (رک) از مرتبهٔ  $p$  را به کار برده‌ایم و به ازای  $h = x_n - x_{n+1}$  به نقطهٔ  $x_n$  رسیده‌ایم. اکنون عمل انتگرالگیری از  $x_n$  تا  $x_{n+1} = x_n + h$  را دوبار، یک بار به ازای نمو فعلی  $h$  و بار دیگر با استفاده از نمو به طول  $h/2$  انجام می‌دهیم. بنابراین در  $x = x_{n+1}$  برای مقدار  $y(x)$  به دو برآورد  $y_h(x_{n+1})$  و  $y_{h/2}(x_{n+1})$  دست می‌یابیم، و مقایسهٔ این دو، برآوردی از خطا را به دست خواهد داد. برای تعیین این برآورد نخست توجه می‌کنیم که روش رونگه-کوئا از مرتبهٔ  $p$ ، بسط مجانبی خطای موضعی به شکل زیر را دارد

$$y_h(x_n + mh) = y(x_n + mh) + C(x_n + mh)h^p + \theta(h^{p+1}) \quad (38.8)$$

در اینجا،  $y_h(x_n + mh)$  معرف تقریبی است برای جواب  $y(x)$  در نقطهٔ  $x = x_n + mh$  که از روش رونگه-کوئا با نمو  $h$  و بسا شروع از مقدار دقیق  $y_n = y(x_n)$  و بعد از  $m$  نمو به دست آمده است. بعلاوه، اگرچه ثابت  $C(x_n + mh)$  به نقطهٔ  $x = x_n + mh$  و به  $f(x, y)$  بستگی دارد ولی مستقل از  $h$  است. بنا بر این

$$y_h(x_{n+1}) = y(x_{n+1}) + C(x_{n+1})h^p + \theta(h^{p+1}) \quad (\text{الف } 39.8)$$

$$y_{h/2}(x_{n+1}) = y(x_{n+1}) + C(x_{n+1})(h/2)^p + \theta(h^{p+1}) \quad (\text{ب } 39.8)$$

با کم کردن رابطهٔ (الف ۳۹.۸) از (ب ۳۹.۸)، ملاحظه می‌کنیم که قسمت اصلی خطا در (ب ۳۹.۸) می‌تواند به صورت زیر برآورد شود

$$C_n \left(\frac{h}{2}\right)^p \approx \frac{y_{h/2}(x_{n+1}) - y_h(x_{n+1})}{1 - 2^p}$$

کمیت

$$D_n = \frac{|y_{h/2}(x_{n+1}) - y_h(x_{n+1})|}{2^p - 1} \quad (40.8)$$

برآورد قابل محاسبه‌ای از خطا را در تقریب  $y_{h/2}(x_{n+1})$  برای ما فراهم می‌سازد و می‌تواند برای تصمیم‌گیری در این باب که آیا نمو به کار رفتهٔ  $h$  دقیقاً مناسب است یا آنکه خیلی بزرگ یا خیلی کوچک است، به کار رود.

اکنون فرض می‌کنیم که یک تحمل خطای موضعی  $\epsilon$  داده شده است و می‌خواهیم خطای برآورد شدهٔ  $D_n$  را کمتر از تحمل خطای موضعی در هر واحد مرحله نگاه داریم،

یعنی می‌خواهیم داشته باشیم

$$D_n \leq \varepsilon h \quad (۴۱۰۸)$$

فرض می‌کنیم که  $y_h(x_{n+1})$ ،  $y_{h/2}(x_{n+1})$  و  $D_n$  را محاسبه کرده‌ایم. اکنون می‌باید تصمیم بگیریم که آیا مقدار  $y_{h/2}(x_{n+1})$  مورد قبول است یا نیست، و برای انتگرالگیری بعدی، نمو  $h$  چه مقدار باید باشد. به روشی که بعداً توضیح داده خواهد شد از روی تحمل خطای مفروض  $\varepsilon$ ، کران پایینی برای خطای  $\varepsilon' < \varepsilon$  به دست می‌آوریم. حالات ممکنه زیر وجود دارند:

$$\varepsilon' < \frac{D_n}{h} < \varepsilon \quad (i)$$

در این حالت، مقدار  $y_{h/2}(x_{n+1})$  را می‌پذیریم، و عمل انتگرالگیری را با همان نمو  $h$ ، از  $x_{n+1}$  ادامه می‌دهیم.

$$\frac{D_n}{h} > \varepsilon \quad (ii)$$

در این حالت، مقدار خطا خیلی بزرگ است و از این رو باید  $h$  را -مثلاً به  $h/2$ - تبدیل و عمل انتگرالگیری را مجدداً از نقطه  $x = x_n$  شروع کنیم.

$$\frac{D_n}{h} < \varepsilon' \quad (iii)$$

در این حالت، دقت بیش از مقدار خواسته شده است. مقدار  $y_{h/2}(x_{n+1})$  را می‌پذیریم و  $2h$  را جایگزین  $h$  می‌کنیم و نقطه  $x_{n+1}$  انتگرالگیری را از سر می‌گیریم.

اگر اندازه نمو را در هر بازه به نصف شدن یا دو برابر شدن محدود کنیم، آنگاه برای یک روش از مرتبه  $p$ ، کران پایین  $\varepsilon'$  می‌تواند برابر با مقدار زیر گذاشته شود

$$\varepsilon' = \varepsilon / 2^{p+1}$$

و این امر بدان جهت است که نصف کردن  $h$ ، خطا را تقریباً به اندازه  $1/2^{p+1}$  تقلیل می‌دهد. برای روش مرتبه ۴ رونگه-کوتسا، داریم  $p=4$ ، از این رو  $\varepsilon' = \varepsilon / 32$ ، در واقع بیشتر اوقات توصیه نمی‌شود که نمو را تغییر دهند و برای اطمینان می‌توان  $\varepsilon' = \varepsilon / 50$  را به کار برده. شکل بسیار پیچیده دیگری از کنترل نمو که  $h$  را به دو برابر شدن یا نصف شدن محدود نمی‌کند به شکل زیر است. با توجه به رابطه (۴۰۰۸) داریم

$$D_n \approx 2C_n \left(\frac{h}{\gamma}\right)^{p+1} \quad (الف ۴۲۰۸)$$

هدف ما آن است که برای مرحله بعدی، نمودی برابر  $\bar{h}$  انتخاب کنیم. از آنجایی که قسمت اصلی خطا در مرحله بعدی برابر با  ${}_{2}C_n(\bar{h}/\psi)^{p+1}$  است،  $h$  باید چنان انتخاب شود که رابطه تحمل خطای (۴۱.۸) برقرار شود، از این رو باید داشته باشیم

$${}_{2}C_n\left(\frac{\bar{h}}{\psi}\right)^{p+1} \leq \varepsilon \bar{h} \quad (\text{ب } ۴۲.۲)$$

و باز با فرض آنکه  $C_n$  خیلی تغییر نمی کند، می توانیم  $C_n$  را از روابط (۴۲.۸ الف) و (۴۲.۸ ب) به گونه زیر حذف کنیم: با توجه به رابطه (۴۲.۸ ب) داریم

$${}_{2}C_n\left(\frac{h}{\psi}\right)^{p+1} \left(\frac{\bar{h}}{\psi}\right)^{p+1} \leq \varepsilon \bar{h} \left(\frac{h}{\psi}\right)^{p+1}$$

$$D_n \frac{(\bar{h})^p}{\psi^{p+1}} \leq \frac{\varepsilon h \cdot h^p}{\psi^{p+1}}$$

$$\bar{h}^p \leq h^p \varepsilon h / D_n$$

$$\bar{h} \approx h(\varepsilon h / D_n)^{1/p} \quad (\text{ب } ۴۲.۸)$$

بنابراین اگر قبلاً به ازای نمود  $h$ ، انتگرالگیری را با موفقیت انجام داده باشیم، نمود انتگرالگیری بعدی را باید برابر با  $\bar{h}$  و یا شاید برای اطمینان بیشتر، کمتر از آن بگیریم. برای مثال فرض کنید که روشی به کار برده ایم که با  $p=4$  داریم  $\varepsilon=10^{-6}$  و  $h=0.1$  و  $D_n$  طبق محاسبه باید برابر با  $10^{-5}$  باشد. پس

$$\bar{h} \approx 0.1(0.01)^{1/4} = 0.01(0.32) = 0.032$$

لذا این شرایط به مقدار خیلی کوچکتر  $h$  نیاز دارد. از طرف دیگر اگر باز  $p=4$  و  $h=0.1$  و  $\varepsilon=10^{-6}$  و  $D_n=10^{-8}$  را حساب کنیم، آنگاه داریم

$$\bar{h} \approx 0.1(10)^{1/4} \approx 0.1(1.78) = 0.178$$

لذا نمودی می تواند تقریباً دو برابر شود. استفاده از نمودهای متغیر به طور قابل ملاحظه ای به پیچیدگی یک برنامه می افزاید و در مجموعه ای از نقاط که به فواصل غیر یکنواخت قرار گرفته اند منجر به نتایجی می شود که برای استفاده کننده ناخوشایند است. نصف کردن و دو برابر کردن بازه ها عموماً برای استفاده کننده قابل قبولتر است. از سوی دیگر، در برنامه هایی که در آنها کنترل نمود به طور خودکار انجام می گیرد برآوردهای خیلی خوبی از دقت فراهم شده است و روی هم رفته کاملاً کارا هستند.

عیب اصلی این روش برای کنترل خطا، کار اضافه ای است که اساساً باید انجام گیرد. در سالیهای اخیر شکل های متفاوتی از روش رونگه-کوتا که برای کنترل نمود مناسب هستند پیشنهاد

شده‌اند که چندتا از این شکل‌های تازه متفاوت با نام‌های مرسن<sup>۱</sup>، ورنر<sup>۲</sup>، و فهلبرگ<sup>۳</sup> همراه هستند، روشی را که توسط فهلبرگ ارائه گردیده‌است به‌اجمال بررسی می‌کنیم و آن را با رکنف ۴۵ [۲۸] نشان می‌دهیم. در این روش، در هر مرحله شش بار محاسبه تابع موردنیاز است ولی برآورد خطای آن به‌صورت خودکار انجام می‌گیرد، و در عین حال دقت آن از روش مرتبه چهارم استانده بیشتر است. فهلبرگ نشان داد که ترکیب مقادیر چهارتا از این توابع با مجموعه‌ای از ضرایب می‌تواند روش مرتبه چهارم را به‌دست دهد، در حالی که ترکیب هر شش مقدار تابع با مجموعه دیگری از ضرایب می‌تواند برای به‌وجود آوردن روش مرتبه پنجم به‌کار رود. سپس مقایسه مقادیر به دست آمده از روش‌های مرتبه‌های چهارم و پنجم منجر به برآورد خطایی می‌شود که می‌تواند برای کنترل  $h$  به‌کار رود.

راهی را که فهلبرگ اختیار کرده خیلی به‌اختصار توضیح می‌دهیم. فرض می‌کنیم که از معادله  $y' = f(x, y)$  تا یک نقطه  $x_n$  به‌ازای نمو  $h_n$  انتگرال گرفته‌ایم و اکنون می‌خواهیم در  $x = x_{n+1}$  برآوردی برای  $y(x)$  به‌دست آوریم. برآورد اول به‌ازای ضرایبی مانند  $c_i$  با فرمول زیر معین می‌شود

$$y_{n+1} = y_n + \sum_{i=1}^6 c_i k_i \quad (۴۳.۸ \text{ الف})$$

و برآورد دوم به‌ازای مجموعه دیگری از ضرایب  $c_i^*$  به‌صورت زیر تعیین می‌شود

$$y_{n+1}^* = y_n + \sum_{i=1}^6 c_i^* k_i \quad (۴۳.۸ \text{ ب})$$

لذا برآورد خطا برای کنترل  $h$  با فرمول زیر محاسبه می‌شود:

$$D_n = \sum_{i=1}^6 (c_i - c_i^*) k_i$$

همچنان که قبلاً توضیح داده شد، این رابطه برای برآورد مناسبی از  $h$  در انتگرال‌گیری بعدی می‌تواند به‌کار رود. توابع  $k_i$  در هر دو فرمول یکی هستند و می‌توانند به‌صورت زیر بیان شوند

$$k_i = h_n f\left(x_n + \alpha_i h_n, y_n + \sum_{j=1}^{i-1} \beta_{ij} k_j\right) \quad i = 1, \dots, 6$$

در مورد ضرایب  $\alpha_i$  و  $\beta_{ij}$  انتخاب‌های ممکنه زیادی وجود دارند که به‌روش مرتبه ۵ رونگه-کوتا منجر می‌شوند. فهلبرگ، مجموعه خاصی از ضرایب را پیشنهاد نموده‌است که مادر اینجا آنها را عرضه نخواهیم کرد. دانشجویان علاقه‌مند، برای اطلاع بیشتر از جزئیات این روش می‌توانند به‌مرجع [۲۸] آخر این کتاب مراجعه کنند.

روش دیگر رونگه-کوتا با کنترل اندازه نمو که توسط ورنر ارائه گردیده، اساس بسیاری از زیر برنامه‌های موفق در حل معادلات دیفرانسیل است که به DVERK موسوم و در کتابهای زیر برنامه‌ای به تعداد زیاد در دسترس هستند. در روش ورنر که با  $rk=56$  نشان داده می‌شود، در هر مرحله به هشت بار محاسبه تابع احتیاج است و از این توابع برای  $y(x)$  دو برآورد به دست می‌آورند که مبنای یکی تقریب مرتبه پنجم است و مبنای دیگری تقریب مرتبه ششم. سپس مقایسه این دو برآورد، مبنایی برای انتخاب نمو فراهم می‌سازد. چند آزمون اولیه در مورد این روش در دانشگاه تورنتو<sup>۱</sup> [۲۹] به عمل آمده است. روش مذکور سپس به صورت زیر برنامه DVERK در آمده و در IMSL، هوستون<sup>۲</sup>، تکزاس<sup>۳</sup>، منتشر شده است. IMSL که معرف کتابهای بین‌المللی ریاضی و آمار<sup>۴</sup> است، مجموعه‌ای از زیر برنامه‌های کاملاً آزمایش شده برای بسیاری از مسائل متنوع آمار و ریاضی را در بردارد. دستیابی به این مجموعه کتاب بر اساس پرداخت حق اشتراک<sup>۵</sup> صورت می‌گیرد و تقریباً برای تمام کامپیوترهای متوسط و بزرگ از جمله IBM، CDC، UNIVAC، Burroughs و Honeywell موجود است. چون بسیاری از مراکز کامپیوتری حق اشتراک برای مجموعه IMSL پرداخت می‌کنند، دیگر ما در اینجا زیر برنامه DVERK را نخواهیم داد. چون در این فصل برای حل چندین مسئله از این زیر برنامه استفاده می‌کنیم، پارامترهای موجود در حکم فراخوانی<sup>۶</sup> و انتخابهای<sup>۷</sup> مختلف را به طور اجمال توضیح می‌دهیم.

در استفاده معمولی در صورت نبود انتخاب<sup>۸</sup> و بعد از تعیین مقادیر اولیه، قسمت اصلی برنامه مربوط به حل معادله دیفرانسیل مرتبه اول  $y' = f(x, y)$  از  $x = XBEGIN$  تا  $x = XM$  شامل يك حلقه DO به صورت زیر است:

```

X = XBEGIN
Y = YBEGIN
DO 10 K=1,M
  XEND = XBEGIN + FLOAT(K)*(XM - XBEGIN)/FLOAT(M)
  CALL DVERK ( N, FCN1, X, Y, XEND, TOL, IND, C, NW, W, IER )
  PRINT 600, XEND, Y(1), C(24)
600  FORMAT(F19.6,E21.8,F15.0)
10  CONTINUE

```

معانی پارامترهای موجود در این زیر برنامه به قرار زیرند:

$N \equiv$  تعداد معادلاتی که باید حل شوند (در اینجا  $N=1$ ).

$FCN1 \equiv$  نام زیر برنامه مربوط به  $f(x, y)$ ، که استفاده کننده باید به صورت يك

برنامه خارجی تهیه کند.

$X \equiv$  مقدار اولیه متغیر مستقل.

1. Toronto
2. Houston
3. Texas
4. International Mathematical and Statistical Library
5. subscription
6. statement
7. options
8. default options

$Y \equiv$  مقدار اولیه متغیر وابسته.

$XEND \equiv$  مقداری از  $x$  که به ازای آن جواب می باید بیرون داده شود.

$TOL \equiv$  میزان تحمل کنترل خطا؛ در عین اینکه خصوصیات تحمل خطا ممکن

است انواع مختلف داشته باشند در انتخاب قراردادی این تحمل سعی

بر این است که خطای کلی<sup>۱</sup> نسبی از تحمل برای کنترل خطا کمتر باشد.

$IND \equiv ۱$  موجب نبودن همه انتخابهایی است که باید استفاده شوند.

$\equiv ۲$  اجازه می دهد که انتخابها صورت گیرند.

$C \equiv$  بردارهای ارتباطی به طول ۲۴؛ اگر  $IND$  روی ۲ گذاشته شده باشد،

برخی از این بردارها را استفاده کننده می تواند تعیین کند؛ این انتخابها

انواع مختلف کنترل خطا، مینیم یا ماکسیمیم نموها، حدود تعداد

ارزیابیهای تابع و غیره را مدنظر قرار می دهند.

$NW \equiv$  اولین بعد<sup>۲</sup> (تعدادسطرهای) ماتریس فضای کار  $W$  باید حداقل به اندازه

$N$  باشد.

$W \equiv$  ماتریس فضای کار که اولین بعد (تعداد سطرهای) آن برابر با  $NW$

است و دومین بعد (تعداد ستونهای) آن باید نا کوچکتر از ۹ باشد.

$IER \equiv$  پرچم<sup>۳</sup>خطا برای نشان دادن انواع مختلف خطاهایی که با آنها مواجه

می شویم.

در حلقه<sup>۴</sup>  $DO$ ، مذکور در بالا، نقاط  $XEND$  مقادیری از  $x$  هستند که به ازای آنها

جواب بیرون داده می شود. در این حالت جواب در  $M$  نقطه<sup>۵</sup> مساوی الفاصله

$XBEGIN + k\Delta X$  بیرون داده خواهد شد که در آن  $\Delta X = (XM - XBEGIN) / M$ .

در خود زیر برنامه،  $DVERK$  نمره<sup>۶</sup>های مناسبی را برای دقت مورد نیاز انجام خواهد داد.

همچنان که مراحل انتگرالگیری پیش می رود، نمو به طور طبیعی تغییر می کند. زیر برنامه

مذکور تعداد محاسبات تابع را هم که برای یافتن جواب در نقطه<sup>۷</sup>  $XEND$  لازم است حفظ

خواهد کرد.  $DVERK$  زیر برنامه ای است با دقت بسیار بالا که در هر مرحله از انتگرالگیری،

حداقل به هشت بار محاسبه تابع احتیاج دارد. تعداد محاسبات تابع که عملاً به کار برده

می شود، در (۲۴)  $C$  ذخیره می گردد و همان گونه که در بالا انجام دادیم می تواند به عنوان

یک خروجی بیرون داده شود.

روش مزبور برای معادله<sup>۸</sup> دیفرانسیل

$$y'(x) = \frac{1}{x^2} - \frac{y(x)}{x} - y^2(x)$$

$$y(1) = -1$$

که در مثال ۲.۸ مورد بررسی قرار گرفت، به کار رفته و برنامه<sup>۹</sup> کامل و نتایج حاصل از آن در

زیر داده شده است

## حل معادلات دیفرانسیل ۴۸۳

```

C USE OF D V E R K TO SOLVE EXAMPLE 8.2 .
  INTEGER IER,IND,K,N,NW
  REAL C(24),TOL,W(1,9),X,XEND,Y(1)
  DATA N , X , Y(1), TOL ,IND,NW
  * / 1 , 1.,-1.,1,E-7, 1 , 1 /
  EXTERNAL FCN1
  DO 10 K=1,4
    XEND = 1. + FLOAT(K)/4.
    CALL DVERK ( N, FCN1, X, Y, XEND, TOL, IND, C, NW, W, IER )
    PRINT 600, XEND,Y(1),C(24)
600  FORMAT(11X,F8.6,5X,E16.8,5X,F11.0)
10  CONTINUE

                                STOP

END
SUBROUTINE FCN1 ( N, X, Y, YPRIME )
  REAL X,Y(1),YPRIME(1)
  YPRIME(1) = (1./X - Y(1))/X - Y(1)*Y(1)
                                RETURN
END

```

## بروندا

X	Y(1)	FCN EVALS
1.25	-0.799999999	16.
1.50	-0.666666664	24.
1.75	-0.57142854	32.
2.00	-0.499999996	40.

نتایج به دست آمده از لحاظ دقت با جوابهای حاصل از به کار بردن روش کلاسیک<sup>۱</sup> مرتبه چهارم با نمو ثابت  $h = 1/32$  قابل مقایسه هستند. از آنجایی که روشهای کلاسیک مرتبه چهارم در هر مرحله به چهار بار محاسبه تابع احتیاج دارند، کلاً ۱۲۸ بار محاسبه تابع لازم خواهد بود تا به دقتی در حدود هفت رقم اعشاری برسیم. در حالی که، برای رسیدن به همان دقت، در زیر برنامه DVERK، تنها به ۴۰ بار محاسبه تابع نیاز است. لازم به تذکر است که برای داشتن خروجی در  $x = 1.75$  به  $x = 1.6$  بار محاسبه تابع نیاز بود که این امر حاکی از بزرگی نمو  $h = 1/4$  بود، و ظاهراً برای رسیدن به دقت  $10^{-7}$  می بایستی این نمو نصف می شد.

در بخش ۱۲.۸ استفاده از زیر برنامه DVERK را برای حل یسک دستگاه معادلات دیفرانسیل مرتبه اول نشان خواهیم داد.

## تمرین

۱-۶.۸ فرض کنید برای حل یسک معادله دیفرانسیل، از روش مرتبه دوم رونگه-کوتا استفاده کرده ایم و کنترل اندازه نمو نیز بر اساس نصف کردن فاصله صورت گرفته است. اگر از نمو  $h = 0.1$  و معیار خطای  $\epsilon = 10^{-6}$  استفاده کنیم و در نقطه  $x = x_n$ ،  $D_n = 10^{-4}$  به دست آید، نمو  $h$  برای مرحله بعدی انتگرالگیری چه مقدار باید باشد؟

۲-۶.۸ برنامه ای برای روش مرتبه دوم رونگه-کوتا که در آن کنترل اندازه نمو به دو برابر-



شدن و نصف شدن محدود می‌شود، بنویسید. این برنامه را برای حل معادله مثال ۲.۸ به ازای  $\epsilon = 10^{-6}$  به کار برید.

