

به نام پروردگار یگانه و یکتا



دانشکده علوم ریاضی

آنالیز عددی پیشرفته

گردآورنده : رضا مختاری

ویرایش ششم : پاییز ۱۴۰۰

فهرست مطالب

۲	فهرست مطالب
۵	۱ مفاهیم مقدماتی
۶	۱.۱ مراجع آنالیز عددی
۶	۲.۱ تاریخچه آنالیز (محاسبات) عددی
۸	۳.۱ مفاهیم اساسی
۱۰	۴.۱ پیش‌نیازهای ریاضی
۱۸	۲ خطاها
۱۸	۱.۲ منابع تولید خطا
۲۱	۲.۲ نمایش اعداد
۲۵	۱.۲.۲ نمایش ۶۴-بیتی ممیز شناور
۲۷	۲.۲.۲ اعداد ماشینی
۳۱	۳.۲ انتشار خطا
۳۵	۱.۳.۲ خطای اعمال ریاضی
۳۷	۲.۳.۲ تقریب توابع یک متغیره
۴۰	۴.۲ بررسی پایداری یک الگوریتم
۵۰	۵.۲ حساب بازه‌ای
۵۲	۶.۲ تقسیم‌بندی مسایل در ریاضیات محاسباتی
۵۷	۳ درونیابی
۵۷	۱.۳ مقدمات
۵۹	۲.۳ مسئله درونیابی کلی
۶۲	۳.۳ درونیابی چندجمله‌ای
۶۴	۱.۳.۳ روش لاگرانژ
۶۶	۲.۳.۳ الگوریتم نویل
۶۸	۳.۳.۳ روش تفاضلات تقسیم‌شده نیوتن
۷۱	۴.۳.۳ روش‌های پیشرو/پسرو نیوتن
۷۸	۵.۳.۳ خطای چندجمله‌ای درون‌یاب

۸۴	درونیابی هریت	۴.۳
۸۵	روش لاگرانژ تعمیم یافته	۱.۴.۳
۸۵	روش تفاضلات تقسیم شده نیوتن تعمیم یافته	۲.۴.۳
۸۷	درونیابی با چند جمله ای های قطعه ای	۵.۳
۸۹	درونیابی هموار اسپلاین	۶.۳
۹۷	B-اسپلاین	۱.۶.۳
۹۹	درونیابی گویا	۷.۳
۱۰۲	روش تفاضلات تقسیم شده وارون	۱.۷.۳
۱۰۴	الگوریتم شبه نویل	۲.۷.۳
۱۰۸	درونیابی مثلثاتی	۸.۳
۱۱۱	تبدیل فوریه سریع	۱.۸.۳
۱۱۴	محاسبه مجموع های سینوسی-کسینوسی	۲.۸.۳
۱۱۶	برونیابی و درونیابی وارون	۹.۳
۱۲۰	نظریه تقریب	۴
۱۲۰	فضای ضرب داخلی	۱.۴
۱۲۹	تقریب کمترین مربعات گسسته	۲.۴
۱۳۱	تقریب کمترین مربعات پیوسته	۳.۴
۱۳۶	درونیابی و تقریب در ابعاد بالاتر	۵
۱۳۷	ضرب تانسوری	۱.۵
۱۴۲	درونیابی با توابع قطعه ای چند جمله ای	۲.۵
۱۴۶	مشتق گیری و انتگرال گیری عددی	۶
۱۴۶	مشتق گیری عددی	۱.۶
۱۴۶	استفاده از چند جمله ای درون یاب	۱.۱.۶
۱۴۹	استفاده از بسط تیلور	۲.۱.۶
۱۴۹	روش گاوس	۳.۱.۶
۱۵۰	فن برون یابی ریچاردسون	۴.۱.۶
۱۵۲	انتگرال گیری عددی	۲.۶
۱۵۴	قواعد نیوتن-کانتس	۱.۲.۶
۱۵۹	استفاده از فن برون یابی	۳.۶
۱۶۳	قواعد گاوس	۴.۶
۱۶۹	مطالب تکمیلی	۵.۶
۱۶۹	انتگرال گیری ضربی	۱.۵.۶
۱۷۰	انتگرال های تکین	۲.۵.۶



فصل ۱

مفاهیم مقدماتی

دی رفت و باز نیاید

فردا را اعتماد نشاید

حال را غنیمت دان

که دیر نپاید

خواجه عبدا... انصاری

در این فصل پس از آرایه مراجع و تاریخچه‌ای از آنالیز عددی، مفاهیم، تعاریف و قضایایی که در ادامه کار به آن‌ها نیاز داریم را مرور می‌کنیم.

۱.۱ مراجع آنالیز عددی

سه برخورد متفاوت با سرفصل آنالیز عددی عبارتند از

• کاربردی و الگوریتمی (روش‌گرایانه)

۱. Introductory Computer Methods & NA, Pennington, 1970

۲. Introduction to NA: An Algorithmic Approach, Deboor, 1972

۳. Numerical Recipes: The Art of Scientific Computing, Press & Teukolsky & Vetterling & Flannery, 2007

۴. NA for Engineers and Scientists, Miller, 2014

• نظری (ایده‌گرایانه)

۱. Theoretical NA, Wendroff, 1967

۲. Theoretical NA, Linz, 1979

۳. A theoretical Introduction to NA, Ryaben'kii & Tsynkov, 2007

۴. Theoretical NA: A Functional Analysis Framework, Atkinson & Han, 2009

• تلفیقی از دو دیدگاه قبلی

۱. Elements of NA, Henrichi, 1964

۲. A First Course in NA, Ralston & Rabinowits, 1978

۳. An Introduction to NA, Atkinson, 1989

۴. NA: Mathematics of Scientific Computing, Kincaid & Cheney, 1991

۵. Numerical Mathematics, Quarteroni & Sacco & Saleri, 2007

۶. Introduction to NA, Stoer & Bulirsch, 1993 & 2002 & 2008

۷. An Introduction to Numerical Methods and Analysis, Epperson, 2013

۸. NA, Burdan & Faires, 2016

۹. A Consice Introduction to NA, Faul, 2016

۲.۱ تاریخچه آنالیز(محاسبات) عددی

آن چه در ادامه می‌آید و از آن به عنوان تاریخچه آنالیز(محاسبات) عددی یاد شده، وقایعی است که بی‌شک با آنالیز عددی در ارتباط بوده و در گسترش آن بی‌تاثیر نبوده‌اند. خواننده علاقمند می‌تواند تاریخچه کامل را در کتاب‌های تاریخ ریاضیات دنبال کند.

۱. تهیه نتایج به صورت اعداد (بشر اولیه!؟)
۲. کارهای بابلی‌ها و مصری‌ها در رابطه با نجوم و مهندسی راه و ساختمان
 - لوح بابلی حدود ۲۰۰۰ سال قبل از میلاد (مجذور اعداد صحیح از ۱ تا ۶۰)
 - ثبت کسوف و خسوف از حدود ۷۵۰ سال قبل از میلاد
 - توجه مصریان به کسر و ابداع روش نابجایی
۳. کارهای ارشمیدس در حدود ۲۲۰ سال قبل از میلاد ($3\frac{1}{7} < \pi < 3\frac{1}{4}$)
۴. ابداع روش تکراری محاسبه \sqrt{a} توسط هرون در حدود ۱۰۰ سال قبل از میلاد
۵. کارهای فیثاغورسیان در رابطه با مجموع عددی سری‌ها
۶. روش دیوفانتوس برای حل معادله درجه دوم در حدود سال ۲۵۰ میلادی
۷. نمادگذاری عددی عربی (هندی)
۸. ساخت جدول‌های مثلثاتی قبل از قرن ۱۰
۹. پیدایش جبر در قرن ۱۶
۱۰. قرن ۱۷
 - انقلاب صنعتی (کارهای نیوتن، اویلر، لاگرانژ، گاوس و بسل)
 - ساخت جدول لگاریتم توسط نپر
 - اختراع خط‌کش محاسبه توسط اوترد
 - اختراع ماشین حساب توسط پاسکال و لایب‌نیتز
 - محاسبات با سری‌های نامتناهی
۱۱. پایه‌ریزی تفاضلات متناهی توسط ژاکوب استرلینگ و بروک تیلور در اوایل قرن ۱۸
۱۲. پیش‌بینی وجود و موضع سیاره نپتون توسط آدامز و لوریه در سال ۱۸۴۵ میلادی
۱۳. پیشرفت و تولید وسیع ماشین حساب در اواخر قرن ۱۹
۱۴. چاپ کتاب آنالیز ریاضی عددی^۱ توسط اسکار بورو [۲۱] در سال ۱۹۳۰ میلادی
۱۵. جنگ جهانی دوم و پیدایش رایانه در سال ۱۹۴۰ میلادی
۱۶. معرفی واژه Numerical Analysis توسط UCLA در سال ۱۹۴۷ میلادی

۳.۱ مفاهیم اساسی

در این بخش پس از بررسی چند تعریف از آنالیز عددی، مفاهیمی اساسی و مرتبط با آنالیز عددی را مطرح می‌کنیم.

تعریف ۱.۱ تعاریف متنوع آنالیز عددی عبارتند از

۱. آنالیز عددی نظریه روش‌های ساخت‌یافته (سازنده)^۲ در آنالیز ریاضی، آنالیز تابعی، جبر، جبرخطی و سایر شاخه‌های ریاضی است [۱۲].
 ۲. آنالیز عددی آنالیز الگوریتم‌های پیوسته است [۲۳].
 ۳. آنالیز عددی مطالعه کمی جواب تقریبی مسایل ریاضی با توجه به خطاها و کران‌های آن‌ها است [۲۵].
 ۴. آنالیز عددی به استنتاج ریاضی، توصیف و تحلیل روش‌هایی می‌پردازد که جواب‌های عددی مسایل ریاضی را تولید می‌کنند [۱].
 ۵. آنالیز عددی مطالعه روش‌های تقریبی و دقت آن‌ها است [۶].
 ۶. آنالیز عددی شامل مطالعه، توسعه و تجزیه و تحلیل الگوریتم‌ها برای به دست آوردن جواب‌های عددی مسایل مختلف ریاضی است. آنالیز عددی ریاضیات محاسبات علمی است [۱۴].
 ۷. آنالیز عددی مطالعه جواب تقریبی مسایل ریاضی با در نظر گرفتن حد و اندازه خطاهای ممکن است [۳].
 ۸. آنالیز عددی مطالعه الگوریتم‌ها برای مسایل ریاضیات پیوسته (متمايز از ریاضیات گسسته) است [۲۴].
- آتکینسون^۳ آنالیز عددی‌دان معاصر، برای شناخت یک موضوع دیدگاه‌های متفاوت زیر را مطرح می‌کند.
۱. (مورد قدیمی) نظری (دیدگاه دانشمندان)
 ۲. (مورد قدیمی) تجربی (دیدگاه مهندسان)
 ۳. (مورد جدید) محاسباتی (دیدگاه متخصصان آنالیز عددی)
- به عبارتی وی علوم محاسباتی را به علوم نظری و علوم تجربی اضافه کرده و معتقد است این علوم به صورت زیر با هم در ارتباط هستند.

علوم نظری \rightleftharpoons علوم محاسباتی \rightleftharpoons علوم تجربی \rightleftharpoons علوم نظری

الگوریتم واژه‌ای است که در آنالیز (محاسبات) عددی بسیار مورد استفاده قرار می‌گیرد و در ادامه به تعریف آن می‌پردازیم.

تعریف ۲.۱ الگوریتم در آنالیز عددی،

- به عنوان توصیفی کامل و بدون ابهام از روش ساختن جواب یک مسئله ریاضی تعریف می‌شود [۱].
- فرآیند بدون ابهامی است که دنباله‌ای متناهی از گام‌هایی که با ترتیب مشخصی اجرا می‌شوند را تشریح می‌کند [۲].

در آنالیز عددی فقط با الگوریتم‌های عددی سرو کار داریم. الگوریتم‌هایی هستند که خروجی آن‌ها جواب عددی یا تقریبی یک مسئله است. الگوریتم عددی همگرا، به آن دسته از الگوریتم‌های عددی گفته می‌شود که در هر تکرار، عناصر دنباله‌ای را تولید می‌کنند که آن دنباله به جواب واقعی مسئله همگرا باشد. باید توجه داشت که ماهیت این دنباله‌ها وابسته به جواب مسئله است و ممکن است عدد، تابع، بردار، ماتریس و

^۲ Constructive methods^۲: روشی که علاوه بر وجود، راه رسیدن به جواب را نیز مشخص می‌کند.

^۳ K.E. Atkinson^۳

غیره باشد. از شبه‌کد^۴ برای توصیف الگوریتم استفاده می‌شود. شبه‌کد را می‌توان زبانی بین زبان‌های سطح بالای رایانه‌ای و زبان محاوره تلقی کرد. هنگام مواجه شدن با یک مسئله حقیقی، سه سوال بحث‌برانگیزی که برای یک متخصص آنالیز عددی مطرح می‌شود، عبارتند از:

۱. آیا جواب مسئله وجود دارد؟

۲. آیا جواب یکتا است؟

۳. حساسیت جواب در برابر اختلالات جزئی (کوچک) داده‌های ورودی در چه حد است؟

تعریف ۳.۱ فرض کنید d و $S(d)$ به ترتیب بیان‌گر داده‌های ورودی و جواب متناظر با داده‌ها، برای مسئله‌ای دلخواه باشند^۵. هم‌چنین فرض کنید $d + \delta d$ داده‌های ورودی اختلال‌یافته (تغییر داده‌شده) و $S(d + \delta d)$ جواب نظیر آن باشد. معیار اندازه‌گیری تغییرات اعمال‌شده در داده‌ها و اختلاف جواب‌های متناظر را به ترتیب با $\|S(d + \delta d) - S(d)\|$ و $\|\delta d\|$ که اعداد نامنفی هستند، نشان می‌دهیم. به ازای داده‌های مفروض d ، مسئله خوش‌طرح (قیافه)^۶ نامیده می‌شود اگر دو شرط زیر برقرار باشد

الف) به ازای تمام داده‌های متعلق به یک همسایگی حول d جواب یکتا موجود باشد، یعنی عددی مانند $\epsilon > 0$ چنان موجود باشد که به ازای هر δd که $\|\delta d\| < \epsilon$ مقدار $S(d + \delta d)$ موجود و یکتا باشد.

ب) جواب یعنی $S(d)$ به طور پیوسته به داده‌ها یعنی d بستگی داشته باشد یعنی اگر $\|\delta d\| \rightarrow 0$ آنگاه داشته باشیم $\|S(d + \delta d) - S(d)\| \rightarrow 0$.

تذکر ۱.۱ در صورتی که به ازای مجموعه‌ای از داده‌ها، بیش از یک جواب موجود باشد، نمی‌توانیم تشخیص دهیم روش عددی ما کدام جواب را تولید می‌کند. ممکن است اطلاعات حقیقی (تجربی) پیرامون مسئله، مشخص کند که جواب مسئله، موجود و یکتا است. اگر متخصص در این مرحله به جواب نرسید، مدل ریاضی مسئله را مورد بررسی قرار داده و برای تعیین وجود و یکتایی جواب در مباحث نظری در جستجوی یافتن قضایایی است که تکلیف وجود و یکتایی جواب را روشن می‌کنند. به طور معمول با افزودن شرایطی، می‌توان جواب‌های اضافی را کنار گذاشت. به عنوان مثال، در ریشه‌یابی یک معادله جبری درجه دو یا بیشتر، می‌توان هدف را تعیین ریشه با بزرگ‌ترین اندازه قرار داد.

تذکر ۲.۱ اگر شرط ب برقرار نباشد، استفاده از یک الگوریتم عددی کار دشواری است، زیرا در این صورت داده‌هایی بسیار نزدیک به d وجود دارند که جواب‌های متناظر آنها متمایز از $S(d)$ هستند.

تعریف ۴.۱ مسئله‌ای که به ازای مجموعه‌ای از داده‌ها خوش‌طرح است، خوش‌حالت (وضع)^۷ نامیده می‌شود اگر هر تغییر جزئی در داده‌های ورودی منجر به تغییر کوچکی در جواب شود و در صورتی که تغییر جواب فاحش باشد، مسئله را بدحالت (وضع)^۸ (مریض حال) نامند.

تذکر ۳.۱ مفاهیم خوش‌طرحی (بدطرحی) و خوش‌حالت (بدحالت) بودن نسبی هستند و نه تنها به مسئله بلکه به داده‌ها نیز بستگی دارند.

تعریف ۵.۱ جواب یک مسئله خوش‌طرح اغلب بستگی لیبشیتس به داده‌ها دارد یعنی اعداد $L > 0$ و $\epsilon > 0$ موجود هستند به طوری که به

ازای هر δd که $\|\delta d\| < \epsilon$ خواهیم داشت

$$\|S(d + \delta d) - S(d)\| \leq L \|\delta d\|.$$

^۴Pseudocode

^۵ماهیت d و $S(d)$ بر حسب نوع مسئله، ممکن است عدد، تابع، بردار، ماتریس و غیره باشد.

^۶Well-posed

^۷Well-conditioned

^۸Ill-conditioned

L به عدد حالت (وضعیت)^۹ معروف است و اگر L خیلی بزرگ نباشد، مسئله خوش‌حالت است. یک مثال جالب در برنامه ill-conditioned.nb بررسی شده است.

تعریف ۶.۱ اگر هنگام به کار بردن یک الگوریتم در حل یک مسئله خوش‌حالت، اختلال کوچکی در داده‌های ورودی، منجر به تولید تغییرات بزرگی در نتایج نهایی گردد، آن الگوریتم ناپایدار^{۱۰} و در غیر این صورت پایدار^{۱۱} نامیده می‌شود. اگر الگوریتمی به ازای انتخاب مشخصی از داده‌ها پایدار گردد، آن را به طور مشروط پایدار^{۱۲} نامند. برنامه stable.nb را ببینید.

۴.۱ پیش‌نیازهای ریاضی

در این بخش به معرفی اصطلاحات و بیان قضایای مورد نیاز می‌پردازیم.

تعریف ۷.۱ با این فرض که $n \in \mathbb{N}$ و $X \subseteq \mathbb{R}$ ، چند فضای پرکاربرد از توابع به صورت زیر تعریف می‌شوند.

$$C(X) = \{f : \Omega \rightarrow \mathbb{R} \mid f \text{ در } \Omega \text{ پیوسته باشد}\}$$

$$C[a, b] = \{f : [a, b] \rightarrow \mathbb{R} \mid f \text{ در } [a, b] \text{ پیوسته باشد}\}$$

$$C^n(X) = \{f : \Omega \rightarrow \mathbb{R} \mid f \text{ در } \Omega \text{ پیوسته باشد و } n \text{ مرتبه مشتق داشته باشد}\}$$

$$C^n[a, b] = \{f : [a, b] \rightarrow \mathbb{R} \mid f \text{ در } [a, b] \text{ پیوسته باشد و } n \text{ مرتبه مشتق داشته باشد}\}$$

$$C^\infty(\Omega) = \{f : X \rightarrow \mathbb{R} \mid f \text{ در } \Omega \text{ پیوسته باشد و از هر مرتبه‌ای مشتق داشته باشد}\}$$

$$C^\infty[a, b] = \{f : [a, b] \rightarrow \mathbb{R} \mid f \text{ در } [a, b] \text{ پیوسته باشد و از هر مرتبه‌ای مشتق داشته باشد}\}$$

در $\Omega \subseteq \mathbb{R}^d$ که در آن $d > 1$ ، تعاریف متناظری تعمیم داده می‌شود.

تذکر ۴.۱ توابع چند جمله‌ای، گویا، مثلثاتی، نمایی و لگاریتمی در رده $C^\infty(\Omega)$ قرار دارند، که در آن دامنه تعریف این توابع است.

قضیه ۱.۱ فرض کنید $f \in C[a, b]$. f در $[a, b]$ ، ماکزیمم و مینیمم مطلق خود را اختیار می‌کند.

قضیه ۲.۱ (مقدار میانی) فرض کنید $f \in C[a, b]$ و $f(a) \neq f(b)$ ، هر مقدار بین $f(a)$ و $f(b)$ را در نقطه‌ای بین a و b به خود می‌گیرد.

قضیه ۳.۱ فرض کنید $f \in C[a, b]$ و f بر (a, b) مشتق‌پذیر باشد. اگر $f(a) = f(b) = k$ ، آن‌گاه عدد x در (a, b) موجود است که $f'(x) = 0$. این قضیه برای $k = 0$ ، به قضیه رُل معروف است.

قضیه ۴.۱ (تعمیم قضیه رُل) فرض کنید $f \in C[a, b]$ و n مشتق مرتبه n ام f بر (a, b) موجود باشد. اگر x_0, \dots, x_n در $[a, b]$ متمایز باشند و $f(x_0) = f(x_1) = \dots = f(x_n) = 0$ ، آن‌گاه عدد x در (a, b) موجود است به طوری که $f^{(n)}(x) = 0$.

Condition-number^۹

Unstable^{۱۰}

Stable^{۱۱}

Conditionally stable^{۱۲}

Rolle^{۱۳}

قضیه ۵.۱ (مقدار میانگین) فرض کنید $f \in C[a, b]$ و f بر (a, b) مشتق‌پذیر باشد. نقطه‌ای مانند c در (a, b) موجود است به طوری که

$$f'(c) = \frac{f(b) - f(a)}{b - a}.$$

قضیه ۶.۱ (کوشی^{۱۴}) فرض کنید $f, g \in C[a, b]$ و f و g بر (a, b) مشتق‌پذیر باشند و g' در (a, b) ناصفر باشد. نقطه‌ای مانند c در (a, b) موجود است به طوری که

$$\frac{f'(c)}{g'(c)} = \frac{f(b) - f(a)}{g(b) - g(a)}.$$

قضیه ۷.۱ (مقدار میانگین برای انتگرال) فرض کنید $f \in C[a, b]$. نقطه‌ای مانند c در $[a, b]$ موجود است به طوری که

$$\int_a^b f(x)dx = (b - a)f(c).$$

قضیه ۸.۱ (مقدار میانگین تعمیم‌یافته برای انتگرال) فرض کنید $f \in C[a, b]$ و g تابعی انتگرال‌پذیر باشد که در $[a, b]$ تغییر علامت ندهد. نقطه‌ای مانند c در $[a, b]$ موجود است به طوری که

$$\int_a^b f(x)g(x)dx = f(c) \int_a^b g(x)dx.$$

قضیه ۹.۱ (اساسی حساب دیفرانسیل و انتگرال) فرض کنید $f \in C[a, b]$ و $F(x) = \int_a^x f(t)dt$ در (a, b) مشتق‌پذیر بوده و $F'(x) = f(x)$.

