

به نام پروردگار یگانه و یکتا



دانشگاه صنعتی اصفهان - دانشکده علوم ریاضی

محاسبات عددی

پدیدآورنده : رضا مختاری

ویرایش ششم : زمستان ۱۳۹۹ و بهار ۱۴۰۰

فهرست مطالب

۱	خطاها
۱	۱.۱ منابع تولید خطا
۳	۲.۱ نمایش اعداد
۴	۳.۱ نمایش اعداد در رایانه
۴	۱.۳.۱ نمایش ۶۴-بیتی ممیز شناور
۶	۲.۳.۱ اعداد ماشینی
۸	۴.۱ انواع خطا
۱۰	۵.۱ خطای محاسبات (فرمول)
۱۲	۱.۵.۱ خطای اعمال ریاضی
۱۵	۲.۵.۱ تقریب توابع یک متغیره
۲۲	۲ ریشه‌یابی (حل معادلات غیرخطی)
۲۲	۱.۲ بررسی کمی ریشه‌ها
۲۴	۲.۲ دنباله‌های همگرا
۲۷	۳.۲ روش‌های عددی
۲۸	۱.۳.۲ روش دوبخشی
۳۰	۲.۳.۲ روش نابجایی
۳۲	۳.۳.۲ روش تکرار ساده
۳۷	۴.۳.۲ روش نیوتن
۴۲	۵.۳.۲ روش وتری
۴۶	۳ حل عددی معادلات دیفرانسیل عادی
۴۶	۱.۳ روش بسط تیلور
۴۸	۲.۳ روش اویلر
۴۹	۳.۳ روش‌های رانگ-کوتا
۵۱	۴.۳ دستگاه معادلات دیفرانسیل مرتبه اول
۵۲	۱.۴.۳ روش اویلر

۵۳	۲.۴.۳ روش اویلر اصلاح شده
۵۵	۳.۴.۳ روش تیلور
۵۵	۴.۴.۳ روش رانگ-کوتای چهار مرحله‌ای
۵۷	۵.۴.۳ معادلات دیفرانسیل عادی مرتبه بالاتر
۵۹	۴ درون‌یابی
۶۰	۱.۴ درون‌یابی
۶۱	۱.۱.۴ روش لاگرانژ
۶۳	۲.۱.۴ روش تفاضلات تقسیم‌شده نیوتن
۶۶	۳.۱.۴ روش‌های پیشرو/پسرو نیوتن
۷۰	۲.۴ خطای چندجمله‌ای درون‌یاب
۷۲	۳.۴ برون‌یابی و درون‌یابی وارون
۷۵	۴.۴ تقریب کم‌ترین مربعات گسسته
۸۱	۵ مشتق‌گیری و انتگرال‌گیری عددی
۸۱	۱.۵ مشتق‌گیری عددی
۸۱	۱.۱.۵ روش مبتنی بر بسط تیلور
۸۴	۲.۵ انتگرال‌گیری عددی
۸۵	۱.۲.۵ قاعده ذوزنقه
۸۷	۲.۲.۵ قاعده سیمسون
۸۹	۳.۲.۵ قاعده نقطه میانی
۹۰	۴.۲.۵ قاعده‌های نیوتن-کاتس
۹۲	۵.۲.۵ کوادراتور گاوس
۹۶	۶.۲.۵ روش رامبرگ
۱۰۰	۶ دستگاه معادلات خطی
۱۰۱	۱.۶ روش‌های مستقیم
۱۰۱	۱.۱.۶ روش حذف گاوسی
۱۰۶	۲.۱.۶ روش حذفی گاوس-جردن
۱۱۲	۳.۱.۶ تجزیه مثلثی
۱۱۷	۲.۶ روش‌های تکراری
۱۱۷	۱.۲.۶ نرم برداری و ماتریسی
۱۱۹	۲.۲.۶ روش‌های مبتنی بر تفکیک ماتریسی
۱۲۳	۳.۶ آنالیز خطا

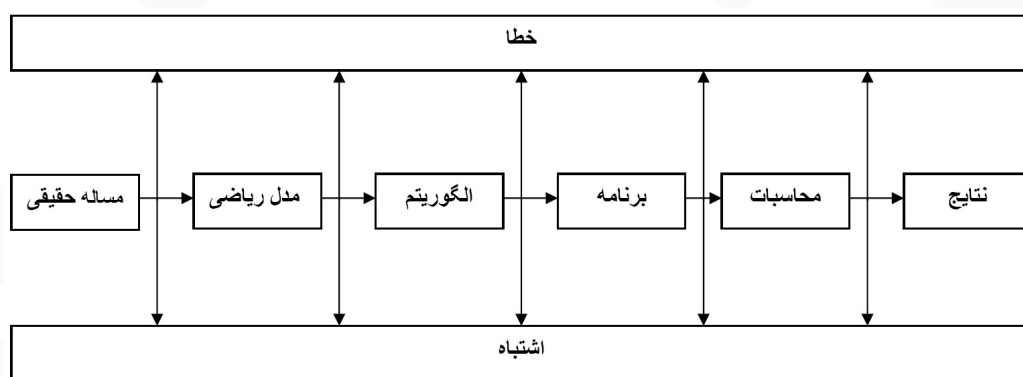
فصل ۱

خطاها

زمانی که به دست آوردن جواب دقیق (واقعی یا تحلیلی) یک مسئله به سادگی امکان پذیر نیست و یا مقرون به صرفه نیست، به کمک روش‌های عددی، یک جواب تقریبی برای مسئله پیدا می‌کنیم. این فرایند به تولید خطا منجر می‌شود. در این فصل قصد داریم منابع تولید خطا و انواع خطا را شناسایی کرده و تا حدی از انتشار خطا جلوگیری کنیم.