۳-۶۰۸ ببینید که آیا مرکز محاسباتی شما مجموعه زیر برنامه‌های IMSL را انجام می‌دهند یا نه. این صورت زیر برنامه DVERK را برای حل معادلات دیفرانسیل زیر به کار برید. در هر حالت TOL را برابر با  $10^{-7}$  قرار دهید و خروجی را در مقادیر XEND به صورت زیر به دست آورید

$$XEND = XO + K(XM - XO) / 10 \quad K = 1, 2, \dots, 10$$

$$y' = x - 1 + y/x \quad (\text{الف})$$

$$XO = 1, XM = 2, y(XO) = 2$$

$$y' = xy^2 \quad (\text{ب})$$

$$XO = 1, XM = 4, y(XO) = 1$$

## ۲.۸ فرمولهای چندمرحله‌ای

الگوریتم مرتبه  $k$  تیلر و روشهای رونگه-کوتا هر دو مثالهایی از روشهای یک مرحله‌ای هستند. این روشها به اطلاعاتی درباره جواب در یک نقطه  $x = x_n$  نیاز دارند تا از روی این جواب، به تعیین  $y$  در نقطه بعدی  $x = x_{n+1}$  بپردازند. در روشهای چندمرحله‌ای، از اطلاعات پیرامون جواب در بیش از یک نقطه استفاده می‌کنند. فرض کنید که برای  $y'$  و  $y$  در چند نقطه با فواصل مساوی مثل  $x_0, x_1, \dots, x_n$  تقریبهایی از قبل به دست آورده باشیم. یک دسته از روشهای چندمرحله‌ای بر پایه اصل انتگرالگیری عددی استوار شده‌اند. اگر از معادله دیفرانسیل  $y' = f(x, y)$  از  $x_n$  تا  $x_{n+1}$  انتگرال بگیریم، خواهیم داشت

$$\int_{x_n}^{x_{n+1}} y' dx = \int_{x_n}^{x_{n+1}} f(x, y(x)) dx$$

یا

$$y_{n+1} = y_n + \int_{x_n}^{x_{n+1}} f(x, y(x)) dx \quad (44.8)$$

برای انجام انتگرالگیری در رابطه (۴۴.۸)،  $f(x, y(x))$  را با یک بسجمله‌ای، که  $f(x, y(x))$  را در  $(m+1)$  نقطه  $x_n, x_{n-1}, x_{n-2}, \dots, x_{n-m}$  در نیایی کند، تقریب می‌زنیم. برای این منظور اگر از قرارداد

$$f(x_k, y(x_k)) = f_k$$

استفاده کنیم، می‌توانیم فرمول پرسو نیوتن (تمرین ۶.۲-۸ را ببینید) از مرتبه  $m$  را به کار گیریم:

$$p_m(x) = \sum_{k=0}^m (-1)^k \binom{-s}{k} \Delta^k f_{n-k} \quad s = \frac{x - x_n}{h}$$

با قراردادن رابطه فوق در (۴۴.۸) و با توجه به تساوی  $dx = h ds$ ، داریم

$$\begin{aligned} y_{n+1} &= y_n + h \int_0^1 \sum_{k=0}^m (-1)^k \binom{-s}{k} \Delta^k f_{n-k} ds \\ &= y_n + h \{ \gamma_0 f_n + \gamma_1 \Delta f_{n-1} + \dots + \gamma_m \Delta^m f_{n-m} \} \end{aligned} \quad (45.8)$$

که در آن

$$\gamma_k = (-1)^k \int_0^1 \binom{-s}{k} ds \quad (الف 45.8)$$

با توجه به تعریف تابع دو جمله‌ای که در فصل ۲ آمده بود، می‌توانیم مقادیر  $\gamma_k$  را به آسانی محاسبه کنیم، که چندتای اول آن به قرار زیرند:

$$\gamma_0 = 1 \quad \gamma_1 = \frac{1}{2} \quad \gamma_2 = \frac{5}{12} \quad \gamma_3 = \frac{3}{8} \quad \gamma_4 = \frac{251}{720}$$

فرمول (۴۵.۸) به روش ادمز-بشفورث معروف است. ساده‌ترین حالت آن، با قراردادن  $m=0$  در فرمول (۴۵.۸)، باردیگر روش اوپلر را به دست می‌دهد. به طور کلی استفاده از فرمول (۴۵.۸) مستلزم داشتن مقدار  $f' = y'$  در  $m+1$  نقطه  $x_n, x_{n-1}, \dots, x_{n-m}$  است. از روی این مقادیر می‌توانیم تفاضلهای

$$\Delta f_{n-1}, \Delta^2 f_{n-2}, \dots, \Delta^m f_{n-m}$$

را تشکیل دهیم و از (۴۵.۸) می‌توانیم  $y_{n+1}$  را محاسبه کنیم؛ و از معادله دیفرانسیل می‌توانیم  $f_{n+1} = f(x_{n+1}, y_{n+1})$  را به دست آوریم. اکنون نقطه  $x_{n+1}$  را نقطه  $x_n$  می‌گیریم و خط جدیدی از تفاضلهای تشکیل می‌دهیم و این روند را تکرار می‌کنیم. به ازای  $m=3$ ، که معمولاً در عمل به کار برده می‌شود، جدول تفاضلهای به صورت زیر خواهد بود

$$\begin{array}{rcccc}
 x_{n-3} & y_{n-3} & f_{n-3} & & \\
 & & & \Delta f_{n-3} & \\
 x_{n-2} & y_{n-2} & f_{n-2} & & \Delta^2 f_{n-2} \\
 & & & \Delta f_{n-2} & \Delta^2 f_{n-3} \\
 x_{n-1} & y_{n-1} & f_{n-1} & & \Delta^2 f_{n-2} \\
 & & & \Delta f_{n-1} & \\
 x_n & y_n & f_n & & 
 \end{array}$$

و فرمول (۴۵.۸) به حالت خاص زیر بدل می شود

$$y_{n+1} = y_n + h \left( f_n + \frac{1}{2} \Delta f_{n-1} + \frac{5}{12} \Delta^2 f_{n-2} + \frac{3}{8} \Delta^3 f_{n-3} \right) \quad (46.8)$$

در عمل از لحاظ محاسباتی آسانتر است که به جای کار کردن با تفاضها با عرضها\* کار کنیم. از تعریف عملگر تفاضل پیشرو  $\Delta$ ، خواهیم داشت

$$\Delta f_{n-1} = f_n - f_{n-1}$$

$$\Delta^2 f_{n-2} = f_n - 2f_{n-1} + f_{n-2}$$

$$\Delta^3 f_{n-3} = f_n - 3f_{n-1} + 3f_{n-2} - f_{n-3}$$

از قراردادن این روابط در (۴۶.۸) و مرتب کردن نتیجه حاصله، خواهیم داشت

$$y_{n+1} = y_n + \frac{h}{24} (55f_n - 59f_{n-1} + 37f_{n-2} - 9f_{n-3}) \quad (47.8)$$

خطای موضعی رابطه (۴۶.۸) را می توان به صورت زیر به دست آورد: با توجه به تمرین ۶.۲-۸ می دانیم که خطا در فرمول پسرو نیوتن به ازای  $n=3$  و  $k=0$  به صورت زیر است

$$h^4 f^{(4)}(\eta) \binom{-5}{4}$$

\* عرضها منظور مقادیر تابع است و نه تفاضهای آنها...م.

بنابراین خطا در فرمول (۴۶.۸) به کمک رابطه

$$E_{AB} = h \int_0^1 h^4 f^{(4)}(\eta) \binom{-s}{4} ds$$

تعیین می‌شود. از آنجا که  $\binom{-s}{4} ds$  در بازه  $[0, 1]$  تغییر علامت نمی‌دهد، نقطه‌ای مانند  $\xi$  بین  $x_{n-3}$  و  $x_{n+1}$  وجود دارد به طوری که

$$E_{AB} = h^5 f^{(4)}(\xi) \int_0^1 \binom{-s}{4} ds$$

$$E_{AB} = h^5 y^{(4)}(\xi) \frac{251}{720} \quad (48.8)$$

برای استفاده از رابطه (۴۷.۸) باید ۴ مقدار آغازین را داشته باشیم. این مقادیر آغازین باید از یک منبع مستقل به دست آمده باشند. برای نشان دادن چگونگی استفاده از رابطه (۴۷.۸) چند مرحله اول انتگرالگیری معادله زیر را به ازای  $h = 0.1$  انجام می‌دهیم

$$y' = -y^2$$

$$y(1) = 1$$

جواب دقیق این مسئله برابر است با  $y = 1/x$ . در جدول زیر چهار مقدار آغازین اول از جواب دقیق به دست آمده‌اند و درایه‌های دیگر از رابطه (۴۷.۸).

$x_n$	$y_n$	$f_n = -y_n^2$	$y(x_n) = 1/x_n$
۱.۰	۱.۰۰۰۰۰۰۰۰	-۱.۰۰۰۰۰۰۰۰	
۱.۱	۰.۹۰۹۰۹۰۹۱	-۰.۸۲۶۴۴۶۲۸	۰.۹۰۹۰۹۰۹۱
۱.۲	۰.۸۳۳۳۳۳۳۳	-۰.۶۹۴۴۴۴۴۴	۰.۸۳۳۳۳۳۳۳
۱.۳	۰.۷۶۹۲۳۰۷۷	-۰.۵۸۹۱۷۱۵۹۸	۰.۷۶۹۲۳۰۷۷
۱.۴	۰.۷۱۴۲۳۶۳۲	-۰.۵۱۰۴۱۹۲۶	۰.۷۱۴۲۸۵۷۱
۱.۵	۰.۶۶۶۶۸۶۰۳۰	-۰.۴۴۲۲۷۰۲۶۶	۰.۶۶۶۶۶۶۶۶۷
۱.۶	۰.۶۲۵۲۴۶۱۳	-۰.۳۹۰۹۳۲۷۲	۰.۶۲۵۰۰۰۰۰

مشاهده می‌شود که مقادیر  $y_n$  که توسط فرمول (۴۷.۸) محاسبه شده‌اند در رقم چهارم خود تقریباً دو واحد خطا دارند. با استفاده از فرمول برآورد خطای موضعی (۴۳.۶) و این امر که داریم

$$|y^v(x)| = \left| \frac{5!}{x^6} \right| \leq 120 \quad 1 < x < 2$$

کران خطا را به دست می‌آوریم

$$|E_{AB}| \leq \frac{251}{720} (120)(10^{-5}) \approx 0.00004$$

این کران در حدود دو برابر خطایی است که از رفتن از یک مرحله به مرحلهٔ دیگر ایجاد می‌شود.

به همین طریق بسا استفاده از انتگرالگیری عددی، فرمولهای دیگری از نوع چندمرحله‌ای می‌توان به دست آورد. در رابطه (۴۳.۸) به جای انتگرالگیری از  $f(x, y)$  از  $x_n$  تا  $x_{n+1}$ ، می‌توانیم مثلاً از  $x_{n-p}$  تا  $x_{n+1}$  به ازای یک عدد صحیح  $p \geq 0$  انتگرالگیری کنیم. اگر با فرمول بسرو نیوتن در  $m+1$  نقطهٔ  $x_n, x_{n-1}, \dots, x_{n-m}$  مجدداً درونیایی کنیم، خواهیم داشت

$$y_{n+1} = y_{n-p} + h \int_{-p}^1 \sum_{k=0}^m (-1)^k \binom{-s}{k} \Delta^k f_{n-k} ds \quad (49.8)$$

در حالت  $p=0$ ، فرمول (۴۴.۸) ادمز-بشفورت به دست می‌آید. متناظر با  $m=1$ ،  $p=3$ ،  $m=3$  و فرمولهای خاص و جالبی از این نوع حاصل می‌شوند. این فرمولها همراه با عبارت خطای موضعی آنها به صورت زیرند

$$y_{n+1} = y_{n-1} + 2hf_n$$

$$E = \frac{h^3}{3} y'''(\xi) \quad (50.8)$$

$$y_{n+1} = y_{n-2} + \frac{2h}{3} (2f_n - f_{n-1} + 2f_{n-2})$$

$$E = \frac{14}{45} h^5 y^{(5)}(\xi)$$

$$(51.8)$$

فرمول (۵۰.۸)، که از لحاظ سادگی با روش اویلر قابل مقایسه است، خطای گسسته‌سازی مطلوبتری دارد. همچنین فرمول (۵۱.۸) که مستلزم داشتن  $f(x, y)$  تنها در سه نقطه است، خطای گسسته‌سازی قابل مقایسه‌ای با روش ادمز-بشفورت (۴۷.۸) دارد. می‌توان نشان داد که به ازای عدد فرد  $m$  و  $m=p$ ، هر فرمولی به صورت (۴۹.۸) دارای این ویژگی است

که ضریب تفاضل  $m$  در آن برابر با صفر می‌شود، و بنا بر این فرمولی حاصل می‌شود که دقت آن بیش از مقداری است که انتظار می‌رود. از سوی دیگر این فرمولها دستخوش ناپایداری بیشتری قرار می‌گیرند، که مفهوم ناپایداری بعداً مورد بحث واقع خواهد شد. يك اشکال اصلی فرمولهای چندمرحله‌ای این است که آنها «خود-آغازگر» نیستند. بنابراین در روش ادمز-بشفورث (۴۷.۸)، و قبل از به‌کار بستن این فرمول باید چهارمقدار متوالی  $f(x, y)$  را در نقاط با فواصل مساوی داشته باشیم. این مقادیر آغازین را باید به کمک روش مستقلی به دست آورد. برای مثال می‌توانیم الگوریتم تیلر یا یکی از روشهای رونگه-کوتا را برای به دست آوردن این مقادیر آغازین به کار ببریم. و نیز برای به دست آوردن دقت مطلوب باید اطمینان حاصل کنیم که این مقادیر آغازین خود به قدر لازم دقیق هستند. اشکال دوم روش ادمز-بشفورث این است که با اینکه خطای موضعی مربوط به گسسته‌سازی برابر با  $O(h^5)$  است، ضریب جمله خطا اندکی بیشتر از همان ضریب برای فرمولهای رونگه-کوتا از مرتبه مشابه است. به همین دلیل روش رونگه-کوتا معمولاً، ولی نه همیشه، از دقت بیشتری برخوردار است. از سوی دیگر فرمولهای چندمرحله‌ای در مقایسه با چهار بار محاسبه کردن در هر مرحله با روش رونگه-کوتا، در هر مرحله تنها به یک بار محاسبه مشتق احتیاج دارند، و بنا بر این این روش به نحو قابل ملاحظه‌ای سریعتر است، و، کارهای محاسباتی کمتری نیاز دارد.

□ مثال ۵.۸: معادله

$$y' = x + y \quad y(0) = 0$$

را از  $x=0$  تا  $x=1$  با استفاده از روش ادمز-بشفورث حل کنید.

**برنامه فورقون** و نتایج حاصله از آن برای این مسئله در زیر داده شده است. جواب دقیق این مسئله چنین است  $y = e^x - x - 1$ . با استفاده از این جواب، چهارمقدار آغازین محاسبه شده‌اند. اولین ستون نتایج مقادیر  $x_n$  را به ازای  $h = 1/32$  می‌دهد، ستون دوم مقادیر  $y_n$  را که به کمک فرمول (۴۷.۸) محاسبه شده‌اند معین می‌کند، و ستون سوم مقدار  $y(x_n)$  را که از روی جواب حساب شده نشان می‌دهد و ستون چهارم مقدار خطای  $e_n = y_n - y(x_n)$  را.

نتایج تقریباً تا شش رقم با معنی دقیق‌اند، که همان چیزی است که از فرمول خطای (۴۸.۸) تقریباً انتظار می‌رفت. از آنجایی که خطای انباشته شده مربوط به گسسته‌سازی، برابر با  $O(h^4)$  است، انتظار می‌رود که اگر اندازه  $h$  نصف شود، این خطا به اندازه  $1/16$  تقلیل یابد.

## برنامه فورترون برای مثال ۵-۸

```

C ADAMS-BASHFORTH METHOD
  INTEGER I,N,NSTEPS
  REAL ERROR,F(4),H,X,BEGIN,XN,YBEGIN,YN
C
  SOLN(X) = EXP(X) - 1. - X
C
C ** INITIALIZE
C PRINT 600
600 FORMAT('ADAMS-BASHFORTH METHOD'/
  * '0',4X,'N',13X,'XN',15X,'YN',13X,'Y(XN)',12X,'ERROR'/)
  NSTEPS = 32
  H = 1./NSTEPS
  YBEGIN = 0.
  XBEGIN = 0.
C
C ** COMPUTE FIRST FOUR POINTS USING EXACT SOLUTION
C F(1) = XBEGIN + YBEGIN
  N = 0
  ERROR = 0.
  PRINT 601, N, XBEGIN, YBEGIN, YBEGIN, ERROR
601 FORMAT(' ',13,4X,4E17.8)
  DO 20 N=1,3
    XN = XBEGIN + N*H
    YN = SOLN(XN)
    F(N+1) = XN + YN
    PRINT 601, N,XN,YN,YN,ERROR
  20 CONTINUE
C
C ** BEGIN ITERATION
C DO 50 N=4,NSTEPS
  YN = YN + (H/24.)*(55.*F(4)-59.*F(3)+37.*F(2)-9.*F(1))
  XN = XBEGIN + N*H
  F(1) = F(2)
  F(2) = F(3)
  F(3) = F(4)
  F(4) = XN + YN
  YOFXN = SOLN(XN)
  ERROR = YN - YOFXN
  PRINT 601, N,XN,YN,YOFXN,ERROR
50 CONTINUE
                                STOP
END

```

## نتایج کامپیوتری برای مثال ۵-۸

N	XN	YN	Y(XN)	ERROR
0	0.	0.	0.	0.
1	0.31250000E-01	0.49340725E-03	0.49340725E-03	0.
2	0.62500000E-01	0.19944459E-02	0.19944459E-02	0.
3	0.93750000E-01	0.45351386E-02	0.45351386E-02	0.
4	0.12500000E-00	0.81484411E-02	0.81484467E-02	-0.55879354E-08
5	0.15625000E-00	0.12868421E-01	0.12868434E-01	-0.12922101E-07
6	0.18750000E-00	0.18730211E-01	0.18730238E-01	-0.26309863E-07
7	0.21875000E-00	0.25770056E-01	0.25770098E-01	-0.41676685E-07
8	0.25000000E-00	0.34025350E-01	0.34025416E-01	-0.65192580E-07
9	0.28125000E-00	0.43534677E-01	0.43534756E-01	-0.78696758E-07
10	0.31250000E-00	0.54337843E-01	0.54337934E-01	-0.90803951E-07
11	0.34375000E-00	0.66475919E-01	0.66476032E-01	-0.11269003E-06
12	0.37500000E-00	0.79991280E-01	0.79991400E-01	-0.12014061E-06
13	0.40625000E-00	0.94927646E-01	0.94927788E-01	-0.14156103E-06

## نتایج کامپیوتری برای مثال ۵.۸ (ادامه)

N	XN	YN	Y(XN)	ERROR
14	0.43750000E-00	0.11133012E-00	0.11133029E-00	-0.16111881E-06
15	0.46875000E-00	0.12924525E-00	0.12924545E-00	-0.19185245E-06
16	0.50000000E 00	0.14872105E-00	0.14872126E-00	-0.21234155E-06
17	0.53125000E 00	0.16980705E 00	0.16980730E-00	-0.24400651E-06
18	0.56250000E 00	0.19255438E-00	0.19255446E-00	-0.26822090E-06
19	0.59375000E 00	0.21701577E-00	0.21701607E-00	-0.29988587E-06
20	0.62500000E 00	0.24324562E-00	0.24324594E-00	-0.31664968E-06
21	0.65625000E 00	0.27130008E-00	0.27130044E-00	-0.34645200E-06
22	0.68750000E 00	0.30123707E-00	0.30123746E-00	-0.39115548E-06
23	0.71875000E 00	0.33311634E-00	0.33311677E-00	-0.42840838E-06
24	0.75000000E 00	0.36699954E-00	0.36700001E-00	-0.46566129E-06
25	0.78125000E 00	0.40295030E-00	0.40295079E-00	-0.49173832E-06
26	0.81250000E 00	0.44103424E-00	0.44103476E-00	-0.52526593E-06
27	0.84375000E 00	0.48131907E-00	0.48131964E-00	-0.56624413E-06
28	0.87500000E 00	0.52387466E 00	0.52387527E 00	-0.61094761E-06
29	0.90625000E 00	0.56877308E 00	0.56877375E 00	-0.66310167E-06
30	0.93750000E 00	0.61608872E 00	0.61608934E 00	-0.71525574E-06
31	0.96875000E 00	0.66589829E 00	0.66589907E 00	-0.77486038E-06
32	0.09999999E 01	0.71828098E 00	0.71828181E 00	-0.82701445E-06

## تکمیل

۷.۰۸-۱ با استفاده از فرمول (۴۵.۸) الف، ضرایب  $\gamma_k$  ( $k=1, \dots, 4$ ) را در فرمول ادمز-بشفور (۴۵.۸) به دست آورید.

۷.۰۸-۲ در فرمول (۴۵.۸)،  $m$  را برابر با ۴ قرار دهید و مانند فرمول (۴۷.۸)، فرمول ادمز-بشفور متناظر با این مقدار  $m$  را بر حسب عرضها به دست آورید.

۷.۰۸-۳ فرمول میلن<sup>۱</sup> (۵۱.۸) و جمله خطای متناظر آن را به دست آورید.

۷.۰۸-۴ برنامه‌ای بنویسید که در آن فرمول میلن برای انتگرالگیری یک معادله دیفرانسیل به ازای نقاط با فواصل مساوی استفاده شده باشد. فرض کنید که سه مقدار آغازین اول معلوم باشند.

۷.۰۸-۵ معادله مثال ۵.۸ را با استفاده از برنامه میلن، و به ازای  $h=1/32$ ، حل، و نتایج به دست آمده را با نتایج داده شده در مثال ۵.۸ مقایسه کنید.

۷.۰۸-۶ با استفاده از روش ادمز-بشفور (۴۷.۸) معادله  $xy' = x - y$  را  $y(2) = 2$  از  $x = 2$  تا  $x = 3$  با به ازای  $h = 0.05$  حل کنید. مقادیر آغازین را از جواب دقیق زیر به دست آورید



$$y(x) = \frac{x}{2} + \frac{2}{x}$$

۷-۷.۸ با استفاده از روش ادمز-بشفورت (۴۷.۸) به ازای  $h = 1/128$  و  $h = 1/64$  معادلهٔ  $y' + y = e^{-x}$  را از  $x = 0$  تا  $x = 1$  حل کنید. دقت نتایج حاصل را برآورد نمایید. مقادیر آغازین از جواب دقیق  $y = xe^{-x}$  به دست می آیند.

۸-۷.۸ با استفاده از فرمول (۴۹.۸) به ازای  $m = 2$  (نه  $m = 3$ ) و فرمول خطا در درونیایی بسجمله‌ای (۱۸.۲)، فرمولهای (۵۱.۸) را به دست آورید. برای انجام این کار، بحث انجام شده در آغاز بخش ۲.۷ مفید خواهد بود.

۹-۷.۸ با بسط  $y_{n+1}$  و  $y_{n-1}$  تا جملات مرتبهٔ سوم پیرامون  $x = x_n$ ، فرمول (۵۰.۸) را تحقیق و فرض کنید که مقادیر آغازین دقیق هستند.

### ۸.۸ روشهای پیشگو-تصحیحی

روشهای چندمرحله‌ای قسمت ۷.۸ با استفاده از بسجمله‌ایهایی به دست آمده بودند که در نقطهٔ  $x_n$  و نقاط ماقبل  $x_{n-1}$  درونیایی شده بودند. این فرمولها اغلب به فرمولهای نوع باز معروفند. فرمولهای نوع بسته ۲ را می توان بر اساس درونیایی بسجمله‌ای در نقطهٔ  $x_{n+1}$  و نیز  $x_n$  و نقاط ماقبل  $x_{n-1}$  به دست آورد. ساده‌ترین فرمول این نوع از تقریب‌زدن انتگرال موجود در رابطهٔ (۴۳.۸) به وسیلهٔ فرمول دوزنقه‌ای (۲۶.۷) به دست می آید. با انجام این عمل فرمول زیر حاصل می شود

$$y_{n+1} = y_n + \frac{h}{2} [f(x_n, y_n) + f(x_{n+1}, y_{n+1})] \quad n = 0, 1, \dots \quad (52.8)$$

خطای این فرمول برابر است با  $-(h^3/12)y'''$ ، و لذا بیانگر بهبودی است در روش اوایلر. ولی (۵۲.۸) يك معادلهٔ ضمنی بر حسب  $y_{n+1}$  است، زیرا  $y_{n+1}$  به عنوان يك شناسه در سمت راست ظاهر می شود.

در حالت کلی اگر  $f(x, y)$  تابعی غیر خطی باشد، نمی توانیم معادلهٔ (۵۲.۸) را نسبت به  $y_{n+1}$  دقیقاً حل کنیم. ولی می توانیم سعی کنیم که  $y_{n+1}$  را از راه بارستی به دست آوریم. بنا بر این با ثابت نگه داشتن  $x_n$ ، نخستین تقریب  $y_{n+1}^{(0)}$  را به وسیلهٔ فرمول اوایلر

$$y_{n+1}^{(0)} = y_n + hf(x_n, y_n) \quad (53.8)$$

برای  $y_{n+1}$  به دست می آوریم، سپس  $f(x_{n+1}, y_{n+1}^{(0)})$  را محاسبه می کنیم و درست‌تر است

معادله (۵۲.۸) قرار می‌دهیم تا تقریب زیر را به دست آوریم

$$y_{n+1}^{(1)} = y_n + \frac{h}{\gamma} [f(x_n, y_n) + f(x_{n+1}, y_{n+1}^{(0)})]$$

بعد  $f(x_{n+1}, y_{n+1}^{(1)})$  را محاسبه می‌کنیم و با استفاده مجدد از معادله (۵۲.۸)، تقریب دیگری به دست می‌آوریم. در حالت کلی، بارست با رابطه

$$y_{n+1}^{(k)} = y_n + \frac{h}{\gamma} [f(x_n, y_n) + f(x_{n+1}, y_{n+1}^{(k-1)})] \quad k = 1, 2, \dots \quad (52.8)$$

تعریف می‌شود. بارست، زمانی به پایان می‌رسد که دو بارست متوالی، با دقت مطلوب متوافق باشند. این بارست، که برای به دست آوردن مقادیر اصلاح شده  $y_{n+1}$  در نقطه ثابت  $x_{n+1}$  به کار می‌رود، اغلب بارست داخلی نامیده می‌شود تا با معادله (۵۲.۸)، که برای تولید مقادیر  $y_n$  به ازای  $n = 0, 1, \dots$  به کار می‌رود، تفاوت داشته باشد. این شیوه عمل را در الگوریتم ۴.۸ خلاصه می‌کنیم.

**الگوریتم ۴.۸:** یک روش پیشگو-تصحیحی مرتبه دوم برای معادله دیفرانسیل  $y(x_0) = y_0, y' = f(x, y)$  با  $h$  داده شده و  $x_n = x_0 + nh$  به ازای جميع مقادیر ثابت  $n = 0, 1, \dots$

۱. با استفاده از فرمول (۵۳.۸)،  $y_{n+1}^{(0)}$  را محاسبه کنید.

۲. با استفاده از فرمول (۵۴.۸)،  $y_{n+1}^{(k)}$ ،  $k = 1, 2, \dots$  را محاسبه کنید، و این عمل را بر مبنای  $k$  بارست انجام دهید تا به ازای مقدار معینی مانند  $\epsilon$  رابطه زیر

$$\frac{|y_{n+1}^{(k)} - y_{n+1}^{(k-1)}|}{|y_{n+1}^{(k)}|} < \epsilon$$

برقرار شود.

برای مشخص کردن  $\epsilon$  در الگوریتم ۴.۸، باید به خاطر داشته باشیم که دقت مورد نظر در هر مرحله به وسیله میزان خطای موجود در فرمول اساسی (۵۲.۸) و نمو  $h$  محدود می‌شود.

برای تطبیق این الگوریتم با حل یک مسئله خاص، باید نکات زیر را مشخص کنیم، (الف) تعداد مراحل لازم،  $N$ . (ب) حداکثر تعداد بارستهای داخلی،  $K$ . (پ) چه باید کرد، در صورت بیشتر شدن  $k$  از  $K$ .

چنین متداول است که یک فرمول صریح مانند فرمول اویلر را فرمول بازگویند، و

يك فرمول ضمنی مانند (۵۲.۸) را فرمول بسته. هنگامی که این فرمولها با هم به صورت يك جفت فرمول به کار برده شوند فرمول نوع باز، پیشگو نیز نامیده می شود، در حالی که فرمول نوع بسته يك تصحیحی خواننده می شود. معمولاً يك فرمول تصحیحی دقیقتر از يك فرمول پیشگو است، حتی زمانی که مرتبهٔ خطای گسسته سازی در هر دو یکی باشد و این امر عمدتاً برای این است که ضریب عبارت خطا کوچکتر است. در رابطه با فرمولهای تصحیحی، طبیعتاً دو سؤال مطرح می شود. اول اینکه «در چه شرایطی بارستهای داخلی روی  $h$ ، همگرا خواهند شد؟»، دوم اینکه «چند بارست برای به دست آوردن دقت مطلوب، لازم است؟» جواب سؤال دوم به چندین عامل بستگی دارد. ولی اگر فرمولهای پیشگو و تصحیحی از مرتبهٔ واحدی باشند، تجربه نشان داده است که اگر نمو  $h$  به طور مناسب انتخاب شده باشد، يك یا دو بار استفاده از فرمول تصحیحی کافی نیست، بهتر است که به جای ادامهٔ بارستها، نمو  $h$  کاهش داده شود. پاسخ سؤال اول در قضیهٔ ۲.۸ مستتر است.