قضیه ۱۰.۱ (تیلور با باقیمانده لاگرانژ) فرض کنید $f \in C^{n+1}[a, b]$ و $f^{(n+1)}$ بر (a, b) موجود باشد و $x_0 \in [a, b]$. در این صورت، به ازای هر $x \in [a, b]$ نقطه‌ای مانند $\xi(x)$ بین x_0 و x وجود دارد که

$$f(x) = P_n(x) + R_n(x, x_0)$$

که در آن

$$P_n(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots + \frac{f^{(n)}(x_0)}{n!}(x - x_0)^n = \sum_{i=0}^n \frac{f^{(i)}(x_0)}{i!}(x - x_0)^i$$

و

$$R_n(x, x_0) = \frac{f^{(n+1)}(\xi(x))}{(n+1)!}(x - x_0)^{n+1}.$$

در این جا $P_n(x)$ چند جمله‌ای تیلور مرتبه n ام f حول x_0 و $R_n(x, x_0)$ جمله باقی‌مانده (یا خطای برش^{۱۵}) متناظر با $P_n(x)$ نامیده می‌شود. اگر $n \rightarrow \infty$ آن‌گاه $P_n(x)$ به یک سری بی‌پایان تبدیل می‌شود که به آن سری تیلور f حول x_0 گویند. در این حالت، شرط بی‌نهایت بار مشتق‌پذیر بودن f در x_0 لازم است.

^{۱۴}Cauchy
^{۱۵}Truncation error

تذکر ۵.۱ در قضیه قبل اگر $x_0 = 0$ آن گاه واژه تیلور به مکولرن تبدیل می شود.

تذکر ۶.۱ مقدار $R_n(x, x_0)$ مقدار خطا در استفاده از P_n به جای f را نشان می دهد. در عمل با یافتن کرانی برای جمله باقی مانده، در واقع برای خطای تقریب f با P_n کرانی پیدا می کنیم.

تذکر ۷.۱ ویژگی مهم چند جمله ای تیلور مرتبه n آن است که P_n و مشتقات تا مرتبه n آن با f و مشتقات تا مرتبه n آن در نقطه x_0 برابر هستند.

تذکر ۸.۱ شکل دیگر (کاربردی) قضیه تیلور به صورت زیر است

$$f(x+h) = f(x) + f'(x)h + \dots + \frac{f^{(n)}(x)}{n!}h^n + \frac{f^{(n+1)}(\xi)}{(n+1)!}h^{n+1} = \sum_{i=0}^n \frac{f^{(i)}(x)}{i!}h^i + \frac{f^{(n+1)}(\xi)}{(n+1)!}h^{n+1}.$$

و با

$$f(x_0) = f(x+h) = f(x) + f'(x)h + \dots + \frac{f^{(n)}(x)}{n!}h^n + \frac{f^{(n+1)}(\xi)}{(n+1)!}h^{n+1} = \sum_{i=0}^n \frac{f^{(i)}(x)}{i!}h^i + \frac{f^{(n+1)}(\xi)}{(n+1)!}h^{n+1},$$

که در آن $h = x_0 - x$.

مثال ۱۰.۱ فرض کنید $f(x) = 2 + 4x - x^3$. آن گاه به وضوح برای $n \geq 3$ داریم

$$P_n(x) = 2 + 4x - x^3, \quad R_n(x, 0) = 0,$$

و برای $n = 2$ خواهیم داشت

$$P_2(x) = 2 + 4x, \quad R_2(x, 0) = -x^3,$$

و برای $n = 1$ می توان نوشت

$$P_1(x) = 2 + 4x, \quad R_1(x, 0) = -2x^2\xi(x),$$

Δ

که در آن $\xi(x)$ نقطه ای بین 0 و x است.

قضیه ۱۱.۱ (تیلور با باقیمانده انتگرالی) اگر $f \in C^{n+1}[a, b]$ آن گاه به ازای هر $x, x_0 \in [a, b]$ داریم

$$f(x) = \sum_{i=0}^n \frac{f^{(i)}(x_0)}{i!}(x-x_0)^i + R_n(x, x_0)$$

که در آن

$$R_n(x, x_0) = \frac{1}{n!} \int_{x_0}^x f^{(n+1)}(t)(x-t)^n dt.$$

برهان. به کمک انتگرال‌گیری جزء به جزء به ازای هر n داریم

$$R_n(x) = -\frac{1}{n!} f^n(x_0)(x - x_0)^n + R_{n-1}(x).$$

با استفاده مکرر از این رابطه خواهیم داشت

$$R_n(x) = -\sum_{i=1}^n \frac{f^{(i)}(x_0)}{i!} (x - x_0)^i + R_0(x).$$

حکم با توجه به رابطه زیر برقرار است

$$R_0(x) = \int_{x_0}^x f'(t) dt = f(x) - f(x_0).$$

□

قضیه ۱۲.۱ (تیلور توابع دومتغیره) اگر $f \in C^{m+1}(\Omega)$ که در آن $\Omega \subseteq \mathbb{R}^r$ ، $\Omega = [a, b] \times [c, d]$ ، آنگاه به ازای هر نقطه $(x+h, y+k)$

در Ω داریم

$$f(x+h, y+k) = \sum_{i=0}^n \frac{1}{i!} \left(h \frac{\partial}{\partial x} + k \frac{\partial}{\partial y} \right)^i f(x, y) + R_n(h, k),$$

که در آن

$$R_n(h, k) = \frac{1}{(n+1)!} \left(h \frac{\partial}{\partial x} + k \frac{\partial}{\partial y} \right)^{n+1} f(x + \theta h, y + \theta k),$$

و θ عددی بین ۰ و ۱ است. هم‌چنین داریم

$$\begin{aligned} \left(h \frac{\partial}{\partial x} + k \frac{\partial}{\partial y} \right)^0 f(x, y) &= f(x, y), \\ \left(h \frac{\partial}{\partial x} + k \frac{\partial}{\partial y} \right)^1 f(x, y) &= \left(h \frac{\partial f}{\partial x} + k \frac{\partial f}{\partial y} \right) (x, y), \\ \left(h \frac{\partial}{\partial x} + k \frac{\partial}{\partial y} \right)^r f(x, y) &= \left(h^r \frac{\partial^r f}{\partial x^r} + r h k \frac{\partial^r f}{\partial x \partial y} + k^r \frac{\partial^r f}{\partial y^r} \right) (x, y), \dots \end{aligned}$$

تذکر ۹.۱ این قضیه برای توابع چندمتغیره به سادگی تعمیم‌پذیر است.

هنگامی که قصد داریم دو دنباله را با هم مقایسه کنیم از تعریف زیر استفاده می‌کنیم.

تعریف ۸.۱ اگر $\{\alpha_n\}$ و $\{\beta_n\}$ دو دنباله همگرا به x باشند، آنگاه بنابر تعریف داریم

$$\forall \epsilon > 0, \exists N_1 : \forall n (n \geq N_1 \rightarrow |\alpha_n - x| < \epsilon),$$

$$\forall \epsilon > 0, \exists N_2 : \forall n (n \geq N_2 \rightarrow |\beta_n - x| < \epsilon).$$

اگر به ازای هر $\epsilon > 0$ داشته باشیم $N_1 < N_2$ ، آنگاه $\{\alpha_n\}$ سریعتر از $\{\beta_n\}$ به x همگرا است.

متأسفانه چون یافتن N به سادگی امکان‌پذیر نیست، این تعریف بیشتر جنبه نظری دارد و در عمل نمی‌توان از آن سود برد. خوشبختانه در عمل تعریف زیر راه‌گشا است.

تعریف ۹.۱ فرض کنید $\{x_n\}$ و $\{\alpha_n\}$ دو دنباله متفاوت باشند. می‌نویسیم

$$x_n = O(\alpha_n),$$

و می‌خوانیم x_n آئی بزرگ α_n است هرگاه ثابت‌های C و N چنان موجود باشند که داشته باشیم

$$|x_n| \leq C |\alpha_n|, \quad \forall n \geq N.$$

اگر برای هر n داشته باشیم $\alpha_n \neq 0$ می‌توان گفت نسبت $|x_n/\alpha_n|$ کران‌دار (با کران C) باقی می‌ماند همان‌طور که $n \rightarrow \infty$. هر چند که در این تعریف $\{\alpha_n\}$ یک دنباله دلخواه است ولی در عمل، بیشتر مواقع $\alpha_n = \frac{1}{n^p}$ اختیار می‌شود که در آن $p > 0$ و در صدد یافتن بزرگ‌ترین مقدار p هستیم به طوری که

$$x_n = O\left(\frac{1}{n^p}\right).$$

گاهی مواقع می‌نویسیم

$$x_n = o(\alpha_n),$$

و می‌خوانیم x_n آئی کوچک α_n است هرگاه داشته باشیم $\lim_{n \rightarrow \infty} (x_n/\alpha_n) = 0$.

تذکر ۱۰.۱ بیشتر مواقع دو دنباله بیان‌شده در تعریف قبل همگرا به صفر هستند و اگر در این صورت داشته باشیم $x_n = O(\alpha_n)$ آن‌گاه می‌توان گفت همگرایی x_n به صفر متناسب با همگرایی α_n به صفر است و یا همگرایی x_n به صفر مانند همگرایی α_n به صفر است در حالی که اگر داشته باشیم $x_n = o(\alpha_n)$ آن‌گاه می‌توان گفت همگرایی x_n به صفر سریع‌تر از همگرایی α_n به صفر است.

مثال ۲.۱ برای $n \in \mathbb{N}$ فرض کنید

$$\alpha_n = \frac{n+1}{n^2} \quad \text{و} \quad \hat{\alpha}_n = \frac{n+3}{n^2}.$$

اگرچه $\lim_{n \rightarrow \infty} \alpha_n = 0 = \lim_{n \rightarrow \infty} \hat{\alpha}_n$ ولی $\{\hat{\alpha}_n\}$ سریع‌تر از $\{\alpha_n\}$ به صفر همگرا است. به جدول زیر نگاه کنید.

n	۱	۲	۳	۴	۵	۱۰۰
α_n	۲/۰۰۰۰۰۰	۰/۷۵۰۰۰	۰/۴۴۴۴۴	۰/۳۱۲۵۰	۰/۲۴۰۰۰	۰/۰۱۰۱
$\hat{\alpha}_n$	۴/۰۰۰۰۰۰	۰/۶۲۵۰۰	۰/۲۲۲۲۲	۰/۱۰۹۳۸	۰/۰۶۴۰۰	۰/۰۰۰۱۰۳

اگر قرار دهیم $\beta_n = 1/n$ و $\hat{\beta}_n = 1/n^2$ آن‌گاه

$$|\alpha_n - 0| = \frac{n+1}{n^2} \leq \frac{n+n}{n^2} = \frac{2}{n} = 2\beta_n,$$

و

$$|\hat{\alpha}_n - 0| = \frac{n+3}{n^2} \leq \frac{n+3n}{n^2} = \frac{4}{n} = 4\hat{\beta}_n,$$

بنابراین

$$\alpha_n = 0 + O\left(\frac{1}{n}\right) \quad \text{و} \quad \hat{\alpha}_n = 0 + O\left(\frac{1}{n^2}\right).$$

Δ

همچنین می‌توان نوشت $\hat{\alpha}_n = o(\alpha_n)$.

مثال ۳.۱ بسط تیلور $\ln(x)$ حول $x_0 = 1$ به صورت زیر است

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

در نتیجه می‌توان نوشت

$$\ln 2 - \sum_{k=1}^{n-1} (-1)^{k-1} \frac{1}{k} = O\left(\frac{1}{n}\right),$$

که نمونه‌ای از همگرایی خیلی کند است. از طرف دیگر به کمک بسط مکلورن e^x داریم

$$e - \sum_{k=0}^{n-1} \frac{1}{k!} = O\left(\frac{1}{n!}\right),$$

Δ

که مثالی از یک همگرایی خیلی تند است.

تعریف ۱.۰.۱ فرض کنید f و g دو تابع دلخواه باشند. می‌نویسیم

$$f(x) = O(g(x)), \quad (x \rightarrow \infty),$$

هرگاه ثابت‌های C و r چنان موجود باشند که داشته باشیم

$$|f(x)| \leq C |g(x)|, \quad \forall x \geq r.$$

به طور کلی می‌نویسیم

$$f(x) = O(g(x)), \quad (x \rightarrow x_0),$$

هرگاه ثابت C و یک همسایگی حول x_0 چنان موجود باشند که به ازای هر x در آن همسایگی داشته باشیم

$$|f(x)| \leq C |g(x)|.$$

همچنین

$$f(x) = o(g(x)), \quad (x \rightarrow x_0),$$

به این معنی است که $\lim_{x \rightarrow x_0} f(x)/g(x) = 0$.

مثال ۴.۱ چون برای $x \geq 1$ داریم $\sqrt{x^2 + 1} \leq \sqrt{2}x$ پس می‌توان نوشت

$$\sqrt{x^2 + 1} = O(x), \quad (x \rightarrow \infty).$$

Δ

مثال ۵.۱ با اعمال قضیه مکورن برای تابع $\cos(x)$ خواهیم داشت

$$\cos h = 1 - \frac{1}{2}h^2 + \frac{1}{24}h^4 \cos \xi(h),$$

که در آن $\xi(h)$ عددی بین 0 و h است. در نتیجه

$$\cos h + \frac{1}{2}h^2 = 1 + \frac{1}{24}h^4 \cos \xi(h),$$

و یا می‌توان نوشت

$$\cos h + \frac{1}{2}h^2 = 1 + O(h^4),$$

زیرا $\frac{1}{24}h^4 \leq |\cos \xi(h)| h^4 \leq \frac{1}{24}h^4$ در نتیجه زمانی که $h \rightarrow 0$ ، $\cos h + \frac{1}{2}h^2$ سریعتر به ۱ میل می‌کند تا h^4 به صفر. Δ

تمرین ۱۰.۱ درستی روابط زیر را بررسی کنید

$$\frac{1}{n \ln n} = o\left(\frac{1}{n}\right),$$

$$\sin h = h - \frac{h^3}{6} + O(h^5), \quad (h \rightarrow 0),$$

$$e^{-1/h} = o(h^2), \quad (h \rightarrow \infty).$$

تعریف ۱۱.۱ فرض کنید $\{x_n\}$ دنباله‌ای همگرا به x باشد و برای هر n داشته باشیم $x_n \neq x$. این دنباله به طور خطی به x همگرا است هرگاه ثابت $C \in (0, 1)$ چنان وجود داشته باشد که

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - x|}{|x_n - x|} = C.$$

در این حالت C به نرخ همگرایی^{۱۶} معروف است. برای $C = 0$ ($C = 1$) همگرایی زیرخطی^{۱۷} (زیرخطی^{۱۸}) نامیده می‌شود. اگر اعداد مثبت و ثابت C و $p > 1$ وجود داشته باشند به طوری که

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1} - x|}{|x_n - x|^p} = C$$

آن‌گاه مرتبه همگرایی^{۱۹} این دنباله p است. C به ثابت خطای مجانبی^{۲۰} معروف است. اگر $p = 2$ ($p = 3$) همگرایی مربعی (مکعبی) نامیده می‌شود ($C < 1$ لزومی ندارد).

^{۱۶}Rate of convergence

^{۱۷}Superlinear

^{۱۸}Sublinear

^{۱۹}Order of convergence

^{۲۰}Asymptotic error constant

مثال ۶.۱ فرض کنید $\{x_n\}$ دنباله‌ای با همگرایی خطی به صفر با

$$\lim_{n \rightarrow \infty} \frac{|x_{n+1}|}{|x_n|} = \circ/5,$$

و $\{\hat{x}_n\}$ دنباله‌ای با همگرایی مربعی به صفر با همان ثابت خطای مجانبی باشد یعنی

$$\lim_{n \rightarrow \infty} \frac{|\hat{x}_{n+1}|}{|\hat{x}_n|^2} = \circ/5.$$

برای سادگی فرض می‌کنیم برای $n \geq N$ داریم

$$\frac{|x_{n+1}|}{|x_n|} \simeq \circ/5, \quad \frac{|\hat{x}_{n+1}|}{|\hat{x}_n|^2} \simeq \circ/5.$$

بنابراین می‌توان نوشت

$$|x_n - \circ| = |x_n| \simeq \circ/5 |x_{n-1}| \simeq (\circ/5)^2 |x_{n-2}| \simeq \dots \simeq (\circ/5)^n |x_0|,$$

و

$$\begin{aligned} |\hat{x}_n - \circ| = |\hat{x}_n| &\simeq \circ/5 |\hat{x}_{n-1}|^2 \simeq (\circ/5)(\circ/5 |\hat{x}_{n-2}|^2)^2 = (\circ/5)^3 |\hat{x}_{n-2}|^4 \\ &\simeq (\circ/5)^2 (\circ/5 |\hat{x}_{n-2}|^2)^4 = (\circ/5)^4 |\hat{x}_{n-2}|^8 \\ &\simeq \dots \simeq (\circ/5)^{2^{n-1}} |\hat{x}_0|^{2^n}. \end{aligned}$$

Δ

ادامه این مثال را در برنامه oc.nb دنبال کنید.

مثال ۷.۱ دنباله $a_n = \frac{1}{\sqrt[n]{n}}$ با نرخ همگرایی $\frac{1}{\sqrt[n]{n}}$ ، همگرایی خطی به صفر دارد. همگرایی دنباله $b_n = \frac{1}{\sqrt[n]{n}}$ به صفر زیرخطی است. در اصل همگرایی این دنباله مربعی است. دنباله $c_n = \frac{1}{n+1}$ همگرایی زیرخطی به صفر دارد.

Δ

تمرین ۲.۱ در مورد مرتبه و نرخ همگرایی دنباله $a_n = \frac{1}{\sqrt[n]{n}}$ چه می‌توان گفت؟

تمرین ۳.۱ نشان دهید در زیر، مرتبه همگرایی دقیقتر تعریف شده است.

فرض کنید $\{x_n\}$ دنباله‌ای همگرا به x باشد و برای هر n داشته باشیم $x_n \neq x$. اگر اعداد مثبت و ثابت C و $p > 1$ وجود داشته باشند به طوری که برای n به اندازه کافی بزرگ داشته باشیم

$$|x_{n+1} - x| \leq C |x_n - x|^p,$$

آنگاه گوئیم مرتبه همگرایی این دنباله p است.

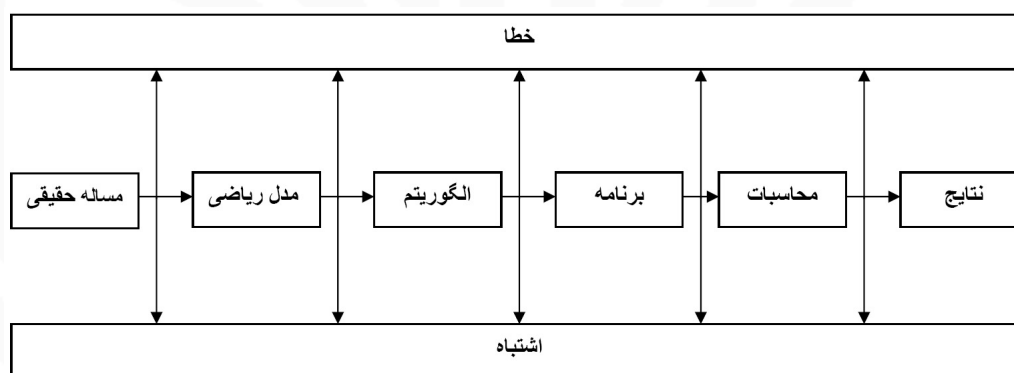
فصل ۲

خطاها

زمانی که به دست آوردن جواب واقعی یک مسئله غیرممکن است و یا مقرون به صرفه نیست، سعی می‌کنیم به کمک روش‌های عددی، یک جواب تقریبی برای مسئله پیدا کنیم. این فرایند منجر به تولید خطا^۱ می‌شود. در این فصل قصد داریم منابع تولید خطا و انواع خطا را شناسایی کرده و تا حدی از انتشار خطا جلوگیری کنیم.

۱.۲ منابع تولید خطا

بیشتر مواقع در عمل با یک مسئله حقیقی (فیزیکی) مواجه هستیم و بنا به دلایلی، جواب عددی (تقریبی) آن را جستجو می‌کنیم. مراحل یافتن جواب عددی چنین مسئله‌ای در روندنمای آمده در شکل ۱.۲ خلاصه می‌شود. این روندنما مکان‌های احتمالی بروز خطا و اشتباه^۲ را نیز نشان می‌دهد. مراحل مختلف این روندنما را در مثال بعد دنبال می‌کنیم.



شکل ۱.۲: فرایند تولید جواب عددی (تقریبی)

مثال ۱.۲ (مسئله حقیقی) می‌خواهیم دوره تناوب حرکت نوسانی و متناوب یک آونگ ساده به جرم m و طول l را به دست آوریم. فرض کنید $\theta(t)$ جابجایی زاویه‌ای آونگ از خط قائم در زمان t باشد. صرف نظر از مقاومت هوا و اصطکاک در لولا، مدل این مسئله به صورت

Error^۱
Bug^۲

$$ml \frac{d^2 \theta}{dt^2} = -mg \sin \theta$$

$$\frac{d^2 \theta}{dt^2} = -\frac{g}{l} \sin \theta, \quad (1.2)$$

در می‌آید که یک معادله دیفرانسیل غیرخطی است. با فرض کوچک بودن θ یعنی

$$\theta = 6^\circ \simeq 0.1047 \text{ rad}, \quad \sin \theta \simeq 0.1045, \quad \theta = 15^\circ \simeq 0.2618 \text{ rad}, \quad \sin \theta \simeq 0.2598,$$

می‌توان فرض کرد $\sin \theta \simeq \theta^{\text{rad}}$ و معادله دیفرانسیل را به صورت $\omega^2 \theta = 0$ نوشت که در آن $\omega^2 = \frac{g}{l}$. این معادله دیفرانسیل خطی

جوابی به صورت

$$\theta(t) = A \sin \omega t + B \cos \omega t,$$

دارد که از آن نتیجه می‌گیریم $T_L = \frac{2\pi}{\omega} = 2\pi \sqrt{\frac{l}{g}}$. از طرف دیگر با حل معادله دیفرانسیل غیرخطی (۱.۲) خواهیم داشت

$$T_N = 2\pi \sqrt{\frac{l}{g}} \left(1 + \frac{1^2}{2^2} \sin^2 \frac{\theta_m}{2} + \frac{1^2 \times 3^2}{2^2 \times 4^2} \sin^4 \frac{\theta_m}{2} + \dots \right),$$

که در آن θ_m جابجایی زاویه‌ای ماکزیمم آونگ است. جالب است توجه داشته باشیم که T_L اولین جمله سری T_N است. در ادامه الگوریتم ۱.۲ را برای این مدل طراحی می‌کنیم.

الگوریتم ۱.۲ الگوریتم یافتن دوره تناوب حرکت یک آونگ ساده.

• ورودی: l ، θ_m و مرز خطا یعنی ϵ

• خروجی: مقدار T_L ، T_N و اختلاف آن‌ها

$$(1) \text{ قرار دهید } g = 9.8, p = 3.14, \text{ و } T_L = 2p \sqrt{\frac{l}{g}}$$

$$(2) \text{ قرار دهید } k = 1, M = \left(\frac{2k-1}{2k} \right)^2, T_1 = T_L, C = T_L \times M \times \sin^{2k} \frac{\theta_m}{2}, \text{ و } T_2 = T_1 + C$$

(3) تا زمانی که $|C| > \epsilon$ گام ۴ را تکرار کنید

$$(4) \text{ قرار دهید } k = k + 1, M = M \times \left(\frac{2k-1}{2k} \right)^2, T_1 = T_2, C = T_L \times M \times \sin^{2k} \frac{\theta_m}{2}, \text{ و } T_2 = T_1 + C$$

(5) قرار دهید $T_N = T_2$ و T_L ، T_N و $T_N - T_L$ را چاپ کرده، متوقف شوید

حال برای الگوریتم ۱.۲ برنامه ۱۲ را با استفاده از نرم‌افزار Mathematica می‌نویسیم. پس از اجرای برنامه ۱۲ در محیط نرم‌افزار Mathematica محاسبات شروع شده و به ازای $\epsilon = 0.0001$ نتایج گزارش شده در جدول ۱.۲ تولید می‌شود. باید توجه داشت که

△

برای $\theta_m < 15^\circ$ اختلاف T_N و T_L از 0.0005 بیشتر نیست.

```

l = Input["طول آونگ را بر حسب متر وارد کنید"];
Thetam = Input["جایجایی زاویه‌ای ماکزیمم را بر حسب درجه وارد کنید"];
Epsilon = Input["مقدار مجاز خطا را وارد کنید"];
st = Sin[Thetam * Pi/۳۶۰]; g = ۹/۸; TL = ۲ * Pi * √(l/g); k = ۱;
M = ((۲k-1)/۲k)²; CC = TL * M * st²k; T۱ = TL; T۲ = T۱ + CC;
While[CC > Epsilon, {
  k++; M = ((۲k-1)/۲k)²; T۱ = T۲; CC = TL * M * st²k; T۲ = T۱ + CC; }];
TN = T۲;
Print["TL = ", TL, " TN = ", TN, " TN - TL = ", TN - TL];

```

برنامه ۱۰۲: برنامه یافتن دوره تناوب حرکت یک آونگ ساده.

با توجه به روندنمای ارایه شده در شکل ۱۰۲ و مثال ۱۰۲، خطاها از نظر منابع تولید به صورت زیر تقسیم‌بندی می‌شوند

$l(cm)$	θ_m	T_L	T_N	$T_N - T_L$	$l(cm)$	θ_m	T_L	T_N	$T_N - T_L$
۱۰	۳	۰/۶۳۴۷۰	۰/۶۳۴۸۱	۰/۰۰۰۱۱	۱۰	۱۲	۰/۶۳۴۷۰	۰/۶۳۶۴۴	۰/۰۰۱۷۴
۲۰	۳	۰/۸۹۷۶۰	۰/۸۹۷۷۵	۰/۰۰۰۱۵	۲۰	۱۲	۰/۸۹۷۶۰	۰/۹۰۰۰۷	۰/۰۰۲۴۷
۳۰	۳	۱/۰۹۹۳۳	۱/۰۹۹۵۲	۰/۰۰۰۱۹	۳۰	۱۲	۱/۰۹۹۳۳	۱/۱۰۲۳۵	۰/۰۰۳۰۲
۱۰	۶	۰/۶۳۴۷۰	۰/۶۳۵۱۳	۰/۰۰۰۴۳	۱۰	۱۵	۰/۶۳۴۷۰	۰/۶۳۷۴۳	۰/۰۰۲۷۳
۲۰	۶	۰/۸۹۷۶۰	۰/۸۹۸۲۱	۰/۰۰۰۶۱	۲۰	۱۵	۰/۸۹۷۶۰	۰/۹۰۱۴۶	۰/۰۰۳۸۶
۳۰	۶	۱/۰۹۹۳۳	۱/۱۰۰۰۸	۰/۰۰۰۷۵	۳۰	۱۵	۱/۰۹۹۳۳	۱/۱۰۴۰۶	۰/۰۰۴۷۳
۱۰	۹	۰/۶۳۴۵۰	۰/۶۳۵۷۰	۰/۰۰۱۲۰	۱۰	۱۸	۰/۶۳۴۷۰	۰/۶۳۸۶۳	۰/۰۰۳۹۳
۲۰	۹	۰/۸۹۷۶۰	۰/۸۹۹۰۰	۰/۰۰۱۴۰	۲۰	۱۸	۰/۸۹۷۶۰	۰/۹۰۳۱۷	۰/۰۰۵۵۷
۳۰	۹	۱/۰۹۹۳۳	۱/۱۰۱۰۳	۰/۰۰۱۷۰	۳۰	۱۸	۱/۰۹۹۳۳	۱/۱۰۶۱۵	۰/۰۰۶۸۲

جدول ۱۰۲: دوره تناوب آونگ ساده با طول‌های مختلف

۱. ذاتی

- مدل (ناشی از صرف‌نظرها، چشم‌پوشی‌ها و ساده‌سازی‌ها مانند فرض $(\sin \theta \simeq \theta^{rad})$)
- داده‌های مدل (ناشی از آزمایشات و اندازه‌گیری‌ها مثل g, l)

۲. محاسباتی

- نمایش اعداد (مانند $\pi = ۳/۱۴$)
- اعمال ریاضی (به عنوان مثال $(\frac{l}{g})$)
- روش‌های (الگوریتم‌های) عددی (محاسباتی) (مثل خطای روش محاسبه $(\sqrt{\frac{l}{g}})$)

تذکر ۱۰۲: خطاهای ذاتی به متخصص آنالیز عددی مربوط نمی‌شود اما برای پرهیز از خطاهای محاسباتی و کنترل آن‌ها باید راه چاره‌ای پیدا

تذکر ۲.۲ اختلاف جواب واقعی و تقریبی از اشتباهات و خطاها ناشی می‌شود. اشتباه را می‌توان برطرف نمود ولی خطا ممکن است اجتناب‌ناپذیر باشد. به عنوان مثال قراردادن ۲۳۳۳ به جای ۳۲۳۲ یک اشتباه است حال آن که عدد π بسط دهدهی نامختوم دارد و در وسایل محاسباتی باید یک بسط دهدهی مختوم برای آن در نظر گرفت که این موجب بروز خطا می‌شود.

تذکر ۳.۲ برای اطلاعات بیشتر در خصوص اثرات مخرب خطاها و اشتباهات به آدرس‌های زیر مراجعه کنید.

<https://www.iro.umontreal.ca/mignotte/IFT2425/Disasters.html>

<https://www.computerworld.com/article/3412197/top-software-failures-in-recent-history.html>

در اینجا فقط به دو مورد زیر اشاره می‌شود.

- عدم موفقیت موشک پاتریوت در جنگ خلیج فارس سال ۱۹۹۱ (۲۸ کشته و ۱۰۰ زخمی) به دلیل وقوع خطای گرد کردن در محاسبات

مسیر

- شکست ماموریت موشک آریان ۵ فرانسه در سال ۱۹۹۶ (۵۰۰ میلیون دلار خسارت مادی) به دلیل وقوع پاریز^۳ در رایانه آن

۲.۲ نمایش اعداد

در این بخش به بررسی نمایش اعداد حقیقی می‌پردازیم. اثبات برخی از قضایا را می‌توان در سایر مراجع آنالیز عددی یافت.

تعریف ۱.۲ هر عدد حقیقی مثبت x نمایشی معروف به نمایش مکانی^۴ به صورت

$$\begin{aligned} x &= a_m \beta^m + a_{m-1} \beta^{m-1} + \dots + a_1 \beta^1 + a_0 \beta^0 + a_{-1} \beta^{-1} + a_{-2} \beta^{-2} + \dots \\ &= (a_m a_{m-1} \dots a_1 a_0 a_{-1} a_{-2} \dots)_\beta, \end{aligned} \quad (2.2)$$

دارد که در آن $\beta > 1$ یک عدد طبیعی معروف به مبنای β ، $m \in \mathbb{Z}$ و $a_i \in \{0, 1, 2, \dots, \beta - 1\}$ برای یکتایی نمایش لازم است که $a_m \neq 0$ و عدد صحیح j چنان وجود داشته باشد که $a_{j-1} = a_j = 0, \dots$ یا به ازای هر N به اندازه کافی بزرگ وجود داشته باشد $N \leq j$ به طوری که $a_{-j} \neq \beta - 1$.

تذکر ۴.۲ رابطه (۲.۲) به نمایش (بسط) عدد x در مبنای β معروف است و اگر $\beta = ۱۰$ اختیار شود به آن، نمایش (بسط) دهدهی (اعشاری) گویند (متداول در زندگی روزمره) و در حالتی که $\beta = ۲$ در نظر گرفته شود به آن، نمایش دودویی (باینری) گفته می‌شود (مبنای کار رایانه).

تذکر ۵.۲ شرط دوم یکتایی نمایش، لازم است چه در غیر این صورت به عنوان مثال می‌توان نوشت

$$\begin{aligned} ۳/۴۷۹۹۹\dots &= ۳/۴۷\bar{9} = ۳ \times ۱۰^0 + ۴ \times ۱۰^{-1} + ۷ \times ۱۰^{-2} + ۹ \times ۱۰^{-3} + ۹ \times ۱۰^{-4} + \dots \\ &= ۳/۴۷ + \frac{۹ \times ۱۰^{-2}}{۱ - ۱۰^{-1}} = ۳/۴۷ + ۰/۰۱ = ۳/۴۸. \end{aligned}$$

قضیه ۱.۲ بسط یک عدد گویا در هر مبنایی یا مختوم است یا نامختوم متناوب.

^۳ Underflow

^۴ Positional representation

نتیجه ۱.۱.۲ بسط عدد گنگ، نامختوم نامتناوب است.

مثال ۲.۲

$$\frac{1}{3} = 0.\overline{3} = (0/1)_2, \quad \frac{3}{8} = 0.375, \quad \frac{1}{10} = 0/1 = (0/\overline{0011})_2, \quad \sqrt{2} = 1.4142\dots, \quad \pi = 3.141592\dots$$

△

قضیه ۲.۲ به ازای هر $\epsilon > 0$ و هر عدد حقیقی x عدد حقیقی y چنان وجود دارد که $|x - y| < \epsilon$.

هنگام کار با رایانه (ماشین حساب) اعداد را در مبنای ۱۰ وارد کرده و انتظار داریم نتایج (خروجی) نیز در همین مبنا نمایش داده شود، ولی این وسایل با مبنای دیگری (امروزه مبنای ۲ و در قدیم مبنای ۱۶ نیز استفاده می‌شده) کار می‌کنند. بنابراین مسئله تغییر مبنا مطرح می‌شود.

پرسش ۱.۲ مبنای ۲ در مقایسه با مبنای ۱۶ چه معایب و مزایایی دارد؟

تغییر مبنا از مبنای ۲ به مبنای ۱۰ به کمک بسط به سادگی امکان‌پذیر است.

مثال ۳.۲

$$\begin{aligned} (1011001/111001)_2 &= 2^6 + 2^4 + 2^3 + 2^0 + 2^{-1} + 2^{-2} + 2^{-3} + 2^{-6} \\ &= (64 + 16 + 8 + 1 + 0.5 + 0.25 + 0.125 + 0.015625) \\ &= 89.890625. \end{aligned}$$

△

برای تغییر مبنا از ۱۰ به ۲، اگر عدد داده شده در مبنای ۱۰، صحیح مثبت باشد، با تقسیم‌های پشت سر هم بر ۲ عدد متناظر در مبنای ۲ به دست می‌آید.

مثال ۴.۲

$$23 = 11 \times 2 + 1, \quad 11 = 5 \times 2 + 1, \quad 5 = 2 \times 2 + 1, \quad 2 = 1 \times 2 + 0.$$

△

$$\text{پس } 23 = (10111)_2$$

اگر عدد داده شده x در مبنای ۱۰ در فاصله $(0, 1)$ باشد و عدد متناظر در مبنای ۲ به صورت $(0/b_1 b_2 b_3 \dots)_2$ فرض شود آن‌گاه

$$x = b_1 \times 2^{-1} + b_2 \times 2^{-2} + b_3 \times 2^{-3} + \dots,$$

و در نتیجه

$$2x = b_1 + b_2 \times 2^{-1} + b_3 \times 2^{-2} + \dots$$

با قرار دادن

$$y = b_2 \times 2^{-1} + b_3 \times 2^{-2} + \dots,$$

خواهیم داشت

$$0 \leq y < 1 \times 2^{-1} + 1 \times 2^{-2} + \dots = \frac{2^{-1}}{1-2^{-1}} = 1.$$

پس $[y] = 0$. از طرف دیگر داریم $2x = b_1 + y$ و بلافاصله از $[2x] = [b_1 + y]$ نتیجه می‌شود $b_1 = [2x]$ زیرا b_1 عددی صحیح است (b_1 یا صفر است یا یک). با قرار دادن $x = y$ و تکرار این روند می‌توان b_2, b_3, \dots را به دست آورد. شرط توقف این روند تکراری زمانی است که x صفر شود یا تناوب به وجود آید و یا ممکن است بسط نامختوم نامتناوب باشد و پس از چند تکرار مجبور شویم روند را متوقف کنیم.

مثال ۵.۲ این روند را برای حالت‌هایی که نمایش، مختوم یا نامختوم متناوب است دنبال می‌کنیم.

با توجه به جدول ۲.۲ می‌توان نوشت $(0.100101)_2 = 0.578125$ و $(0.00011)_2 = 0.1$.

i	x	$2x$	b_i
۱	۰/۵۷۸۱۲۵	۱/۱۵۶۲۵	۱
۲	۰/۱۵۶۲۵	۰/۳۱۲۵	۰
۳	۰/۳۱۲۵	۰/۶۲۵	۰
۴	۰/۶۲۵	۱/۲۵	۱
۵	۰/۲۵	۰/۵	۰
۶	۰/۵	۱/۰	۱
	۰		

i	x	$2x$	b_i
۱	۰/۱	۰/۲	۰
۲	۰/۲	۰/۴	۰
۳	۰/۴	۰/۸	۰
۴	۰/۸	۱/۶	۱
۵	۰/۶	۱/۲	۱
۶	۰/۲	۰/۴	۰
	۰/۴		

جدول ۲.۲: نمایش دودویی اعداد کوچک‌تر از واحد

تمرین ۱.۲ برنامه‌ای بنویسید که عدد حقیقی x و عدد طبیعی n را از ورودی گرفته و نمایش دودویی عدد x را تولید کند و اگر نمایش نامختوم نامتناوب است حداکثر تا n رقم نمایش، چاپ شود.

تعریف ۲.۲ برای نمایش اعداد در ماشین، ابتدا نمایشی معروف به نمایش ممیز ثابت^۵ در نظر گرفته شد که در آن نظیر برخی از اعداد حقیقی مثبت x می‌توان یک عدد ماشینی به صورت زیر در نظر گرفت

$$\bar{x} = (a_{n-1} \dots a_1 a_0 / a_{-1} a_{-2} \dots a_{-t})_\beta,$$

که در آن $\beta > 1$ یک عدد طبیعی معروف به مبنا، $a_i \in \{0, 1, 2, \dots, \beta - 1\}$ و t و n اعداد طبیعی مشخص و ثابتی هستند. در اصل در این نمایش مکان ممیز مشخص و ثابت است و فرض می‌شود a_i ها همگی صفر نیستند.

مجموعه تمام اعداد ماشینی با نمایش ممیز ثابت با $\mathbb{F}_\circ(\beta, n, t) = \mathbb{F}_\circ$ نمایش داده می‌شود و برخی از نکات مهم در این نمایش عبارتند از

$$\text{Card}(\mathbb{F}_\circ) = 2 \times \beta^n \times \beta^t - 1.$$

• با این فرض که $\gamma = \beta - 1$ ، بزرگترین و کوچکترین عدد عبارتند از

$$x_{\max} = (\gamma \dots \gamma / \gamma \dots \gamma)_\beta = \gamma \sum_{i=-t}^{n-1} \beta^i = \beta^n - \beta^{-t} \simeq \beta^n, \quad x_{\min} = (0 \dots 0 / 0 \dots 0 1)_\beta = \beta^{-t},$$

- اعداد در مجموعه $[-x_{\max}, -x_{\min}] \cup [x_{\min}, x_{\max}]$ هم فاصله و با فاصله β^{-t} قرار دارند،
- برای نمایش اعداد بسیار بزرگ (کوچک) در این نمایش با مشکل مواجه می‌شویم و در نتیجه این نمایش برای محاسبات علمی مناسب نیست ولی برای بسیاری از کاربردها مانند حسابداری، این نمایش سودمند است و هنوز هم ماشین‌هایی بر این اساس ساخته می‌شوند.

تعریف ۳.۲ هر عدد حقیقی مثبت x نمایشی معروف به نمایش علمی^۶ به صورت

$$x = ({}^{\circ} / a_1 \dots a_t a_{t+1} \dots)_{\beta} \times \beta^e,$$

دارد که در آن $\beta > 1$ یک عدد طبیعی معروف به مبنا و $a_i \in \{0, 1, 2, \dots, \beta - 1\}$ ، $e \in \mathbb{Z}$ به توان (نما) معروف است و $m = {}^{\circ} / a_1 \dots a_t a_{t+1} \dots$ ماننیس (جزء کسری) نام دارد و $1/\beta \leq m < 1$ و برای یکتایی نمایش لازم است که $a_1 \neq 0$ (در این حالت گفته می‌شود نمایش نرمال شده است).

نمایش مکانی نقطه شروعی برای تعریف نمایش ممیز ثابت بود حال آنکه برای تعریف نمایش پر کاربرد ممیز شناور^۷ از نمایش علمی ایده گرفته شد. در این نمایش اعداد بسیار بزرگ (کوچک) به سادگی قابل نمایش هستند و یک روش متداول برای نمایش اعداد در رایانه است که از بدو پیدایش مورد توجه سازندگان سخت‌افزار رایانه قرار گرفته است و چنین به نظر می‌رسد که تا حدودی به طور سلیقه‌ای با آن برخورد شده است.

تعریف ۴.۲ نظیر برخی از اعداد حقیقی مثبت x نمایشی معروف به نمایش ممیز شناور به صورت

$$\hat{x} = ({}^{\circ} / a_1 \dots a_t)_{\beta} \times \beta^e = m \times \beta^{e-t},$$

تعریف می‌شود که در آن $\beta > 1$ یک عدد طبیعی معروف به مبنا، e یعنی توان به زیرمجموعه‌ای از \mathbb{Z} تعلق دارد، $t \in \mathbb{N}$ تعداد ارقام بامعناى ماننیس یعنی ${}^{\circ} / a_1 \dots a_t$ را نشان می‌دهد، $a_i \in \{0, 1, 2, \dots, \beta - 1\}$ و برای یکتایی نمایش لازم است که $a_1 \neq 0$ (نمایش نرمال شده). به وضوح برای $m = (a_1 \dots a_t)_{\beta}$ داریم $\beta^{t-1} - 1 \leq m \leq \beta^t - 1$ و اعداد صحیح L و U چنان وجود دارند که $L \leq e \leq U$. همچنین a_1 به باارزش‌ترین رقم و a_t به کم‌ارزش‌ترین رقم معروف است.

مجموعه تمام اعداد ماشینی با نمایش ممیز شناور همراه با صفر با $\mathbb{F} = \mathbb{F}(\beta, t, L, U)$ نمایش داده می‌شود و برخی از نکات مهم در این نمایش عبارتند از

$$\text{Card}(\mathbb{F}) = 2 \times (\beta - 1) \times \beta^{t-1}(U - L + 1) + 1 \quad \bullet$$

- با این فرض که $\gamma = \beta - 1$ ، بزرگترین و کوچکترین عدد عبارتند از

$$x_{\max} = ({}^{\circ} / \gamma \dots \gamma)_{\beta} \times \beta^U = \gamma \beta^U \sum_{i=1}^t \beta^{-i} = \beta^U (1 - \beta^{-t}) \simeq \beta^U, \quad x_{\min} = ({}^{\circ} / 1 0 \dots 0)_{\beta} \times \beta^L = \beta^{L-1},$$

- اعداد در مجموعه \mathbb{F} هم فاصله نیستند و فاصله آنها به توان بستگی دارد. هرچه e بزرگتر (کوچکتر) باشد فاصله اعداد بیشتر (کمتر) است.

باید توجه داشت که از یک طرف صفر به مجموعه اعداد ممیز شناور اضافه شده است و از طرف دیگر در یک همسایگی حول صفر هیچ عددی قابل نمایش نیست. به منظور برطرف کردن این مشکلات، تعریف زیر مطرح می‌شود.