قضیهٔ ۲.۸ اگر  $f(x, y)$  و  $\partial f / \partial y$  بر حسب  $x$  و  $y$  در بازهٔ بسته  $[a, b]$  پیوسته باشند، در صورتی که  $h$  چنان کوچک انتخاب شده باشد که به ازای  $x = x_n$  و کلیهٔ  $y$ ها با شرط

$$|y - y_{n+1}| \leq |y_{n+1}^{(0)} - y_{n+1}|$$

$$\left| \frac{\partial f}{\partial y} \right| h < 2 \quad (55.8)$$

برقرار باشد، آنگاه بارستهای داخلی که به وسیلهٔ (۵۴.۸) تعریف شده اند همگرا خواهند بود. برای اثبات این مطلب، اولاً ملاحظه می کنیم که در فرمول بارست (۵۴.۸)،  $x_n$  ثابت است. بنا بر این، اگر قرار دهیم  $Y^{(k)} = y_{n+1}^{(k)}$  آنگاه، می توانیم فرمول (۵۴.۸) را به صورت

$$Y^{(k)} = F(Y^{(k-1)})$$

بنویسیم که در آن داریم

$$F(Y) = \frac{h}{\tau} f(x_{n+1}, Y) + C$$

و در این رابطه  $C$  به  $h$  بستگی دارد، نه به  $Y$ . می توان این مطلب را به مثابهٔ موردی از بارست نقطهٔ ثابت که در قسمت ۳.۳ بررسی شد در نظر گرفت. در فرع قضیهٔ ۱.۳، ثابت کردیم که چنین بارستی همگرا خواهد بود، به شرط آنکه  $F'(Y)$  پیوسته باشد و به ازای جميع مقادیر  $Y$  با شرط  $|Y - y_{n+1}| \leq |Y^{(0)} - y_{n+1}|$ ،  $F(Y)$  ثابت  $F(Y)$  است، در رابطه

$$|F'(Y)| < 1$$

صدق کند. از آنجا که داریم  $F'(Y) = (h/2) \partial f / \partial y$  و چون  $\partial f / \partial y$  کراندار و بنا بر فرض ناصفر است، فرمول بارست (۵۴.۸) همگرا خواهد شد اگر

$$F'(Y) = \left| \frac{h}{2} \frac{\partial f}{\partial y} \right| < 1$$

یعنی اگر

$$h < \frac{2}{|\partial f / \partial y|}$$

و چون تساوی  $F'(Y) = (h/2) \partial f / \partial y$  برقرار است، قضیه ثابت می‌شود.

□ مثال ۶۰۸: معادله

$$y' = x - \frac{1}{y} \quad y(0) = 1$$

را از  $x = 0$  تا  $x = 0.2$ ، با استفاده از الگوریتم  $4.8$  و به ازای  $h = 0.1$  حل کنید  
از آنجا که خطا در فرمول (۵۴.۸) برابر است با  $(h^3/12)y'''$  و چون با مشتقگیری از معادله بالا خواهیم داشت  $2 \approx y'''(0)$ ، بنابراین خطا تقریباً برابر با  $0.0002$  خواهد بود. بنا بر این نمی‌توانیم در نتایج، دقتی بسا بیش از سه رقم اعشار انتظار داشته باشیم

مرحله ۱

$$y_0^{(0)} = 0.9 \quad \text{باروش اولی:}$$

$$y_1^{(1)} = 0.8994 \quad \text{با فرمول (۵۴.۸):}$$

$$y_2^{(2)} = 0.8994$$

و چون  $y_1^{(1)}$  و  $y_2^{(2)}$  تا چهار رقم اعشار متوافق‌اند، این جواب را می‌پذیریم و مقدار  $y_1$  را برابر با  $1.00118 = f(x_1, y_1) = y_1'$ ، به دست می‌آوریم

مرحله ۲، با روش اولی داریم

$$y_4^{(0)} = 0.8994 + 0.1(-1.00118) = 0.79922$$

و با فرمول (۵۴.۸) داریم

$$y_4^{(1)} = 0.8994 + 0.05 \left[ -1.0118 + \left( 0.2 - \frac{1}{0.79922} \right) \right] = 0.79962$$

$$y_4^{(2)} = 0.8994 + 0.05 \left[ -1.0118 + \left( 0.2 - \frac{1}{0.7962} \right) \right] = 0.7960$$

$$y_4^{(3)} = 0.7960$$

جواب  $y_4 = 0.7960$  را می‌پذیریم و  $y_4$  را محاسبه نموده و به مرحله بعدی می‌پردازیم. در جریان محاسبات، انتظار داریم که دقت به تدریج کاهش یابد. ظاهراً در اینجا به ازای  $h = 0.1$  به دو یا سه بار استفاده از فرمول «تصحیحی» نیاز داریم. علت آن عمدتاً این است که مرتبه پیشگویی که مورد استفاده واقع شده از مرتبه تصحیحی کمتر است. برای تحقیق همگرایی بارستهای داخلی به ازای  $h = 0.1$  در این مثال  $\partial f / \partial y = 1/y^2$  را محاسبه می‌کنیم، و به موجب قضیه ۲.۸ می‌خواهیم که  $h$  کمتر از  $2y^2$  باشد. جواب  $y$  را نداریم، اما از مراحل بالا آشکار است که در بازه  $[0.5, 0.2]$  داریم  $y > 0.7$ . بنابراین اگر داشته باشیم  $0.98 = (0.7)^2 < 2h$ ، بارستهای داخلی همگرا می‌شوند. □

### تمرین

۱-۸.۸ برای معادله خاص  $y' = Ay$  و  $y(0) = 1$  نشان دهید که اگر داشته باشیم  $|Ah/2| < 1$ ، فرمول تصحیحی دوزنقه‌ای (۵۲.۸) به معادله تفاضلی منجر می‌شود که جواب آن برابر است با

$$y_n = \left[ \frac{(1 + Ah/2)}{(1 - Ah/2)} \right]^n$$

۲-۸.۸ برای جوابی که در تمرین ۱-۸.۸ به دست آمد، نشان دهید که به ازای مقدار ثابت  $x = x_n = nh$  داریم

$$\lim_{h \rightarrow 0} y_n = e^{Ax_n}$$

۳-۸.۸ معادله  $y' = x^2 + y$  و  $y(0) = 1$  را از  $x = 0$  تا  $x = 0.5$  با استفاده از روش اویلر به عنوان پیشگو و فرمول (۵۴.۸) به عنوان تصحیحی، حل کنید. نمو  $h$  را طوری تعیین کنید که در  $x = 0.5$  دقتی تا چهار رقم اعشار به دست آید. عملیات را با نمو  $h = 0.05$  آغاز کنید.

### ۹.۸ روش آدمز-مولتن<sup>۱</sup>

فرمول تصحیحی از مرتبه بالاتر را می‌توان با استفاده از یک بسجمله‌ای که عمل درونیایی

را در نقاط  $x_{n-m}, \dots, x_n, x_{n+1}$  (به ازای عدد صحیح  $m > 0$ ) انجام می‌دهد، به دست آورد. فرمول پسونیوتن که عمل درونیایی را در  $m+2$  نقطه بر حسب  $s = (x - x_n)/h$  انجام می‌دهد چنین است

$$P_{m+1}(s) = \sum_{k=0}^{m+1} (-1)^k \binom{m+1-s}{k} \Delta^k f_{n+1-k} \quad (56.8)$$

این تفاضلهای بر اساس مقادیر  $f_{n-m}, \dots, f_n, f_{n+1}$  به دست آمده‌اند. اگر از فرمول (56.8) از  $x_n$  تا  $x_{n+1}$  انتگرال بگیریم و فرمول (43.8) را به کار ببریم، خواهیم داشت

$$y_{n+1} = y_n + h(\gamma'_0 f_{n+1} + \gamma'_1 \Delta f_n + \dots + \gamma'_{m+1} \Delta^{m+1} f_{n-m}) \quad (57.8)$$

که در آن

$$\gamma'_k = (-1)^k \int_0^1 \binom{m+1-s}{k} ds \quad k = 0, 1, \dots, m+1$$

چند مقدار اول  $\gamma'_k$  عبارت‌اند از

$$\gamma'_0 = 1 \quad \gamma'_1 = -\frac{1}{2} \quad \gamma'_2 = -\frac{1}{12} \quad \gamma'_3 = -\frac{1}{24} \quad \gamma'_4 = -\frac{1}{720}$$

خطا در فرمول (57.8) که از خطا در بسجمله‌ای درونیاب ناشی شده چنین است

$$E = \gamma'_{m+2} h^{m+3} y^{(m+3)}(\xi) \quad (58.8)$$

حالت  $m=2$  حالتی است که اغلب به کار برده می‌شود. اگر در (57.8) به ازای  $m=2$  تفاضلهای بر حسب عرضها بیان شوند، خواهیم داشت

$$y_{n+1} = y_n + \frac{h}{24} (9f_{n+1} + 19f_n - 5f_{n-1} + f_{n-2}) \quad (59.8)$$

و خطا برابر است با

$$E_{AM} = -\frac{19}{720} h^5 y^{(5)}(\xi) \quad (60.8)$$

فرمول (57.8) به فرمول ادمز-مولتن معروف است. واضح است که فرمول ادمز-مولتن از مرتبه چهارم (59.8)، یک فرمول مصحح از نوع بسته است، زیرا  $f_{n+1} = f(x_{n+1}, y_{n+1})$  متضمن کمیت مجهول  $y_{n+1}$  است. می‌توان نشان داد که اگر  $h$  به اندازه کافی کوچک و شرط

$|< \partial f / \partial y| < 9h / 24$  برقرار باشد، آنگاه فرایند بارست که بر اساس فرمول (۵۹.۸) به دست آمده، همگرا خواهد شد. یک فرمول پیشگوی ساده، فرمول مرتبه چهارم ادمز-بشفورت (۴۷.۸) است که برای استفاده با این فرمول مصحح مناسب خواهد بود. در این حالت، پیشگو و مصحح هم مرتبه هستند. اگر  $h$  به طور مناسب انتخاب شود، آنگاه بایکبار استفاده از فرمول مصحح، بهبودی مهمی در دقت حاصل می شود.

الگوریتم ۵۰.۸: روش پیشگو-تصحیحی ادمز-مولتن برای معادله دیفرانسیل  $y' = f(x, y)$  با  $h$  ثابت و  $x_n = x_0 + nh$  و مفروضات  $(y_0, f_0), (y_1, f_1), (y_2, f_2), (y_3, f_3)$  به ازای مقادیر ثابت  $n = 3, 4, \dots$

۱. با استفاده از فرمول زیر،  $y_{n+1}^{(0)}$  را محاسبه می کنیم

$$y_{n+1}^{(0)} = y_n + \frac{h}{24} (55f_n - 59f_{n-1} + 37f_{n-2} - 9f_{n-3})$$

۲.  $f_{n+1}^{(0)} = f(x_{n+1}, y_{n+1}^{(0)})$  را محاسبه می کنیم.

۳. مقدار

$$y_{n+1}^{(k)} = y_n + \frac{h}{24} [9f(x_{n+1}, y_{n+1}^{(k-1)}) + 19f_n - 5f_{n-1} + f_{n-2}]$$

$$k = 1, 2, \dots$$

را محاسبه می کنیم.

۴. عمل بارستن روی  $k$  را تا حصول رابطه

$$\frac{|y_{n+1}^{(k)} - y_{n+1}^{(k-1)}|}{|y_{n+1}^{(k)}|} < \varepsilon$$

به ازای  $\varepsilon$  مفروض ادامه می دهیم.

هنوز این الگوریتم کامل نیست مگر آنکه مشخص نماییم که در صورت نبود همگرایی در مرحله چهارم چه باید بکنیم. زیر برنامه ای مانند DVERK از مشخصه کاملتری برای زیر برنامه ای که در همه جا منظور ما را در حل معادلات دیفرانسیل بر آورده سازد، بر خوردار است.

فرمول تصحیحی علاوه بر ایجاد دقت بهتر، کار مفید دیگری نیز انجام می دهد. این فرمول بر آوردی از خطای گسسته سازی محلی به دست می دهد که برای تصمیم گیری در مورد اینکه آیا نمو  $h$  برای ایجاد دقت خواسته شده کافی هست یا نه، به کار می رود. برای آنکه

این روش برآورد خطا را برای روش پیشگو-تصحیحی که مشتمل بر فرمولهای ادمز-بشفورت و ادمز-مولتن است آزمایش کنیم، برآورد خطای موضعی را برای هر یک از آنها می نویسیم

$$E_{AB} = \frac{251}{720} h^5 y''(\xi_1)$$

$$E_{AM} = -\frac{19}{720} h^5 y''(\xi_2) \quad (۶۱.۸)$$

گیریم  $y_{n+1}^{(0)}$  معرف مقدار  $y_{n+1}$  حاصل از فرمول (۴۷.۸)، و  $y_{n+1}^{(1)}$  مقدار حاصل از یکبار استفاده از الگوریتم ۵.۸ باشد. اگر فرض شود که مقادیر  $f$  در تمام نقاط تا  $x_n$  و منجمله خود آن دقیق باشند و اگر  $y(x_{n+1})$  معرف مقدار دقیق  $y$  در نقطه  $x_{n+1}$  باشد، آنگاه از رابطه (۶۱.۸)، برآوردهای خطا را به دست خواهیم آورد:

$$y(x_{n+1}) - y_{n+1}^{(0)} = \frac{251}{720} h^5 y''(\xi_1) \quad (الف ۶۲.۸)$$

$$y(x_{n+1}) - y_{n+1}^{(1)} = -\frac{19}{720} h^5 y''(\xi_2) \quad (ب ۶۲.۸)$$

در حالت کلی،  $\xi_1 \neq \xi_2$ . ولی اگر در بازه مورد نظر، فرض شود که  $y''(x)$  تقریباً ثابت است، آنگاه با تقریب (ب ۶۲.۸) از (الف ۶۲.۸)، برآورد زیر را برای  $y'$  به دست خواهیم آورد:

$$h^5 y'' = \frac{720}{270} (y_{n+1}^{(1)} - y_{n+1}^{(0)})$$

با گذاردن این مقدار در رابطه (ب ۶۲.۸)، خواهیم داشت

$$\begin{aligned} y(x_{n+1}) - y_{n+1}^{(1)} &= -\frac{19}{720} (y_{n+1}^{(1)} - y_{n+1}^{(0)}) \\ &\approx -\frac{1}{14} (y_{n+1}^{(1)} - y_{n+1}^{(0)}) = D_{n+1} \end{aligned} \quad (۶۳.۸)$$

بنابراین خطا در مقدار تصحیح شده، تقریباً برابر با  $-\frac{1}{14}$  تفاضل بین مقدار تصحیح

شده و مقدار پیشگو شده است.

همان گونه که قبلاً اشاره شد، بهتر است که فرمول تصحیحی را تنها یکبار به کار برد. اگر میزان دقت حاصل از رابطه (۶۳.۸) کافی نباشد، بهتر است که به جای استفاده از فرمول مذکور برای بیش از یکبار، نمو را کاهش داد.



دریک برنامه همه منظوره برای حل معادلات دیفرانسیل، از برآورد خطا به طریق زیر استفاده می‌شود: فرض کنید که بخواهیم مانند رابطه (۴۱.۸)، خطای موضعی دریک مرحله واحد\* کراندار باشد به طوری که داشته باشیم

$$E_1 \leq \frac{|D_{n+1}|}{h} \leq E_2$$

و مقادیر آغازین نیز معین شده باشند. مراحل انجام کار چنین است.

۱. با استفاده از فرمول (۴۷.۸)،  $y_{n+1}^{(0)}$  را به دست می‌آوریم.  $f_{n+1}^{(0)}$  را محاسبه می‌کنیم.

۲. با استفاده از فرمول (۵۹.۸)،  $y_{n+1}^{(1)}$  را به دست می‌آوریم.  $f_{n+1}^{(1)}$  را محاسبه می‌کنیم.

۳. از روی رابطه (۶۳.۸)،  $|D_{n+1}|$  را محاسبه می‌کنیم.

۴. اگر  $E_1 \leq |D_{n+1}|/h \leq E_2$ ، با استفاده از همان مقدار  $h$ ، مرحله انتگرالگیری بعدی را انجام می‌دهیم.

۵. اگر  $|D_{n+1}|/h > E_2$ ، آنگاه نمو  $h$  بسیار بزرگ است و می‌باید کاهش یابد. در این حالت، معمول است که  $h/2$  را به جای  $h$  می‌گذارند و چهارمقدار آغازین را مجدداً محاسبه می‌کنند و سپس به مرحله ۱ بازمی‌گردند.

۶. اگر  $|D_{n+1}|/h < E_1$ ، آنگاه دقت به دست آمده بیش از اندازه لازم است. از این رو با قراردادن  $2h$  به جای  $h$  در وقت مصرفی کامپیوتر می‌توانیم صرفه‌جویی کنیم و در بازه‌ای به طول  $2h$ ، چهارمقدار آغازین جدید را مجدداً محاسبه کنیم و به مرحله ۱ برگردیم.

چنانچه در بالا ذکر شد در استفاده از روشهای پیشگو-تصحیحی با نمو متغیر، لازم است که اولاً، در ابتدای کار روشی برای به دست آوردن مقادیر آغازین لازم داشته باشیم؛ ثانیاً، روشی برای به دست آوردن مقادیر لازم در مراحل نصف را، وقتی بازه نصف می‌شود، داشته باشیم؛ ثالثاً، روشی برای به دست آوردن مقادیر لازم  $y$  برای هنگامی که بازه دو برابر می‌شود داشته باشیم. برای هر یک از این سه وضعیت، فرمولهای خاصی می‌توانند به دست آیند. این فرمولها پیچیدگی یک برنامه را سخت افزایش می‌دهند. اما یک ترکیب نسبتاً آرمانی از این فرمولها، به کارگیری روش مرتبه چهارم، رونگه-کوتا (۳۷.۸)، همراه بسایک زوج پیشگو-تصحیحی مرتبه چهارم، از قبیل فرمولهای (۴۷.۸) و (۵۹.۸) است. بنابراین روش رونگه-کوتا می‌تواند در ابتدای کار برای شروع مسئله، هنگام نصف کردن بازه و یا دو برابر کردن آن به کار رود، در حالی که وقتی نمو ثابت است، زوج پیشگو-تصحیحی می‌تواند برای ادامه طبیعی آن به کار رود.

پیش از اتمام این قسمت، باید یادآور شویم که فرمولهای پیشگو-تصحیحی زیساد

\* منظور از یک مرحله واحد، یک مرحله کامل به ازای یک  $h$  معین است. م.م.

دیگری نیز وجود دارند، و به ویژه فرمولهای زیر که منسوب به میلن<sup>۱</sup> هستند، غالباً به کار برده می‌شوند

$$y_{n+1}^{(0)} = y_{n-2} + \frac{2h}{3} (2f_n - f_{n-1} + 2f_{n-2}) \quad E_M^0 = \frac{2\lambda}{90} h^5 y^{(5)}(\xi_1) \quad (۶۴.۸ \text{ الف})$$

$$y_{n+1}^{(1)} = y_{n-1} + \frac{h}{3} (f_{n+1}^{(0)} + 4f_n + f_{n-1}) \quad E_M^1 = -\frac{1}{90} h^5 y^{(5)}(\xi_2) \quad (۶۴.۸ \text{ ب})$$

معادله (۶۴.۸ الف)، در قسمت ۶.۸ به دست آمده بود و معادله (۶۴.۸ ب) بر اساس قاعده سیمپسن برای انتگرالگیری عددی مبتنی است. همان گونه که برای فرمولهای ادمز-مولتن عمل کردیم، می‌توانیم نشان دهیم که برآورد خطای موضعی با رابطه زیر معین می‌شود

$$D_{n+1} = -\frac{1}{24} (y_{n+1}^{(1)} - y_{n+1}^{(0)}) \quad (۶۵.۸)$$

به نظر می‌رسد که برآورد خطا در روش میلن تا اندازه‌ای مطلوبتر از روش ادمز-مولتن باشد، ولی چنانچه خواهیم دید در برخی موارد رابطه (۶۴.۸ ب) دستخوش ناپایداری قرار می‌گیرد.

درحالی که برای انتگرالگیری از معادلات دیفرانسیل، در کتابهای مختلف روشهای زیادی ذکر شده‌اند، روش مرتبه چهارم رونگه-کوتا و روشهای پیشگوه-تصحیحی نظیر ادمز-مولتن یا میلن (۶۴.۸)، متداولترین روشها در ایالات متحده هستند. اگر چه برای همه مسائل، هیچ روشی کلاً بهتر از روش دیگر نیست، ولی بجاست که در حالت کلی به مزایا و معایب هر یک از این روشها اشاره کنیم.

امتیاز مهم روشهای رونگه-کوتا در خود-آغاز گر بودن آنهاست. علاوه بر آن پایدارند، دقت خوبی به دست می‌دهند، و به عنوان یک برنامه کامپیوتری جای نسبتاً کمی از حافظه را اشغال می‌کنند. روشهای RK (رونگه-کوتا) ای استانده، خطاهای موضعی را برآورد نمی‌کنند، لذا استفاده کننده راهی برای درک مناسب بودن یا نبودن نمو  $h$  به کار رفته، ندارد. البته می‌توان روشهای کنترل نمو را، که در قسمت ۶.۸ بیان شد، به کار برد ولی این امر از نظر وقت کامپیوتری پرهزینه خواهد بود. اشکال عمده دوم آن است که روش رونگه-کوتا در مقایسه با روشهای پیشگوه-تصحیحی مرتبه چهارم که تنها به دو بار محاسبه تابع احتیاج دارند، در هر مرحله از انتگرالگیری به چهار بار محاسبه تابع نیاز دارد. در بعضی از مسائل زمان محاسبه برای روشهای رونگه-کوتا تقریباً دو برابر است.

روشهای پیشگو-تصحیحی در هر مرحله به طور خودکار بر آورد خطا را به ما می دهند، و بدین طریق به برنامه اجازه می دهند که برای يك دقت مطلوب مقدار بهینه‌ای از  $h$  را انتخاب نماید. ضمناً روشهای مزبور سریع نیز هستند، زیرا در هر مرحله تنها به دو بار محاسبهٔ تابع احتیاج دارند. از سوی دیگر، نوشتن زیر برنامه‌های پیشگو-تصحیحی بسیار پیچیده است و برای شروع و دو برابر کردن و نصف کردن نمو  $h$  به تکنیکهای خاصی نیاز دارند و این روشها ممکن است دستخوش ناپایداری عددی قرار گیرند (قسمت ۱۱.۸ را ببینید). سالیهای زیادی بود که روشهای رونگه-کوتا برای انجام کارهای همه‌منظوره تقریباً به طور انحصاری در امریکا به کار می رفت، و اخیراً روشهای پیشگو-تصحیحی مقبولیت یافته اند.

در چند سال گذشته، روشهای همه‌منظورهٔ پیچیده‌ای که در آنها هم مرتبه‌های متغیر و هم نمو‌های متغیر به کار برده می شوند توسعه پیدا کرده اند. روشهای ادمز، که قبلاً شرح دادیم، از جملهٔ متداولترین روشهای مرتبه-متغیر و نمو-متغیر هستند. هدف از این روشها انتخاب مرتبهٔ مناسب و نمو مناسب به طور خودکار است، و این انتخابها مقدار کار لازم برای دستیابی به دقتی مشخص برای مسئله‌ای معین را حداقل می سازند. مزیت مهم دیگر این روشها خود-آغازگر بودن آنهاست، زیرا در شروع کار می توانند يك روش پایین-مرتبه را به کار برند و زمانی که نمو تغییر می کند برای جبران از دست رفتن مقادیر، می توانند به آسانی اصلاح شوند. توضیح کامل در مورد زیر برنامهٔ DIFSUB که بر اساس روش ادمز از نوع مرتبه-متغیر و نمو-متغیر تنظیم شده در کتاب گیر<sup>۱</sup> [صص ۱۶۷-۱۵۸] آمده است. زیر برنامه‌ای به نام DVOGER وجود دارد که آن نیز بر اساس روش گیر در برنامه‌های IMSL تنظیم شده و برای اجرا روی جدیدترین کامپیوترها مورد استفاده قرار گرفته است.

□ مثال ۷.۸: با استفاده از فرمولهای پیشگو-تصحیحی ادمز-مولتن و به ازای  $h = 1/32$ ، مثال ۵.۸ را حل کنید، نتایج به دست آمده را با نتایج حاصله در آن مثال مقایسه کنید. برنامه و نتایج به دست آمده از کامپیوتر در زیر داده شده اند. در این حالت، (مقدار تصحیح شده)  $x_n$  و  $y_n$ ، بر آورد خطای موضعی  $D_n$ ، و خطای واقعی  $e_n$  را مشخص کرده ایم. از مقایسهٔ این نتایج با نتایج به دست آمده در مثال ۵.۸، بهبودی قطعی در دقت را ملاحظه می کنیم، به ویژه زمانی که  $x$  به سمت ۱ نزدیک می شود، که در این حالت نتایج تا هفت یا هشت رقم با معنی صحیح اند. به نظر می رسد که بر آورد خطای موضعی نسبتاً ثابت و به طور کلی از خطای واقعی  $e_n$  اندکی کمتر است. ولی با بررسی دقیقتر درمی یابیم که در مراحل ۵ تا ۱۳، نتایج تنها تا شش رقم با معنا صحیح اند. توضیح این امر چنین است که برای این مراحل، مقادیر  $y_n$  از لحاظ مرتبهٔ مقدار کوچکتر از مقادیر آن در حالت  $x \rightarrow 1$  می باشند. چون  $D_n$  يك آزمون خطای مطلق است، تعداد ارقام با معنی دقیق در نتیجه را مشخص نمی سازد. وقتی که محاسبات در حساب با ممیز شناور انجام می گیرند، این امری عادی است.