تعریف ۵.۲ اگر در تعریف اعداد با نمایش ممیز شناور $a_1 = 0$ فقط برای $e = L$ مجاز باشد، مجموعه بزرگتری از اعداد ماشینی معروف به اعداد با نمایش غیرنرمال شده ممیز شناور^۸ به دست می آید که با $\mathbb{F}_D = \mathbb{F}_D(\beta, t, L, U)$ نمایش داده می شود.

باید توجه داشت که صفر در این نمایش وجود دارد، یکنایی نمایش حفظ شده است و به علاوه $\mathbb{F} \subset \mathbb{F}_D$. برخی از نکات مهم در این نمایش عبارتند از

$$\text{Card}(\mathbb{F}_D) = \text{Card}(\mathbb{F}) + 2(\beta^{t-1} - 1) \bullet$$

$$x_{\min} = \beta^{L-t} \bullet$$

• می توان اعدادی با ماتیس بین ۱ (متناظر با حالت $0 \dots 01$) و $1 - \beta^{t-1}$ (متناظر با حالت $0\gamma \dots 0\gamma$) که به بازه $(-\beta^{L-1}, \beta^{L-1})$ تعلق دارند نیز نمایش داد.

مثال ۶.۲ مجموعه اعداد $\mathbb{F}_D(2, 3, -1, 2)$ را مشخص کنید.

در مجموعه $\mathbb{F}(2, 3, -1, 2)$ ، $16 = (2-1) \times 2^{2-1} \times (2+1+1)$ عدد مثبت وجود دارد که عبارتند از

$$x_{\max} = (0/111)_2 \times 2^2 = \frac{7}{4}, \quad (0/110)_2 \times 2^2 = 3, \quad (0/101)_2 \times 2^2 = \frac{5}{2}, \quad (0/100)_2 \times 2^2 = 2, \\ \frac{7}{4}, \quad \frac{3}{2}, \quad \frac{5}{4}, \quad 1, \quad \frac{7}{8}, \quad \frac{3}{4}, \quad \frac{5}{8}, \quad \frac{1}{2}, \quad \frac{7}{16}, \quad \frac{3}{8}, \quad \frac{5}{16}, \quad x_{\min} = (0/100)_2 \times 2^{-1} = \frac{1}{4}.$$

اعداد غیرنرمال شده نامنفی عبارتند از

$$(0/011)_2 \times 2^{-1} = \frac{3}{16}, \quad (0/010)_2 \times 2^{-1} = \frac{1}{8}, \quad (0/001)_2 \times 2^{-1} = \frac{1}{16}, \quad (0/000)_2 \times 2^{-1} = 0.$$

تعریف ۶.۲ Δ فاصله عدد ۱ و عدد قابل نمایش بعد از آن به اپسیلون ماشین معروف است و با ϵ_M نمایش داده می شود و به عنوان مثال در $\mathbb{F}(\beta, t, L, U)$ مقدار آن به صورت زیر به دست می آید

$$(0/10 \dots 01)_\beta \times \beta^1 - (0/10 \dots 0)_\beta \times \beta^1 = (0/0 \dots 01)_\beta \times \beta^1 = \beta^{1-t}.$$

تنوع انتخاب کمیت های β, t, L و U موجب شد که ماشین های محاسباتی (رایانه های) گوناگونی تولید شوند که در برخی موارد نتایج اجرای یک برنامه بر روی ماشین های مختلف، متفاوت بود. به همین دلیل IEEE^۹ در سال ۱۹۸۵ استاندارد ۷۵۴ را وضع و سازندگان این وسایل را به پیروی از آن ترغیب کرد. این استاندارد در سال ۱۹۸۹ توسط IEC^{۱۰} تأیید شد و به صورت استاندارد IEC559 مطرح گردید. در این استاندارد دو قالب^{۱۱} موسوم به دقت ساده^{۱۲} و دقت دو برابر^{۱۳} برای اعداد با نمایش ممیز شناور پیشنهاد شده است.

۱.۲.۲ نمایش ۶۴-بیتی ممیز شناور

این نمایش پیش از این به دقت دو برابر (مضاعف) معروف بوده و متناظر با نوع double در زبان C است. در این نمایش برای هر عدد در مبنای ۲، ابتدا یک ساختار به طول ۶۴ بیت در نظر گرفته می شود. اولین بیت به علامت معروف است و با s نمایش داده می شود و بلافاصله بعد از

Floating point de-normalized^۸
 Institute of Electrical and Electronics Engineers^۹
 International Electrical Commission^{۱۰}
 Format^{۱۱}
 Single precision^{۱۲}
 Double precision^{۱۳}

آن ۱۱ بیت برای مشخصه^{۱۴} در نظر گرفته می‌شود و با c نمایش داده می‌شود و ۵۲ بیت باقی‌مانده برای مانتیس کنار گذاشته می‌شود و آن را با f نشان می‌دهند.

سپس برای هر عدد نمایشی به صورت $(-1)^s \times 2^e \times (1 + f)$ در نظر گرفته می‌شود و توان به کمک رابطه $e = c - 1023$ به دست می‌آید. با توجه به محدودیت $2047 = 2^{11} - 1 = (110001)_2 = 1 \leq c \leq (110001)_2 = 2047$ داریم $0 \leq c \leq 1024$ و $-1023 \leq e \leq 1023$. از طرف دیگر محدودیت $(0/110001)_2 \leq f \leq (0/110001)_2$ باعث می‌گردد که $2 < (1/110001)_2 \leq (1 + f)_2 \leq 0.1$ به نظر می‌رسد برای مانتیس ۵۲ بیت در اختیار داریم ولی در واقع ساختار $1 + f$ موجب می‌شود که ۵۳ بیت داشته باشیم که آن یک بیت اضافه به بیت پنهان معروف است و برای یکنایی نمایش لازم است. اگر کوچک‌ترین و بزرگ‌ترین عدد مثبت قابل نمایش را به ترتیب با mN و MN نشان دهیم آن‌گاه

$$mN = 2^{-1022} \times (1/00000)_2 \simeq 2/22507 \times 10^{-208},$$

$$MN = 2^{1024} \times (1/00000)_2 \simeq 1/79769 \times 10^{208}.$$

تذکر ۶.۲ در این نمایش $\epsilon_M = 2^{-52} \simeq 2/220446 \times 10^{-16}$ (همان eps در نرم‌افزار MATLAB) و از توان -1023 برای نمایش صفر و پیام پاریز^{۱۵} (زمانی که اعدادی با اندازه خیلی کوچک تولید می‌شوند) و از توان 1024 با مانتیس مثبت برای نمایش ∞ (inf) در MATLAB، صور مبهم (NaN در MATLAB) و پیام سرریز^{۱۶} (زمانی که اعدادی با اندازه خیلی بزرگ تولید می‌شوند) استفاده می‌شود.

در ادامه به دنبال بزرگ‌ترین عدد صحیح مثبت M هستیم که هر عدد صحیح x با شرط $0 \leq x \leq M$ در این نمایش به طور دقیق قابل نمایش باشد. به وضوح، تمام اعداد صحیح نامنفی که بزرگتر از $2^{52} \times (1/110001)_2 = 2^{52} - 1$ نباشند به طور دقیق قابل نمایش هستند و به علاوه 2^{52} نیز به صورت $2^{52} \times (1/00000)_2$ قابل نمایش است. اما تعداد ارقام در مانتیس جهت نمایش $1 + 2^{52}$ کافی نیست (۵۳ رقم در مانتیس لازم است). بنابراین $M = 2^{52} \simeq 9/007199 \times 10^{15}$ و در نتیجه تمام اعداد صحیح ۱۵ رقمی و بسیاری از اعداد ۱۶ رقمی در این نمایش به طور دقیق قابل نمایش هستند (گفته می‌شود دقت در نمایش ۶۴-بیتی ۱۵ الی ۱۶ رقم است). البته باید توجه داشت که بسیاری از توان‌های ۲ نیز به طور دقیق قابل نمایش هستند.

تذکر ۷.۲ محور اعداد با ممیز شناور بر خلاف محور اعداد حقیقی، متناهی و گسسته است. بعضی از نرم‌افزارها مانند Mathematica، محدودیت‌های سخت‌افزار را به کمک برنامه‌های نرم‌افزاری برطرف می‌کنند و به اصطلاح دقت را بالا می‌برند که ممکن است کمی سرعت انجام محاسبات پایین بیاید. برنامه uoflow.nb را ببینید.

تمرین ۲.۲ آیا $\mathbb{F}(2, 53, -1023, 1024)$ بیانگر مجموعه اعداد ممیز شناور ۶۴-بیتی است؟

تمرین ۳.۲ در نمایش ۳۲-بیتی (دقت ساده متناظر با نوع float در زبان C) هشت بیت برای مشخصه منظور می‌شود. تمام کمیت‌های معرفی‌شده در این بخش را، در این نمایش تعیین کنید. آیا $\mathbb{F}(2, 24, -127, 128)$ بیانگر مجموعه اعداد ممیز شناور ۳۲-بیتی است؟

برای جزئیات بیشتر به مرجع [۱۷] و برای جزئیات کاملتر به استاندارد بیان‌شده مراجعه کنید.

^{۱۴}Characteristic

^{۱۵}Underflow

^{۱۶}Overflow

۲.۲.۲ اعداد ماشینی

اعداد $\mathbb{F} = \mathbb{F}(\beta, t, L, U)$ به اعداد ماشینی t -رقمی نیز معروف هستند. به وضوح $\mathbb{F} \subset \mathbb{R}$ و اگر $x \in \mathbb{R}$ و $x \notin \mathbb{F}$ ، باید نگاهی مانند $r: \mathbb{R} \rightarrow \mathbb{F}$ ساخت به گونه‌ای که

$$|x - r(x)| \leq |x - y|, \quad \forall y \in \mathbb{F}.$$

اگر نمایش علمی x به صورت $x = (\circ/a_1 \cdots a_t a_{t+1} \cdots)_\beta \times \beta^e$ باشد، عدد ماشینی t -رقمی نظیر x با نگه داشتن t رقم از مانتیس و کنار گذاشتن بقیه ارقام به دست می‌آید (البته با این فرض که برای نمایش e مشکلی نداشته باشیم) و برای این کار روش‌های زیر موجود است

۱. روش قطع کردن (برش)^{۱۷}

۲. روش گرد کردن معمولی^{۱۸}

۳. روش گرد کردن به زوج^{۱۹}

در روش قطع کردن، t رقم از مانتیس حفظ شده و بقیه ارقام کنار گذاشته می‌شود در حالی که در روش گرد کردن معمولی، ابتدا روش قطع کردن اعمال شده، سپس اگر $a_{t+1} \geq \frac{\beta}{4}$ یک واحد به a_t اضافه می‌گردد. اما در روش گرد کردن به زوج، ابتدا روش قطع کردن اعمال شده و در هر یک از حالت‌های زیر یک واحد به a_t اضافه می‌گردد

$$a_{t+1} > \frac{\beta}{4} \bullet$$

$$a_{t+1} = \frac{\beta}{4} \bullet \text{ و رقم مخالف صفری در سمت راست } a_{t+1} \text{ مشاهده شود،}$$

$$a_{t+1} = \frac{\beta}{4} \bullet \text{ و رقم مخالف صفری در سمت راست } a_{t+1} \text{ مشاهده نشود و } a_t \text{ فرد باشد.}$$

امروزه روش‌های قطع کردن (به دلیل خطای زیاد) و گرد کردن به زوج (به دلیل عدم توجیحات قابل قبول) منسوخ شده‌اند (اگرچه زمانی رایانه‌هایی بر مبنای آنها ساخته شد). بنابراین می‌توان نگاهی را به صورت زیر در نظر گرفت

$$r(x) = \begin{cases} (\circ/a_1 \cdots a_t)_\beta \times \beta^e, & a_{t+1} < \frac{\beta}{4}, \\ (\circ/a_1 \cdots a_t + \circ/0 \cdots 01)_\beta \times \beta^e, & a_{t+1} \geq \frac{\beta}{4}. \end{cases}$$

پرسش ۲.۲ تفاوت 47^{km} با 47000^m یا تفاوت $3/7$ با $3/70$ یا $3/700$ در چیست؟

تعریف ۷.۲ منظور از ارقام بامعنای یک عدد مخالف صفر، ارقام مخالف صفر، صفرهای بین دو رقم مخالف صفر و صفرهایی است که در سمت راست عدد به منظور نشان دادن نوعی دقت قرار داده می‌شوند (تمام ارقام مانتیس).

△

مثال ۷.۲ عدد $0/000704500$ حداقل ۴ رقم بامعنا و حداکثر ۶ رقم بامعنا دارد.

^{۱۷}Chopping

^{۱۸}Rounding

^{۱۹}Rounding to even

تعریف ۸.۲ اگر \hat{x} تقریبی از x باشد، $e_a(x) = |x - \hat{x}|$ را خطای مطلق \hat{x} نسبت به x نامند. e_a منحصر به فرد است و در عمل بیشتر مواقع قابل تعیین نیست و به جای آن از هر عدد b_a استفاده می‌شود که کمتر از $e_a(x)$ نباشد. b_a یکتا نیست و به کران خطای مطلق معروف است. بنابراین $e_a(x) \leq b_a$ و در نتیجه $\hat{x} - b_a \leq x \leq \hat{x} + b_a$ بعضی مواقع از نمایش $x = \hat{x} \pm b_a$ استفاده می‌شود.

مثال ۸.۲ اگر عدد $\hat{x} = ۱/۷۳۲$ به عنوان تقریبی از $x = \sqrt{۳}$ باشد، آن‌گاه

$$e_a(x) = \left| \sqrt{۳} - ۱/۷۳۲ \right| = ۱/۷۳۲۰۵۰۸۰۷۵\dots - ۱/۷۳۲ = ۰/۰۰۰۰۵۰۸۰۷۵\dots$$

از طرفی می‌دانیم $۱/۷۳۲۱ < \sqrt{۳} < ۱/۷۳۲۰$ ، بنابراین $۱/۷۳۲۰ < \sqrt{۳} - ۱/۷۳۲ < ۰$ در نتیجه $b_a = ۰/۰۰۰۰۱$ معیاری برای نزدیکی $۱/۷۳۲$ به $\sqrt{۳}$ است. \triangle

پرسش ۳.۲ آیا خطای مطلق معیار مناسبی برای مقایسه خطاها است؟

پاسخ. خیر. به عنوان مثال خطای مطلق یک صندوق‌دار بانک، تاییست و دروازه‌بان را در نظر بگیرید.

تعریف ۹.۲ اگر \hat{x} تقریبی از $x \neq ۰$ باشد، $e_r(x) = \frac{|x - \hat{x}|}{|x|} = \frac{e_a(x)}{|x|}$ خطای نسبی \hat{x} نسبت به x نامیده می‌شود و مشابه خطای مطلق منحصر به فرد است و بیشتر مواقع در عمل قابل تعیین نیست و از کران خطای نسبی استفاده می‌شود. $۱۰۰ \times e_r(x)$ به درصد خطا معروف است.

قضیه ۳.۲ اگر \hat{x} تقریبی از x باشد و b_a یک کران خطای مطلق برای این تقریب باشد، آن‌گاه

$$e_r(x) \leq \frac{b_a}{|\hat{x}| - b_a}.$$

به علاوه اگر b_a نسبت به $|\hat{x}|$ خیلی کوچک باشد داریم

$$e_r(x) \leq \frac{b_a}{|\hat{x}|}.$$

□

برهان. به کمک تعریف b_a و با استفاده از خواص نابرابری‌ها نشان دهید.

مثال ۹.۲ عدد $\hat{x} = ۱/۷۳۲$ را به عنوان تقریبی از $x = \sqrt{۳}$ در نظر گرفته، داریم

$$e_r(x) = \frac{|\sqrt{۳} - ۱/۷۳۲|}{\sqrt{۳}} = \frac{۱/۷۳۲۰۵۰۸۰۷۵\dots - ۱/۷۳۲}{۱/۷۳۲۰۵۰۸۰۷۵\dots} = \frac{۰/۰۰۰۰۵۰۸۰۷۵\dots}{۱/۷۳۲۰۵۰۸۰۷۵\dots}$$

پس

$$e_r(x) = ۰/۰۰۰۰۰۲۹۳۳۳۷\dots < ۰/۰۰۰۰۰۳.$$

اما با توجه به قضیه ۳.۲ و $b_a = ۰/۰۰۰۰۱$ می‌توان نوشت

$$e_r(x) \leq \frac{۰/۰۰۰۰۱}{۱/۷۳۲ - ۰/۰۰۰۰۱} = \frac{۰/۰۰۰۰۱}{۱/۷۳۱۹} = ۰/۰۰۰۰۰۵۷۷۴۰۰۵\dots < ۰/۰۰۰۰۰۶$$

و یا

$$e_r(x) \leq \frac{۰/۰۰۰۰۱}{۱/۷۳۲} = ۰/۰۰۰۰۰۵۷۷۳۶۷۲\dots < ۰/۰۰۰۰۰۶.$$

Δ

تعریف ۱۰.۲ اگر $\hat{x} = (0/a_1 \dots 0/a_t)_\beta \times \beta^e$ تقریبی برای x باشد، بزرگترین عدد طبیعی n (حداکثر t) که در نابرابری

$$e_a(x) \leq 0.5 \times \beta^{e-n},$$

صدق کند، تعداد ارقام بامعناى درست \hat{x} نسبت به x نامیده می‌شود.

پرسش ۴.۲ در نظر بگیرید \hat{x} تقریبی از x باشد. آیا هر چه تعداد ارقام بامعناى \hat{x} بیشتر باشد، می‌توان گفت \hat{x} تقریب بهتری است؟ پاسخ. خیر. با توجه به تعریف، هر چه تعداد ارقام بامعناى درست \hat{x} نسبت به x بیشتر باشد، خطای \hat{x} کمتر است.

مثال ۱۰.۲ از اعداد $\hat{x} = 99/96$ و $\tilde{x} = 100/7$ کدام یک تقریب بهتری برای $x = 100$ است؟

$$\hat{x} = 99/96 \rightarrow e = 2, \quad |x - \hat{x}| = 0.04 \leq 0.5 \times 10^{2-n} \rightarrow n = 3,$$

$$\tilde{x} = 100/7 \rightarrow e = 3, \quad |x - \tilde{x}| = 0.7 \leq 0.5 \times 10^{3-n} \rightarrow n = 2.$$

Δ

بنابراین \hat{x} تقریب بهتری برای x است.

مثال ۱۱.۲ تعداد ارقام بامعناى درست $\hat{x} = 100/31$ را نسبت به $x = 100/3104$ مشخص کنید.

$$\hat{x} = 100/31 \rightarrow e = 3, \quad |x - \hat{x}| = 0.0004 \leq 0.5 \times 10^{3-n} \rightarrow n = 6,$$

Δ

و چون \hat{x} فقط پنج رقم بامعنا دارد پس $n = 5$.

قضیه ۴.۲ اگر \hat{x} گردشده x تا n رقم بامعنا باشد، آن‌گاه \hat{x} دارای n رقم بامعناى درست است.

قضیه ۵.۲ اگر کمیت تقریبی و دقیق را در (به) توانی از مبنا ضرب (تقسیم) کنیم، خطای نسبی تغییر نمی‌کند، ولی خطای مطلق در (به) آن توان ضرب (تقسیم) می‌شود. همچنین تعداد ارقام بامعناى درست نیز تغییر نمی‌کند.

ارتباط خطای نسبی (دقت) با تعداد ارقام بامعناى درست، در قضیه بعد بیان شده است.

قضیه ۶.۲ اگر کمیت تقریبی \hat{x} دارای n رقم بامعناى درست نسبت به کمیت دقیق x باشد، آن‌گاه $e_r(x) < 0.5 \times \beta^{1-n}$ به شرط آن که ارقام بامعناى درست \hat{x} از یک رقم ۱ و $n-1$ رقم صفر در جلوی آن تشکیل نشده باشد. برعکس اگر $e_r(x) \leq 0.5 \times \beta^{-n}$ ، آن‌گاه \hat{x} نسبت به x ، دارای دست کم n رقم بامعناى درست است.

مثال ۱۲.۲ تقریبی از $\sqrt{3}$ ارایه دهید که خطای نسبی آن از 10^{-4} کمتر باشد.

بنابر قضیه ۶.۲ اگر \hat{x} تقریبی از $x = \sqrt{3}$ باشد که 5 رقم بامعناى درست داشته باشد آن‌گاه $10^{-4} < 5 \times 10^{-5} < e_r(x)$. از این رو، با توجه به قضیه ۴.۲، کافی است \hat{x} گردشده $\sqrt{3}$ تا 5 رقم بامعنا باشد، یعنی $\hat{x} = 1/7221$.

Δ

تمرین ۴.۲ نشان دهید

$$e_a(r(x)) \leq 0.5 \times \beta^{e-t}, \quad e_r(r(x)) \leq 0.5 \times \beta^{1-t} = 0.5 \times \epsilon_M. \quad (۳.۲)$$

اگر $r(x)$ از روش قطع کردن به دست آمده باشد، این کران‌های خطا چگونه تغییر می‌کنند؟

تعریف ۱۱.۲ کمیت $0.5 \times \epsilon_M$ به r واحد 2^0 یا واحد گرد کردن 2^1 یا دقت ماشین 2^2 معروف است. به نظر می‌رسد دلیل این نام‌گذاری آن است که اگر $\epsilon \leq 0.5 \times \epsilon_M$ آن‌گاه $r(1 + \epsilon) = 1$.

تذکر ۸.۲ گاهی مواقع اتفاق می‌افتد که $r(x) \notin \mathbb{F}$. برای نمونه به موارد زیر در $\mathbb{F} = \mathbb{F}(10, 3, -9, 9)$ توجه کنید.

$$\begin{aligned} r(0.1234 \times 10^{12}) &= 0.123 \times 10^{12} \notin \mathbb{F}, \\ r(0.9998 \times 10^9) &= 0.100 \times 10^9 \notin \mathbb{F}, \\ r(0.0123 \times 10^{-9}) &= 0.123 \times 10^{-10} \notin \mathbb{F}, \\ r(0.1234 \times 10^{-11}) &= 0.123 \times 10^{-11} \notin \mathbb{F}. \end{aligned}$$

در دو حالت اول سرریز و در دو حالت دوم پاریز اتفاق افتاده است. دو مورد آخر به صورت زیر قابل برطرف کردن هستند

$$\begin{aligned} r(0.0123 \times 10^{-9}) &= 0.012 \times 10^{-9} \in \mathbb{F}, \\ r(0.1234 \times 10^{-11}) &= 0 \in \mathbb{F}. \end{aligned}$$

ولی متأسفانه دیگر روابط (۳.۲) برقرار نیستند.

رایانه‌ها با پدیده سرریز (پاریز) به عنوان نوعی بی‌نظمی در محاسبات برخورد می‌کنند. وقوع این پدیده‌ها ممکن است به توقف محاسبات منجر شود. با مقیاس کردن مناسب داده‌های ورودی و بررسی نتایج میانی محاسبات، تا حدودی می‌توان از وقوع این پدیده‌ها جلوگیری کرد. به عنوان نمونه برای محاسبه $\|x\| = \sqrt{x_1^2 + x_2^2}$ عبارت $\|x\| = |x_1| \sqrt{1 + (x_2/x_1)^2}$ یا $\|x\| = |x_2| \sqrt{1 + (x_1/x_2)^2}$ مناسب‌تر است. چون در برخورد با این پدیده‌ها برای هر روش عددی راه‌کارهای حفاظتی مخصوصی نیاز است و این پدیده‌ها به ندرت اتفاق می‌افتند، فرض می‌کنیم هیچ محدودیتی روی e (توان در نمایش ممیز شناور) نداشته باشیم و از تعریف زیر استفاده می‌کنیم.

تعریف ۱۲.۲ برای $x \in \mathbb{R}$ به گونه‌ای که $x_{\min} \leq |x| \leq x_{\max}$ تعریف می‌کنیم

$$r(x) = x(1 + \epsilon), \quad |\epsilon| \leq 0.5 \times \epsilon_M.$$

بنابراین $x \in \epsilon$ خطای مطلق تقریب x با $r(x)$ و ϵ خطای نسبی این تقریب است.

تذکر ۹.۲ به نظر می‌رسد ایده این تعریف از خطای نسبی آمده در (۳.۲) ناشی شده است. زیرا داریم

$$\frac{|x - r(x)|}{|x|} \leq 0.5 \times \epsilon_M,$$

و یا

$$-0.5 \times \epsilon_M \leq \frac{x - r(x)}{x} \leq 0.5 \times \epsilon_M,$$

واز آنجا

$$(-0.5 \times \epsilon_M)x \leq x - r(x) \leq (0.5 \times \epsilon_M)x,$$

در نتیجه

$$x(1 - 0.5 \times \epsilon_M) \leq r(x) \leq x(1 + 0.5 \times \epsilon_M).$$

تعریف ۱۳.۲ اعمال جمع، منها، ضرب و تقسیم در یک ماشین محاسباتی به صورت اعمال ممیز شناور به کار می‌روند که برای $x, y \in \mathbb{F}$ به صورت زیر تعریف می‌شوند

$$\begin{cases} r(x + y) = (x + y)(1 + \epsilon_1), \\ r(x - y) = (x - y)(1 + \epsilon_2), \\ r(x * y) = (x * y)(1 + \epsilon_3), \\ r(x/y) = (x/y)(1 + \epsilon_4), \end{cases} \quad |\epsilon_i| \leq 0.5 \times \epsilon_M, \quad i = 1, 2, 3, 4.$$

تذکر ۱۰.۲ اگر E یک عبارت شامل تعدادی اعمال ریاضی باشد، منظور از $r(E)$ آن است که تمام اعمال در آن عبارت باید به صورت اعمال ممیز شناور در نظر گرفته شوند. به عنوان نمونه داریم

$$r(x * (y + z)) = (x * ((y + z)(1 + \epsilon_1)))(1 + \epsilon_2), \quad \forall x, y, z \in \mathbb{F},$$

که در آن $|\epsilon_1|, |\epsilon_2| \leq 0.5 \times \epsilon_M$. همچنین منظور از $r(\sqrt{x})$ ، $r(\sin x)$ و سایر اعمال مشابه، اعمال r بر روی نتیجه توابع تقریبی نظیر است که در ماشین استفاده می‌شود.

۳.۲ انتشار خطا

بسیاری از اصول موضوعه میدان اعداد حقیقی مانند یکتایی عضو خنثی، شرکت‌پذیری و غیره، در مجموعه اعداد ماشینی برقرار نیستند. به عبارتی روابطی که از نظر ریاضی معادل هستند از نظر محاسباتی نتایج متفاوتی تولید می‌کنند. بنابراین از نظر عددی مهم است که بین طرح‌های ارزیابی مختلف حتی اگر از دیدگاه نظری معادل باشند تفاوت قائل باشیم. به مثال زیر از مرجع [۲۲] توجه کنید.

مثال ۱۳.۲ فرض کنید می‌خواهیم در یک ماشین با $t = 8$ و $\beta = 10$ ، مجموع سه عدد زیر را به دست آوریم.

$$a = 0.23371258 \times 10^{-4}, \quad b = 0.33678429 \times 10^2, \quad c = -0.33677811 \times 10^2.$$

با یک مقایسه ساده با جواب دقیق یعنی 10^{-3} $\times 0.641371258$ ، متوجه می‌شویم که مقدار

$$r(a + (b + c)) = r(0.23371258 \times 10^{-4} + 0.61800000 \times 10^{-3}) = 0.64137126 \times 10^{-3},$$

دقیق‌تر از مقدار

$$r((a + b) + c) = r(0.33678452 \times 10^2 - 0.33677811 \times 10^2) = 0.64100000 \times 10^{-3},$$

است. در واقع خطای ناشی از رابطه اخیر به علت پدیده از بین رفتن ارقام با معنا^{۲۳} است. این پدیده در تفاضل دو عدد هم علامت نزدیک به هم اتفاق می‌افتد. Δ

تعریف ۱۴.۲ اعمال جمع، منها، ضرب و تقسیم و همچنین اعمال $\sqrt{\cdot}$ ، \sin ، و غیره به اعمال مقدماتی^{۲۴} معروف هستند.

حال قصد داریم مفهوم الگوریتم را قدری دقیق‌تر بیان کنیم و به همین منظور فرض کنید در یک مسئله نتایج y_1, \dots, y_m از داده‌های ورودی x_1, \dots, x_n به دست می‌آیند. با تعریف

$$x = [x_1 \cdots x_n]^T, \quad y = [y_1 \cdots y_m]^T,$$

حل مسئله یعنی تعیین $y = \phi(x)$ که در آن $\phi: \Omega \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^m$ یک تابع برداری چندمتغیره است، یعنی

$$y_j = \phi_j(x_1, \dots, x_n), \quad j = 1, \dots, m.$$

فرض کنید ϕ_j و مشتقات جزئی مرتبه اول آن برای $j = 1, \dots, m$ بر Ω پیوسته باشند، \hat{x} تقریبی از x باشد،

$$e_a(x) = x - \hat{x}, \quad e_a(x_i) = x_i - \hat{x}_i, \quad i = 1, \dots, n,$$

بیانگر خطای مطلق و

$$e_r(x_i) = \frac{x_i - \hat{x}_i}{x_i}, \quad x_i \neq 0, \quad i = 1, \dots, n,$$

نشان‌دهنده خطای نسبی باشد. به کمک بسط تیلور توابع چندمتغیره داریم

$$e_a(y_j) = y_j - \hat{y}_j = \phi_j(x) - \phi_j(\hat{x}) \doteq \sum_{i=1}^n (x_i - \hat{x}_i) \frac{\partial \phi_j(\hat{x})}{\partial x_i} = \sum_{i=1}^n e_a(x_i) \frac{\partial \phi_j(\hat{x})}{\partial x_i}, \quad j = 1, \dots, m,$$

که می‌توان آن را به شکل فشرده $e_a(y) \doteq J e_a(x)$ نوشت که در آن

$$J = \left[\frac{\partial \phi_j(\hat{x})}{\partial x_i} \right]_{m \times n},$$

^{۲۳}Cancellation

^{۲۴}Elementary operations

ماتریس ژاکوبین تابع ϕ و منظور از \hat{x} تقریب مرتبه اول است. کمیت $\frac{\partial \phi_j(\hat{x})}{\partial x_i}$ نشان‌دهنده میزان حساسیت y_j نسبت به اختلال x_i از $e_a(x_i)$ است. حال اگر $x_i \neq 0$ و $y_j \neq 0$ آن‌گاه می‌توان نوشت

$$e_r(y_j) \doteq \sum_{i=1}^n \frac{x_i}{\phi_j(\hat{x})} \frac{\partial \phi_j(\hat{x})}{\partial x_i} e_r(x_i), \quad j = 1, \dots, m,$$

کمیت $\frac{x_i}{\phi_j(\hat{x})} \frac{\partial \phi_j(\hat{x})}{\partial x_i}$ نشان‌دهنده آن است که خطای نسبی $e_r(x_i)$ به چه اندازه روی خطای نسبی $e_r(y_j)$ اثرگذار است. این ضرایب بزرگ‌نمایی نسبت به تغییرات $x_i \rightarrow \lambda x_i$ و $y_j \rightarrow \mu y_j$ ثابت باقی می‌مانند و به اعداد حالت (وضعیت)^{۲۵} معروف هستند. اگر تمام اعداد حالت کوچک باشند، مسئله $y = \phi(x)$ خوش‌حالت و در غیر این صورت بدحالت است. این آنالیز خطا به آنالیز خطای دیفرانسیلی^{۲۶} یا آنالیز خطای پیشرو^{۲۷} معروف است و در ارتباط با آن باید به نکات زیر توجه داشت

- در مقابل این آنالیز خطا، یک آنالیز خطای قدیمی موسوم به آنالیز خطای پسرو^{۲۸} در برخی مراجع یافت می‌شود که نخستین بار ویلیکینسون^{۲۹} در مباحث جبرخطی مطرح کرد. در این آنالیز خطا سعی بر آن است که $e_a(x)$ به گونه‌ای یافت شود که $y + e_a(y) = \phi(x + e_a(x))$ ، این طور نیست که یک مسئله یا خوش‌حالت باشد یا بدحالت، بلکه بسته به داده‌های مسئله، ممکن است مسئله خوش‌حالت یا بدحالت باشد.
- برای یک مسئله بدحالت، تغییرات کوچک روی داده‌ها به خطای بزرگی روی جواب منجر می‌شود.
- این تعریف خوش‌حالت بودن خیلی کلی است و دو ایراد دارد

۱. باید $x_i \neq 0$ و $y_j \neq 0$ ،

۲. باید m, n عدد بررسی شوند که برای برخی مسایل مانند مسئله حل دستگاه خطی ممکن است بررسی آنها زمان‌گیر باشد،

به همین دلیل بهتر است با توجه به مسئله، حالت آن بررسی شود. به عنوان مثال در مسئله حل دستگاه خطی $Ax = b$ کمیت $\kappa = \|A^{-1}\| \|A\|$ عدد وضعیت مسئله است و با توجه به بزرگی یا کوچکی بودن آن، مسئله بدحالت یا خوش‌حالت خواهد بود.

به کمک آنالیز خطای پیشرو می‌توان به دو پرسش زیر پاسخ داد.

۱. (مستقیم) اگر خطای مطلق (نسبی) x_i ها معلوم باشد، خطای مطلق (نسبی) y_j ها چقدر خواهد بود؟

۲. (وارون) برای آنکه خطای مطلق (نسبی) y_j ها از ϵ کمتر باشد، حداکثر خطای مطلق (نسبی) x_i ها چقدر می‌تواند باشد؟

مثال ۱۴.۲ یک استوانه به شعاع قاعده $\frac{r}{4}$ و ارتفاع $\sqrt{2}r$ را در نظر بگیرید. اگر شعاع و ارتفاع استوانه و عدد π را با دقت $5S^\circ$ وارد محاسبات کنیم، حجم این استوانه با چه خطایی به دست می‌آید؟

می‌دانیم حجم یک استوانه به کمک قاعده $\pi r^2 h$ تعیین می‌شود که در آن r شعاع قاعده و h ارتفاع استوانه است. اگر تعریف کنیم

$$z = V(p, r, h) = \pi r^2 h$$

^{۲۵}Condition numbers

^{۲۶}Differential error analysis

^{۲۷}Forward error analysis

^{۲۸}Backward error analysis

^{۲۹}Wilkinson

^{۳۰}منظور S رقم بامعنا است. حرف نخست واژه significant است.

آن‌گاه باید V به ازای $p = \frac{4}{3}$, $r = \sqrt{2}$ و $h = \sqrt{2}$ ارزیابی شود. چون دقت $5S$ مدنظر است

$$\hat{p} = 3/1416, \quad \hat{r} = 1/3333, \quad \hat{h} = 1/4142,$$

به طوری که $e_a(p), e_a(r), e_a(h) \leq 0/5 \times 10^{-4}$. بنابراین

$$\hat{z} = V(\hat{p}, \hat{r}, \hat{h}) = 3/4116 \times 1/3333^2 \times 1/4142 = 7/8980.$$

از طرفی

$$e_a(z) \doteq e_a(p)\hat{r}^2\hat{h} + 2e_a(r)\hat{p}\hat{r}\hat{h} + e_a(h)\hat{p}\hat{r}^2 \leq (\hat{r}^2\hat{h} + 2\hat{p}\hat{r}\hat{h} + \hat{p}\hat{r}^2) \times 0/5 \times 10^{-4} \lesssim 0/0010.$$

در نتیجه حجم استوانه با خطایی کمتر از $0/0010$ برابر است با $7/8980$. به کمک یک ماشین حساب 10^0 رقمی نتیجه $z = 7/898458555$ به دست می‌آید و بنابراین خطای واقعی عبارت است از

$$e_a(z) = |z - \hat{z}| = 4/585554 \times 10^{-4} < 0/0010.$$

Δ

تمرین ۵.۲ یک استوانه به شعاع قاعده $\frac{4}{3}$ و ارتفاع $\sqrt{2}$ در نظر بگیرید. اگر بخواهیم حجم این استوانه را با دقت $4S$ به دست آوریم، شعاع قاعده و ارتفاع استوانه و حتی عدد π را با چه دقتی وارد محاسبات کنیم؟

مثال ۱۵.۲ اعداد $x = \sqrt{5}$ و $y = \frac{\pi}{11}$ را با چه دقتی در نظر بگیریم تا مقدار $6x^2(\ln x + \sin 2y)$ با خطایی حداکثر به اندازه $0/5 \times 10^{-2}$ حساب شود.

اگر فرض کنیم

$$z = \phi(x, y) = 6x^2(\ln x + \sin 2y),$$

آن‌گاه

$$\frac{\partial \phi}{\partial x} = 12x(\ln x + \sin 2y) + 6x, \quad \frac{\partial \phi}{\partial y} = 12x^2 \cos 2y,$$

و در نتیجه با فرض $x = 2/2$ و $y = 0/3$ یا $1S$ یا $2S$ در نظر می‌گیریم داریم

$$e_a(z) \doteq \left| e_a(x) \frac{\partial f(2/2, 0/3)}{\partial x} + e_a(y) \frac{\partial f(2/2, 0/3)}{\partial y} \right| \simeq |48/9e_a(x) + 47/9e_a(y)|,$$

و برای برقراری نابرابری $e_a(z) \leq 0/5 \times 10^{-2}$ باید داشته باشیم

$$|48/9e_a(x) + 47/9e_a(y)| \leq 0/5 \times 10^{-2}.$$

یک جواب نامعادله اخیر عبارت است از

$$e_a(x) \leq 0.5 \times 10^{-4}, \quad e_a(y) \leq 0.5 \times 10^{-4}.$$

Δ

مثال ۱۶.۲ برای محاسبه $y = \phi(a, b, c) = a^2 b^3 c^4$ داریم $e_r(y) \simeq 2e_r(a) + 3e_r(b) + 4e_r(c)$.

Δ

به دو مثال بعد از مرجع [۲۲] توجه کنید.

مثال ۱۷.۲ مسئله $y = \phi(a, b, c) = a + b + c$ با توجه به

$$e_r(y) \doteq \frac{a}{a+b+c} e_r(a) + \frac{b}{a+b+c} e_r(b) + \frac{c}{a+b+c} e_r(c),$$

خوش حالت است اگر جمعوندهای a, b, c در مقایسه با $a + b + c$ کوچک باشند. به عنوان نمونه اگر $a = 10^n$ و $b = -a$ و $c = 10^{-n}$

Δ

آن‌گاه هرچه n بزرگتر باشد، مسئله بدحالت‌تر خواهد بود.

مثال ۱۸.۲ به منظور یافتن ریشه با اندازه کوچکتر معادله $y^2 + 2py - q = 0$ ، مسئله زیر را در نظر می‌گیریم

$$y = \phi(p, q) = -p + \sqrt{p^2 + q}.$$

$$\text{چون } \frac{\partial \phi}{\partial q} = \frac{1}{2\sqrt{p^2 + q}} \text{ و } \frac{\partial \phi}{\partial p} = -1 + \frac{p}{\sqrt{p^2 + q}} = \frac{-y}{\sqrt{p^2 + q}} \text{ داریم}$$

$$e_r(y) \doteq \frac{-p}{\sqrt{p^2 + q}} e_r(p) + \frac{q}{2y\sqrt{p^2 + q}} e_r(q) = \frac{-p}{\sqrt{p^2 + q}} e_r(p) + \frac{p + \sqrt{p^2 + q}}{2\sqrt{p^2 + q}} e_r(q).$$

برای $q \geq 0$ ، اندازه ضرایب بزرگ‌نمایی (اعداد حالت) حداکثر یک است و در نتیجه مسئله خوش حالت است ولی برای $q < 0$ ، هر چه q به $-p^2$

Δ

نزدیکتر باشد، مسئله بدحالت‌تر خواهد بود.

ماتریس هیلبرت (در مسایل تقریب ظاهر می‌شود)، ماتریس واندرموند (در درونبایی چندجمله‌ای به دست می‌آید) و چندجمله‌ای خائن ویلکینسون (در ریشه‌یابی به آن برخورد می‌کنیم) نمونه‌هایی از معروفترین مسایل هستند که از خود بدحالتی آشکاری نشان می‌دهند.

۱.۳.۲ خطای اعمال ریاضی

به کمک آنالیز خطای پیشرو به سادگی می‌توان تخمین خطای چهار عمل اصلی را به صورت زیر به دست آورد

$$y = \phi(x_1, x_2) = x_1 \pm x_2, \quad e_a(y) \doteq e_a(x_1) \pm e_a(x_2), \quad e_r(y) \doteq \frac{x_1}{x_1 \pm x_2} e_r(x_1) \pm \frac{x_2}{x_1 \pm x_2} e_r(x_2),$$

$$y = \phi(x_1, x_2) = x_1 \times x_2, \quad e_a(y) \doteq x_2 e_a(x_1) + x_1 e_a(x_2), \quad e_r(y) \doteq e_r(x_1) + e_r(x_2),$$

$$y = \phi(x_1, x_2) = x_1/x_2, \quad e_a(y) \doteq \frac{x_2 e_a(x_1) - x_1 e_a(x_2)}{x_2^2}, \quad e_r(y) \doteq e_r(x_1) - e_r(x_2).$$

مثال ۱۹.۲ اگر اعداد مثبت a و b هر یک n رقم بامعنای درست نسبت به اعداد A و B داشته باشند، حداقل تعداد ارقام بامعنای درست $a + b$ و ab را تعیین کنید.

بنابر قضیه ۶.۲، $e_r(a) \leq 5 \times 10^{-n}$ و $e_r(b) \leq 5 \times 10^{-n}$. با توجه به کران خطای نسبی عمل جمع داریم

$$e_r(a + b) \leq \max\{e_r(a), e_r(b)\} \leq 5 \times 10^{-n} = 0.5 \times 10^{-(n-1)}.$$

پس $a + b$ حداقل $n - 1$ رقم بامعنای درست دارد. به علاوه می‌توان نوشت

$$e_r(ab) \leq e_r(a) + e_r(b) = 5 \times 10^{-n} + 5 \times 10^{-n} = 10 \times 10^{-n} = 10^{-(n-2)-1} < 0.5 \times 10^{-(n-2)}.$$

△

بنابراین ab دست کم $n - 2$ رقم بامعنای درست دارد.

تذکر ۱۱.۲ با توجه به کران خطای نسبی عمل منها، باید از تفاضل دو عدد هم‌علامت نزدیک به هم جلوگیری کرد زیرا موجب می‌شود ارقام بامعنای نتیجه از بین برود. به وضوح ضرایب بزرگ‌نمایی (اعداد حالت) بزرگ هستند و تولید چنین خطایی ممکن است زبانبار نباشد ولی انتشار آن می‌تواند خطرناک (به این مفهوم که رشد خطای زیادی به همراه دارد) باشد. باید توجه داشت که با یک بررسی ساده مشخص می‌شود که تفاضل دو عدد هم‌علامت که نزدیک به هم نباشند خطرناک نیست. بعضی مواقع می‌توان جهت اجتناب از عمل تفاضل در محاسبات، از اتحادها کمک گرفت. به عنوان مثال

$$e^{x-y} = \frac{e^x}{e^y}, \quad 1 - \cos x = 2 \sin^2 \frac{x}{2}, \quad \ln x - \ln y = \ln \frac{x}{y}.$$

اگر عمل منها اجتناب‌پذیر نباشد باید سعی کرد عمل منها با دقت دو برابر (یا بیشتر) انجام شود.

مثال ۲۰.۲ در صورتی که $g(x) = x \left(\sqrt{1 + \frac{1}{x}} - 1 \right)$ مقدار $g(10^9)$ را با دقت 10^{-5} به دست آورید. به سادگی داریم

$$1 + \frac{1}{x} = 1.000000001, \quad \sqrt{1 + \frac{1}{x}} = 1.0000000005.$$

بنابراین ماشین حساب^{۳۱} نتیجه زیر را به دست می‌دهد

$$g(10^9) = 10^9 (1.0000000005 - 1) = 0.0000000005.$$

به این ترتیب با یک ماشین حساب 10^{-5} رقمی $g(10^9) = 0.0000000005$ که نتیجه‌ای نادرست است. در حالی که اگر از یک رایانه استفاده شود، نتیجه درست $g(10^9) = 0.3333333333$ به دست می‌آید. دلیل این خطای فاحش، تفاضل دو عدد نزدیک به هم در محاسبات است که به از بین رفتن ارقام بامعنا منجر می‌شود. برای به دست آوردن تقریب بهتر، روش محاسبه را به صورت زیر تغییر می‌دهیم

$$g(x) = x \left(\sqrt{1 + \frac{1}{x}} - 1 \right) \times \frac{\sqrt{(1 + \frac{1}{x})^2} + \sqrt{1 + \frac{1}{x}} + 1}{\sqrt{(1 + \frac{1}{x})^2} + \sqrt{1 + \frac{1}{x}} + 1} = \frac{1}{\sqrt{(1 + \frac{1}{x})^2} + \sqrt{1 + \frac{1}{x}} + 1}.$$

^{۳۱} ماشین حساب‌های جدید مقدار متفاوتی به دست می‌دهند.