## حل معادلات دیفرانسیل ۵۰۳

زمانی که با اعداد خیلی بزرگ یا خیلی کوچک در مقایسه بسا ۱ کار می‌شود، برای تعیین تعداد ارقام با معنی صحیح، آزمون نسبی نسبت به آزمون مطلق، شاخص بهتری است. برای مثال، آزمون خطای نسبی برای فرمول ادمز-مولتن به جای فرمول (۶۳.۸) فرمول زیر خواهد بود

$$\bar{D}_{n+1} = \frac{1}{14} \frac{|y_{n+1}^{(0)} - y_{n+1}^{(1)}|}{|y_{n+1}^{(1)}|}$$

## برنامه فورترن برای مثال ۷.۸

```

C ADAMS-MOULTON METHOD
  INTEGER I,N,NSTEPS
  REAL ERROR,F(4),H,XBEGIN,XN,YBEGIN,YN
C
  SOLN(X) = EXP(X) - 1. - X
C
C ** INITIALIZE
  PRINT 600
600 FORMAT('!ADAMS-MOULTON METHOD'/
*'0',3X,'N',14X,'XN',15X,'YN',9X,'DN = YN - YNP',8X,'ERROR'/)
  NSTEPS = 32
  H = 1./NSTEPS
  YBEGIN = 0.
  XBEGIN = 0.
C
C ** COMPUTE FIRST FOUR POINTS USING EXACT SOLUTION
  F(1) = XBEGIN + YBEGIN
  N = 0
  ERROR = 0.
  DIFF = 0.
  PRINT 601,N,XBEGIN,YBEGIN,DIFF,ERROR
601 FORMAT(' ',13,4X,4E17.8)
  DO 20 N=1,3
    XN = XBEGIN + N*H
    YN = SOLN(XN)
    F(N+1) = XN + YN
    PRINT 601, N,XN,YN,DIFF,ERROR
  20 CONTINUE
C
C ** BEGIN ITERATION
  DO 50 N=4,NSTEPS
  C PREDICT USING ADAMS-BASHFORTH FORMULA
    YNPRED = YN + (H/24.)*(55.*F(4)-59.*F(3)+37.*F(2)-9.*F(1))
    XN = XBEGIN + N*H
    FNPRED = XN + YNPRED
  C CORRECT USING ADAMS-MOULTON FORMULA
    YN = YN + (H/24.)*(9.*FNPRED + 19.*F(4) - 5.*F(3) + F(2))
    DIFF = (YN - YNPRED)/14.
    F(1) = F(2)
    F(2) = F(3)
    F(3) = F(4)
    F(4) = XN + YN
    YOFXN = SOLN(XN)
    ERROR = YN - YOFXN
    PRINT 601, N,XN,YN,DIFF,ERROR
  50 CONTINUE
  END
  STOP

```

## نتایج کامپیوتری برای مثال ۷.۸

N	XN	YN	DN	ERROR
0	0.	0.	0.	0.
1	0.31250000E-01	0.49340725E-03	0.	0.
2	0.62500000E-01	0.19944459E-02	0.	0.
3	0.93750000E-01	0.45351386E-02	0.	0.
4	0.12500000E-00	0.81484520E-02	0.78164571E-09	0.53551048E-08
5	0.15625000E-00	0.12868445E-01	0.90637643E-09	0.11408702E-07
6	0.18750000E-00	0.18730249E-01	0.88143028E-09	0.11175871E-07
7	0.21875000E-00	0.25770108E-01	0.91469178E-09	0.10011718E-07
8	0.25000000E-00	0.34025417E-01	0.93132257E-09	0.13969839E-08
9	0.28125000E-00	0.43534759E-01	0.96458407E-09	0.37252903E-08
10	0.31250000E-00	0.54337924E-01	0.99784564E-09	0.83819032E-08
11	0.34375000E-00	0.66476036E-01	0.99784564E-09	0.37252903E-08
12	0.37500000E-00	0.79991416E-01	0.10643686E-08	0.15832484E-07
13	0.40625000E-00	0.94927801E-01	0.11308917E-08	0.13969839E-07
14	0.43750000E-00	0.11133030E-00	0.11308917E-08	0.14901161E-07
15	0.46875000E-00	0.12924545E-00	0.10643686E-08	0.55879354E-08
16	0.50000000E 00	0.14872127E-00	0.11974147E-08	0.74505806E-08
17	0.53125000E 00	0.16980730E-00	0.13304609E-08	0.18626451E-08
18	0.56250000E 00	0.19255465E-00	0.13304609E-08	0.37252903E-08
19	0.59375000E 00	0.21701607E-00	0.13304609E-08	0.
20	0.62500000E 00	0.24324595E-00	0.13304609E-08	0.11175871E-07
21	0.65625000E 00	0.27130044E-00	0.13304609E-08	0.11175871E-07
22	0.68750000E 00	0.30123746E-00	0.13304609E-08	0.
23	0.71875000E 00	0.33311676E-00	0.13304609E-08	-0.37252903E-08
24	0.75000000E 00	0.36700001E-00	0.15965530E-08	-0.74505806E-08
25	0.78125000E 00	0.40295079E-00	0.15965530E-08	0.37252903E-08
26	0.81250000E 00	0.44103477E-00	0.15965530E-08	0.74505806E-08
27	0.84375000E 00	0.48131964E-00	0.18626451E-08	0.74505806E-08
28	0.87500000E 00	0.52387527E 00	0.15965530E-08	0.74505806E-08
29	0.90625000E 00	0.56877375E 00	0.21287372E-08	0.
30	0.93750000E 00	0.61608943E 00	0.21287372E-08	0.
31	0.96875000E 00	0.66589906E 00	0.21287372E-08	-0.74505806E-08
32	0.99999999E 01	0.71828180E 00	0.21287372E-08	-0.74505806E-08

## تمرین

۹.۸-۱ نشان دهید که با فرمول

$$y_{n+1}^{(k)} = y_n + \frac{h}{24} [9f(x_{n+1}, y_{n+1}^{(k-1)}) + 19f_n - 5f_{n-1} + f_{n-2}]$$

$$k = 1, 2, \dots$$

نقطه  $x_n$

تعریف می شود همگرا خواهد بود، به شرط آنکه داشته باشیم  $|(9h/24)(\partial f/\partial y)| < 1$  (قسمت ۸.۸ را ببینید).

۲=۹۰۸ با استفاده از روابط (۵۷.۸) و (۵۸.۸)، فرمول خطا را برای روش ادمز-مولتن (۶۰.۸) به دست آورید.

۳=۹۰۸ فرمول برآورد خطای محلی (۶۵.۸) را برای فرمولهای پیشگو-تصحیحی میلن (۶۴.۸) به دست آورید.

۴=۹۰۸ با استفاده از فرمولهای پیشگو-تصحیحی معادله  $y' = y + x^2$ ،  $y(0) = 1$  را از  $x=0$  تا  $x=2$  با  $h=0.1$  حل کنید. مقادیر آغازین صحیح تا ۶ رقم اعشاری به قرار زیرند

$$\begin{aligned}y(0) &= 1.000000 \\y(0.1) &= 1.105513 \\y(0.2) &= 1.222208 \\y(0.3) &= 1.359576\end{aligned}$$

$D_{n+1}$  را حساب و خطا را در  $x=2$  برآورد کنید.

### ۱۰.۸° پایداری روشهای عددی

نخستین باری که کامپیوترها برای حل معادلات دیفرانسیل به طور وسیع به کار گرفته شدند، ملاحظه شد که خطاهای ناشی از برخی از فرمولهای متداول برای انتگرالگیری، مانند فرمولهای میلن (۶۴.۸) در جواب، در مقایسه با خطایی که تنها از خطای گسسته سازی انتظار می رود، خیلی بیشترند. علاوه، زمانی که نمو کوچکتر می شد، این خطاها به ازای مقدار ثابتی از  $x$ ، به جای آنکه کوچکتر شوند، عملاً بزرگتر می شدند. برای تجسم این رفتار، روشی را که در قسمت ۷.۸ به دست آورده بودیم، بررسی می کنیم

$$y_{n+1} = y_{n-1} + 2hf_n \quad (66.8)$$

که در آن خطای گسسته سازی برابر با  $\frac{1}{6}h^3 y'''(\xi)$  است. انتظار می رفت که به ازای  $h$  ثابت، این روش نسبت به روش اویلر که دارای خطایی برابر با  $O(h^2)$  است، نتایج دقیقتری به دست دهد. ولی مسئله ساده

$$y' = -2y + 1 \quad y(0) = 1 \quad (67.8)$$

را که جواب دقیق آن برابر است با  $y = \frac{1}{4}e^{-2x} + \frac{1}{4}$  در نظر می گیریم.

نتایجی که در جدول ۱۰.۸ آورده شده اند توسط کامپیوتر و با کارگیری نمو  $h = 1/32$  به دست آمده اند. ستون شامل مقادیر انتخابی است از  $x$  که در این مقادیر جوابها چاپ می شوند،  $Y(N)$  معرف جواب دقیق است،  $Y_1(N)$  معرف جوابی است که با روش

## جدول ۱.۸

X(N)	Y(N)	Y1(N)	Y2(N)	E1(N)	E2(N)
0.	1.000000	1.000000	1.000000	0.	0.
0.0312500	0.9697065	0.9687500	0.9697065	-0.0009565	-0.0000000
0.5000000	0.6839397	0.6780370	0.6840817	-0.0059027	0.0001420
1.0000000	0.5676676	0.5633943	0.5678247	-0.0042733	0.0001571
1.5000000	0.5248935	0.5225730	0.5251328	-0.0023205	0.0002392
2.0000000	0.5091578	0.5080376	0.5097007	-0.0011202	0.0005429
2.2500000	0.5055545	0.5047962	0.5064264	-0.0007583	0.0008719
2.5000000	0.5033690	0.5028620	0.5047904	-0.0005070	0.0014214
3.0000000	0.5012394	0.5010190	0.5050759	-0.0002203	0.0038365
3.5000000	0.5004559	0.5003628	0.5108669	-0.0000931	0.0104110
3.7500000	0.5002765	0.5002165	0.5174337	-0.0000601	0.0171571
3.7812500	0.5002598	0.5002029	0.4819995	-0.0000568	-0.0182603
3.8125000	0.5002440	0.5001903	0.5196837	-0.0000538	0.0194397
3.8437500	0.5002293	0.5001784	0.4795391	-0.0000509	-0.0206902
3.8750000	0.5002154	0.5001672	0.5222413	-0.0000482	0.0220260
3.9062500	0.5002023	0.5001568	0.4767589	-0.0000456	-0.0234434
3.9375000	0.5001901	0.5001470	0.5251465	-0.0000431	0.0249564
3.9687500	0.5001785	0.5001378	0.4736156	-0.0000408	-0.0265630
4.0000000	0.5001677	0.5001292	0.5284445	-0.0000386	0.0282768

اولر به دست آمده است،  $Y_2(N)$  معرف جوابی ناشی از فرمول (۶۶.۸) و  $E_1(N)$  و  $E_2(N)$  خطاهای متناظر آنها. روش (۶۶.۸) به مقادیر آغازین  $y_0$  و  $y_1$  احتیاج دارد. مقدار

دقیقی را که از جواب دقیق یعنی  $\frac{1}{4} + \frac{1}{4}e^{-2x}$  به دست می آید، برای تعیین  $y_1$  به کار می بریم.

ستونهای خطا نشان می دهند که  $E_2(N)$  در چند مرحله اول به طور قابل ملاحظه ای از  $E_1(N)$  کوچکتر است، اما بعد سریعاً افزایش می یابد، به طوری که در  $x = 2.25$ ،  $E_2(N)$  از  $E_1(N)$  بزرگتر می شود. به ازای  $x \rightarrow 4$ ، جواب به طرف حالت ماندگار  $1/2$  نزدیک می شود. روش اولر در عمل با خطای یکنوا نزولی به این جواب حالت ماندگار نزدیک می شود. در حالی که در روش (۶۶.۸) خطا به طور نمایی بزرگ می شود. بعلاوه، به طوری که چند مرحله آخر (که نتایج در هر مرحله از انتگرالگیری چاپ شده اند) نشان می دهند، خطاهای  $E_2(N)$  علامت نوسانی\* دارند. بعد از  $x = 4$ ،  $Y_2(N)$  دارای ارقام دقیق با معنای صحیحی نخواهد بود. پدیده ای که در این مثال نشان داده شد به ناپایداری عددی معروف است.

برای کمک به درک این رفتار، معادلهٔ تفاضلی (۶۶.۸) را عمیقتر بررسی می کنیم. در مورد مثالی که بررسی شد، داریم  $f(x) = 1 - 2y$ ، و از این رو معادله (۶۶.۸) به صورت زیر درمی آید

## 1. steady-state

\* منظور آن است که به ترتیب یکبار مثبت و یکبار منفی می شوند. س. م.

$$y_{n+1} + 2hy_n - y_{n-1} = 2h \quad y_0 = 1 \quad (۶۸.۸)$$

می‌توانیم با استفاده از روشهای مذکور در قسمت ۲.۸، این معادله تفاضلی را صریحاً حل کنیم. جواب عمومی معادله (۶۸.۸) برابر است با

$$y_n = C_1 \beta_1^n + C_2 \beta_2^n + \frac{1}{2} \quad (۶۹.۸)$$

که در آن  $\beta_1$  و  $\beta_2$  ریشه‌های معادله مشخصه

$$\beta^2 + 2h\beta - 1 = 0$$

هستند، این ریشه‌ها برابرند با

$$\beta_{1,2} = -2h \pm \sqrt{1 + 4h^2}$$

اگر  $\sqrt{1 + 4h^2}$  را به سری تیلر تا جملات خطی بسط دهیم، این ریشه‌ها را می‌توانیم به شکل زیر بیان کنیم

$$\beta_1 = 1 - 2h + \mathcal{O}(h^2)$$

$$\beta_2 = -(1 + 2h) + \mathcal{O}(h^2)$$

با قراردادن این مقادیر در معادله (۶۹.۸) خواهیم داشت

$$y_n = C_1 (1 - 2h + \mathcal{O}(h^2))^n + C_2 (-1)^n (1 + 2h + \mathcal{O}(h^2))^n + \frac{1}{2} \quad (۷۰.۸)$$

در مبحث حساب دیفرانسیل و انتگرال نشان داده‌ایم که

$$\lim_{\varepsilon \rightarrow 0} (1 + \varepsilon)^{1/\varepsilon} = e$$

با استفاده از این حد و رابطه  $n = x_n/h$ ، به‌ازای  $x_n$  ثابت، نتیجه می‌شود

$$\lim_{h \rightarrow 0} (1 + 2h)^n = \lim_{h \rightarrow 0} (1 + 2h)^{(1/2h)(2x_n)} = e^{2x_n}$$

و نیز خواهیم داشت

$$\lim_{h \rightarrow 0} (1 - 2h)^n = e^{-2x_n}$$

از این رو به‌ازای  $h \rightarrow 0$ ، جواب (۷۰.۸) به‌جواب

$$y_n = \left( C_1 e^{-2x_n} + \frac{1}{2} \right) + C_2 (-1)^n e^{2x_n} \quad (۷۱.۸)$$



نزدیک می‌شود. بنابراین جملهٔ اول به سمت جواب حقیقی معادلهٔ دیفرانسیل میل می‌کند. جملهٔ دوم اضافی (غیر مربوط) است و تنها ناشی از این امر است که یک معادلهٔ تفاضلی مرتبهٔ دوم را به جای معادلهٔ تفاضلی مرتبهٔ اول قراردادیم. اگر تمام عملیات حسابی دقیق باشند، با استفاده از شرایط اولیه نتیجه می‌گیریم  $C_1 = 0$ ، و لذا جواب دقیق از معادلهٔ (۷۱.۸) به دست می‌آید. ولی در عمل مقداری خطا وارد می‌شود که عمدتاً به علت گرد کردن یا به علت مقادیر آغازین غیردقیقی است، و از این رو  $C_1$  دقیقاً صفر نخواهد شد. بنابراین در هر مرحله از انتگرالگیری، خطای کوچکی داخل می‌شود، که مقدار آن بر اثر ضرب شدن در عامل  $e^{2.5n}$  (۱-) که به طور نمایی افزایش می‌یابد، بزرگ می‌شود. چون قسمت اصلی جواب حقیقی به طور نمایی نزول می‌کند، خطای حاصل از جواب اضافی، سرانجام بر جواب حقیقی مسلط شده منجر به نتایج کاملاً غیردقیقی می‌شود.

با بیانی خالی از دقت می‌توانیم روشی را ناپایدار گوئیم، که خطاهای وارد در آن هنگام محاسبات، به نسبت نمایی افزایش یابند.

روشهای یک مرحله‌ای مانند روشهایی از نوع رونگه-کوتا به ازای  $h$  به اندازهٔ کافی کوچک، ناپایداری عددی نشان نمی‌دهند. روشهای چندمرحله‌ای در بعضی حالات ممکن است به ازای تمام مقادیر  $h$  و در برخی دیگر به ازای یک رشته از مقادیر  $h$  ناپایدار باشند. برای اینکه ببینیم آیا یک روش چندمرحله‌ای پایدار است یا نیست، می‌توانیم چنین عمل کنیم: اگر روش چندمرحله‌ای به یک معادلهٔ تفاضلی از مرتبهٔ  $k$  منجر شود، ریشه‌های معادلهٔ مشخصهٔ متناظر با معادلهٔ تفاضلی همگن را به دست می‌آوریم و آنها را  $\beta_i$ ،  $(i = 1, \dots, k)$ ، می‌نامیم. در این صورت جواب عمومی معادلهٔ تفاضلی همگن چنین خواهد شد

$$y_n = c_1 \beta_1^n + c_2 \beta_2^n + \dots + c_k \beta_k^n \quad (72.8)$$

یکی از این جوابها، مثلاً  $\beta_1^n$ ، هنگامی که  $h \rightarrow 0$ ، به سمت جواب دقیق معادلهٔ تفاضلی میل می‌کند. همهٔ جوابهای دیگر خرابی هستند. یک روش چندمرحله‌ای را پایدار قوی گویند اگر  $h \rightarrow 0$ ، ریشه‌های اضافی در شرط

$$|\beta_i| < 1 \quad i = 2, 3, \dots, k$$

صدق کنند. در این شرایط هر خطایی که در محاسبه وارد شود، هنگام افزایش  $n$  از بین می‌رود، در حالی که اگر یکی از ریشه‌های خارجی  $\beta_i$  از لحاظ قدر مطلق بیشتر از یک باشد، خطاها به طور نمایی بزرگ می‌شوند.

در معادلهٔ تفاضلی کلی  $y' = f(x, y)$ ، به دست آوردن ریشه‌های معادلهٔ مشخصه،  $\beta_i$ ، غیرممکن خواهد بود. ولی برای آنکه معیاری از ناپایداری یک روش به دست آید، معمولاً بررسی معادلهٔ خاص  $y' = \lambda y$ ، که در آن  $\lambda$  یک عدد ثابت است، کفایت می‌کند.

در ابتدا روش مرتبهٔ چهارم ادمز-بشفورت را بررسی می‌کنیم. اگر در معادلهٔ (۷۲.۸)

## 1. strongly stable

قرار دهیم  $f(x, y) = \lambda y$ ، خواهیم داشت

$$y_{n+1} - y_n - \frac{h\lambda}{24} (55y_n - 59y_{n-1} + 37y_{n-2} - 9y_{n-3}) = 0 \quad (73.8)$$

معادله مشخصه این معادله تفاضلی چنین است:

$$\beta^4 - \beta^3 - \frac{h\lambda}{24} (55\beta^3 - 59\beta^2 + 37\beta - 9) = 0$$

البته ریشه‌های این معادله توابعی از  $h\lambda$  هستند. رسم بر این است که معادله مشخصه را به شکل

$$\rho(\beta) + h\lambda\sigma(\beta) = 0 \quad (74.8)$$

می‌نویسند، که در آن  $\rho(\beta)$  و  $\sigma(\beta)$  بسجمله‌ایهایی هستند که به صورت زیر تعریف می‌شوند

$$\rho(\beta) = \beta^4 - \beta^3$$

$$\sigma(\beta) = -\frac{1}{24} (55\beta^3 - 59\beta^2 + 37\beta - 9)$$

دیده می‌شود که به ازای  $h \rightarrow 0$ ، معادله (74.8) به  $\rho(\beta) = 0$  تبدیل می‌شود که ریشه‌های آن برابرند با  $\beta_1 = 1$ ،  $\beta_2 = \beta_3 = \beta_4 = 0$ . به ازای  $h \neq 0$ ، جواب عمومی معادله (73.8) به شکل

$$y_n = c_1 \beta_1^n + c_2 \beta_2^n + c_3 \beta_3^n + c_4 \beta_4^n$$

است که در آن  $\beta_i$ ها جوابهای معادله (74.8) هستند. می‌توان نشان داد که وقتی  $h \rightarrow 0$ ،  $\beta_i^n$  به جواب مطلوب معادله  $y' = \lambda y$  نزدیک می‌شود، در حالی که سایر ریشه‌ها متناظر با جوابهای خارجی هستند. چون ریشه‌های معادله (74.8) توابع پیوسته‌ای از  $h$  هستند، از آنجا نتیجه می‌شود که وقتی  $h$  به اندازه کافی کوچک باشد و به ازای  $i = 2, 3, 4$ ، شرط  $|\beta_i| < 1$  برقرار است. و از این رو بنا بر تعریف پایداری، روش ادمز-بشفورت قویاً پایدار است. همه روشهای چندمرحله‌ای به یک معادله مشخصه به شکل (74.8) منجر می‌شوند که گاهی سمت چپ آنها را بسجمله‌ای پایداری نامند. تعریف پایداری می‌تواند بر اساس بسجمله‌ای پایداری مطرح شود. روشی را پایدار قوی گویند، که در آن همه ریشه‌های  $\rho(\beta) = 0$  ریشه ساده  $\beta = 1$ ، از لحاظ قدر مطلق کمتر از یک باشند.

حال خواص پایداری روش میلن (64.8) را که به صورت

$$y_{n+1} = y_{n-1} + \frac{h}{3} (f_{n+1} + 4f_n + f_{n-1}) \quad (75.8)$$

بیان شده است، بررسی می‌کنیم. بازم از قرارداد  $f(x, y) = \lambda y$ ، خواهیم داشت

$$y_{n+1} - y_{n-1} - \frac{h\lambda}{3} (y_{n+1} + 4y_n + y_{n-1}) = 0$$

که معادلهٔ مشخصهٔ آن چنین می‌شود:

$$\rho(\beta) + h\lambda\sigma(\beta) = 0 \quad (۷۶.۸)$$

که در آن

$$\rho(\beta) = \beta^2 - 1$$

$$\sigma(\beta) = \beta^2 + 4\beta + 1$$

در اینجا معادلهٔ  $\rho(\beta) = 0$  ریشه‌هایی به صورت  $\beta_1 = 1$  و  $\beta_2 = -1$  دارد و از این رو بنا بر تعریف بالا، روش میلن پایدار قوی نیست. برای آنکه مفهوم ضمنی این مطلب را درک کنیم، ریشه‌های بسجمله‌ای پایداری (۷۶.۸) را محاسبه می‌کنیم. به‌ازای مقدار بسیار کوچک  $h$  داریم

$$\beta_1 = 1 + \lambda h + \mathcal{O}(h^2)$$

$$\beta_2 = -(1 - \lambda h/3) + \mathcal{O}(h^2) \quad (۷۷.۸)$$

از این رو جواب عمومی معادلهٔ (۷۵.۸) چنین است

$$y_n = c_1 (1 + \lambda h + \mathcal{O}(h^2))^n + c_2 (-1)^n (1 - \lambda h/3 + \mathcal{O}(h^2))^n$$

اگر  $n$  را مساوی  $x_n/h$  بگذاریم، و  $h$  را به سمت صفر میل دهیم این جواب به‌جواب

$$y_n = c_1 e^{\lambda x_n} + c_2 (-1)^n e^{-\lambda x_n/3} \quad (۷۸.۸)$$

نزدیک می‌شود. در این حالت پایداری به‌علامت  $\lambda$  بستگی پیدا می‌کند. اگر  $\lambda > 0$  و لذا جواب مطلوب به‌طور نمایی صعودی باشد، واضح است که جواب خارجی به‌طور نمایی نزولی خواهد بود و لذا روش میلن پایدار می‌شود. از سوی دیگر اگر  $\lambda < 0$ ، آنگاه روش میلن ناپایدار خواهد بود زیرا جواب اضافی به‌طور نمایی افزایش می‌یابد و سرانجام جواب مطلوب را محو می‌کند. روشهایی از این نوع را که پایداریشان برای معادلهٔ آزمون  $y' = \lambda y$  بستگی به‌علامت  $\lambda$  دارد، روشهای پایدار ضعیف نامند، برای معادلهٔ کلیتر  $y' = f(x, y)$ ، هرگاه در بازهٔ انتگرالگیری  $\partial f / \partial y < 0$  از روش میلن می‌توان پایداری ضعیف را انتظار داشت.