در این صورت با توجه به

$$\left(1 + \frac{1}{x}\right)^2 = 1/000000002, \quad \sqrt{\left(1 + \frac{1}{x}\right)^2} = 1/000000000,$$

نتیجه زیر به دست می‌آید

$$g(10^9) = \frac{1}{1/000000000 + 1/000000000 + 1} = 0/333333333.$$

△

این مثال نه تنها نشان می‌دهد که ممکن است یک ماشین نتایج نادرستی تولید کند بلکه به خوبی نشان می‌دهد که به دلیل خطای گرد کردن، ممکن است محاسبه دو عبارت مختلف که از نظر ریاضی هم‌ارز هستند به نتایج متفاوتی منجر شود. از این رو باید از نظر عددی بین الگوریتم‌هایی که از نظر ریاضی هم‌ارز هستند تفاوت قائل شویم.

تذکر ۱۲.۲ با توجه به خطای نسبی اعمال ضرب و تقسیم، باید در محاسبات از ضرب (تقسیم) اعداد بزرگ (کوچک) در (به) اعداد تقریبی پرهیز کرد.

مثال ۲۱.۲ در محاسبه 10000π با توجه به

$$\begin{aligned} \pi = 3/14 &\rightarrow 10000 \times \pi = 31400, \\ \pi = 3/142 &\rightarrow 10000 \times \pi = 31420, \\ \pi = 3/1416 &\rightarrow 10000 \times \pi = 31416, \end{aligned}$$

درمی‌یابیم که در تقریب اول خطایی به اندازه ۱۶ واحد، در تقریب دوم خطایی نزدیک به ۴ واحد و در تقریب سوم خطایی کمتر از ۱/۰ مرتکب شده‌ایم.

△

تذکر ۱۳.۲ چون هر عمل محاسباتی خطایی به همراه دارد، یک قاعده کلی دیگر آن است که از حجم محاسبات تا آنجا که ممکن است کاسته شود.

مثال ۲۲.۲ به جای عبارت $a_0 + a_1x + \dots + a_nx^n$ از عبارت $a_0 + (a_1x + a_2x^2 + \dots + a_nx^n)x$ استفاده شود که تعداد ضرب کمتری دارد. این راه‌کار به قاعده هورنر^{۳۲} معروف است.

△

۲.۳.۲ تقریب توابع یک متغیره

به کمک آنالیز خطای پیشرو می‌توان کران خطای برخی از اعمال را به سادگی به دست آورد. به عنوان نمونه در محاسبه جذر یک عدد داریم

$$y = \phi(x) = \sqrt{x}, \quad e_a(y) \doteq \frac{1}{2\sqrt{x}} e_a(x), \quad e_r(y) \doteq \frac{1}{2} e_r(x).$$

^{۳۲} Horner's rule

بنابر قضیه تیلور، می‌توان به جای کار کردن با یک تابع پیچیده از چند جمله‌ای تیلور نظیر آن استفاده کرد. مثال‌هایی که در ادامه خواهند آمد چگونگی این تقریب را نشان می‌دهند.

مثال ۲۳.۲ مطلوب است محاسبه مقدار $e^{\frac{\pi}{10}}$ با دقت $3S$.

به کمک قضیه تیلور داریم

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots + \frac{x^n}{n!} + \frac{x^{n+1}}{(n+1)!} e^\xi,$$

که در آن $x < \xi < 0$ در نتیجه

$$e^{\frac{\pi}{10}} = 1 + \frac{\pi}{10} + \frac{\left(\frac{\pi}{10}\right)^2}{2!} + \dots + \frac{\left(\frac{\pi}{10}\right)^n}{n!} + \frac{\left(\frac{\pi}{10}\right)^{n+1}}{(n+1)!} e^\xi,$$

که در آن $1 < \frac{\pi}{10} < \xi < 0$. چون e^x تابعی صعودی است پس $3 < e^1 < e^\xi < e^0 = 1$. از طرفی $0.315 < \frac{\pi}{10}$ و بنابراین

$$\text{خطا} \leq \left| \frac{\left(\frac{\pi}{10}\right)^{n+1}}{(n+1)!} e^\xi \right| < \frac{3 \times (0.315)^{n+1}}{(n+1)!}.$$

حال باید داشته باشیم $\frac{3 \times (0.315)^{n+1}}{(n+1)!} < 0.5 \times 10^{-2}$ که نتیجه می‌دهد $n \geq 3$.

$$e^{\frac{\pi}{10}} \simeq 1 + 0.315 + \frac{(0.315)^2}{2!} + \frac{(0.315)^3}{3!} = 1.37.$$

△

مثال ۲۴.۲ (همگرایی سریع) می‌خواهیم تابع $\cos x$ را به ازای مقادیر $|x| < \frac{\pi}{4}$ با دقت $5D$ ارزیابی کنیم. با توجه به سری تیلور

$$\cos x = \sum_{i=0}^{\infty} (-1)^i \frac{x^{2i}}{(2i)!},$$

و این که

$$\left| \frac{x^{2n}}{(2n)!} \right| < \frac{\left(\frac{\pi}{4}\right)^{2n}}{(2n)!} < \frac{1/6^{2n}}{(2n)!},$$

باید داشته باشیم

$$\frac{1/6^{2n}}{(2n)!} < 0.5 \times 10^{-5},$$

که نتیجه می‌دهد $n \geq 6$.

△

مثال ۲۵.۲ (همگرایی کند) با توجه به

$$\frac{1}{1+t} = 1 - t + t^2 - t^3 + \dots,$$

داریم

$$\int_0^x \frac{dt}{1+t} = \int_0^x (1 - t + t^2 - t^3 + \dots) dt,$$

و با

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

برای ارزیابی $\ln(1+x)$ با دقت $5D$ ، باید داشته باشیم

$$\left| \frac{x^n}{n} \right| < 0.5 \times 10^{-5},$$

Δ

که برای $x = 0.99$ نتیجه می‌دهد $n \geq 582$.

مثال ۲۶.۲ تابع $Si(x) = \int_0^x \frac{\sin t}{t} dt$ را که در آن $x \geq 0$ ، در نظر بگیرید. دست کم چند جمله از بسط مک‌لورن تابع $f(t) = \sin t$ لازم است تا $Si(1)$ با دقت $6D$ مشخص شود؟
با اعمال قضیه تیلور برای تابع f داریم

$$Si(1) = \int_0^1 \frac{1}{t} \left(f(0) + \frac{f'(0)}{1!}t + \frac{f''(0)}{2!}t^2 + \dots + \frac{f^{(k)}(0)}{k!}t^k + \frac{f^{(k+1)}(\xi(t))}{(k+1)!}t^{k+1} \right) dt,$$

و یا

$$Si(1) = \int_0^1 \frac{1}{t} \left(t - \frac{t^3}{3!} + \frac{t^5}{5!} + \dots + (-1)^n \frac{t^{2n+1}}{(2n+1)!} + (-1)^{n+1} \frac{t^{2n+2}}{(2n+2)!} \sin(\xi(t)) \right) dt.$$

در نتیجه

$$Si(1) = \int_0^1 \left(1 - \frac{t^2}{3!} + \frac{t^4}{5!} + \dots + (-1)^n \frac{t^{2n}}{(2n+1)!} \right) dt + \frac{(-1)^{n+1}}{(2n+2)!} \int_0^1 t^{2n+1} \sin(\xi(t)) dt.$$

بنابراین

$$Si(1) = 1 - \frac{1}{3 \times 3!} + \frac{1}{5 \times 5!} + \dots + (-1)^n \frac{1}{(2n+1) \times (2n+1)!} + E,$$

که در آن

$$E = \frac{(-1)^{n+1}}{(2n+2)!} \int_0^1 t^{2n+1} \sin(\xi(t)) dt.$$

حال می‌توان نوشت

$$|E| = \frac{1}{(2n+2)!} \left| \int_0^1 t^{2n+1} \sin(\xi(t)) dt \right| \leq \frac{1}{(2n+2)!} \int_0^1 t^{2n+1} |\sin(\xi(t))| dt.$$

پس

$$|E| \leq \frac{1}{(2n+2)!} \int_0^1 t^{2n+1} dt = \frac{1}{(2n+2) \times (2n+2)!}.$$

از $\frac{1}{(2n+2) \times (2n+2)!} < 0.5 \times 10^{-6}$ نتیجه می‌شود $n \geq 4$ بنابراین

$$Si(1) \simeq 1 - \frac{1}{3 \times 3!} + \frac{1}{5 \times 5!} - \frac{1}{7 \times 7!} + \frac{1}{9 \times 9!},$$

و یا

$$Si(1) \simeq 1 - 0.05555556 + 0.00166667 - 0.0000283 + 0.00000002 = 0.9460831.$$

Δ

تمرین ۶.۲ مطلوب است تعیین تقریبی از عدد π با دقت ۳S به کمک بسط مکلورن تابع $f(x) = \tan^{-1} x$.

۴.۲ بررسی پایداری یک الگوریتم

این بخش را با تعریف دیگری از الگوریتم آغاز می‌کنیم.

تعریف ۱۵.۲ یک دنباله متناهی از اعمال مقدماتی که تعیین می‌کنند چگونه جواب یک مسئله به کمک داده‌های ورودی حساب می‌شود، الگوریتم (عددی) نامیده می‌شود.

در ادامه قصد داریم یک الگوریتم را دقیق‌تر بررسی کنیم. در هر مرحله از محاسبات یک مجموعه عملوندی^{۳۳} از اعداد وجود دارد که یا همان داده‌های ورودی (x_i) یا نتایج میانی تولید شده هستند. یک عمل مقدماتی بر روی یک یا دو عنصر از یک مجموعه عملوندی اثر کرده و یک عدد جدید تولید می‌کند که یا یک نتیجه نهایی (y_j) یا یک نتیجه میانی است. عناصری از مجموعه عملوندی که در ادامه محاسبات به آنها نیازی نیست، از مجموعه عملوندی کنار گذاشته می‌شوند. اولین مجموعه عملوندی فقط شامل داده‌های ورودی یعنی x_1, \dots, x_n و آخرین مجموعه عملوندی فقط شامل داده‌های خروجی یعنی y_1, \dots, y_m است. اگرچه در یک مجموعه عملوندی ترتیب (و حتی تکرار) مهم نیست، یک مجموعه عملوندی را با یک بردار نمایش می‌دهیم

$$x^{(0)} = x, \quad x^{(i)} = \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_{n_i}^{(i)} \end{bmatrix} \in \mathbb{R}^{n_i}.$$

هر عمل مقدماتی متناظر با یک نگاشت بر روی یک مجموعه عملوندی است و یک عمل مقدماتی را با یک نگاشت مقدماتی به صورت زیر نظیر می‌کنیم

$$\begin{cases} \phi^{(i)} : D_i \subseteq \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{n_{i+1}}, \\ \phi^{(i)}(x^{(i)}) = x^{(i+1)}. \end{cases}$$

صرف نظر از جایگشت‌های بی‌اهمیت $x^{(i)}$ و $x^{(i+1)}$ ، $\phi^{(i)}$ به طور یکتا تعریف می‌شود. اعمال مقدماتی متنوع در یک الگوریتم موجب می‌شود که ϕ به یک دنباله از نگاشت‌های مقدماتی تجزیه شود

$$\begin{cases} \phi = \phi^{(r)} \circ \phi^{(r-1)} \circ \dots \circ \phi^{(0)}, \\ \phi^{(i)} : D_i \subseteq \mathbb{R}^{n_i} \rightarrow D_{i+1} \subseteq \mathbb{R}^{n_{i+1}}, \end{cases}$$

که در آن $n_0 = n$ و $n_{r+1} = m$.

مثال ۲۷.۲ برای مسئله $y = \phi(a, b, c) = a + b + c$ متناظر با الگوریتم $(a + b) + c$ نگاشت‌های

$$\phi^{(0)}(a, b, c) = \begin{bmatrix} a + b \\ c \end{bmatrix}, \quad \phi^{(1)}(u, v) = u + v,$$

و متناظر با الگوریتم $a + (b + c)$ نگاشت‌های

$$\phi^{(0)}(a, b, c) = \begin{bmatrix} a \\ b + c \end{bmatrix}, \quad \phi^{(1)}(u, v) = u + v,$$

Δ

در نظر گرفته می‌شوند.

مثال ۲۸.۲ برای مسئله $y = \phi(a, b) = a^2 - b^2$ متناظر با الگوریتم $a^2 - b^2$ نگاشت‌های

$$\phi^{(0)}(a, b) = \begin{bmatrix} a^2 \\ b \end{bmatrix}, \quad \phi^{(1)}(u, b) = \begin{bmatrix} u \\ b^2 \end{bmatrix}, \quad \phi^{(2)}(v, w) = v - w,$$

و متناظر با الگوریتم $(a + b)(a - b)$ نگاشت‌های

$$\phi^{(0)}(a, b) = \begin{bmatrix} a \\ b \\ a + b \end{bmatrix}, \quad \phi^{(1)}(a, b, u) = \begin{bmatrix} u \\ a - b \end{bmatrix}, \quad \phi^{(2)}(v, w) = vw,$$

در نظر گرفته می‌شوند که می‌توان آنها را به صورت زیر خلاصه‌تر کرد

$$\phi^{(0)}(a, b) = \begin{bmatrix} a^2 \\ b^2 \end{bmatrix}, \quad \phi^{(1)}(u, v) = u - v, \quad \phi^{(0)}(a, b) = \begin{bmatrix} a + b \\ a - b \end{bmatrix}, \quad \phi^{(1)}(u, v) = uv.$$

باید توجه داشت که $\phi^{(0)}$ ‌هایی که در روابط اخیر تعریف شده‌اند، مقدماتی نیستند ولی چون متناظر با اعمال مقدماتی مستقل از هم هستند، به کار بردن آنها به منظور خلاصه‌تر شدن الگوریتم‌ها مانعی ندارد.

Δ

با این امید که بتوان معیاری برای بررسی کیفیت یک الگوریتم به دست آورد، ابتدا باید دید چرا الگوریتم‌های مختلف در حل یک مسئله، نتایج متفاوتی تولید می‌کنند. به بیان ساده شاید بتوان گفت که انتشار خطا برای الگوریتم‌های مختلف، نقش متفاوتی بازی می‌کند. به مثال زیر از مرجع [۲۲] توجه کنید.

مثال ۲۹.۲ در ادامه مثال ۱۳.۲ برای $|\epsilon_1|, |\epsilon_2| \leq 0.5 \times \epsilon_M$ داریم

$$\hat{y} = r((a + b) + c) = ((a + b)(1 + \epsilon_1) + c)(1 + \epsilon_2) = (a + b + c)\left(1 + \frac{a + b}{a + b + c}\epsilon_1(1 + \epsilon_2) + \epsilon_2\right),$$

و اگر $e_r(y) = \frac{\hat{y} - y}{y}$ ، به وضوح

$$e_r(y) \doteq \frac{a + b}{a + b + c}\epsilon_1 + \epsilon_2.$$

برای یک الگوریتم ضرایب $\frac{a + b}{a + b + c}$ و ۱ و برای الگوریتم دیگر ضرایب ۱ و $\frac{b + c}{a + b + c}$ ، موجب رشد خطاهای گرد کردن میانی می‌شوند. برای اعداد در نظر گرفته شده داریم

$$\frac{a + b}{a + b + c} \simeq 0.5 \times 10^5, \quad \frac{b + c}{a + b + c} \simeq 0.97,$$

که نشان می‌دهد چرا الگوریتم $a + (b + c)$ به نتیجه دقیق‌تری منجر شده است. Δ

در ادامه قصد داریم انتشار خطاهای گرد کردن میانی را در یک الگوریتم بررسی کنیم. همان طور که بیان شد، یک الگوریتم در حل مسئله $y = \phi(x)$ متناظر با دنباله‌ای از نگاشت‌های مقدماتی است، یعنی

$$x = x^{(0)} \rightarrow x^{(1)} = \phi^{(0)}(x^{(0)}) \rightarrow \dots \rightarrow y = x^{(r+1)} = \phi^{(r)}(x^{(r)}).$$

با این فرض که $\phi^{(i)}$ روی D_i به طور پیوسته مشتق‌پذیر باشد، تعریف می‌کنیم

$$\psi^{(i)} : D_i \rightarrow \mathbb{R}^m, \quad \psi^{(i)} = \phi^{(r)} \circ \phi^{(r-1)} \circ \dots \circ \phi^{(i)}.$$

$\psi^{(i)}$ به نگاشت باقی‌مانده مرحله نام معروف است و $\psi^{(0)} = \phi$. اگر $J\psi^{(i)}$ و $J\phi^{(i)}$ به ترتیب بیانگر ژاکوبین $\phi^{(i)}$ و $\psi^{(i)}$ باشند، به کمک قاعده زنجیره‌ای داریم

$$J\phi(x) = J\phi^{(r)}(x^{(r)}) \cdot J\phi^{(r-1)}(x^{(r-1)}) \dots J\phi^{(0)}(x^{(0)}),$$

$$J\psi^{(i)}(x) = J\phi^{(r)}(x^{(r)}) \cdot J\phi^{(r-1)}(x^{(r-1)}) \dots J\phi^{(i)}(x^{(i)}).$$

انتظار داریم در مرحله نام، $x^{(i+1)}$ به دست آید ولی با در نظر گرفتن خطای داده‌های ورودی و خطای گرد کردن میانی، $\hat{x}^{(i+1)} = r(\phi^{(i)}(\hat{x}^{(i)}))$ تولید می‌شود. بنابراین

$$e_a(x^{(i+1)}) := r(\phi^{(i)}(\hat{x}^{(i)})) - x^{(i+1)} = (r(\phi^{(i)}(\hat{x}^{(i)})) - \phi^{(i)}(\hat{x}^{(i)})) + (\phi^{(i)}(\hat{x}^{(i)}) - \phi^{(i)}(x^{(i)})). \quad (۴.۲)$$

بنابراین باید دو نوع خطا انتظار داشته باشیم که به خطای (گرد کردن) تولیدشده^{۳۴} (متناظر با عبارت اول در سمت راست (۴.۲)) و خطای منتشرشده^{۳۵} (متناظر با عبارت دوم در سمت راست (۴.۲)) معروف هستند. مشابه آنالیز خطای پیشرو، به سادگی برای خطای منتشرشده داریم

$$\phi^{(i)}(\hat{x}^{(i)}) - \phi^{(i)}(x^{(i)}) \doteq J\phi^{(i)}(\hat{x}^{(i)})e_a(x^{(i)}),$$

و برای محاسبه خطای تولیدشده، ابتدا می‌دانیم برای $j = 1, \dots, n_{i+1}$ اگر $\phi_j^{(i)} : D_i \rightarrow \mathbb{R}$ یک نگاشت مقدماتی باشد، آنگاه خواهیم داشت

$$r(\phi_j^{(i)}(u)) = (1 + \epsilon_j)\phi_j^{(i)}(u), \quad |\epsilon_j| \leq 0.5 \times \epsilon_M.$$

حال برای

$$\phi^{(i)}(u) = [\phi_1^{(i)}(u) \dots \phi_{n_{i+1}}^{(i)}(u)]^T,$$

که یا فقط یک نگاشت مقدماتی یا تعدادی نگاشت مقدماتی مستقل از هم را شامل می‌شود، می‌توان نوشت

$$r(\phi^{(i)}(u)) = (I + E_{i+1})\phi^{(i)}(u),$$

که در آن I یک ماتریس همانی $n_{i+1} \times n_{i+1}$ و $E_{i+1} = \text{diag}[\epsilon_1, \dots, \epsilon_{n_{i+1}}]$ یک ماتریس قطری $n_{i+1} \times n_{i+1}$ است. بنابراین

$$r(\phi^{(i)}(\hat{x}^{(i)})) - \phi^{(i)}(\hat{x}^{(i)}) = E_{i+1}\phi^{(i)}(\hat{x}^{(i)}),$$

و چون $E_{i+1}\phi^{(i)}(\hat{x}^{(i)}) \doteq E_{i+1}\phi^{(i)}(x^{(i)})$ (زیرا جملات خطای تقریب $\phi^{(i)}(x^{(i)})$ با $\phi^{(i)}(\hat{x}^{(i)})$ در عناصر قطری E_{i+1} ضرب شده و خطاهایی با مرتبه بالاتر تولید می‌کنند که قابل چشم‌پوشی هستند)، می‌توان نوشت

$$r(\phi^{(i)}(\hat{x}^{(i)})) - \phi^{(i)}(\hat{x}^{(i)}) \doteq E_{i+1}\phi^{(i)}(x^{(i)}) = E_{i+1}x^{(i+1)} := e_{i+1}.$$

در این رابطه e_{i+1} نقش خطای مطلق و E_{i+1} نقش خطای نسبی را بازی می‌کند. در نتیجه خطای مرحله i ام عبارت است از

$$e_a(x^{(i+1)}) \doteq e_{i+1} + J\phi^{(i)}(\hat{x}^{(i)})e_a(x^{(i)}) = E_{i+1}x^{(i+1)} + J\phi^{(i)}(\hat{x}^{(i)})e_a(x^{(i)}), \quad i \geq 0.$$

بنابراین با توجه به $e_a(x^{(0)}) = e_a(x)$ داریم

$$e_a(x^{(1)}) \doteq J\phi^{(0)}(\hat{x}^{(0)})e_a(x^{(0)}) + e_1,$$

$$e_a(x^{(2)}) \doteq J\phi^{(1)}(\hat{x}^{(1)})(J\phi^{(0)}(\hat{x}^{(0)})e_a(x^{(0)}) + e_1) + e_2,$$

\vdots ,

$$e_a(x^{(r+1)}) \doteq J\phi^{(r)}(\hat{x}^{(r)}) \dots J\phi^{(0)}(\hat{x}^{(0)})e_a(x^{(0)}) + J\phi^{(r)}(\hat{x}^{(r)}) \dots J\phi^{(1)}(\hat{x}^{(1)})e_1 + \dots + J\phi^{(r)}(\hat{x}^{(r)})e_r + e_{r+1}.$$

حال به کمک نگاشت‌های باقی‌مانده، می‌توان اثر خطای داده‌های ورودی و $e_a(x)$ و خطای گرد کردن میانی e_i ، روی نتیجه نهایی $y = x^{(r+1)}$ را در رابطه مهم زیر خلاصه کرد

$$\begin{aligned} e_a(y) &\doteq J\phi(\hat{x})e_a(x) + J\psi^{(1)}(\hat{x}^{(1)})e_1 + \dots + J\psi^{(r)}(\hat{x}^{(r)})e_r + e_{r+1} \\ &= J\phi(\hat{x})e_a(x) + J\psi^{(1)}(\hat{x}^{(1)})E_1x^{(1)} + \dots + J\psi^{(r)}(\hat{x}^{(r)})E_rx^{(r)} + E_{r+1}y. \end{aligned}$$

بنابراین بزرگی اندازه ماتریس ژاکوبین $J\psi^{(i)}(\hat{x}^{(i)})$ نشان‌دهنده اثر خطای گرد کردن میانی e_i روی جواب است.

مثال ۳۰.۲ در مثال ۲۸.۲ دو الگوریتم برای حل مسئله $y = \phi(a, b) = a^2 - b^2$ ارائه شده است. برای الگوریتم $a^2 - b^2$ داریم

$$x = x^{(0)} = \begin{bmatrix} a \\ b \end{bmatrix}, \quad x^{(1)} = \begin{bmatrix} a^2 \\ b \end{bmatrix}, \quad x^{(2)} = \begin{bmatrix} a^2 \\ b^2 \end{bmatrix}, \quad x^{(2)} = y = a^2 - b^2.$$

به سادگی خواهیم داشت

$$\psi^{(1)}(u, v) = u - v^2, \quad \psi^{(2)}(u, v) = u - v,$$

واز آنجا

$$J\phi(x) = [2a \quad -2b]^T, \quad J\psi^{(1)}(x^{(1)}) = [1 \quad -2b]^T, \quad J\psi^{(2)}(x^{(2)}) = [1 \quad -1]^T.$$

بنابراین

$$E_1 = \begin{bmatrix} \epsilon_1 & \circ \\ \circ & \circ \end{bmatrix}, \quad e_1 = E_1 x^{(1)} = \begin{bmatrix} \epsilon_1 a^2 \\ \circ \end{bmatrix},$$

$$E_r = \begin{bmatrix} \circ & \circ \\ \circ & \epsilon_r \end{bmatrix}, \quad e_r = E_r x^{(r)} = \begin{bmatrix} \circ \\ \epsilon_r b^2 \end{bmatrix},$$

و $e_r = \epsilon_r(a^2 - b^2)$ به طوری که $|\epsilon_1|, |\epsilon_r| \leq \circ/5 \times \epsilon_M$ در نتیجه چون $e_a(x) = [e_a(a) \ e_a(b)]^T$ خواهیم داشت

$$e_a(y) \doteq 2ae_a(a) - 2be_a(b) + a^2\epsilon_1 - b^2\epsilon_r + (a^2 - b^2)\epsilon_r.$$

برای الگوریتم $(a-b)(a+b)$ داریم

$$x = x^{(0)} = \begin{bmatrix} a \\ b \end{bmatrix}, \quad x^{(1)} = \begin{bmatrix} a+b \\ a-b \end{bmatrix}, \quad x^{(r)} = y = a^2 - b^2.$$

به سادگی خواهیم داشت

$$\psi^{(1)}(u, v) = uv, \quad J\psi^{(1)}(x^{(1)}) = [a-b \ a+b]^T, \quad J\phi(x) = [2a \ -2b]^T.$$

بنابراین

$$E_1 = \begin{bmatrix} \epsilon_1 & \circ \\ \circ & \epsilon_r \end{bmatrix}, \quad e_1 = E_1 x^{(1)} = \begin{bmatrix} \epsilon_1(a+b) \\ \epsilon_r(a-b) \end{bmatrix},$$

و $e_r = \epsilon_r(a^2 - b^2)$ به طوری که $|\epsilon_1|, |\epsilon_r| \leq \circ/5 \times \epsilon_M$ در نتیجه خواهیم داشت

$$e_a(y) \doteq 2ae_a(a) - 2be_a(b) + (a^2 - b^2)(\epsilon_1 + \epsilon_r + \epsilon_r).$$

△

تعریف ۱۶.۲ برای هر الگوریتم کمیت $J\psi^{(1)}(\hat{x}^{(1)})e_1 + \dots + J\psi^{(r)}(\hat{x}^{(r)})e_r + e_{r+1}$ به اثر کلی خطای گرد کردن 36 آن الگوریتم معروف است. از دو الگوریتم در حل یک مسئله، آن الگوریتمی که اثر کلی خطای گرد کردن کوچک‌تری داشته باشد از نظر عددی قابل اعتمادتر 37 است، البته به شرط آن که داده‌های ورودی هر دو الگوریتم یکسان باشند.

مثال ۳۱.۲ در ادامه مثال ۳۰.۲، اثر کلی خطای گرد کردن برای الگوریتم اول حداکثر $\frac{\epsilon_M}{4}$ و برای الگوریتم دوم حداکثر $3|a^2 - b^2| \frac{\epsilon_M}{4}$ است. بنابراین اگر

$$3|a^2 - b^2| \frac{\epsilon_M}{4} < (a^2 + b^2 + |a^2 - b^2|) \frac{\epsilon_M}{4},$$

³⁶Total effect of rounding
³⁷Numerically more trustworthy

که هم‌ارز است با

$$2|a^2 - b^2| < a^2 + b^2,$$

و یا

$$-a^2 - b^2 < 2(a^2 - b^2) < a^2 + b^2,$$

که نتیجه می‌دهد $\frac{1}{4} < \left(\frac{a}{b}\right)^2 < 3$. در این صورت الگوریتم دوم از الگوریتم اول و در غیر این صورت الگوریتم اول از الگوریتم دوم قابل‌اعتمادتر است. برای $a = 0.3227$ و $b = 0.3134$ در یک ماشین با $\beta = 10$ و $t = 4$ داریم

$$r(a^2 - b^2) = 0.6580 \times 10^{-2}, \quad r((a+b)(a-b)) = 0.6562 \times 10^{-2},$$

و چون $\left(\frac{a}{b}\right)^2 \simeq 1.07$ پس الگوریتم دوم از نظر عددی قابل‌اعتمادتر است. این مطلب با توجه به مقدار دقیق‌تر $y = 0.656213 \times 10^{-2}$ تأیید می‌شود. Δ

صرف نظر از اینکه چه الگوریتمی در حل مسئله استفاده می‌شود، کران خطای زیر را داریم

$$|E_{r+1}y| \leq |y| \frac{\epsilon_M}{\gamma}, \quad |y| = [|y_1| \cdots |y_m|].$$

از طرفی چون ممکن است داده‌های ورودی اعداد ماشینی (قابل‌نمایش به طور دقیق) نباشند، کران خطای زیر را می‌توان در نظر گرفت

$$|e_a(x)| \leq |x| \frac{\epsilon_M}{\gamma}.$$

بنابراین بدون توجه به اینکه چه الگوریتمی برای حل مسئله در نظر گرفته می‌شود، باید همواره خطای

$$e_a^\circ(y) := (|J\phi(x)| |x| + |y|) \frac{\epsilon_M}{\gamma},$$

معروف به خطای ذاتی^{۳۸} را از هر الگوریتمی انتظار داشته باشیم.

تعریف ۱۷.۲ خطای گرد کردن e_i بی‌ضرر^{۳۹} است، هرگاه سهم آن در محاسبه $e_a(y)$ ، حداکثر به اندازه خطای ذاتی باشد، یعنی

$$|J\psi^{(i)}(\hat{x}^{(i)})e_i| \simeq e_a^\circ(y),$$

و اگر e_i ها همه بی‌ضرر باشند، الگوریتم از نظر عددی پایدار (خوش‌رفتار)^{۴۰} نامیده می‌شود.

مثال ۳۲.۲ در ادامه مثال ۳۱.۲، با توجه به خطای ذاتی، یعنی

$$e_a^\circ(y) = (2|a| + 2|b|) \begin{bmatrix} |a| \\ |b| \end{bmatrix} + |a^2 - b^2| \frac{\epsilon_M}{\gamma} = (2a^2 + 2b^2 + |a^2 - b^2|) \frac{\epsilon_M}{\gamma},$$

Inherent error^{۳۸}

Harmless^{۳۹}

Numerically stable(well-behaved)^{۴۰}

△

نتیجه می‌گیریم هر دو الگوریتم از نظر عددی پایدار هستند.

تذکر ۱۴.۲ جهت آشنایی با سایر مفاهیم مرتبط مانند گراف نظیر یک الگوریتم، الگوریتم خوش‌خیم^{۴۱} و غیره به مقاله [۴] مراجعه کنید.

مثال‌های زیر از مرجع [۲۲] را به دقت بررسی کنید.

مثال ۳۳.۲ می‌خواهیم ریشه با اندازه کوچکتر معادله $y^2 + 2py - q = 0$ را بیابیم که در آن $p, q > 0$ و $p \gg q$. یعنی قصد داریم مسئله خوش‌حالت $y = \phi(p, q) = -p + \sqrt{p^2 + q}$ را حل کنیم. ابتدا الگوریتم زیر به ذهن می‌رسد

$$x = x^{(0)} = \begin{bmatrix} p \\ q \end{bmatrix}, x^{(1)} = \begin{bmatrix} p \\ q \\ p^2 \end{bmatrix}, x^{(2)} = \begin{bmatrix} p \\ p^2 + q \end{bmatrix}, x^{(3)} = \begin{bmatrix} p \\ \sqrt{p^2 + q} \end{bmatrix}, x^{(4)} = y = -p + \sqrt{p^2 + q},$$

و از آنجا داریم

$$\psi^{(1)}(u, v, w) = -u + \sqrt{w + v}, \quad \psi^{(2)}(u, v) = -u + \sqrt{v}, \quad \psi^{(3)}(u, v) = v - u,$$

و بنابراین

$$J\psi^{(1)}(x^{(1)}) = [-1 \quad \frac{1}{2\sqrt{p^2+q}} \quad \frac{1}{2\sqrt{p^2+q}}]^T, \quad J\psi^{(2)}(x^{(2)}) = [-1 \quad \frac{1}{2\sqrt{p^2+q}}]^T, \quad J\psi^{(3)}(x^{(3)}) = [-1 \quad 1]^T,$$

پس

$$e_1 = \begin{bmatrix} \circ \\ \circ \\ \epsilon_1 p^2 \end{bmatrix}, \quad e_2 = \begin{bmatrix} \circ \\ \epsilon_2 (p^2 + q) \end{bmatrix}, \quad e_3 = \begin{bmatrix} \circ \\ \epsilon_3 \sqrt{p^2 + q} \end{bmatrix}, \quad e_4 = \epsilon_4 y.$$

در نتیجه

$$e_a(y) \doteq \frac{-y}{\sqrt{p^2+q}} e_a(p) + \frac{1}{2\sqrt{p^2+q}} e_a(q) + \frac{p^2}{2\sqrt{p^2+q}} \epsilon_1 + \frac{\sqrt{p^2+q}}{2} \epsilon_2 + \sqrt{p^2+q} \epsilon_3 + e_4.$$

از طرفی چون

$$e_a^\circ(y) = \left(\frac{py}{\sqrt{p^2+q}} + \frac{q}{2\sqrt{p^2+q}} + y \right) \frac{\epsilon_M}{2},$$

پس داریم

$$\frac{e_a^\circ(y)}{y} = \left(\frac{p}{\sqrt{p^2+q}} + \frac{p + \sqrt{p^2+q}}{2\sqrt{p^2+q}} + 1 \right) \frac{\epsilon_M}{2},$$

و بنابراین

$$\frac{\epsilon_M}{2} < \frac{e_a^\circ(y)}{y} < 3 \frac{\epsilon_M}{2}.$$

متناظر با عمل $\sqrt{\cdot}$ خطای $\frac{\sqrt{p^2+q}}{-p+\sqrt{p^2+q}} \epsilon_2 = k \epsilon_2$ تولید می‌شود که در آن $k = \frac{p\sqrt{p^2+q+p^2+q}}{q} \simeq \frac{2p^2}{q}$ و چون $p \gg q$ به وضوح $k \gg 3$ و این حاکی از مضر بودن خطای e_2 و در نتیجه ناپایداری عددی الگوریتم است. حال با داده ورودی $x = x^{(0)} = [p \ q]^T$ الگوریتم زیر را در

نظر می‌گیریم

$$x^{(1)} = \begin{bmatrix} p \\ q \\ p^r \end{bmatrix}, \quad x^{(r)} = \begin{bmatrix} p \\ q \\ p^r + q \end{bmatrix}, \quad x^{(r)} = \begin{bmatrix} p \\ q \\ \sqrt{p^r + q} \end{bmatrix}, \quad x^{(r)} = \begin{bmatrix} p \\ q \\ p + \sqrt{p^r + q} \end{bmatrix}, \quad y = \frac{q}{p + \sqrt{p^r + q}},$$

و از آنجا داریم

$$\psi^{(r)}(u, v, w) = \frac{v}{u + w}, \quad J\psi^{(r)}(x^{(r)}) = \left[\frac{-q}{(p + \sqrt{p^r + q})^2} \quad \frac{1}{p + \sqrt{p^r + q}} \quad \frac{-q}{(p + \sqrt{p^r + q})^2} \right]^T,$$

و چون $e_r = [0 \ 0 \ \epsilon_r \sqrt{p^r + q}]^T$ ، متناظر با عمل $\sqrt{\cdot}$ خطای $\epsilon_r = k\epsilon_r$ ، $k = -\frac{\sqrt{p^r + q}}{p + \sqrt{p^r + q}}$ به وضوح $|k| < 1$ حاکی از بی‌ضرر بودن خطای گرد کردن نظیر است و با بررسی اثر سایر e_i ها می‌توان نتیجه گرفت این الگوریتم از نظر عددی پایدار است. برنامه `sqr.nb` را ببینید. Δ

مثال ۳۴.۲ برای $k \in \mathbb{N}$ و $x \in \mathbb{R}$ داده‌شده، مقدار $\cos(kx)$ را از رابطه زیر به دست می‌آوریم

$$\cos(m+1)x = 2 \cos x \cos(mx) - \cos(m-1)x, \quad m = 1, \dots, k-1.$$

اگر در محاسبه $\cos x$ خطایی داشته باشیم، چقدر در محاسبه $\cos(kx)$ اثر دارد؟

برای مسئله $y = \phi(x) = \cos(kx)$ ، با تعریف $x^{(0)} = x$ و $\phi^{(0)}(x) = \cos x := c$ به وضوح داریم

$$\psi^{(1)}(c) = \cos(kx) = \cos(k \cos^{-1} c), \quad J\psi^{(1)}(c)e_1 = \frac{k \sin(kx)}{\sin x} \cos x e_1 = k \cot x \sin(kx) e_1, \quad |e_1| \leq \frac{\epsilon_M}{2}.$$

با توجه به $\frac{\epsilon_M}{2} = (k|x \sin(kx)| + |\cos(kx)|) \frac{\epsilon_M}{2}$ ، واضح است که برای $|x|$ کوچک و $|kx| \simeq 1$ ، خطای e_1 مضر و در نتیجه الگوریتم ناپایدار است. Δ

مثال ۳۵.۲ برای $k \in \mathbb{N}$ و $x \in \mathbb{R}$ داده‌شده، مقدار $\cos(kx)$ و $\sin(kx)$ را از روابط بازگشتی زیر به دست می‌آوریم

$$\begin{cases} \cos(mx) = \cos x \cos(m-1)x - \sin x \sin(m-1)x, \\ \sin(mx) = \sin x \cos(m-1)x + \cos x \sin(m-1)x, \end{cases} \quad m = 1, \dots, k. \quad (5.2)$$

اگر در محاسبه $\cos x$ و $\sin x$ خطایی داشته باشیم، چقدر در محاسبه $\cos(kx)$ و $\sin(kx)$ اثر دارد؟

یک ماتریس دوران به صورت زیر در نظر بگیرید

$$U = \begin{bmatrix} c & -s \\ s & c \end{bmatrix},$$

که در آن $c = c_1 = \cos x$ و $s = s_1 = \sin x$ ، U یک ماتریس متعامد^{۴۲} است، یعنی $U^{-1} = U^T$ و به کمک آن می‌توان روابط بازگشتی

(۵.۲) را به صورت زیر خلاصه کرد

$$\begin{bmatrix} c_m \\ s_m \end{bmatrix} = U \begin{bmatrix} c_{m-1} \\ s_{m-1} \end{bmatrix}, \quad m = 1, \dots, k,$$

که در آن $c_m = \cos(mx)$ و $s_m = \sin(mx)$ بنابراین

$$\begin{bmatrix} c_k \\ s_k \end{bmatrix} = U^k \begin{bmatrix} c_0 \\ s_0 \end{bmatrix} = U^k \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

از طرفی بنابر

$$\frac{\partial U}{\partial c} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I, \quad \frac{\partial U}{\partial s} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} := A.$$

داریم

$$\frac{\partial U^k}{\partial s} = \frac{\partial}{\partial s}(UU^{k-1}) = \frac{\partial U}{\partial s}U^{k-1} + U \frac{\partial U^{k-1}}{\partial s} = AU^{k-1} + U \frac{\partial}{\partial s}(UU^{k-2}) = AU^{k-1} + U \left(\frac{\partial U}{\partial s}U^{k-2} + U \frac{\partial U^{k-2}}{\partial s} \right).$$

با تکرار این فرایند خواهیم داشت

$$\frac{\partial U^k}{\partial s} = AU^{k-1} + UAU^{k-2} + \dots + U^{k-1}A = kAU^{k-1},$$

زیرا $AU = UA$ و چون $\frac{\partial U}{\partial c} = I$ داریم $\frac{\partial U^k}{\partial c} = kU^{k-1}$ بنابراین

$$\frac{\partial U^k}{\partial c} = k \begin{bmatrix} c_{k-1} & -s_{k-1} \\ s_{k-1} & c_{k-1} \end{bmatrix}, \quad \frac{\partial U^k}{\partial s} = k \begin{bmatrix} -s_{k-1} & -c_{k-1} \\ c_{k-1} & -s_{k-1} \end{bmatrix}.$$

برای مسئله $y = \phi(x) = \begin{bmatrix} \cos(kx) \\ \sin(kx) \end{bmatrix}$ با تعریف $x^{(c)} = x$ و $\phi^{(c)}(x) = \begin{bmatrix} \cos x \\ \sin x \end{bmatrix}$ به وضوح داریم

$$\psi^{(1)} \begin{bmatrix} c \\ s \end{bmatrix} = U^k \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} c_k \\ s_k \end{bmatrix},$$

و از آنجا می‌توان نوشت

$$J\psi^{(1)}(x^{(1)})e_1 = \frac{\partial U^k}{\partial c} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \cos x \epsilon_1 + \frac{\partial U^k}{\partial s} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \sin x \epsilon_1 = \epsilon_1 k \cos x \begin{bmatrix} c_{k-1} \\ s_{k-1} \end{bmatrix} + \epsilon_1 k \sin x \begin{bmatrix} -s_{k-1} \\ c_{k-1} \end{bmatrix}.$$

با توجه به

$$e_a^\circ(y) = \begin{bmatrix} k|x \sin(kx)| + |\cos(kx)| \\ k|x \cos(kx)| + |\sin(kx)| \end{bmatrix} \frac{\epsilon_M}{2},$$

واضح است که برای $|x|$ کوچک و $|kx| \simeq 1$ داریم $|k \sin x| \simeq 1$ و $|k \cos x| \simeq k$ و بنابراین خطای محاسبه $\sin x$ بی‌ضرر ولی خطای محاسبه $\cos x$ مضر و در نتیجه الگوریتم ناپایدار است. به عنوان تمرین بررسی کنید که محاسبه $\cos(kx)$ با این الگوریتم برای $|x|$ کوچک و $|kx| \simeq 1$ نسبت به الگوریتم مثال قبل از نظر عددی قابل اعتمادتر است. Δ

مثال ۳۶.۲ برای $k \in \mathbb{N}$ و $x \in \mathbb{R}^+$ داده‌شده، به منظور محاسبه مقادیر $\cos(kx)$ و $\sin(kx)$ به کمک روابط بازگشتی

$$\begin{cases} \cos(m+1)x = \cos x \cos(mx) - \sin x \sin(mx), \\ \sin(m+1)x = \sin x \cos(mx) + \cos x \sin(mx), \end{cases} \quad m = 1, \dots, k-1,$$

و روابط

$$\begin{cases} dc_{m+1} := \cos(m+1)x - \cos(mx), \\ = 2(\cos x - 1)\cos(mx) - \sin x \sin(mx) - \cos x \cos(mx) + \cos(mx), \\ = -4 \sin^2\left(\frac{x}{2}\right) \cos(mx) + (\cos(mx) - \cos(m-1)x), \end{cases}$$

و

$$\begin{cases} ds_{m+1} := \sin(m+1)x - \sin(mx), \\ = 2(\cos x - 1)\sin(mx) + \sin x \cos(mx) - \cos x \sin(mx) + \sin(mx), \\ = -4 \sin^2\left(\frac{x}{2}\right) \sin(mx) + (\sin(mx) - \sin(m-1)x), \end{cases}$$

روابط بازگشتی زیر معرفی می‌شوند

$$\begin{cases} c_m := c_{m-1} + dc_m, & dc_{m+1} := t c_m + dc_m, \\ s_m := s_{m-1} + ds_m, & ds_{m+1} := t s_m + ds_m, \end{cases} \quad m = 1, \dots, k,$$

که در آن

$$c_0 = 1, s_0 = 0, t = 2dc_1, dc_1 = -2 \sin^2\left(\frac{x}{2}\right), ds_1 = \sin x = 2 \sin\left(\frac{x}{2}\right) \cos\left(\frac{x}{2}\right) = \sqrt{-dc_1(2 + dc_1)}.$$

برای مسئله $y = \phi(x) = \begin{bmatrix} \cos(kx) \\ \sin(kx) \end{bmatrix}$ با تعریف $x^{(\circ)} = x$ و $s = \phi^{(\circ)}(x) = \sin\left(\frac{x}{2}\right)$ ، اگر در محاسبه $\sin\left(\frac{x}{2}\right)$ خطایی داشته باشیم، چقدر در محاسبه $\cos(kx)$ و $\sin(kx)$ اثر دارد؟

به وضوح داریم

$$\psi^{(1)}(s) = \begin{bmatrix} c_k \\ s_k \end{bmatrix},$$

و از آنجا می توان نوشت

$$J\psi^{(1)}(s)e_1 = \begin{bmatrix} \frac{\partial c_k}{\partial s} \\ \frac{\partial s_k}{\partial s} \\ \frac{\partial s}{\partial s} \end{bmatrix} s \epsilon = \begin{bmatrix} \frac{\partial c_k}{\partial x} \frac{\partial x}{\partial s} \\ \frac{\partial s_k}{\partial x} \frac{\partial x}{\partial s} \\ \frac{\partial x}{\partial x} \frac{\partial x}{\partial s} \end{bmatrix} s \epsilon = \sqrt{k} \tan\left(\frac{x}{\sqrt{k}}\right) \begin{bmatrix} -s_k \\ c_k \end{bmatrix} \epsilon.$$

واضح است که برای x کوچک و $kx \simeq 1$ ، خطای محاسبه s بی ضرر بوده و در نتیجه الگوریتم پایدار است. به عنوان تمرین بررسی کنید که محاسبه $\cos(kx)$ با این الگوریتم برای x کوچک و $kx \simeq 1$ نسبت به الگوریتم مثال قبل از نظر عددی قابل اعتمادتر است. Δ

برنامه `cskx.nb` را بررسی کنید. مثال 5 صفحه 25 مرجع [۲۲] را ببینید.