در عمل همهٔ روشهای چندمرحله‌ای در رشته‌ای از مقادیر برای نمو  $h$ ، ناپایداری

نشان می‌دهند مثلاً، روش مرتبهٔ دوم ادمز-بشفورث را که به‌صورت

$$y_{n+1} = y_n + \frac{h}{4} \{3f_n - f_{n-1}\}$$

تعریف می‌شود در نظر می‌گیریم. اگر این روش برای معادلهٔ آزمون  $y' = \lambda y$  به کار گرفته شود، معادلهٔ تفاضلی

$$y_{n+1} - y_n - \frac{h\lambda}{4} \{3y_n - y_{n-1}\} = 0$$

به دست خواهد آمد، و از روی آن بسجمله‌ای پایداری

$$\beta^2 - \beta - \frac{h\lambda}{4} \{3\beta - 1\}$$

یا معادله

$$\beta^2 - \left(1 + \frac{3h\lambda}{4}\right)\beta + \frac{h\lambda}{4} = 0$$

به دست می‌آید. اگر  $\lambda < 0$ ، هر دو ریشهٔ این معادلهٔ درجهٔ دوم از لحاظ قدرمطلق کوچکتر از یک هستند به شرط آنکه  $0 < h\lambda < 1$  - . در این حالت چون خطاهای ناشی از جواب خارجی بزرگ نخواهند شد، پایداری مطلق به دست خواهیم آورد. ولی اگر  $|h\lambda| > 1$ ، آنگاه یکی از این ریشه‌ها از لحاظ قدرمطلق از یک بزرگتر می‌شود و تاحدی به ناپایداری بر می‌خوریم. شرط  $0 < h\lambda < 1$  -، نمو  $h$  را که برای این روش می‌تواند به کار رود به طور مؤثری محدود می‌کند. برای مثال اگر  $\lambda = -100$  آنگاه، برای اطمینان از پایداری باید  $h < 0.01$  انتخاب شود. یک روش چندمرحله‌ای را به ازای مقادیری از  $h\lambda$  که به ازای آنها ریشه‌های بسجمله‌ای پایداری آن (۷۴.۸) از لحاظ قدرمطلق کمتر از یک باشند، روش پایدار مطلق نامند. روشهای مختلف، نواحی مختلفی برای پایداری مطلق دارند. عموماً ما آن روشهایی را ترجیح می‌دهیم که بزرگترین ناحیه را برای پایداری مطلق دارند. برای مثال، می‌توان نشان داد که نواحی پایداری در روشهای ضمنی ادمز-مولتن، ۱۰ برابر بزرگتر از نواحی پایداری روشهای ادمز-شفورت با همان مرتبه است. مخصوصاً روش مرتبهٔ دوم ادمز-مولتن که با رابطهٔ

$$y_{n+1} = y_n + h \left( f_{n+1} - \frac{1}{4} f_n + \frac{1}{4} f_{n-1} \right)$$

معین می‌شود، برای معادلهٔ آزمون  $y' = \lambda y$  و به ازای  $\lambda < 0$  و  $-\infty < h\lambda < 0$ ، پایداری مطلق است.

برای معادلاتی به شکل  $y' = \lambda y$  و با  $\lambda > 0$ ، جواب مطلوب مانند  $e^{\lambda x}$ ، به طور نمایی بزرگ خواهد شد. هر روش چندمرحله‌ای باید یک ریشه، ریشهٔ اصلی، داشته باشد که این ریشه تقریبی است برای جواب مطلوب. سپس همهٔ ریشه‌های خارجی دیگر باید از لحاظ قدرمطلق کمتر از این ریشهٔ اصلی باشند. روشی با این خصوصیت را که همهٔ ریشه‌های

خارجی بسجمله‌ای پایداری آن از لحاظ قدرمطلق کمتر از ریشهٔ اصلی باشند روش پایدار نسبی<sup>۱</sup> نامند. نواحی پایداری برای روشهای چندمرحله‌ای در کتاب گیر [۳۰] به‌طور مبسوط مورد بحث قرار گرفته‌اند.

### تمرین

۱-۱۰۰۸ نشان دهید که فرمول تصحیحی بر اساس قاعدهٔ دوزنقه‌ای (۵۲.۸)، برای معادلاتی به شکل  $y' = \lambda y$  پایدار است (تمرین ۸.۸-۱ را ببینید).

۲-۱۰۰۸ نشان دهید که ریشه‌های معادلهٔ مشخصهٔ (۷۶.۸)، وقتی  $h \rightarrow 0$  می‌تواند به شکل (۷۷.۸) بیان شود، و نشان دهید که جواب معادلهٔ تفاضلی (۷۵.۸) وقتی  $h \rightarrow 0$  به جواب (۷۸.۸) نزدیک می‌شود.

۳-۱۰۰۸ يك برنامهٔ کامپیوتری بنویسید که ریشه‌های معادلهٔ مشخصهٔ (۷۳.۸) را برای فرمول ادمز-بشفورث به دست دهد. برای این برنامه قرار دهید  $\lambda = -1$  و  $h = 0.1$ . مقدار تقریبی  $\bar{h}$  را، که بیشتر از آن يك یا چند ریشهٔ این معادله از لحاظ قدرمطلق بزرگتر از يك می‌شوند، تعیین نمایید. و بدین طریق کران بالایی برای  $h$  به دست آورید که به ازای مقدار  $h$  بیشتر از آن کران، روش ادمز-بشفورث ناپایدار شود.

۴-۱۰۰۸ به ازای  $h = 1/2$  و از  $x = 0$  تا  $x = 6$ ، معادلهٔ (۶۷.۸) را با روش میلن (۶۴.۸) حل کنید. مقادیر آغازین را از جدول ۱.۸ به دست آورید. به نتیجهٔ ناپایداری روی جواب توجه نمایید.

### ۱۱.۸. کنترل و پخش خطای گرد کردن

در قسمت ۴.۸، خطای گسسته‌سازی را به صورت

$$e_n = y(x_n) - y_n$$

تعریف کردیم که در آن  $y(x_n)$  جواب حقیقی معادلهٔ دیفرانسیل است و  $y_n$  جواب دقیق معادلهٔ تفاضلی تقریبی برای معادلهٔ دیفرانسیل. در عمل چون در کامپیوترها طول کلمه<sup>۲</sup> محدود است، يك مقدار  $\bar{y}_n$  به دست می‌آوریم که، به علت خطاهای گرد کردن، با  $y_n$  تفاوت دارد. این تفاوت که با

$$r_n = y_n - \bar{y}_n$$

نشان داده می‌شود خطای گرد کردن<sup>۳</sup> نام دارد، یعنی اختلاف بین جواب دقیق معادلهٔ

1. relatively stable
2. word lengths
3. accumulated round-off-error

تفاضلی و جوابی که از کامپیوتر در نقطه  $x = x_n$  به دست می آید. در هر مرحله از انتگرالگیری یک خطای گرد کردن حاصل می شود که آن را خطای گرد کردن موضعی می نامیم و با  $\varepsilon_n$  مشخص می کنیم. برای مثال، در روش اویلر،  $\varepsilon_n$  به صورت زیر تعریف می شود

$$\tilde{y}_{n+1} = \tilde{y}_n + hf(x_n, \tilde{y}_n) + \varepsilon_n$$

خطای گرد کردن انباشته، اصلاً برابر با مجموع خطاهای گرد کردن موضعی نیست، زیرا هر یک از خطاهای موضعی پخش شده و ممکن است همچنان که محاسبات پیش می روند، هر یک از این خطاها افزایش یا کاهش یابد. در حالت کلی، موضوع پخش خطای گرد کردن خوب توجیه نشده است و نتایج نظری بسیار کمی درباره این موضوع در دست است. خطای گرد کردن انباشته به عوامل زیادی بستگی دارد، از جمله، (۱) نوع حسابی که در کامپیوتر به کار برده می شود، یعنی حساب با ممیز ثابت یا با ممیز شناور، (۲) روشی که کامپیوتر برای گرد کردن به کار می برد، (۳) ترتیبی که عملیات حسابی انجام می گیرند، (۴) شیوه عددی که به کار برده می شود.

چنان که در بخش ۱۰.۸، در بحث ناپایداری عددی نشان داده شد، اثر پخش گرد کردن می تواند خطرناک باشد. اما حتی در روشهای پایدار، از دست دادن دقت بر اثر خطاهای گرد کردن اجتناب ناپذیر است. این امر در فصل ۷ که قانون دزونز نقه ای برای محاسبه انتگرال به کار برده شده بود ملاحظه شد. از دست دادن دقت در بازه بزرگتر ممکن است آنقدر جدی باشد که نتایج را به طور کلی بی اعتبار سازد.

از خطای گرد کردن انباشته، با برخی فرضیات آماری درباره توزیع خطاهای گرد کردن موضعی، می توان برآوردهایی به دست آورد. در اینجا این گونه مسائل دنبال نخواهند شد، بلکه به دنبال یک شیوه ساده و مؤثری هستیم که هنگام حل معادلات دیفرانسیل، از دست دادن دقت بر اثر خطاهای گرد کردن را کاهش دهد.

بسیاری از فرمولهای مورد بحث برای حل معادلات دیفرانسیل در این فصل رامی توان

به شکل

$$y_{n+1} = y_n + h \Delta y_n$$

نوشت، کسه در آن  $h \Delta y_n$  معرف نموی است متضمن ترکیبات  $f(x, y)$  در نقاط انتخابی. این نمودر مقایسه با خود  $y_n$ ، معمولاً نموی کوچک است. در تشکیل حاصل جمع  $y_n + h \Delta y_n$  در حساب با ممیز شناور، کامپیوتر  $h \Delta y_n$  را به طرف راست منتقل می کند تا توان  $h \Delta y_n$  با توان  $y_n$  توافق نماید و با انجام این کار، بیتهایی از سمت راست حذف می شوند. سپس عمل جمع صورت می گیرد، ولی بر اثر حذف بیتها، خطای گرد کردن پدید می آید. برای اینکه این مطلب را روشنتر ببینیم، بیاییم دو عدد با ممیز شناور  $(10^4)$   $(0.54772)$  و  $(10^2)$   $(0.3856)$  را با هم جمع کنیم. با این فرض که طول کلمه را چهار رقم اعشاری

بگیریم. اگر عدد دوم را دو رقم به سمت راست انتقال دهیم، دو رقم آخر آن حذف می‌شوند و با افزودن آن به اولین عدد، مقدار  $(10^4)(0.5510)$  حاصل می‌شود، در حالی که اگر گرد کردن را به طور صحیح انجام می‌دادیم نتیجه می‌باید برابر با  $(10^4)(0.5511)$  می‌شد. البته این مثال مثالی است بیش از حد واضح، زیرا کامپیوترها با بیت‌های دودویی و طول کلمهٔ بزرگتر کار می‌کنند، اما با این وجود اثرات انباشتگی<sup>۱</sup> می‌توانند جدی باشند.

اکنون شیوهٔ ساده‌ای را ذکر می‌کنیم که خطاهایی از این نوع را به طور قابل توجهی کاهش می‌دهد. در ابتدا هر مقدار محاسبه شده‌ای از  $y_n$  به صورت دقت مضاعف ذخیره، و سپس  $y_n \Delta h$  با دقت ساده محاسبه می‌شود، و برای این کار تنها آن قسمت از  $y_n$  که متضمن دقت ساده و برای تشکیل  $h \Delta y_n$  لازم است به کار می‌رود، مجموع  $y_n = h \Delta y_n$  به صورت دقت مضاعف به دست می‌آید، و  $y_{n+1} = y_n + h \Delta y_n$  با دقت مضاعف ذخیره می‌شود. این شیوه را انباشتگی دقت مضاعف جزئی<sup>۲</sup> نامند. در بعضی\* از کامپیوترها، انجام محاسبات با دقت مضاعف به صورت یک دستورالعمل<sup>۳</sup> موجود است، ولی حتی اگر این نیز در کامپیوتر موجود نباشد، در هر مرحله از انتگرالگیری، تنها یک عمل جمع با دقت مضاعف می‌باید انجام گیرد. قسمت اصلی محاسبه، عبارت است از تعیین  $h \Delta y_n$ ، و این کار با دقت ساده انجام می‌پذیرد. کار اضافی و نیز حافظهٔ اضافی در رابطه با این امر بسیار اندک است. از طرف دیگر با انجام این کار به دست آوردن دقت ممکن، می‌تواند بسیار قابل توجه باشد، مخصوصاً زمانی که دقت زیاد در یک بازهٔ بزرگ مطلوب باشد. در واقع این شیوهٔ عمل در امر کاهش انباشتگی خطای گرد کردن آنقدر مؤثر است که هیچ برنامه‌ای از برنامه‌های همه منظورهٔ پیش ساخته برای حل معادلات دیفرانسیل نیست که برای ایجاد شکلی از این انباشتگی دقت مضاعف جزئی، امکاناتی فراهم نکرده باشد.

در اینجا بی‌مناسبت نیست که چند کلمه‌ای به عنوان توضیح بیاوریم. دقت یک انتگرالگیری عددی به خطای گسسته‌سازی و خطای گرد کردن انباشته بستگی دارد. برای آنکه خطای گسسته‌سازی کم باشد، طبیعتاً نمو  $h$  کوچک انتخاب می‌شود. از طرف دیگر، هرچه  $h$  کوچکتر گرفته شود، مراحل انتگرالگیری بیشتری باید پیموده شود که در پی آن احتمالاً خطای گرد کردن بیشتر می‌شود. بنابراین مقدار بهینه‌ای از نمو  $h$  وجود دارد که این مقدار بهینه برای یک کامپیوتر مفروض و یک مسئلهٔ مفروض، بهترین دقت را به دست می‌دهد. پیدا کردن این مقدار بهینه در عمل بدون اینکه مقدار زیادی از وقت کامپیوتر را بگیرد بسیار مشکل است. ولی وجود این مقدار بهینه نشان می‌دهد که اگر نمو را خیلی کوچک بگیریم، خطراتی را در پی خواهد داشت.

### 1. cumulative                      2. partial double-precision accumulation

\* در متن اصلی کتاب، کلمهٔ «بعضی» از کامپیوترها آمده است، ولی انجام محاسبات با دقت مضاعف در بیشتر کامپیوترهای همه‌منظوره موجود است. — م.

### 3. instruction

□ مثال ۸.۸: معادله

$$y' = \frac{1}{x^2} - \frac{y}{x} - y^2 \quad y(1) = -1$$

را با استفاده از روش ادمز-بشفورت، به ازای  $h = 1/256$  از  $x = 1$  تا  $x = 3$  حل کنید. این کار را یک بار با انباشتگی دقت مضاعف جزئی و یک بار بدون آن انجام دهید. نتایج به دست آمده از کامپیوتر در زیر داده شده اند. نمو عمداً آنقدر کوچک انتخاب شده که خطای گسسته سازی قابل صرف نظر باشد. نتایج بعد از هر ۱۶ مرحله چاپ شده اند. جواب دقیقی این مسئله برابر با  $y = -1/x$  است. بنابراین دقت می تواند به آسانی

نتایج کامپیوتری برای مثال ۸.۸

X	SINGLE PRECISION	PARTIAL DOUBLE PRECISION
0.99999999	-0.99999999	-0.99999999
1.06250000	-0.94117642	-0.94117647
1.12500000	-0.88888878	-0.88888889
1.18750000	-0.84210509	-0.84210526
1.24999990	-0.79999977	-0.80000000
1.31249990	-0.76190444	-0.76190476
1.37500000	-0.72727232	-0.72727273
1.42750000	-0.69565168	-0.69565218
1.50000000	-0.66666608	-0.66666667
1.56249990	-0.63999934	-0.64000001
1.62499990	-0.61538386	-0.61538462
1.68750000	-0.59259175	-0.59259260
1.75000000	-0.57142763	-0.57142858
1.81250000	-0.55172310	-0.55172415
1.87499990	-0.53333220	-0.53333335
1.93749990	-0.51612781	-0.51612905
2.00000000	-0.49999869	-0.50000001
2.06250000	-0.48484711	-0.48484850
2.12500000	-0.47058678	-0.47058825
2.18749990	-0.45714134	-0.45714287
2.24999990	-0.44444284	-0.44444446
2.31250000	-0.43243076	-0.43243245
2.37500000	-0.42105088	-0.42105265
2.43750000	-0.41025458	-0.41024643
2.49999990	-0.39999810	-0.40000002
2.56249990	-0.39024193	-0.39024393
2.62500000	-0.38095033	-0.38095240
2.68750000	-0.37209089	-0.37209304
2.75000000	-0.36363416	-0.36363639
2.81249990	-0.35555328	-0.35555558
2.87499990	-0.34782372	-0.34782612
2.92750000	-0.34042308	-0.34042556
3.00000000	-0.33333080	-0.33333336



آزمایش شود. در  $x = 3$  نتایج با دقت مضاعف جزئی، تا سه واحد در رقم هشتم دقیق هستند و با دقت ساده تا  $253$  واحد در رقم هشتم. از آنجایی که این خطا از گرد کردن ناشی شده است، این مثال به وضوح مؤثر بودن دقت مضاعف جزئی را در امر کاهش انباشتگی خطای گرد کردن نشان می‌دهد.

### تمرین

۱-۱۱۰۸ برنامه‌ای بر اساس روش ادمز-بشفورت بسازید که هم از انباشتگی دقت مضاعف جزئی و هم از دقت ساده در آن استفاده شود.

۲-۱۱۰۸ برنامهٔ تمرین ۱-۱۱۰۸ را برای حل معادلهٔ زیر از  $x = 0$  تا  $x = 2$  با استفاده از نمو ثابت  $h = 0.01$  به کار ببرید

$$y' = -2y \quad y(0) = 1$$

مقادیر آغازین می‌توانند از جواب دقیق  $h = e^{-2x}$  به دست آیند. خطای ناشی از گرد کردن چقدر است؟

۳-۱۱۰۸ برنامه‌ای برای روش کلاسیک مرتبهٔ چهارم رونگه-کوتا بنویسید که هم از انباشتگی دقت مضاعف و هم از دقت ساده در آن استفاده شود. این برنامه را برای حل معادلهٔ تمرین ۲-۱۱۰۸ به ازای همان مقدار از نمو به کار گیرید.

### ۱۲.۸\* دستگاه‌های معادلات دیفرانسیل

در بسیاری از زیر برنامه‌های معادلات دیفرانسیل همه منظوره، فرض بر آن است که یک معادلهٔ دیفرانسیل مرتبهٔ  $N$  به صورت یک دستگاه از  $N$  معادلهٔ مرتبهٔ اول بیان شده است. برای یک معادلهٔ مرتبهٔ  $N$  که به شکل

$$y^{(N)} = f(x, y(x), y'(x), \dots, y^{(N-1)}(x)) \quad (79.8)$$

می‌باشد، این تبدیل همیشه به صورت زیر انجام می‌پذیرد: با انتخاب  $y_1 = y$ ، داریم

$$y_1' = y_2$$

$$y_2' = y_3$$

$$y_3' = y_4$$

.....

$$y_{N-1}' = y_N$$

$$y_N' = f(x, y_1, y_2, \dots, y_N) \quad (80.8)$$



$$k_{\psi} = hf\left(x_n + \frac{h}{\psi}, y_n + \frac{k_{\psi}}{\psi}, z_n + \frac{l_{\psi}}{\psi}\right)$$

$$l_{\psi} = hg\left(x_n + \frac{h}{\psi}, y_n + \frac{k_{\psi}}{\psi}, z_n + \frac{l_{\psi}}{\psi}\right)$$

$$k_{\psi} = hf(x_n + h, y_n + k_{\psi}, z_n + l_{\psi})$$

$$l_{\psi} = hg(x_n + h, y_n + k_{\psi}, z_n + l_{\psi})$$

بدیهی است که این کار را برای یک دستگاه معادلات نیز می‌توان تعمیم داد. باید توجه داشت که همهٔ نموها با زیرنماهای کوچکتر را می‌بایسد قبل از نموهای بعدی با زیرنماهای بالاتر محاسبه نمود.

فرمولهای ادمز-مولتن که برای معادله‌های (۸۲.۸) به کار رفته‌اند به صورت زیر عمل می‌کنند:

$$y_{n+1}^{(0)} = y_n + \frac{h}{\psi\psi} [\Delta\Delta f(x_n, y_n, z_n) - \Delta\Delta f(x_{n-1}, y_{n-1}, z_{n-1}) \\ + 3\Delta f(x_{n-2}, y_{n-2}, z_{n-2}) - \Delta f(x_{n-3}, y_{n-3}, z_{n-3})]$$

$$z_{n+1}^{(0)} = z_n + \frac{h}{\psi\psi} [\Delta\Delta g(x_n, y_n, z_n) - \Delta\Delta g(x_{n-1}, y_{n-1}, z_{n-1}) \\ + 3\Delta g(x_{n-2}, y_{n-2}, z_{n-2}) - \Delta g(x_{n-3}, y_{n-3}, z_{n-3})]$$

$$y_{n+1}^{(1)} = y_n + \frac{h}{\psi\psi} [\Delta f(x_{n+1}, y_{n+1}^{(0)}, z_{n+1}^{(0)}) + \Delta f(x_n, y_n, z_n) \\ - \Delta f(x_{n-1}, y_{n-1}, z_{n-1}) + \Delta f(x_{n-2}, y_{n-2}, z_{n-2})]$$

$$z_{n+1}^{(1)} = z_n + \frac{h}{\psi\psi} [\Delta g(x_{n+1}, y_{n+1}^{(0)}, z_{n+1}^{(0)}) + \Delta g(x_n, y_n, z_n) \\ - \Delta g(x_{n-1}, y_{n-1}, z_{n-1}) + \Delta g(x_{n-2}, y_{n-2}, z_{n-2})] \quad (84.8)$$

در بخش ۶.۸، برای حل يك معادلهٔ دیفرانسیل، از زیر برنامه‌ای به نام DVERK از مجموعهٔ برنامه‌های IMSL استفاده کردیم. در اینجا، از این زیر برنامه برای حل يك دستگاه معادلات دیفرانسیل مرتبهٔ اول استفاده می‌کنیم. در زیر برنامهٔ DVERK، معرف متغیر مستقل است، در حالی که  $K = 1, \dots, N$ ،  $Y(K)$  برای تعیین برداری از متغیرهای غیر مستقل به طول  $N$  به کار می‌رود و فرض بر این است که  $N$  معادلهٔ مرتبهٔ اول به شکل (۸۱.۸)

وجود دارد.  $YPRIME(K)$  با شرط  $K = 1, 2, \dots, N$ ، برای تعیین بردار تابعهای  $f_1, \dots, f_N$  در سمت راست (۸۱۰۸) به کار می رود. زیر برنامه FCN برای تعریف  $YPRIME(K)$  به کار رفته است. استفاده از زیر برنامه DVERK برای یک دستگاه معادلات، از جهاتی مشابه با به کار بردن آن برای یک معادله تنهاست. مثال زیر این مورد استفاده را روشن می سازد.

□ مثال ۹۰۸: دستگاه معادلات زیر را به صورت یک دستگاه معادلات مرتبه اول بیان کنید و آن را با استفاده از زیر برنامه DVERK از  $x=0$  تا  $x=1$  حل نمایید:

$$\frac{d^2 z}{dx^2} = z^2 - y + e^x$$

$$\frac{d^2 y}{dx^2} = z - y^2 - e^x$$

$$z(0) = z'(0) = 0 \quad y(0) = 1 \quad y'(0) = -2 \quad (85.8)$$

در این مثال  $x$  متغیر مستقل است، در حالی که  $z(x)$  و  $y(x)$  متغیرهای وابسته می باشند. برای آنکه این دستگاه را به صورت یک دستگاه مرتبه اول بیان کنیم، قرار می دهیم  $z(x) = y_1(x)$  و  $y(x) = y_2(x)$  در این صورت، این دستگاه مرتبه اول همراه با شرایط اولیه اش به صورت زیر درمی آید

$$y_1'(x) = y_2(x) \quad y_1(0) = 0.0$$

$$y_2'(x) = y_4(x) \quad y_2(0) = 1.0$$

$$y_3'(x) = y_1^2(x) - y_2(x) + e^x \quad y_3(0) = 0.0$$

$$y_4'(x) = y_1(x) - y_2^2(x) - e^x \quad y_4(0) = -2.0 \quad (86.8)$$

برنامه فورترن و نتایج جزئی آن در زیر داده شده اند. مقادیر به دست آمده حداقل تا هشت رقم با معنا صحیح اند. به نظر می رسد که برای دستیابی به این دقت به ازای هر خروجی، ۱۶ بار محاسبه تابع مورد نیاز باشد. این امر مستلزم این است که یک مرحله داخلی یا نمونه تقریبی  $h = 0.05$  به کار گرفته شود.

```
C PROGRAM TO SOLVE EXAMPLE 8.9 USING DVERK (IMSL) .
  INTEGER IER,IND,K,N,NW
  REAL C(24),TOL,W(5,9),X,XEND,Y(4)
  DATA N, X, Y, TOL,IND,NW
  / 4, 0., 0., 1., 0., -2., 1.E-9, 1, 5 /
  EXTERNAL FCN2
  DO 12 K=1,10
    XEND = FLOAT(K)/10.
    CALL DVERK ( N, FCN2, X, Y, XEND, TOL, IND, C, NW, W, IER )
    PRINT 600, XEND,Y(1),Y(2),C(24)
```

```

600   FORMAT(3X,F3.1,3X,2(E16.8,3X),F4.0)
12   CONTINUE
                                STOP
END
SUBROUTINE FCN2 ( N, X, Y, YPRIME )
INTEGER N
REAL X, Y(N), YPRIME(N)
YPRIME(1) = Y(3)
YPRIME(2) = Y(4)
YPRIME(3) = Y(1)**2 - Y(2) + EXP(X)
YPRIME(4) = Y(1) - Y(2)**2 - EXP(X)
                                RETURN
END

```

### نتایج کامپیوتری برای مثال ۹.۸

X	Y(1)	Y(2)	FCN EVALS
.1	5.12342280E - 04	7.90476884E - 01	16
.2	4.19528369E - 03	5.63595308E - 01	32
.3	1.44796017E - 02	3.21283135E - 01	48
.4	3.50756908E - 02	6.44861308E - 02	64
.5	6.99842327E - 02	-2.07035152E - 01	80
.6	1.23532042E - 01	-4.94906488E - 01	96
.7	2.00446026E - 01	-8.02372169E - 01	112
.8	3.05983760E - 01	-1.13460479E + 00	128
.9	4.46147292E - 01	-1.49915828E + 00	144
1.0	6.28019076E - 01	-1.90666076E + 00	160

□

همان گونه که این مثال نشان می دهد، DVERK زیر بر نامه ای است که به کار بردن خجیلی ساده است و هنگامی که به دقت بالایی نیاز باشد، کارایی آن خجیلی زیاد است.

### تمرین

۱-۱۲.۸ معادله مرتبه دوم

$$y''(x) = 2(e^{2x} - y^2)^{1/2}$$

$$y(0) = 0 \quad y'(0) = 1$$

را به صورت دستگامی از معادلات مرتبه اول بنویسید و با استفاده از روش کلاسیک مرتبه چهارم رونگه-کوتا به ازای نمره های ثابت  $h = 1/64$  و  $h = 1/128$  آن را از  $x = 0$  تا  $x = 1$  حل نمایید. دقت نتایج حاصل را برآورد کنید.

۲-۱۲.۸ معادله مرتبه دوم زیر را با استفاده از فرمولهای ادمز-مولتن (۸۴.۸)، به ازای نمره ثابت  $h = 0.1$ ، از  $x = 1$  تا  $x = 2$  حل کنید:

$$y''(x) = 2y^3$$

$$y(1) = 1 \quad y'(1) = -1$$

چهار مقدار آغازین برای  $y(x)$  و  $f(x, y) = 2y^3$  لازم است. این مقادیر را از جواب دقیق  $y(x) = 1/x$  به دست آورده سپس نتایج خود را با جواب دقیق مقایسه کنید.

۳-۱۲.۸ بررسی کنید که آیا مرکز محاسباتی شما برنامه‌های IMSL را آّبونه شده است یا نه. اگر چنین است، معادلهٔ تمرین ۱-۱۲.۸ با استفاده از DVERK و به‌ازای مقادیر  $K/10$ ،  $K/100$ ،  $K/1000$ ، حل نمایید.

### ۱۳.۸° معادلات دیفرانسیل سخت‌۱

در عده‌ای از زمینه‌های مهم، از جمله واکنشهای شیمیایی، سیستمهای کنترل و شبکه‌های الکترونیکی، کار بردها به‌دستگاههایی از معادلات دیفرانسیل منجر می‌شوند که حل آنها دشوار است، زیرا فرایندهای مختلف در دستگاه، با مقیاسهای زمانی کاملاً متفاوتی رفتار می‌کنند. مثلاً اگر جواب يك معادلهٔ دیفرانسیل به‌صورت  $y(x) = C_1 e^{-x} + C_2 e^{-1000x}$  داده شده باشد، مؤلفهٔ دوم این جواب در مقایسه با مؤلفهٔ اول آن با افزایش  $x$  بسیار سریعتر تحلیل می‌رود. بیشتر روشهایی که برای حل معادلات دیفرانسیل بیان کرده‌ایم، وقتی برای حل مسائلی با جوابهایی از این نوع به‌کار روند، ناپایداری زیادی از خود نشان می‌دهند. مسائلی را که مؤلفه‌های جواب آنها مقیاسهای زمانی سخت متفاوتی دارند مسائل سخت نامند.

برای مثال، معادلهٔ مرتبهٔ دوم

$$\frac{d^2 y}{dx^2} + 1001 \frac{dy}{dx} + 1000y = 0 \quad (۸۷.۸)$$

را در نظر می‌گیریم. جواب عمومی معادلهٔ (۸۷.۸) برابر است با

$$y(x) = Ae^{-x} + Be^{-1000x}$$

اگر شرایط اولیه را  $y(0) = 1$  و  $y'(0) = -1$  بگیریم، جواب دقیق خواهد شد:

$$y(x) = e^{-x}$$

اکنون با استفاده از روش مرتبهٔ چهارم رونگه-کوتسا، سعی می‌کنیم معادلهٔ (۸۷.۸) را به‌ازای این شرایط اولیه حل کنیم. وقتی دستگاه فوق را به‌صورت يك دستگاه مرتبهٔ اول (بخش ۱۲.۸ را ببینید) بنویسیم، چنین می‌شود

$$\begin{aligned} \frac{dy_1}{dx} &= y_2 & y_1(0) &= 1 \\ \frac{dy_2}{dx} &= -1001y_2 - 1000y_1 & y_2(0) &= -1 \end{aligned} \quad (۸۸.۸)$$

به ازای نمره‌های  $0.0002 < h$ ، روش رونگه-کو تا جوابهایی به دست می‌دهد که جواب  $e^{-x}$  را به خوبی تقریب می‌زند. ولی معنی  $h = 0.0002$  این است که در هر بازه، ۵۰۰ مرحله انتگرالگیری می‌باید وجود داشته باشد. از آنجایی که جواب مطلوب برابر است با  $y(x) = e^{-x}$  به نظر می‌رسد که انتخاب خیلی بزرگتر نمو  $h$  مانعی نداشته باشد. ولی اگر  $h = 0.0003$ ، این مقدار باز خیلی کوچک است، جواب عددی به سمت  $\infty$  جهش می‌یابد. توضیح این رفتار به عواملی که برای شرایط پایداری روش اختیار شده بستگی دارد. برای روش رونگه-کو تا، ناحیه پایداری چنان است که می‌باید داشته باشیم

$$1000h < 2.78$$

یا  $h < 0.00278$ . یعنی برای مسئله بالا، نمو  $h$  به علت سرعت تغییر مؤلفه جواب، یعنی  $e^{-1000x}$ ، به دلایل پایداری محدود می‌شود. روش ادمز-مولتن و سایر روشهای چندمرحله‌ای استاندارد نیز به همین صورت نمو  $h$  را محدود می‌کنند.

بررسی گسترده برای یافتن روشهای مناسب در حل معادلات دیفرانسیل سخت هنوز ادامه دارد. موفقترین روشها ظاهراً روشهای ضمنی هستند. برای مثال، کاربرد روش ذوزنقه‌ای (۵۲.۸) تاحدی با موفقیت همراه بوده است. برای این روش تمام نیمصفحه منفی ناحیه پایداری است، لذا  $h$  به شرایط پایداری محدود نمی‌شود (کتاب گیر [۳۰] را ببینید). با اعمال این روش برای يك دستگاه دو معادله دو مجهولی به شکل

$$y_1' = f_1(x, y_1, y_2)$$

$$y_2' = f_2(x, y_1, y_2)$$

روش ذوزنقه‌ای چنین خواهد شد

$$y_{1, n+1} = y_{1n} + \frac{h}{\gamma} \{ f_1(x_n, y_{1n}, y_{2n}) + f_1(x_{n+1}, y_{1, n+1}, y_{2, n+1}) \}$$

$$y_{2, n+1} = y_{2n} + \frac{h}{\gamma} \{ f_2(x_n, y_{1n}, y_{2n}) + f_2(x_{n+1}, y_{1, n+1}, y_{2, n+1}) \}$$

(۸۹.۸)

با تخصیص این معادلات به دستگاه خطی (۸۸.۸)، خواهیم داشت

$$y_{1, n+1} = y_{1n} + \frac{h}{\gamma} \{ y_{2n} + y_{2, n+1} \}$$

## 1. explodes

\* منظور آن است که با انتخاب  $h = 0.0003$ ، نا پایداری عددی ایجاد می‌شود یعنی مقادیر متغیرهای غیرمستقل مرتباً افزایش می‌یابند. م.

$$y_{2,n+1} = y_{2n} + \frac{h}{4} \{-1001y_{2n} - 1000y_{1n} \\ + (-1001y_{2,n+1} - 1000y_{1,n+1})\}$$

در حالت معمولی، این معادلات با روش بارستی حل می‌شوند، اما به علت خطی بودن معادلات، يك دستگاه ضمنی بر حسب مجهولات  $y_{1,n+1}$  و  $y_{2,n+1}$  به صورت زیر می‌توانیم به دست آوریم:

$$y_{1,n+1} - \frac{h}{4} y_{2,n+1} = y_{1n} + \frac{h}{4} y_{2n} \\ \frac{h}{4} (1000)y_{1,n+1} + \left(1 + \frac{1001h}{4}\right) y_{2,n+1} = y_{2n} - \frac{h}{4} \{1001y_{2n} \\ + 1000y_{1n}\} \quad (90.8)$$

اکنون  $h$  را برابر با ۱٫۰ انتخاب می‌کنیم، لذا دستگاه (۹۰٫۸) چنین خواهد شد

$$y_{1,n+1} - 0.25y_{2,n+1} = y_{1n} + 0.25y_{2n} \\ 50y_{1,n+1} + 51.25y_{2,n+1} = -49.25y_{2n} - 50y_{1n} \quad (91.8)$$

به ازای  $n=0$  داریم،  $y_{10}=1$ ،  $y_{20}=-1$  و با توجه به دستگاه (۹۱٫۸) خواهیم داشت

$$y_{11} = 0.904762 \quad y_{21} = -0.904762$$

که با رعایت نمو بزرگی که به کار رفته، تقریب مناسبی برای جواب دقیق

$$y(0.1) = e^{-0.1} = 0.904837$$

است. بعد از ۱۰ مرحله، به ازای  $h=0.1$  خواهیم داشت  $y_1(1.0) \approx 0.367573$  که در مقایسه با جواب دقیق  $y(1.0) = e^{-1.0} = 0.367879$  جواب مطلوبی است. ولی در استفاده از روش دوزنقه‌ای برای مسائل غیر خطی سخت، يك اصلاح اساسی وجود دارد که می‌باید به کار رود. برای يك معادلهٔ تنهای  $y' = f(x, y)$ ، روش دوزنقه‌ای ضمنی است و به وسیلهٔ رابطهٔ زیر تعریف می‌شود

$$y_{n+1} = y_n + \frac{h}{4} \{f(x_n, y_n) + f(x_{n+1}, y_{n+1})\} \quad (92.8)$$

این معادله، معادله‌ای است ضمنی با  $n$  ثابت که باید با يك روش بارستی نسبت به  $y_{n+1}$



حل شود. در حالت معمولی، بارست نقطه ثابت که به صورت زیر تعریف می شود، مورد استفاده قرار می گیرد

$$y_{n+1}^{(m+1)} = y_n + \frac{h}{\gamma} \{f(x_n, y_n) + f(x_{n+1}, y_{n+1}^{(m)})\} \quad m = 0, 1, \dots$$

که در آن  $y_{n+1}^{(0)}$  تقریبی است برای  $y_{n+1}$  که با روشی غیر از روشی نظیر اوپلر به دست می آید. همان گونه که در قسمت ۸.۸ نشان داده شد، اگر  $\left| \frac{h}{\gamma} \frac{\partial f}{\partial y} \right| < 1$ ، این بارست نقطه ثابت همگرا خواهد شد و چون در مسائل سخت،  $|\partial f / \partial y|$  بسیار بزرگ است، این روش بارستی برای همگرایی به نمودهای خیلی کوچک نیاز پیدا می کند. ولی می توانیم معادله (۹۲.۸) را با روش بارستی نیوتن نسبت به  $y_{n+1}$  به صورت زیر حل کنیم: قرار می دهیم  $\bar{y} = y_{n+1}$ ، و معادله (۹۲.۸) را دوباره به شکل

$$F(\bar{y}) = \bar{y} - y_n - \frac{h}{\gamma} f(x_n, y_n) - \frac{h}{\gamma} f(x_{n+1}, \bar{y}) = 0 \quad (93.8)$$

می نویسیم. اگر  $\bar{y}^{(0)} = y_{n+1}^{(0)}$  تقریب اولیه ای برای  $\bar{y}$  باشد، آنگاه تقریبهایی متوالی طبق روش نیوتن و از راه بارستی

$$\bar{y}^{(m+1)} = \bar{y}^{(m)} + \frac{F(\bar{y}^{(m)})}{F'(\bar{y}^{(m)})} \quad m = 0, 1, \dots$$

به دست می آید، که با توجه به فرمول (۹۳.۸) نتیجه می شود:

$$F'(\bar{y}) = 1 - \frac{h}{\gamma} \frac{\partial f}{\partial \bar{y}}(x_{n+1}, \bar{y})$$

در این حالت، وقتی  $\partial f / \partial \bar{y}$  بزرگ و منفی باشد، که ویژگی عادی مسائل سخت است، در رابطه با همگرایی مشکلی وجود نخواهد داشت. ولی روش نیوتن برای يك معادله تنها به محاسبه  $\partial f / \partial \bar{y}$  نیاز دارد، و برای يك دستگاه معادلات به محاسبه عناصر ماتریس ژاکوبی. بنا بر این می توان انتظار داشت که زیر برنامه های لازم برای حل معادلات دیفرانسیل سخت تاحدی پیچیده باشند.

برای يك دستگاه معادلات خطی به صورت

$$y' = Ay$$

که در آن  $A$  ماتریس ثابتی است، سختی مسئله به توسط ویژه مقادیرهای ماتریس  $A$  معین می شود. اگر ویژه مقادیرهای ماتریس  $A$  منفی و خیلی از هم فاصله داشته باشند، آنگاه دستگاه سخت بوده و می توان انتظار داشت که حل آن باروشهای معمولی مشکل باشد. برای مثال دستگاه (۸۸.۸)، ماتریس  $A$  به صورت

$$A = \begin{bmatrix} 0 & 1 \\ -1000 & -1001 \end{bmatrix}$$

است و ویژه مقادیرهای آن برابر با  $-1$  و  $-1000$  است.  
برای دستگاههای غیرخطی کلیتر به شکل

$$y' = f(x, y)$$

سختی دستگاه با ویژه مقادیرهای ماتریس ژاکوبی

$$\left(\frac{\partial f}{\partial y}\right)$$

معین می شود. برای بحث کاملتر در مورد مسائل سخت و حل آنها به روشهای دیگر، خوانند  
به کتاب گیر [۳۰] ارجاع داده می شود.

### تمرین

۱-۱۳۰۸ سعی کنید با استفاده از روش مرتبه ۴ رونگه-کوتا به ازای نموداری

$$h = 0.0001, 0.0002, 0.0003, 0.001$$

دستگاه (۸۸.۸) را از  $x = 0$  تا  $x = 2$  حل کنید. تحقیق کنید که به ازای  $h = 0.0003$  و  $h = 0.001$ ، جواب جهش می یابد درحالی که به ازای  $h = 0.0002$  و  $h = 0.0001$  جواب دقیق  $y = e^{-x}$  تقریب مطلوبی به دست می آید.

۲-۱۳۰۸ برای دستگاه

$$y_1' = y_2$$

$$y_2' = -200y_1 - 102y_2$$

نشان دهید که ویژه مقادیرهای ماتریس ضرایب  $-2$  و  $-100$  هستند و از این رو جواب  
عمومی به صورت زیر است

$$y_1(x) = y_2(x) = Ae^{-2x} + Be^{-100x}$$

در شرایط  $y(0) = 1$  و  $y'(0) = -2$  که متناظر با  $y_1(0) = 1$  و  $y_2(0) = -2$  هستند،  
جواب دقیق برابر است با  $y(x) = e^{-2x}$ . با استفاده از روش ذوزنقه ای و به ازای  $h = 0.001$   
این دستگاه را از  $x = 0$  تا  $x = 1$  حل و نتایج خود را با جواب دقیق مقایسه کنید.

## ۹

مسائل با مقدار مرزی<sup>۱</sup> در معادلات دیفرانسیل معمولی

در فصل هشتم روشهای عددی برای حل مسائل با مقدار اولیه را بررسی کردیم. در این گونه مسائل کلیه شرایط اولیه در یک نقطه داده شده‌اند. در این فصل مسائلی را بررسی می‌کنیم که در آنها، شرایط در بیش از یک نقطه مشخص شده‌اند. یک مسئله ساده از مسئله با مقدار مرزی مرتبه دوم به صورت

$$y''(x) = y(x) \quad y(0) = 0 \quad y(1) = 1 \quad (1.9)$$

و یک مثال از مسئله با مقدار مرزی مرتبه چهارم به صورت

$$y^{(4)}(x) + k y(x) = q \quad (2.9 \text{ الف})$$

$$y(0) = y'(0) = 0 \quad (2.9 \text{ ب})$$

$$y(L) = y''(L) = 0 \quad (2.9 \text{ پ})$$

است. در اینجا  $L$  می‌تواند معرف انتهای میله‌ای به طول  $L$  باشد که بار یکنواخت  $q$  بر آن وارد شده‌است. شرط (۲.۹ ب) گویای آن است که انتهای میله در  $x = 0$  تثیب شده‌است و در حالی که (۲.۹ پ) حکایت از این دارد که میله در  $x = L$  فقط تکیه داده شده‌است.

برای حل مسائلی از این قبیل سه روش را بررسی می‌کنیم: روش تفاضلهای منتهای ۱ و روشهای «نشانه‌گیری» که از روشهای فصل هشتم اقتباس شده‌اند و روش پهلوی هم‌گذاری ۲.

## ۱.۹ روشهای تفاضل منتهای

فرض می‌کنیم که یک معادلهٔ دیفرانسیل خطی از مرتبهٔ بالاتر از یک با شرایطی که در دوسر بازهٔ  $[a, b]$  مشخص شده، داده شده است. بازهٔ  $[a, b]$  را به  $N$  قسمت مساوی با پهنای  $h$  تقسیم می‌کنیم. قرار می‌دهیم  $x_0 = a$  و  $x_N = b$  و نقاط داخلی شبکهٔ ۴ را به صورت زیر تعریف می‌کنیم

$$x_n = x_0 + nh \quad n = 1, 2, \dots, N-1$$

مقادیر متناظر  $y$  در این نقاط شبکه، به صورت زیر نشان داده می‌شوند:

$$y_n = y(x_0 + nh) \quad n = 0, 1, \dots, N$$

گاهی با نقاطی بیرون از بازهٔ  $[a, b]$  سروکار خواهیم داشت. این نقاط، نقاط خارجی شبکه نامیده می‌شوند. نقاطی که درست چپ  $x_0$  قرار دارند با علامت  $-h$  علامت  $x_{-1} = x_0 - h$  و  $x_{-2} = x_0 - 2h$  و غیره نشان داده می‌شوند و آنهایی که درست راست  $x_N$  قرار دارند با  $x_{N+1} = x_N + h$ ،  $x_{N+2} = x_N + 2h$  و غیره. بدیهی است که مقادیر متناظر  $y$  در نقاط خارجی شبکه با  $y_{N+1}$ ،  $y_{N+2}$ ،  $y_{-1}$ ،  $y_{-2}$  و غیره نشان داده می‌شوند.

برای حل یک مسئلهٔ با مقدار مرزی با روش تفاضلهای منتهای، به جای همهٔ مشتقهایی که در معادله و نیز در شرایط مرزی ظاهر می‌شوند، تقریبی مناسب از تفاضلهای گذشته می‌شود. معمولاً تفاضلهای مرکزی، به علت دقت بیشتری که عاید می‌سازند، ترجیح داده می‌شوند. برخی نمونه‌های تقریبی تفاضل مرکزی در زیر داده شده‌اند (فصل ۷ را ببینید):

$$y'(x_n) \approx \frac{y_{n+1} - y_{n-1}}{2h}$$

$$y''(x_n) \approx \frac{y_{n+1} - 2y_n + y_{n-1}}{h^2}$$

$$y^{(4)}(x_n) \approx \frac{y_{n+2} - 4y_{n+1} + 6y_n - 4y_{n-1} + y_{n-2}}{h^4} \quad (3.9)$$

در هر مورد نمایش تفاضلهای محدود، یک تقریب  $\theta(h^2)$  برای مشتق مربوطه است. برای

- |                         |                         |                |
|-------------------------|-------------------------|----------------|
| 1. finite differences   | 2. shooting             | 3. collocation |
| 4. interior mesh points | 5. exterior mesh points |                |

نشان دادن این روش، معادلهٔ دیفرانسیل خطی مرتبهٔ ۲

$$y''(x) + f(x)y' + g(x)y = q(x) \quad (۲.۹)$$

را در شرایط مرزی

$$y(x_0) = \alpha \quad (۵.۹)$$

$$y(x_N) = \beta \quad (۶.۹)$$

در نظر می‌گیریم. تقریب تفاضل متناهی برای معادلهٔ (۲.۹) به صورت زیر است

$$\frac{y_{n-1} - 2y_n + y_{n+1}}{h^2} + \frac{f(x_n)(y_{n+1} - y_{n-1})}{2h} + g(x_n)y_n = q(x_n)$$

$$n = 1, 2, \dots, N-1$$

از ضرب طرفین معادله در  $h^2$  و قراردادن  $f(x_n) = f_n$  به ازای  $n = 1, \dots, N-1$  گروه‌بندی جملات، خواهیم داشت:

$$\left(1 - \frac{h}{\gamma} f_n\right) y_{n-1} + (-2 + h^2 g_n) y_n + \left(1 + \frac{h}{\gamma} f_n\right) y_{n+1} = h^2 q_n$$

$$n = 1, 2, \dots, N-1 \quad (۷.۹)$$

از آنجایی که  $y_0$  و  $y_N$  با شرایط (۵.۹) و (۶.۹) مشخص شده‌اند، بنابراین معادلهٔ (۷.۹) یک دستگاه خطی است با  $N-1$  معادله و  $N-1$  مجهول  $y_n$  ( $n = 1, 2, \dots, N-1$ ). از قراردادن  $\alpha$  به جای  $y_0$  و  $\beta$  به جای  $y_N$  در معادلهٔ (۷.۹)، دستگاه معادلات فوق به صورت زیرین درمی‌آید

$$(-2 + h^2 g_1) y_1 + \left(1 + \frac{h}{\gamma} f_1\right) y_2 = h^2 q_1 - \left(1 - \frac{h}{\gamma} f_1\right) \alpha$$

$$\left(1 - \frac{h}{\gamma} f_2\right) y_1 + (-2 + h^2 g_2) y_2 + \left(1 + \frac{h}{\gamma} f_2\right) y_3 = h^2 q_2$$

$$\left(1 - \frac{h}{\gamma} f_3\right) y_2 + (-2 + h^2 g_3) y_3 + \left(1 + \frac{h}{\gamma} f_3\right) y_4 = h^2 q_3$$

.....

$$\left(1 - \frac{h}{\gamma} f_{N-2}\right) y_{N-2} + (-2 + h^2 g_{N-2}) y_{N-2} + \left(1 + \frac{h}{\gamma} f_{N-2}\right) y_{N-1}$$

$$= h^2 q_{N-2}$$



اکنون اگر معادله (۷.۹) را به ازای  $n = 1$  بنویسیم، سپس به جای  $y_0/(1 - \gamma h)$ ،  $y_0$  را قرار دهیم خواهیم داشت

$$\left[ (-2 + h^2 g_1) + \frac{1 - (h/2)f_1}{1 - \gamma h} \right] y_1 + \left( 1 + \frac{h}{2} f_1 \right) y_2 = h^2 q_1 \quad (10.9)$$

حال در دستگاه (۸.۹) به جای اولین معادله می توان معادله (۱۰.۹) را گذارد. تمامی معادلات دیگر از دستگاه (۸.۹) بدون تغییر باقی می ماند، و با استفاده از الگوریتم ۳.۴ می توان دستگاه حاصل را دوباره حل کرد. ولی باید توجه داشت که رابطه (۹.۹) تنها یک تقریب  $\theta(h)$  برای شرط مرزی (۹.۹) است (به قسمت ۱.۷، مراجعه کنید).  
در این صورت دقت جواب نیز از مرتبه  $h$  خواهد بود. برای به دست آوردن جوابی که در همه جا از مرتبه  $h^2$  باشد، به جای رابطه (۹.۹) تقریب

$$\frac{y(x_0 + h) - y(x_0 - h)}{2h} + \gamma y(x_0) = 0$$

یا مرتب شده آن

$$y_1 - y_{-1} + 2h\gamma y_0 = 0 \quad (11.9)$$

را قرار می دهیم. از آنجا که یک نقطه خارجی  $y_{-1}$  در رابطه وارد کرده ایم، اکنون باید  $y_0$  و نیز  $y_1, y_2, \dots, y_{N-1}$  را به عنوان مجهول در نظر بگیریم. اما چون  $N$  مجهول داریم می باید  $N$  معادله نیز داشته باشیم. با در نظر گرفتن  $n = 0$  در (۷.۹) می توان یک معادله اضافی به دست آورد. سپس اگر با استفاده از رابطه (۱۱.۹)،  $y_{-1}$  را حذف کنیم، برای دو معادله اول خواهیم داشت

$$\left[ 2h\gamma \left( 1 - \frac{h}{2} f_0 \right) + (-2 + h^2 g_0) \right] y_0 + 2y_1 = h^2 q_0 \quad n = 0$$

$$\left( 1 - \frac{h}{2} f_1 \right) y_0 + (-2 + h^2 g_1) y_1 + \left( 1 + \frac{h}{2} f_1 \right) y_2 = h^2 q_1 \quad n = 1$$

بقیه معادلات همان معادلاتی هستند که در دستگاه (۸.۹) ظاهر شده اند. این دستگاه باز سه قطری ولی حالا از مرتبه  $N$  است و می تواند صریحاً با استفاده از الگوریتم ۳.۴، دوباره حل شود.

به روشنی دیده می شود که دقت حاصله از روش تفاضلهای متناهی به کوچک بودن شبکه و مرتبه تقریب تفاضل محدود، بستگی خواهد داشت. هر قدر شبکه کوچکتر شود، به همان نسبت تعداد معادلاتی که باید حل شوند افزایش می یابند. در نتیجه مدت لازم برای کامپیوتر ممکن است زیاد شود و مشکل بتوان دقت خوبی به دست آورد. برای شبکه هایی

با همان اندازه، استفاده از تقریبهای مرتبه بالاتر، موجب دقت بیشتری می‌شود، اما به پیچیدگی قابل ملاحظه‌ای، بخصوص نزدیک به نقاط انتهایی بازه مورد نظر که مقادیر خارجی در آن نقاط نامعلوم هستند، منجر می‌شود.