۵.۲ حساب بازه‌ای

به کمک آنالیز خطای پیشرو به سادگی می توان اثر خطاهای گرد کردن میانی یک الگوریتم را در حل یک مسئله نه چندان پیچیده مشخص کرد. تعیین اثر کلی خطای گرد کردن میانی الگوریتم‌هایی که از پیچیدگی بالایی برخوردار باشند به عبارتی الگوریتم‌هایی که شامل تعداد زیادی خطاهای گرد کردن میانی هستند به سادگی امکان پذیر نیست. به کمک حساب بازه‌ای^{۴۳} می توان با در نظر گرفتن خطای داده‌های ورودی و خطاهای گرد کردن میانی، برای خطای مطلق یک الگوریتم یک کران بالا به دست آورد. حساب بازه‌ای بر اساس این واقعیت است که هر عدد حقیقی $a \in \mathbb{R}$ که یا یک داده ورودی است یا یک نتیجه میانی یا نهایی، مقداری نامعلوم است و تنها اطلاعی که از آن داریم یک بازه کوچک شامل آن عدد است. بنابراین در حساب بازه‌ای متناظر با عدد a بازه $\tilde{a} = [a', a'']$ در نظر گرفته می شود که در آن a' و a'' اعداد ماشینی هستند و $\tilde{a} \in a$. متناظر با هر عملی روی اعداد باید عملی روی بازه‌ها تعریف کرد. به منظور تعریف عمل دوتایی $\circ \in \{+, -, \times, /\}$ روی بازه‌های \tilde{a} و \tilde{b} قرار می دهیم $\tilde{a} \circ \tilde{b} := \tilde{c}$ به طوری که $\tilde{c} \subset a \circ b$.

مثال ۳۷.۲ برای عمل جمع داریم $[c', c''] := [a', a''] + [b', b'']$ که در آن

$$c' = \max\{\gamma' \in \mathbb{F} | \gamma' \leq a' + b'\}, \quad c'' = \min\{\gamma'' \in \mathbb{F} | \gamma'' \geq a'' + b''\},$$

و برای عمل ضرب داریم $[c', c''] := [a', a''] \times [b', b'']$ که در آن برای $a', b' > 0$

$$c' = \max\{\gamma' \in \mathbb{F} | \gamma' \leq a' \times b'\}, \quad c'' = \min\{\gamma'' \in \mathbb{F} | \gamma'' \geq a'' \times b''\}.$$

Δ

سایر اعمال به طریق مشابه تعریف می شوند.

تذکر ۱۵.۲ بیشتر مواقع می توان از تعاریف زیر استفاده کرد

$$\left\{ \begin{array}{l} [a', a''] + [b', b''] = [a' + b', a'' + b''], \\ [a', a''] - [b', b''] = [a' - b'', a'' - b'], \\ [a', a''] \times [b', b''] = [\min\{a' \times b', a' \times b'', a'' \times b', a'' \times b''\}, \max\{a' \times b', a' \times b'', a'' \times b', a'' \times b''\}], \\ [a', a''] / [b', b''] = [a', a''] \times \left[\frac{1}{b''}, \frac{1}{b'}\right], \quad \text{if } 0 \notin [b', b'']. \end{array} \right.$$

زمانی که یک الگوریتم بر اساس حساب بازه‌ای دنبال شود در پایان یک بازه به دست می‌آید که نه تنها جواب مسئله به آن بازه تعلق دارد بلکه می‌توان به کمک آن بازه کرانی برای خطای مطلق جواب نیز تعیین کرد. اگرچه این کران خطا قابل اعتماد است ولی در بیشتر حالات ممکن است بدبینانه^{۴۴} باشد. بنابراین کافی نیست که فقط اعمال حسابی با اعمال بازه‌ای نظیر جایگزین شوند بلکه باید ملاحظاتی نیز در نظر گرفت و راه‌کارهایی به کار گرفت.

مثال ۳۸.۲ برای ارزیابی تابع $y = \phi(x) = x^3 - 3x^2 + 3x$ ابتدا بهتر است به کمک قاعده هورنر آن را به صورت $y = \phi(x) = ((x - 3)x + 3)x$

$$u = x - 3, \quad v = u \times x, \quad w = v + 3, \quad y = w \times x,$$

بازنویسی کرد. بازه نظیر اعداد دقیق مانند ۳ یا به صورت $[3, 3]$ یا با انتخاب مناسبی از $\epsilon > 0$ به صورت $[3 - \epsilon, 3 + \epsilon]$ در نظر گرفته می‌شود. حال برای ارزیابی ϕ در نقطه $x \in \tilde{x} = [0/9, 1/1]$ داریم

$$\begin{cases} \tilde{u} = \tilde{x} - [3, 3] = [-2/1, -1/9], & \tilde{v} = \tilde{u} \times \tilde{x} = [-2/31, -1/71], \\ \tilde{w} = \tilde{v} + [3, 3] = [0/69, 1/29], & \tilde{y} = \tilde{w} \times \tilde{x} = [0/621, 1/419]. \end{cases}$$

اگرچه بازه به دست آمده شامل جواب مسئله است ولی بزرگ بودن این بازه دلالت بر خطای زیادی دارد. از طرفی چون $\phi'(x) = 3(x-1)^2 \geq 0$ پس ϕ یک تابع صعودی است و باید داشته باشیم

$$\{\phi(x) | x \in \tilde{x}\} = [\phi(0/9), \phi(1/1)] = [0/999, 1/001],$$

یعنی $[0/999, 1/001] \rightarrow [0/9, 1/1] : \phi$. اگر با اعداد ماشینی دو رقمی کار کنیم، به ازای $x = 0/9$ داریم

$$u = -2/1, \quad v = -1/9, \quad w = 1/1, \quad y = 0/99,$$

و به ازای $x = 1/1$ خواهیم داشت

$$u = -1/9, \quad v = -2/1, \quad w = 0/9, \quad y = 0/99.$$

حال با بازنویسی تابع ϕ به صورت $y = \phi(x) = 1 + (x - 1)^3$ به ازای $x \in \tilde{x} = [0/9, 1/1]$ داریم

$$\begin{cases} \tilde{u} = \tilde{x} - [1, 1] = [-0/1, 0/1], & \tilde{v} = \tilde{u} \times \tilde{u} = [-0/01, 0/01], \\ \tilde{w} = \tilde{v} \times \tilde{u} = [-0/001, 0/001], & \tilde{y} = \tilde{w} + [1, 1] = [0/999, 1/001]. \end{cases}$$

اگر این بار نیز با اعداد ماشینی دو رقمی کار کنیم، به ازای $x = 0/9$ داریم

$$u = -0/1, \quad v = 0/01, \quad w = -0/001, \quad y = 1/0,$$

و به ازای $x = ۱/۱$ خواهیم داشت

$$u = ۰/۱, \quad v = ۰/۰۱, \quad w = ۰/۰۰۱, \quad y = ۱/۰.$$

△

تذکر ۱۶.۲ همان طور که در مثال قبل دیده شد، برای دنبال کردن یک الگوریتم در حساب بازه‌ای نه تنها باید اعمال حسابی را با اعمال بازه‌ای نظیر جایگزین کرد بلکه باید راه‌کارهای مناسبی نیز به کار گرفت. برای اطلاعات بیشتر در مورد حساب بازه‌ای، به مرجع [۱۸] مراجعه کنید.

پروژه ۱.۲ بازنویسی الگوریتم‌های عددی کارشناسی با حساب بازه‌ای.

سمینار ۱.۲ حساب بازه‌ای گذشته و حال با تکیه بر کاربردهای جدید.

سمینار ۲.۲ تخمین خطای گرد کردن به صورت آماری.

۶.۲ تقسیم‌بندی مسایل در ریاضیات محاسباتی

بسیاری از مسایل مطرح‌شده در ریاضیات کاربردی و به ویژه آنالیز عددی به صورت زیر مدل می‌شوند

$$T : X \rightarrow Y, \quad Tx = y, \quad (۶.۲)$$

که در آن X و Y دو فضای برداری نرم‌دار مناسب هستند که بر اساس ورودی و خروجی مسئله تعریف می‌شوند و T عملگری است که متناظر با مسئله تعریف می‌شود. مسایلی که در عمل با آنها برخورد می‌کنیم به سه دسته اساسی زیر تقسیم می‌شوند

۱. مسایل مستقیم^{۴۵}: در این مسایل T و x معلوم هستند و y مجهول است. ارزیابی توابع یک یا چندمتغیره عددی یا برداری یا ماتریسی، محاسبه مشتق یا انتگرال یک تابع یک یا چندمتغیره از این نوع به شمار می‌روند،
۲. مسایل وارون^{۴۶}: در این مسایل T و y معلوم بوده و x مجهول است. انواع معادلات، دستگاه معادلات خطی/غیرخطی، معادلات دیفرانسیل عادی یا جزئی، معادلات انتگرال، مسایل ویژه‌مقدار و مقدار تکین جزء مسایل وارون هستند،
۳. مسایل شناسایی^{۴۷}: در این مسایل x و y معلوم و T مجهول است. مسئله درونبایی معروفترین مسئله از نوع شناسایی است. به بیان مهندسی سیستم، چون ورودی و خروجی معلوم هستند، این مسئله به مهندسی وارون^{۴۸} یا جعبه سیاه^{۴۹} نیز معروف است. برخی از مسایل کمینه‌سازی تابعی^{۵۰}، برنامه‌ریزی ریاضی و آمار به طور غیرمستقیم در زمره مسایل شناسایی به حساب می‌آیند.

^{۴۵} Direct problems

^{۴۶} Inverse problems

^{۴۷} Identification problems

^{۴۸} Inverse engineering

^{۴۹} Black box

^{۵۰} Functional minimization

مسائل شناسایی از پیچیدگی بالایی برخوردارند در حالی که مسائل مستقیم ساده بوده و بسیاری از ابهامات و پرسش‌های مربوط به آنها تاکنون بررسی شده‌اند ولی با این حال سرعت پیشرفت پژوهش در این راستا کند است. بسیاری از مسائلی که امروزه با آنها مواجه می‌شویم از نوع مسائل وارون هستند زیرا این مسائل از اهمیت بالایی برخوردارند و یک جایگاه مرکزی در آنالیز عددی برای خود باز کرده‌اند. نوع خطی این مسائل هم از دیدگاه نظری هم از دیدگاه عددی به خوبی مطالعه و بررسی شده است، ولی هنوز مسائلی در این راستا جهت مطالعه و تحقیق وجود دارد. امروزه مسائل وارون غیرخطی راستای بیشتر پژوهش‌ها را به خود منعطف کرده است زیرا این مسائل پیچیده از جذابیت خاصی برخوردارند. باید توجه داشت که گاهی مواقع ممکن است حل یک مسئله از یک نوع به سادگی امکان‌پذیر نباشد و در این صورت سعی می‌کنیم آن مسئله را به یک مسئله از همان نوع یا حتی نوع دیگر تبدیل کرده و به حل آن مبادرت ورزیم.

مثال ۳۹.۲ مسئله مقدار اولیه

$$\begin{cases} \frac{dx(t)}{dt} = f(t)x(t), & t > 0, \\ x(0) = 1, \end{cases}$$

که از نوع وارون است به کمک روش جداسازی متغیرها به مسئله مستقیم زیر تبدیل می‌شود

$$x(t) = e^{\int_0^t f(s)ds}.$$

△

این در حالی است که مسئله جدید با مشتق‌گیری به همان مسئله مقدار اولیه تبدیل می‌شود.

مثال ۴۰.۲ مسئله مقدار اولیه

$$\begin{cases} \frac{dx(t)}{dt} = f(t, x(t)), & t > 0, \\ x(0) = 1, \end{cases}$$

که از نوع وارون است به کمک انتگرال‌گیری به مسئله وارون زیر تبدیل می‌شود که یک معادله انتگرال است.

$$x(t) = 1 + \int_0^t f(s, x(s))ds$$

△

مسئله جدید با مشتق‌گیری به همان مسئله مقدار اولیه تبدیل می‌شود.

تعریف ۱۸.۲ اگر X و Y دو فضای برداری باشند، فضای حاصل ضرب $X \times Y$ به صورت زیر تعریف می‌شود

$$X \times Y := \{(x, y) | x \in X, y \in Y\}.$$

فضای $X \times Y$ همراه با جمع برداری

$$(x_1, y_1) + (x_2, y_2) = (x_1 + x_2, y_1 + y_2), \quad \forall x_1, x_2 \in X, \quad \forall y_1, y_2 \in Y,$$

و ضرب اسکالر

$$\alpha(x, y) = (\alpha x, \alpha y), \quad \forall \alpha \in \mathbb{R}, \quad \forall x \in X, \quad \forall y \in Y,$$

یک فضای برداری است.

مثال ۴۱.۲ مسئله مقدار اولیه

$$\begin{cases} \frac{dx(t)}{dt} = f(t, x(t)), & 0 < t < 1, \\ x(0) = \alpha, \end{cases}$$

را به صورت عملگری بنویسید.

با معرفی $X = C^1[0, 1]$ ، $Y = C[0, 1] \times \mathbb{R}$ و $Tx(t) = \left(\frac{dx(t)}{dt} - f(t, x(t)), x(0) \right)$ به وضوح $T: X \rightarrow Y$ و شکل عملگری مسئله به صورت $Tx(t) = (0, \alpha)$ خلاصه می‌شود.

Δ

تمرین ۷.۲ شکل عملگری مسایل مقدار مرزی یا مقدار اولیه-مرزی

$$\begin{cases} -\Delta u(x, y) = f(x, y), & (x, y) \in \Omega, \\ u(x, y) = g(x, y), & (x, y) \in \partial\Omega, \end{cases} \quad \begin{cases} u_t(x, y, t) = \Delta u(x, y, t) + f(x, y, t), & (x, y, t) \in \Omega \times \mathbb{R}^+, \\ u(x, y, 0) = u_0(x, y), & (x, y) \in \Omega, \\ u(x, y, t) = g(x, y, t), & (x, y, t) \in \partial\Omega, \end{cases}$$

را بنویسید که در آن $\Omega \in \mathbb{R}^2$ یک مجموعه باز و همبند و $\partial\Omega$ مرز Ω است و $\Delta = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2}$. این مسایل را برای حالتی که $\Omega \in \mathbb{R}^1$ یا $\Omega \in \mathbb{R}^2$ نیز بازنویسی کرده و شکل عملگری آنها را بنویسید.

تمرین ۸.۲ مسئله ویژه مقدار را به صورت عملگری بنویسید.

بیشتر مواقع فضاها X و Y نامتناهی‌البعد بوده و به جای آنها از فضاها (زیرفضاهای) متناهی‌البعد X_n و Y_n استفاده می‌کنیم و در واقع مسئله عملگری (۶.۲) را به مسئله ساده‌تر

$$T_n: X_n \rightarrow Y_n, \quad T_n x_n = y_n,$$

گسسته‌سازی می‌کنیم.

مثال ۴۲.۲ مسئله انتگرال‌گیری معین $y = Tx$ که در آن عملگری از $X = C[a, b]$ به توی $Y = \mathbb{R}$ است و $Tx = \int_a^b x(t) dt$ را در نظر بگیرید. این مسئله به کمک یک قاعده انتگرال‌گیری به صورت زیر گسسته‌سازی می‌شود

$$X_n = \mathbb{R}^n, \quad Y_n = \mathbb{R}, \quad x_n = [x(t_1) \cdots x(t_n)]^T, \quad y_n = T_n x_n = \sum_{i=1}^n a_i x(t_i).$$

Δ

تعریف ۱۹.۲ اگر در حل یک مسئله وارون داشته باشیم $\|x - x_n\| \leq B(n)$ ، تابع $B(n)$ به کران بالای خطا^{۵۱} یا تخمین اولیه (پیشین)^{۵۲} معروف است. بیشتر مواقع داریم $B(n) = cn^{-p}$ و در این حالت می‌نویسیم $x = x_n + O(n^{-p})$ و می‌گوییم روش خطایی از مرتبه p دارد.

مثال ۴۳.۲ در روش دویخشی در یافتن ریشه یکتای معادله $x(t) = 0$ متعلق به بازه $[a, b]$ داریم $|t - t_n| \leq \frac{b-a}{\gamma^n}$. همچنین در روش

انتگرال‌گیری دوزنقه یعنی

$$T(h) = \frac{h}{\gamma} \left(x(a) + 2 \sum_{i=1}^{n-1} x(t_i) + x(b) \right),$$

^{۵۱} Upper error bound

^{۵۲} A priori estimate

که در آن

$$h = (b - a)/n, \quad t_i = a + ih, \quad i = 0, \dots, n,$$

داریم

$$\left| \int_a^b x(t) dt - T(h) \right| \leq \frac{(b-a)^2}{12n^2} \max_{t \in [a,b]} |x''(t)|.$$

بنابراین اگر $x \in C^2[a, b]$ روش دوزنقه خطایی از مرتبه دو دارد. Δ

تذکر ۱۷.۲ بعضی مواقع داریم $B(n) = B(n, x)$. این کران جنبه نظری داشته و در عمل کاربرد ندارد ولی اگر $B(n) = B(n, x_n)$ در این حالت به کران (تخمین) ثانویه (پسین)^{۵۳} معروف است و به کمک اطلاعاتی از مسئله اصلی یا اعمال فرضیاتی، به یک کران قابل محاسبه تبدیل می‌شود.

مثال ۴۴.۲ در روش نقطه ثابت برای یافتن تنها ریشه $x(t) = 0$ در بازه $[a, b]$ ، ابتدا معادله را به صورت $t = z(t)$ بازنویسی کرده و تلاش می‌کنیم نقطه ثابت z یعنی همان ریشه x را بیابیم. قضیه نقطه ثابت را جهت یادآوری بیان می‌کنیم.

قضیه ۷.۲ اگر $T: [a, b] \rightarrow [a, b]$ و $z \in C[a, b]$ و z در $[a, b]$ دست کم یک نقطه ثابت دارد؛

ب) هم‌چنین اگر $z'(x)$ بر (a, b) موجود باشد و عدد ثابت و مثبت L چنان وجود داشته باشد که

$$\forall t \in (a, b), \quad |z'(t)| \leq L < 1,$$

آن‌گاه نقطه ثابت z در $[a, b]$ یکتا است؛

پ) به علاوه دنباله $\{t_n\}_{n \geq 1}$ تولید شده توسط فرایند $t_n = z(t_{n-1})$ به ازای هر t_0 در $[a, b]$ به نقطه ثابت z یعنی t همگرا است و

$$\begin{cases} |t - t_n| \leq L^n |t - t_0|, \\ |t - t_n| \leq \frac{L^n}{1-L} |t_1 - t_0|, \\ |t - t_n| \leq \frac{L}{1-L} |t_n - t_{n-1}|. \end{cases}$$

اولین نابرابری دلالت بر کران بالای خطا دارد و به کمک آن همگرایی ثابت می‌شود، نابرابری دوم یک تخمین اولیه و نابرابری سوم یک تخمین ثانویه است. Δ

تذکر ۱۸.۲ اگر کران بالای خطا با خطای واقعی فاصله زیادی داشته باشد، بدبینانه^{۵۴} و در غیر این صورت واقع‌بینانه^{۵۵} است.

تعریف ۲۰.۲ اگر داشته باشیم $x - x_n = e_n + g_n$ که در آن $\|g_n\| / \|e_n\| = 0$ در آن‌گاه e_n به تخمین خطای مجانبی^{۵۶} معروف است.

A posteriori estimate^{۵۳}
Pessimistic^{۵۴}
Realistic^{۵۵}
Asymptotic error estimate^{۵۶}

مثال ۴۵.۲ به کمک بسط تیلور داریم

$$x''(t_0) = \frac{x(t_0 - h) - 2x(t_0) + x(t_0 + h)}{h^2} + C_4 h^2 x^{(4)}(t_0) + C_6 h^4 x^{(6)}(t_0) + \dots$$

با معرفی $e_h = C_4 h^2 x^{(4)}(t_0) + \dots$ و $g_h = C_6 h^4 x^{(6)}(t_0) + \dots$ واضح است که $|g_h/e_h| \rightarrow 0$ هرگاه $h \rightarrow 0$ و x به اندازه کافی هموار باشد. بنابراین e_h تخمین خطای مجانی و از مرتبه دو است.

△