در عمل بهتر است که دستگاه معادلات به ازای چندین مقدار مختلف  $h$  حل شوند. در این صورت مقایسه جوابها در نقاط واحدی از شبکه، دقت حاصل را نشان می‌دهند. بعلاوه از روند برونمایی که در بخش ۵.۷ شرح داده شد، می‌توان معمولاً برای کسب جوابهای بهبود یافته استفاده کرد. به موجب توافقی که برای حل دستگاه معادلات تفاضلهای منتهی شده است، در برونمایی به سمت حد به صورت زیر عمل می‌شود. گیریم  $y_h(x)$  معرف جواب تقریبی مسئله با مقدار مرزی بر روی تعداد  $N = (b-a)/h$  زیر تقسیم در بازه  $[a, b]$  در یکی از نقاط  $x$  از شبکه باشد. گیریم  $y_{h/2}(x)$  جواب تقریبی همان مسئله روی تعداد  $2N = (b-a)/(h/2)$  زیر تقسیم از بازه  $[a, b]$  باشد. اکنون در  $N-1$  نقطه  $x_1 = a+h, x_2 = a+2h, \dots, x_{N-1} = a+(N-1)h$ ، دو جواب تقریبی  $y_h(x_n)$  و  $y_{h/2}(x_n)$  را در دست داریم. با استفاده از روش برونمایی در این نقاط خواهیم داشت

$$y'(x_n) \approx y_n^{(1)} := \frac{2y_{h/2}(x_n) - y_h(x_n)}{h} \quad n = 1, 2, \dots, N-1$$

معمولاً برونمایی فوق، بهبودی قابل توجهی در تقریب ایجاد می‌کند.

□ مثال ۱.۹: با استفاده از روش تفاضلهای منتهی، مسئله با مقدار مرزی (۱.۹) را حل کنید.

هر گاه بگیریم  $f(x) = 0, g(x) = -1, q(x) = 0$ ، از قرارداد  $y_0 = 0$  و  $y_N = 1$  در دستگاه (۱.۹)، دستگاه زیر را به دست می‌آوریم

$$(-2-h^2)y_1 + y_2 = 0$$

$$y_{n-1} + (-2-h^2)y_n + y_{n+1} = 0 \quad n = 2, 3, \dots, N-2$$

$$y_{N-2} + (-2-h^2)y_{N-1} = -1$$

این دستگاهی است از  $N-1$  معادله با  $N-1$  مجهول:  $y_1, y_2, \dots, y_{N-1}$ . این دستگاه، با استفاده از زیر برنامه‌ای بر اساس الگوریتم ۳.۴ و به ازای مقادیر  $h = 0.1$  و  $h = 0.05$  روی کامپیوتر IBM 7090 حل شده بود. نتایج حاصله در زیر آمده است.



## نتایج کامپیوتری برای مثال ۱۰۹

XN	YN(H = 0.05)	YN(H = 0.10)	YN(1)	Y(XN)
0	0	0	0	0
0.05	.04256502			.04256363
0.10	.08523646	.08524469	.08523372	.08523369
0.15	.12812098			.12811689
0.20	.17132582	.17134184	.17132048	.17132045
0.25	.21495896			.21495239
0.30	.25912950	.25915240	.25912187	.25912183
0.35	.30394787			.30393920
0.40	.34952610	.34955449	.34951663	.34951659
0.45	.39597815			.39596794
0.50	.44342014	.44345213	.44340946	.44340942
0.55	.49197068			.49195965
0.60	.54175115	.54178427	.54174010	.54174004
0.65	.59288599			.59287506
0.70	.64550304	.64553425	.64549263	.64549258
0.75	.69973386			.69972418
0.80	.75571401	.75573958	.75570550	.75570543
0.85	.81358345			.81357635
0.90	.87348684	.87350228	.87348166	.87348163
0.95	.93557395			.93557107
1.00	1	1	1	1

ستون چهارم، مقادیر برونیایی شده‌ای را نشان می‌دهد که در بازه‌هایی به طول ۰٫۱ از فرمول

$$y_n^{(1)} = \frac{4y_{h/2}(x_n) - y_h(x_n)}{3}$$

به دست آمده‌اند. مقادیر داده شده در ستون آخر از جواب دقیق مسئله

$$y(x) = \frac{\sinh x}{\sinh 1}$$

به دست آمده‌اند. این نتایج نشان می‌دهند که به ازای  $h = 0.1$ ، جواب از سه تا چهار رقم بامعنی و به ازای  $h = 0.05$  از چهار تا پنج رقم بامعنی صحیح است. در حالی که جواب به دست آمده از برونیاسابی در حدود هفت رقم بامعنی صحیح دارد. برای به دست آوردن دقتی برابر با هفت رقم بامعنی، بدون استفاده از برونیایی لازم است که بازه  $[0, 1]$  تقریباً به یک شبکه ۱۰۰ نقطه‌ای ( $h = 0.01$ ) تقسیم شود. □

## تمرین

۱۰۹-۱ مسئله با مقدار مرزی

مسائل با مقدار مرزی در ... ۵۳۳

$$\frac{d^2 y}{dx^2} + y = 0 \quad y(0) = 0 \quad y(1) = 1$$

را از طریق روشهای تفاضلی حل کنید.  $h$  را برابر  $1/4$  بگیرید و دستگاه حاصل را بسا استفاده از ماشین حساب جیبی حل نمایید.

جواب:  $y_1 = 0.2943$ ,  $y_2 = 0.5702$ ,  $y_3 = 0.8104$

این جواب را با جواب دقیق  $y = (\sin x)/(\sin 1)$  مقایسه کنید.

۳-۱۰۹ بسا استفاده از يك شبکه  $h = 0.1$  مسئله با مقدار مرزی (۱۰۹) را وقتی شرط  $y'(0) + y(0) = 0$  جایگزین شرط  $y(0) = 0$  شده باشد حل کنید.

۳-۱۰۹ يك دستگاه  $\theta(h^2)$  تفاضل منتهای برای تقریب زدن جواب مسئله بسا مقدار مرزی زیر بنویسید

$$y'' + xy' + y = 2x \quad y(0) = 1 \quad y(1) = 0$$

$h$  را برابر  $0.1$  بگیرید و دستگاه را به شکل ماتریسی بنویسید و سپس با استفاده از يك برنامه کامپیوتری بر اساس الگوریتم ۳.۴، این دستگاه را حل کنید.

۴-۱۰۹ نشان دهید که می توان روش بارستی گاوس-زایدل را نیز برای حل دستگاه مثال ۱۰۹ به کار برد و این جواب حاصل از راه روش بارستی را با دقتی برابر با چهار رقم بامعنی به دست آورد. برای این مسئله، آیا روش مستقیم کاراتر از روش بارستی است؟

۵-۱۰۹ مسئله با مقدار مرزی

$$y'' + 2y' + y = x \quad y(0) = 0 \quad y(1) = 0$$

را به ازای  $h = 1/8$  و  $h = 1/16$  با استفاده از روشهای تفاضل حل کنید، و نتیجه را با استفاده از برونیایی بهبود بخشید.

## ۲.۹ روشهای نشانه گیری

در مسائل خطی با مقدار مرزی، از چندین روش می توان استفاده کرد. روش تفاضلهای که در بخش قبل بیان شد، در این گونه موارد روش مناسبی است. در روشهای دیگر سعی بر این است که جوابهای مستقل خطی برای معادلات دیفرانسیل به دست آورند و آنها را طوری ترکیب کنند که در شرایط مرزی صدق نمایند. برای معادلات غیر خطی روش اخیر رانمی توان به کار برد. روش تفاضلهای را می توان برای مسائل غیر خطی تطبیق داد، اما این امر مستلزم تخمین زدن جوابهای آزمایشی و سپس بهتر ساختن آنها به وسیله روندهای بارستی است. در این کار علاوه بر آنکه برنامه نویسی پیچیده ای لازم است، تضمینی هم برای همگرایی بارستها

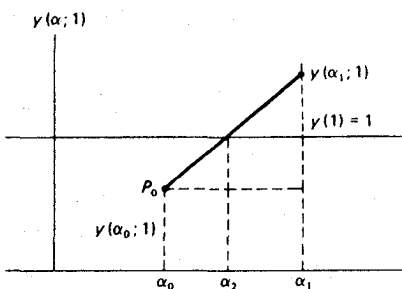
وجود ندارد. روش نشانه‌گیری که در این بخش توضیح داده می‌شود هم برای مسائل خطی و هم برای مسائل غیرخطی قابل استفاده است. در اینجا نیز هیچ تضمینی برای همگرایی وجود ندارد، اما استفاده از این روش بسیار ساده است و زمانی که این روش همگرا باشد، معمولاً کاراتر از روشهای دیگر است.

باردیگر مسئله‌ای را که در (۱۰۹) داده شده بود در نظر می‌گیریم. می‌خواهیم روشهای با مقدار اولیه را که در فصل ۸ مورد بحث قرار گرفته بود به کار ببریم، اما برای انجام این کار می‌باید هم مقدار  $y(0)$  در دست باشد و هم مقدار  $y'(0)$ . از آنجایی که  $y'(0)$  قبلاً معین شده است، آن را به صورت پارامتر مجهول، مثلاً  $\alpha$ ، که باید چنان تعیین شود که جواب حاصل، مقدار  $y(1)$  را با دقت مطلوبی به دست دهد، در نظر می‌گیریم. بنا بر این شیب اولیه را تخمین می‌زنیم و یک روند بارستی چنان بنا می‌کنیم که به سمت شیب دقیق همگرا شود. بگیریم  $\alpha_0$  و  $\alpha_1$  دو تخمین برای شیب اولیه  $y'(0)$  باشند، و فرض می‌کنیم  $y(\alpha_0; 1)$  و  $y(\alpha_1; 1)$  مقادیری از  $y$  باشند که از انتگرال‌گیری معادلهٔ دیفرانسیل در  $x = 1$  به دست آمده‌اند. از لحاظ نموداری می‌توان این وضعیت را به‌طور ترسیمی با شکل‌های ۱۰۹ و ۲۰۹ نمایش داد.

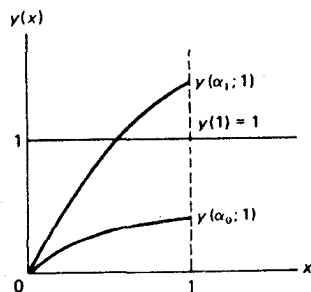
در شکل ۱۰۹، جوابهای با مسائل مقدار اولیه ترسیم شده‌اند، در حالی که در شکل ۲۰۹،  $y(\alpha; 1)$  به صورت تابعی از  $\alpha$  ترسیم شده است. اما یک تقریب معمولاً بهتر برای  $\alpha$  را می‌توان از راه درونیابی خطی به دست آورد. محل تقاطع خط واصل بین نقاط  $P_0$  و  $P_1$  با خط  $y(1) = 1$  مختص  $\alpha$  می‌دارد که با رابطه

$$\alpha_2 = \alpha_0 + (\alpha_1 - \alpha_0) \frac{y(1) - y(\alpha_0; 1)}{y(\alpha_1; 1) - y(\alpha_0; 1)} \quad (12.9)$$

معین می‌شود. حال با استفاده از مقادیر اولیه  $y(0) = 0$  و  $y'(0) = \alpha_2$ ، از معادلهٔ دیفرانسیل انتگرال می‌گیریم تا  $y(\alpha_2; 1)$  به دست آید. باردیگر با استفاده از درونیابی خطی بر پایهٔ



شکل ۲۰۹



شکل ۱۰۹

$\alpha_1$  و  $\alpha_2$ ، تقریب دیگر  $\alpha_3$  را می‌توانیم به‌دست آوریم. این روند تا زمانی که همگرایی به‌دست آید، یعنی تاجایی که  $y(\alpha_i; 1) = 1$  یا  $y(1) = 1$  تا تعداد ارقام مطلوب اعشاری تطابق نماید، ادامه می‌یابد. هیچ تضمینی برای همگرایی این روند با رستی وجود ندارد. واضح است که سرعت همگرایی به این امر که تا چه اندازه برآوردهای اولیه خوب بوده بستگی دارد. این برآوردها گاهی با بررسیهای فیزیکی مسئله و زمانی با نمایش ترمیمی ساده جواب، به‌دست می‌آیند.

روند بالا برای يك مسئله کلی مرتبه دوم با مقدار مرزی

$$y'' = f(x, y, y') \quad y(0) = y_0 \quad y(b) = y_b \quad (13.9)$$

در الگوریتم ۱۰.۹ خلاصه شده است.

**الگوریتم ۱۰.۹:** روش نشانه‌گیری برای مسائل با مقدار مرزی مرتبه دوم

۱. گیریم  $\alpha_k$  تقریبی برای شیب اولیه مجهول  $\alpha = y'(0)$  باشد. (دو مقدار اولیه  $\alpha_0$  و  $\alpha_1$  را با استفاده از شهود فیزیکی انتخاب می‌کنیم).
۲. با استفاده از یکی از روشهای فصل ۸، مسئله با مقدار اولیه

$$y'' = f(x, y, y') \quad y(0) = y_0 \quad y'(0) = \alpha_k$$

- از  $x = 0$  تا  $x = b$  حل می‌کنیم. جواب در نقطه  $x = b$  را  $y(\alpha_k; b)$  می‌نامیم.
۳. تقریب بعدی را از درونبایی خطی زیر به‌دست می‌آوریم

$$\alpha_{k+1} = \alpha_{k-1} + (\alpha_k - \alpha_{k-1}) \frac{y_b - y(\alpha_{k-1}; b)}{y(\alpha_k; b) - y(\alpha_{k-1}; b)}$$

$$k = 1, 2, \dots$$

۴. مراحل ۲ و ۳ را تا برقراری نامساوی  $|y(\alpha_k; b) - y_b| < \varepsilon$  به‌ازای يك  $\varepsilon$  که قبلاً معین شده، ادامه می‌دهیم.

روش با رستی که در الگوریتم ۱۰.۹ به‌کار رفته کار بردی است از روش خط قاطع که در فصل ۳ توضیح داده شده بود.

این روش در مورد دستگاه معادلات مرتبه بالاتر پیچیده‌تر و به‌دست آوردن همگرایی دشوارتر می‌شود. وضعیت کلی برای دستگاه غیر خطی را می‌توان به‌صورت زیر نمایش داد. يك دستگاه چهار معادله چهار مجهولی:

$$x' = f(x, y, z, w, t)$$

$$y' = g(x, y, z, w, t)$$

$$z' = h(x, y, z, w, t)$$

$$w' = l(x, y, z, w, t)$$

(۱۴.۰۹)

را که  $t$  در آن معرف متغیر مستقلی است، در نظر می‌گیریم. در  $t=0$  دوشروط، مثلاً

$$x(0) = x_0$$

$$y(0) = y_0$$

و در  $t=T$ ، دوشروط، مثلاً

$$z(T) = z_T$$

$$w(T) = w_T$$

داده شده‌اند. گیریم  $z(0) = \alpha$  و  $w(0) = \beta$  مقادیر اولیه و دقیق برای  $z(0)$  و  $w(0)$  باشند، و فرض می‌کنیم  $\alpha_0$  و  $\beta_0$  تخمینهایی برای مقادیر اولیه باشند. حال از دستگاه (۱۴.۹) انتگرال می‌گیریم و مقادیر  $z$  و  $w$  حاصل در  $t=T$  را با  $z(\alpha_0, \beta_0; T)$  و  $w(\alpha_0, \beta_0; T)$  نشان می‌دهیم.

چون  $z$  و  $w$  در  $t=T$  آشکارا تساوی از  $\alpha$  و  $\beta$  هستند، می‌توانیم  $z(\alpha, \beta; T)$  و  $w(\alpha, \beta; T)$  را نسبت به این دو متغیر با جملات خطی، به‌سری تیلر بسط دهیم

$$z(\alpha, \beta; T) = z(\alpha_0, \beta_0; T) + (\alpha - \alpha_0) \frac{\partial z}{\partial \alpha}(\alpha_0, \beta_0; T)$$

$$+ (\beta - \beta_0) \frac{\partial z}{\partial \beta}(\alpha_0, \beta_0; T)$$

$$w(\alpha, \beta; T) = w(\alpha_0, \beta_0; T) + (\alpha - \alpha_0) \frac{\partial w}{\partial \alpha}(\alpha_0, \beta_0; T)$$

$$+ (\beta - \beta_0) \frac{\partial w}{\partial \beta}(\alpha_0, \beta_0; T) \quad (15.9)$$

می‌توانیم  $z(\alpha, \beta; T)$  و  $w(\alpha, \beta; T)$  را برابر مقادیر مطلوبشان  $z_T$  و  $w_T$  قرار دهیم، اما پیش از اینکه بتوانیم (۱۵.۹) را نسبت به مقادیر تصحیحی  $\alpha - \alpha_0$  و  $\beta - \beta_0$  حل کنیم، باید مشتقات جزئی در (۱۵.۹) را به‌دست آوریم. از آنجایی که جوابهای  $z$  و  $w$  معلوم نیستند، لذا نمی‌توانیم به‌طور تحلیلی، این مشتقات جزئی را به‌دست آوریم. ولی، می‌توانیم مقادیر عددی تقریبی آنها را معین کنیم. برای این کار دستگاه (۱۴.۹) را یکبار با شرایط اولیه  $x_0, y_0, \alpha_0, \beta_0$ ، و یکبار با شرایط  $x_0, y_0, \alpha_0 + \Delta\alpha_0, \beta_0$ ، و سپس با شرایط  $x_0, y_0, \alpha_0, \beta_0 + \Delta\beta_0$  که در آن  $\Delta\alpha_0$  و  $\Delta\beta_0$  نمو‌های کوچکی هستند حل می‌کنیم. از حذف متغیرهای  $x_0$  و  $y_0$  که در طی عملیات ثابت می‌مانند، خارج قسمتهای تفاضلی زیر را تشکیل می‌دهیم

$$\frac{z(\alpha_0, \beta_0 + \Delta\beta_0; T) - z(\alpha_0, \beta_0; T)}{\Delta\beta_0} \approx \frac{\partial z}{\partial \beta}(\alpha_0, \beta_0; T)$$

$$\frac{w(\alpha_0, \beta_0 + \Delta\beta_0; T) - w(\alpha_0, \beta_0; T)}{\Delta\beta_0} \approx \frac{\partial w}{\partial \beta}(\alpha_0, \beta_0; T)$$

$$\frac{z(\alpha_0 + \Delta\alpha_0, \beta_0; T) - z(\alpha_0, \beta_0; T)}{\Delta\alpha_0} \approx \frac{\partial z}{\partial \alpha}(\alpha_0, \beta_0; T)$$

$$\frac{w(\alpha_0 + \Delta\alpha_0, \beta_0; T) - w(\alpha_0, \beta_0; T)}{\Delta\alpha_0} \approx \frac{\partial w}{\partial \alpha}(\alpha_0, \beta_0; T)$$

بعد از اینکه  $z_T$  را جایگزین  $z(\alpha, \beta, T)$  و  $w_T$  را جایگزین  $w(\alpha, \beta, T)$  کردیم، می‌توانیم (۱۵.۹) را نسبت به مقادیر تصحیحی  $\delta\alpha_0 = \alpha - \alpha_0$  و  $\delta\beta_0 = \beta - \beta_0$  حل کنیم و برآوردهای جدید  $\alpha_1 = \alpha_0 + \delta\alpha_0$  و  $\beta_1 = \beta_0 + \delta\beta_0$  را برای پارامترهای  $\alpha$  و  $\beta$  به دست آوریم. اکنون تمامی روند فوق با شروع از مقادیر  $x_0, y_0, \alpha_0, \beta_0$  به عنوان شرایط اولیه تکرار شده است.

بنابراین هر بارست مشتمل بر سه بار حل دستگاه (۱۴.۹) است. در حالت کلی اگر تعداد  $n$  پارامتر اولیه مجهول وجود داشته باشد، هر بارست مستلزم  $n+1$  جواب دستگاه اصلی است. روشی که در اینجا به کار گرفته شده، باروش اصلاحی نیوتن برای یافتن ریشه‌های معادلات چند متغیری هم‌ارز است (به بخش ۲.۵ نگاه کنید).

در حل با کامپیوتر، مسائل با مقدار مرزی یکی از دشوارترین رده‌های مسائل است. در این مورد هیچ تضمینی برای همگرایی وجود ندارد، تخمینهای اولیه خوبی باید در دسترس باشند، و معمولاً آزمون و خطای فراوان و نیز مقدار زیادی وقت برای کامپیوتر لازم است.

□ مثال ۲.۹: با استفاده از روش نشانه‌گیری، مسئله (۱.۹) را حل کنید. این کار را با

تقریبهای اولیه  $\alpha_0 = 0.3$  و  $\alpha_1 = 0.4$  برای  $y'(0)$  و  $h = 0.1$  شروع کنید.

با استفاده از زیر برنامه استاندارد RK4\* برای حل معادلات دیفرانسیل که در فصل ۸

توضیح داده شد و همراه با درون‌یابی خطی بر اساس رابطه (۱۲.۹)، جوابهای زیر به دست

آمده‌اند. عمل بارستی در شرط  $|\alpha_{k+1} - \alpha_k| < 1 \times 10^{-6}$  متوقف شده است.

$k$	$\alpha_k$	$y(\alpha_k; 1)$
۰	۰٫۳۰۰۰۰۰۰۰۰	۰٫۳۵۲۵۶۰۷۷
۱	۰٫۴۰۰۰۰۰۰۰۰	۰٫۴۷۰۰۰۸۱۰۳
۲	۰٫۸۵۰۹۱۷۱۲	۰٫۹۹۹۹۹۹۹۹
۳	۰٫۸۵۰۹۱۷۱۲	۰٫۹۹۹۹۹۹۹۹

مقدار دقیق  $y'$  در  $x=0$  برابر است با  $\sinh^{-1} 1 = 0.85091813$  برای این مسئله همگرایی بسیار سریع است. با توجه به نمو تخمینی که در نظر گرفته شده است، دقت به دست آمده به طور استثنایی خوب است. با استفاده از روشهای تفاضل متناهی بخش ۱۰.۹، برای به دست آوردن دقت قابل مقایسه، نمو  $h=0.01$  لازم خواهد بود. مع هذا روش تفاضل متناهی باز از لحاظ محاسبه کارتر است.  $\square$

$\square$  مثال ۳.۹: با استفاده از روش نشانه گیری، مسئله با مقدار مرزی غیرخطی زیر را حل کنید

$$yy'' + 1 + y'^2 = 0 \quad y(0) = 1 \quad y(1) = 2 \quad (16.9)$$

حل: گیریم  $\alpha_0 = 0.05$ ،  $\alpha_1 = 1.0$  دو تقریب برای شیب مجهول  $y'(0)$  باشند. باز با استفاده از زیر برنامه RK4 برای حل معادلات دیفرانسیل و درونیابی خطی به ازای نمو  $h=1/64$  نتایج زیر به دست آمده اند:

$\alpha_i$	$y(\alpha_i; 1)$
۰٫۰۵۰۰۰۰۰۰	۰٫۹۹۹۹۹۹۹۹
۰٫۹۹۹۹۹۹۹۹	۱٫۴۱۴۲۱۳۳
۱٫۰۷۰۷۱۰۷۱	۱٫۸۴۷۷۵۸۲
۱٫۹۵۵۴۱۱۸	۱٫۹۷۷۵۷۸۶
۱٫۹۹۹۸۲۹۶۸	۱٫۹۹۹۹۱۴۶۳
۱٫۹۹۹۹۹۹۴۰	۱٫۹۹۹۹۹۹۵۲
۲٫۰۰۰۰۰۰۳۵	۲٫۰۰۰۰۰۰۰۰

شیب دقیق دره  $x = 2$  برابر است با  $y'(0) = 2$ . بعد از هفت بارست مشاهده می‌شود که شیب اولیه تا شش رقم بامعنی صحیح است، در حالی که مقدار  $y$  در  $x = 1$  حداقل تا هفت رقم بامعنی صحیح است. با استفاده از درون‌یابی درجه دوم، بعد از سه بارست اول همگرایی می‌توانست تسریع شود.

آشکار است که تعداد بارست لازم، به انتخاب تقریبهای اولیه  $\alpha_1, \alpha_0$  بستگی دارد. این تقریها را می‌توان گاهی از ملاحظات نموداری یا فیزیکی به دست آورد. □

## تمرین

۲۰۹-۱ جواب عددی معادله زیر را به دست آورید

$$2yy'' - y'^2 + 4y^2 = 0 \quad y\left(\frac{\pi}{6}\right) = \frac{1}{4} \quad y\left(\frac{\pi}{2}\right) = 1$$

به عنوان تقریبهای اولیه برای  $y'(\pi/6)$ ، بگیریید  $\alpha_0 = 0.85$  و  $\alpha_1 = 0.8$ ، و بارستها را تا جایی ادامه دهید که شرط در  $x = \pi/2$  تا پنج رقم بامعنی صحیح باشد.

جواب:  $y = (\sin x)^2$ ؛ شیب اولیه برابر است با  $y'(\pi/6) = \sqrt{3}/2$ .

۲۰۹-۲ در مثال ۳.۹ با استفاده از درون‌یابی درجه دوم بر اساس  $\alpha_0, \alpha_1, \alpha_2$  تقریب بعدی را به دست آورید. چند بارست صرفه جویی شده است؟

۲۰۹-۳ با استفاده از روش نشانه گیری، مسائل زیر را حل کنید:

$$y'' = 2y^3, \quad y(1) = 1, \quad y(2) = \frac{1}{4} \quad (\text{الف})$$

$y'(1) = 0$  را به عنوان اولین تخمین بگیریید. (جواب دقیق  $y = 1/x$ ).

$$y'' = e^y, \quad y(0) = y(1) = 0 \quad (\text{ب})$$

$y'(0) = 0$  را به عنوان اولین تخمین بگیریید.

## ۳.۹ روشهای پهلوی هم گذاری

درسالهای اخیر برای حل مسائل با مقدار مرزی، هم در حالت یک بعدی و هم در مورد ابعاد بالاتر، توجه زیادی بر روشهای تقریبی متمرکز شده است. در این روشهای تقریبی، به جای آنکه جواب در یک مجموعه از نقاط گسسته جستجو شود، سعی بر آن است که یک ترکیب خطی از توابع مستقل خطی به دست آورند که تقریبی را برای جواب ارائه دهد. در حقیقت



این اندیشه اصلی ریشه بسیار قدیمی دارد و به گیلر کین<sup>۱</sup> و ریتس<sup>۲</sup> [۳۱] برمی گردد، اما در این اواخر این اندیشه‌ها شکل جدیدی، تحت عنوان «عنصر منتهای»<sup>۳</sup> به خود گرفته اند (به کتاب استرننگ<sup>۴</sup> و فیکس<sup>۵</sup> [۳۱] مراجعه کنید)، و از نظر قالب و محتوی به آن حدی رسیده اند که اکنون با روشهای تفاضل منتهای کاملاً رقابت می کنند.

ما به طور خیلی اجمالی به فکر موجود در پشت این روشهای تقریبی اشاره می کنیم و توجه خود را بیشتر روش موسوم به روش پهلوی هم گذاری متمرکز می سازیم. (به کتاب استرننگ و فیکس مراجعه کنید [۳۱]).

برای ساده شدن کار، فرض می کنیم که يك مسئله با مقدار مرزی مرتبه دوم خطی داریم که به فرم زیر نوشته ایم

$$Ly \equiv -y'' + p(x)y' + q(x)y = r(x) \quad a \leq x \leq b \quad (۱۷.۹ \text{ الف})$$

$$a_0 y(a) - a_1 y'(a) = \alpha \quad (۱۷.۹ \text{ ب})$$

$$b_0 y(b) + b_1 y'(b) = \beta \quad |a_0| + |b_0| \neq 0$$

گیریم  $\{\psi_j(x)\}$ ،  $j = 1, \dots, N$ ، مجموعه ای از توابع مستقل خطی باشند، که باید به گونه ای که بعداً توضیح داده خواهد شد، انتخاب شوند. در این صورت سعی می کنیم جواب تقریبی برای رابطه (۱۷.۹) به شکل

$$U_N(x) = \sum_{j=1}^N c_j \psi_j(x) \quad (۱۸.۹)$$

پیدا کنیم. ضرایب  $\{c_j\}$  در بسط فوق چنان انتخاب می شوند که اندازه خطای صادق در مسئله با مقدار مرزی را به حداقل برساند. بر حسب اینکه اندازه خطا چگونه تعریف شده باشد، روشهای مختلفی پدید می آیند.

در روش پهلوی هم گذاری ضرایب چنان انتخاب می شوند که  $U_N(x)$  در شرایط (۱۷.۹ ب) و در معادله (الف) دقیقاً در نقاط درون بازه  $[a, b]$  صدق کنند. از این رو  $\{c_j\}$  ها در معادلات زیر صدق کنند

$$a_0 U_N(a) - a_1 U_N'(a) = \alpha$$

$$b_0 U_N(b) + b_1 U_N'(b) = \beta$$

$$LU_N(x_i) - r(x_i) = 0 \quad i = 1, \dots, N-2 \quad (۱۹.۹)$$

که در اینجا  $x_i$  ها مجموعه ای از نقاط متمایز در بازه  $[a, b]$  هستند. زمانی که (۱۹.۹)

به تفصیل نوشته شود، به صورت یک دستگاه  $N$  معادله  $N$  مجهولی خطی بسرحسب  $\{c_j\}$  درخواست آمد. وقتی (۱۹.۹)، مثلاً با یکی از روشهای فصل ۴، حل شود، جواب آن یعنی  $\{c_j\}$ ، در (۱۸.۹) گذاشته می شود تا جواب مطلوب به دست آید. تحلیل خطا برای این روش بسیار پیچیده و بیرون از ظرفیت این کتاب است. در عمل با افزایش تعداد توابع پایه یعنی  $N$ ، می توان یک رشته از تقریبهایی را به دست آورد. سپس از مقایسه این جوابهای تقریبی در مجموعه ای از نقاط ثابت در بازه  $[a, b]$ ، بر آوردی برای دقت می توان تعیین کرد. اکنون به تأمل در انتخاب توابع پایه  $\{\psi_j(x)\}$  می پردازیم. این توابع معمولاً طوری انتخاب می شوند که دارای یک یا چندتا از ویژگیهای زیر باشند:

$$(i) \quad \psi_j(x) \text{ها در بازه } [a, b] \text{ به طور پیوسته مشتق پذیر باشند}$$

$$(ii) \quad \psi_j(x) \text{ها در بازه } [a, b] \text{ متعامد باشند، یعنی}$$

$$\int_a^b \psi_j(x) \psi_k(x) dx = 0 \quad j \neq k \text{ به ازای}$$

$$(iii) \quad \psi_j(x) \text{ها توابعی «ساده» مانند بسجمله ایها یا توابع مثلثاتی باشند،}$$

$$(iv) \quad \psi_j(x) \text{ها در آن شرایط مرزی (در صورت وجود) که همگن باشند، صدق کنند.}$$

یکی از توابع پایه، که معمولاً به کار برده می شود، مجموعه  $\{\sin j\pi x\}$  است.  $j = 1, \dots, N$  است که در بازه  $[0, 1]$  قائم است. باید توجه داشت که به ازای تمامی  $j$ ها در  $x=0$  و در  $x=1$  تساوی  $\sin j\pi x = 0$  برقرار است. یک مجموعه پایه مهم دیگر، مجموعه  $\{P_j(x)\}$  است که در آن  $j = 0, \dots, N$  است که در آن  $P_j(x)$  بسجمله ایهای نژاندر هستند که در فصل ۶ تشریح شدند. این بسجمله ایها در بازه  $[-1, 1]$  متعامدند، بالاخره  $\psi_j(x)$ ها را می توان چنان انتخاب کرد که بسجمله ایهای درجه سوم-تکه ای باشند (به فصل ۶ نگاه کنید).

به عنوان مثال، روش پهلوی هم گذاری را برای معادله (۱۰.۹) که به صورت زیر نوشته

می شود به کار می بریم

$$U'''(x) - U(x) = 0 \quad (۲۰.۹ \text{ الف})$$

$$U(0) = 0 \quad (۲۰.۹ \text{ ب})$$

$$U(1) = 1$$

برای توابع پایه بسجمله ایهایی انتخاب می کنیم و به دنبال جواب تقریبی  $U_N(x)$  به شکل

$$U_N(x) = c_1 x + c_2 x^2 + c_3 x^3 \quad (۲۱.۹)$$

می رویم. ملاحظه می کنیم که بدون توجه به انتخاب  $c_j$ ها، تساوی  $U_N(0) = 0$  برقرار است.

## 1. piecewise-cubic polynomials

چون سه ضریب وجود دارد، باید سه شرط روی  $U_N(x)$  اعمال کنیم. یکی از شروط این است که  $U_N(x)$  باید در شرط مرزی  $x=1$  صدق کند، بنابراین یک معادله برای  $c_3$  ها به صورت زیر داریم

$$U_N(1) = c_1 + c_2 + c_3 = 1 \quad (22.9)$$

با تکیه بر اینکه  $U_N(x)$  دقیقاً در دو نقطه از بازه  $[0, 1]$  در معادله (۲۰.۹) الف) صدق می کند، می توانیم دو شرط اضافی اعمال کنیم. نقاط  $x_0 = 1/4$  و  $x_1 = 3/4$  را به دلخواه انتخاب می کنیم، مستقیماً از محاسبه نتیجه می شود که

$$U_N''(x) - U_N(x) = -c_1 x + (2 - x^2)c_2 + (6x - x^3)c_3$$

و بنابراین داریم

$$U_N''\left(\frac{1}{4}\right) - U_N\left(\frac{1}{4}\right) = -\frac{1}{4}c_1 + \frac{31}{16}c_2 + \frac{95}{64}c_3 = 0$$

$$U_N''\left(\frac{3}{4}\right) - U_N\left(\frac{3}{4}\right) = -\frac{3}{4}c_1 + \frac{23}{16}c_2 + \frac{261}{64}c_3 = 0 \quad (23.9)$$

دستگاه معادلات (۲۲.۹) و (۲۳.۹) را می توان مستقیماً حل کرد و جوابهای زیر را به دست آورد

$$c_1 = 0.85223700 \quad c_2 = -0.013852700 \quad c_3 = 0.16161600$$

اگر این مقادیر را در (۲۱.۹) قرار دهیم جواب تقریبی زیر به دست می دهد

$$U_N(x) = 0.852237x - 0.0138527x^2 + 0.161616x^3 \quad (24.9)$$

اکنون این جواب تقریبی را برای پیدا کردن مقدار تقریبی  $U(x)$  یا حتی  $U'(x)$  در هر نقطه از بازه  $[0, 1]$  می توان به کار برد.

برای اینکه ببینیم تقریب به دست آمده  $U_N(x)$  تا چه اندازه به مقدار دقیق

$$U(x) = \sinh x / \sinh 1$$

نزدیک است، در زیر فهرست چند مقدار تطبیقی را می آوریم (به جدول ۱.۹ نگاه کنید). بنابراین به نظر می رسد که در بدترین مقادیر نزدیک به نقطهٔ میانی بازه، دو تا از مقادیر تسا سه رقم با هم توافق دارند. با توجه به تعداد کم توابع پایه که برای تقریب  $U_N(x)$  به کار برده شده، به نظر می رسد که نتیجه کاملاً مطلوب است. برای به دست آوردن نتایج دقیقتر، فقط باید تعداد توابع پایه را افزایش دهیم.

$x$	$U_N(x)$	$U(x)$
۰٫۱۰	۰٫۰۸۵۲۴۷	۰٫۰۸۵۳۳۷
۰٫۲۵	۰٫۲۱۴۷۱۹	۰٫۲۱۴۹۵۲
۰٫۵۰	۰٫۴۲۲۶۷۵	۰٫۴۲۳۳۰۹
۰٫۷۵	۰٫۶۹۹۵۶۷	۰٫۶۹۹۷۲۴
۰٫۹۰	۰٫۸۷۳۶۱۱	۰٫۸۷۳۴۸۱

## تمرین

۳۰۹-۱ مسئله با مقدار مرزی

$$U'''(x) - U(x) = x \quad U(0) = 0 \quad U(1) = 1$$

را با روش پهلوی هم گذاری حل کنید. برای تابع آزمون<sup>۱</sup> از بسجمله‌ای پایه

$$U_N(x) = c_0 x + c_1 x^2 + c_2 x^3 + \dots + c_N x^N$$

استفاده کنید. ابتدا به ازای  $N=3$  و سپس به ازای  $N=4$  مسئله را حل و جوابها را در نقاط انتخاب شده از بازه مورد نظر مقایسه کنید. همچنین نتایج تقریبی را با جواب دقیق زیر مقایسه کنید

$$U(x) = \frac{\sinh x}{\sinh 1} - x$$

۳۰۹-۲ با استفاده از روش پهلوی هم گذاری سعی کنید که مسئله با مقدار مرزی

$$U''(x) + U(x) = x$$

$$U(0) = 0 \quad U(1) = 1$$

را حل کنید. عملیات را با تابع آزمون

$$U_N(x) = x + \sum_{j=1}^N c_j \sin j\pi x$$

که خود به خود به ازای جميع مقادیر  $c_j$  در شرایط مرزی صدق می کند، شروع کنید. مسئله را به ازای  $N=2$  و  $N=4$  حل، و نتایج را مقایسه کنید.

## ضمیمه

## کتابهای زیر برنامه‌ای

فهرست زیر شرح مختصری است از برخی از مهمترین مجموعه‌های نرم افزارها که شامل زیر برنامه‌های آزمون شده برای حل کلیه مسائل مهم این کتاب است. اطلاعات بیشتر در باب امکان دستیابی به آنها را می‌توانید از منابعی که نام برده‌ایم به دست آورید.

## ۱. IMSL (کتابهای بین‌المللی ریاضی و آمار)

احتمالاً این کاملترین مجموعه‌ای است که می‌توان از بازار تهیه کرد. این مجموعه شامل تعدادی در حدود ۲۳۵۵ زیر برنامه است که برای حل کلیه مسائل این کتاب و زمینه‌های دیگری مانند محاسبات آماری و نیز بهینه‌سازی مقید کاربرد دارد. تمامی این برنامه‌ها به زبان ANSI FORTRAN نوشته و طوری تطبیق داده شده‌اند که روی کلیه کامپیوترهای بزرگ قابل اجرا هستند.

Source :IMSL, Inc. GNB Building, 7500 Bellaire Blvd.  
Houston, Texas 77036.

## ۲. PORT

یک مجموعه نسبتاً کامل و آزمایش شده‌ای است از زیر برنامه‌ها که برای تمامی مسائل متداول در زمینه آنالیز عددی کاربرد دارد. این مجموعه، به زبان PFORT که زیر مجموعه‌ای قابل

- 
1. software package
  2. International Mathematical and Statistical Library

حمل از زبان ANSI FORTRAN است، نوشته و طوری طراحی شده است که به آسانی می توان آن را از يك کامپیوتر به کامپیوتر دیگر منتقل کرد.

Source: Bell Telephone Laboratories, Murray Hill, New Jersey.

### EISPACK .۳

مجموعه ای است برای حل مسائل استاندارد مربوط به ویژه مقدارها و ویژه بردارها. این مجموعه به زبان ANSI FORTRAN کد گذاری شده و شکلی کاملاً مستقل از کامپیوتر اجرا کننده دارد. این مجموعه يك مجموعه نرم افزاری است با کیفیت بسیار بالا و بسیار قابل اطمینان و شامل تعداد زیادی امکانات کمکی برای استفاده کننده (به [۳۲] نگاه کنید).

Source: National Energy Software Center, Argonne National Laboratories 9700 S. Cass Ave, Argonne Illinois 60439.

### LINPACK .۴

يك مجموعه نرم افزاری برای حل دستگاه معادلات خطی و نیز مسائل مربوط به کوچکترین توانهای دوم. این بسته به زبان ANSI FORTRAN نوشته شده و مستقل از کامپیوتر، در زمینه حساب حقیقی، حساب مختلط و حساب با دقت مضاعف در دسترس است. این بسته به نحوی گسترده در چندین مرکز محاسباتی، مورد آزمایشهای زیادی قرار گرفته است.

Source: National Energy Software Center, Argonne National Laboratories, 9700 S. Cass Ave., Argonne, Illinois 60439.

1. Hamming, R. W.: *Numerical Methods for Scientists and Engineers*, McGraw-Hill, New York, 1962.
2. Henrici, P. K.: *Elements of Numerical Analysis*, John Wiley, New York, 1964.
3. Traub, J. F.: *Iterative Methods for the Solution of Equations*, Prentice-Hall, New Jersey, 1963.
4. Scarborough, J. B.: *Numerical Mathematical Analysis*, Johns Hopkins, Baltimore, 1958.
5. Hildebrand, F. B.: *Introduction to Numerical Analysis*, McGraw-Hill, New York, 1956.
6. Müller, D. E.: "A method of solving algebraic equations using an automatic computer," *Mathematical Tables and Other Aids to Computation (MTAC)*, vol. 10, 1956, pp. 208–215.
7. Hastings, C. Jr.: *Approximations for Digital Computers*, Princeton University Press, New Jersey, 1955.
8. Milne, W. E.: *Numerical calculus*, Princeton University Press, New Jersey, 1949.
9. Lanczos, C.: *Applied Analysis*, Prentice-Hall, New Jersey, 1956.
10. Householder, A. S.: *Principles of Numerical Analysis*, McGraw-Hill, New York, 1953.
11. Faddeev, D. K., and V. H. Faddeeva: *Computational Methods of Linear Algebra*, Freeman, San Francisco, 1963.
12. Carnahan, B., et al.: *Applied Numerical Methods*, John Wiley, New York, 1964.
13. *Modern Computing Methods*, Philosophical Library, New York, 1961.
14. McCracken, D., and W. S. Dorn: *Numerical Methods and Fortran Programming*, John Wiley, New York, 1964.
15. Henrici, P. K.: *Discrete Variable Methods for Ordinary Differential Equations*, John Wiley, New York, 1962.
16. Hamming, R. W.: "Stable Predictor-Corrector Methods for Ordinary Differential Equations," *Journal of the Association for Computing Machinery (JACM)*, vol. 6, no. 1, 1959, pp. 37–47.
17. Rice, J. R.: *The Approximation of Functions*, vols. 1 and 2, Addison-Wesley, Reading, Mass., 1964.
18. Forsythe, G., and C. B. Moler; *Computer Solution of Linear Algebraic Systems*, Prentice-Hall, New Jersey, 1967.
19. Isaacson, E., and H. Keller: *Analysis of Numerical Methods*, John Wiley, New York, 1966.
20. Stroud, A. H., and D. Secrest: *Gaussian Quadrature Formulas*, Prentice-Hall, New Jersey, 1966.
21. Johnson, L. W., and R. D. Riess: *Numerical Analysis*, Addison-Wesley, Reading, Mass., 1977.
22. Forsythe, G. E., M. A. Malcolm, and C. D. Moler: *Computer Methods for Mathematical Computations*, Prentice-Hall, New Jersey, 1977.
23. Stewart, G. W., *Introduction to Matrix Computation*, Academic Press, New York, 1973.
24. Wilkinson, J. H.: *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford, 1965.
25. Ralston, A.: *A First Course in Numerical Analysis*, McGraw-Hill, New York, 1965.

26. Shampine, L. and R. Allen: *Numerical Computing*, Saunders, Philadelphia, 1973.
27. Gautschi, W.: "On the Construction of Gaussian Quadrature Rules from Modified Moments," *Math. Comp.*, vol. 24, 1970, pp. 245-260.
28. Fehlberg, E.: "Klassische Runge-Kutta-Formeln vierter und niedriger Ordnung mit Schrittweitenkontrolle und ihre Anwendung auf Wärmeleitungsprobleme," *Computing*, vol. 6, 1970, pp. 61-71.
29. Hull, T. E., W. H. Enright, and R. K. Jackson: *User's Guide for DVERK—A Subroutine for Solving Non-Stiff ODE's*, TR 100, Department of Computer Science, University of Toronto, October, 1976.
30. Gear, C. W.: *Numerical Initial Value Problems in Ordinary Differential Equations*, Prentice-Hall, New Jersey, 1971.
31. Strang, G., and G. Fix: *An Analysis of the Finite Element Method*, Prentice-Hall, New Jersey, 1973.
32. Smith, B. T., J. M. Boyle, J. J. Dongarra, B. S. Garbow, Y. Ikebe, V. C. Klema, and C. B. Moler: "Matrix Eigensystem routines—EISPACK Guide," *Lecture Notes in Computer Science*, vol. 6, Springer-Verlag, Heidelberg, 1976.
33. Ortega, J. M., and W. C. Rheinboldt: *Iterative Solution of Nonlinear Equations in Several Variables*, Academic Press, New York, 1970.
34. Robinson, S. R.: "Quadratic Interpolation Is Risky," *SIAM J. Numer. Analysis*, vol. 16, 1979, pp. 377-379.
35. Rivlin, T. J.: *An Introduction to the Approximation of Functions*, Blaisdell, Waltham, Mass., 1969.
36. Winograd, S.: "On Computing the Discrete Fourier Transform," *Math. Comp.*, vol. 32, 1978, pp. 175-199.
37. Cooley, J. W., and J. W. Tukey: "An Algorithm for the Machine Calculation of Complex Fourier Series," *Math. Comp.*, vol. 19, 1965, pp. 297-301.
38. Ehlich, H., and K. Zeller: "Auswertung der Normen von Interpolationsoperatoren," *Math. Annalen*, vol. 164, 1966, pp. 105-112.
39. de Boor, C., and A. Pinkus: "Proof of the Conjectures of Bernstein and Erdős," *J. Approximation Theory*, vol. 24, 1978, pp. 289-303.
40. de Boor, C.: *A Practical Guide to Splines*, Springer-Verlag, New York, 1978.
41. Wendroff, B.: *Theoretical Numerical Analysis*, Academic Press, New York, 1966.
42. Wilkinson, J. H.: *Rounding Errors in Algebraic Processes*, Prentice-Hall, New Jersey, 1963.



## فهرست اسامی خاص

Adams-Bashforth	ادمز-بشفورت
Adams-Moulton	ادمز-مولتن
Cramer	کرامر
Descartes	دکارت
Euler	اویلر
Jacobi	ژاکوبی
Lagrange	لاگرانژ
Laguerre	لاگر
Lebesgue	لیبگ
Milne	میلن
Muller	مولر
Murnaghan-Wrench	مورناهان-رنج
Newton	نیوتن
Parseval	پارسوال
Rayleigh	ریلی
Remez	رهمه
Rolle	رول
Romberg	رامبرگ
Runge-Kutta	رونکه-کوتا
Steffensen	استفنسن
Taylor	تیلر
Vandermonde	واندرموند

## واژه‌نامه

subscription	آبونمان
free-end	آزاد-انتهای
trial and error	آزمون و خطا
hyperplane	اُبرصفحه
standard	استانده
adjoint	الحاقتی
compact scheme	الگوی فشرده
accumulation	انباشتگی
option	انتخاب
standard deviation	انحراف معیار
step size	اندازهٔ نمو
matrix updating	باب روز کردن ماتریس
iteration	بارست
inner iteration	-داخلی
inverse iteration	-معکوس
fixed point iteration	-نقطهٔ ثابت
track	باریکه
interval	بازه
ill-conditional	بدشروط
ill-conditioning	بدشروطی
body program	بدنهٔ برنامه

data fitting	برازاندن داده‌ها
curve fitting	برازاندن منحنی
subroutine	برنامه فرعی
extrapolation	برونبایی
superposition	برهم‌نهی
software package	بسته نرم‌افزاری
characteristic polynomial	بسجمله‌ای مشخصه
well-conditioned	به‌شرط
bit	بیت
overrelaxation	بیش‌واهاش
flag	پرچم
back-substitution	پسجایگذاری
collocation	پهلوی هم‌گذاری
underflow	پی‌ریز
forward substitution	پیشجایگذاری
biasing	پیشقداری
iteration function	تابع بارست
real valued function	تابع حقیقی مقدار
integrand	تابع زیر علامت انتگرال
univariate function	تابعی از یک متغیر تصادفی
discrete transform	تبدیل گسسته
factorization	تجزیه به عوامل
binary search	تجسس دودویی
backward error analysis	تحلیل خطای پسرو
pivoting strategy	تدبیر لولا‌گزینی
transposed	ترانهاده
transposition	ترانهش
stochastic	تصادفی
filtering	تصفیه
forward shift	تعویض پیشرو
difference	تفاضل
forward difference	پیشرو

divided difference	منقسم
finite difference	های متناهی
synthetic division	تقسیم ترکیبی
monomial	تکجمله‌ای
multiplicity	تکرر
piecewise	تکه‌ای
singularity	تکینگی
steepest descent	تندترین کاهش
sparse	تنک
register	ثبتات
permutation	جایگشت
summand	جمع‌وند
jump	جهش
nontrivial solution	جواب نابدیهی
weighted sum	حاصلجمع وزین
steady-state	حالت ماندگار
interval arithmetic	حساب بازه‌ای
call-statement	حکم فراخوانی
loop	حلقه
error	خطا
residual error	بی‌باقیمانده
round-off error	بی‌گرد کردن
line search	خط‌جویی
contour lines	خطوط مرزی
self-starting	خود آغازگر
well-behaved	خوش رفتار
entry	درایه
interpolant	درون‌یاب
interpolation	درون‌یابی

iterative interpolation	- بارستی
osculatory interpolation	- بوسانی
linear interpolation	- خطی
spline interpolation	- قلمی
single precision	دقت ساده
double precision	- مضاعف
plotter	رسام
method	روش
secant method	- خط متقاطع
chopping method	- قطع کردن
damped method	- میرا
undamped method	روش غیر میرا
truncating method	روش قطع کردن
rounding method	روش گرد کردن
flowchart	روند نما
phase angle	زاویه فاز
superdiagonal	زبرقطری
superscript	زیرنمایه
subdiagonal	زیرقطری
subscript	زیرنمایه
structured	ساختیافته
consistent	سازگار
bookkeeping	ساماندهی
overflow	سرریز
truncated series	سریهای ناتمام
pseudorandom	شبه تصادفی
perturbed argument	شناسه آشفته
shifted power form	صورت توانی منتقله

nested form	صورت تودرتو
premultiplication	ضرب ازچپ
postmultiplication	ضرب از راست
pulse	ضربه
leading coefficient	ضریب جمله پیشرو
condition number	عدد شرط
transcendental number	عدد غیر جبری
complex number	عدد همثافت
operand	عملوند
finite elements	عنصرهای منتهای
angular frequency	فراوانی زاویه‌ای
shifted power form	فرم توانی تغییر مکان یافته
corrector formula	فرمول تصحیحی
central difference-formula	فرمول تفاضل-مرکزی
compact	فشرده
panel	قطعه
peak	قله
comment card	کارت توضیحی
minor	کهپاد
small-oh	کهصفر
discretization	گسسته‌سازی
matrix	ماتریس
upper triangular matrix	- بالامثلثی
lower triangular matrix	- پایین‌مثلثی
dense matrix	- چگال
tridiagonal matrix	- سه قطری
off-diagonal matrix	- غیر قطری
mantissa	مانتیس

orthogonal	متعامد
boundary_value problem	مسئلهٔ مقدار مرزی
circular	مستدیر
pivotal equation	معادلهٔ لولایی
intermediate_value	مقدار میانی
time scale	مقیاس زمانی
floating_point	ممیز شناور
plausible	موجه نما
big_oh	مهضفر
damping	میرایی
instability	ناپایداری
annual region	ناحیهٔ طوقی
improper	ناسره
imbedding	نشانیدن
breakpoint	نقطهٔ انفصال
exterior mesh_points	نقطه‌های خارج شبکه
increment	نمو
sample	نمونه
harmonic oscillation	نوسانهای همساز
noise	نوفه
noisy	نوفه‌ای
variance	واریانس
waltz	والس
relaxation	واهلش
existence and uniqueness	وجود و یکتایی
compilation	همگردانی
univariate	یک متغیری
monotone	یکنوا