۱.۱ منابع تولید خطا

بیشتر مواقع در عمل با یک مسئله حقیقی (فیزیکی) مواجه هستیم و بنابر دلایلی، جواب تقریبی (عددی) آن را جستجو می‌کنیم. مراحل یافتن جواب تقریبی یک مسئله حقیقی در روندنمای آمده در شکل ۱.۱ خلاصه می‌شود. این روندنما مکان‌های احتمالی بروز خطا و همچنین اشتباهات را نیز نشان می‌دهد.



شکل ۱.۱: فرایند تولید جواب عددی (تقریبی)

تذکر ۱.۱ اختلافی که بین جواب دقیق و تقریبی وجود دارد ممکن است از اشتباهات و خطاها ناشی شده باشد. اشتباه را می‌توان برطرف کرد ولی خطا بیشتر اوقات اجتناب ناپذیر است. به عنوان مثال قراردادن 2322 به جای 2232 یک اشتباه است و استفاده از $3/14$ به جای عدد π که بسط دهدهی نامختوم دارد، موجب بروز خطا می‌شود.

مراحل مختلف روندنمای آمده در شکل ۱.۱ را در مثال بعد دنبال می‌کنیم.

مثال ۱.۱ (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 در زمان t باشد. به کمک برخی از قوانین و اصول فیزیک و ریاضیات، صرف نظر از مقاومت هوا و اصطکاک در لولا، مدل این مسئله به صورت

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

یا $\frac{d^2 \theta}{dt^2} = -\frac{g}{l} \sin \theta$ به دست می‌آید که یک معادله دیفرانسیل عادی غیرخطی است. با فرض کوچک بودن θ یعنی

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

$$\theta = 15^\circ \simeq 0.262 \text{ rad}, \quad \sin \theta \simeq 0.259,$$

می‌توان فرض کرد $\sin \theta \simeq \theta^{\text{rad}}$ و معادله دیفرانسیل غیرخطی را به صورت زیر نوشت

$$\frac{d^2 \theta}{dt^2} + \omega^2 \theta = 0, \quad \omega^2 = \frac{g}{l}.$$

این معادله دیفرانسیل خطی جوابی متناوب به صورت $\theta(t) = c_1 \sin \omega t + c_2 \cos \omega t$ دارد و بنابراین دوره تناوب آونگ ساده عبارت است از $T = \frac{2\pi}{\omega} = 2\pi \sqrt{\frac{l}{g}}$.

△

با توجه به این مثال، می‌توان خطاها را از نقطه نظر منبع تولید به صورت زیر تقسیم‌بندی کرد.

انواع خطا

۱. ذاتی

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

۲. محاسباتی

- نمایش اعداد (مانند $\pi = 3.14$)
- اعمال ریاضی (به عنوان مثال $\frac{l}{g}$)
- روش‌های (الگوریتم‌های) عددی (محاسباتی) (مثل خطای روش محاسبه $(\sqrt{\frac{l}{g}})$)

تذکر ۲.۱ خطاهای ذاتی به محاسبات عددی مربوط نمی‌شود اما برای پرهیز از خطاهای محاسباتی و کنترل آن‌ها باید راه چاره‌ای پیدا کرد.

تذکر ۳.۱ 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^۱ در رایانه آن

۲.۱ نمایش اعداد

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

قضیه ۱.۱ هر عدد حقیقی مثبت x نمایشی به صورت

$$\begin{aligned} x &= a_m \beta^m + a_{m-1} \beta^{m-1} + \dots + a_1 \beta^1 + a_0 \beta^0 + a_{-1} \beta^{-1} + a_{-2} \beta^{-2} + \dots \\ &= (a_m a_{m-1} \dots a_1 a_0 / a_{-1} a_{-2} \dots)_\beta, \end{aligned} \quad (1.1)$$

دارد که در آن $a_m \neq 0$ و $a_i \in \{0, 1, 2, \dots, \beta - 1\}$, $m \in \mathbb{Z}$.

تذکر ۴.۱ رابطه (۱.۱) 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).

قضیه ۲.۱ نمایش یک عدد گویا (در هر مبنایی) یا مختوم است یا نامختوم متناوب.

نتیجه ۱.۲.۱ بسط یک عدد گنگ، نامختوم نامتناوب است.

مثال ۲.۱ به موارد زیر در مبناهای متفاوت توجه کنید

$$\begin{aligned} \frac{2}{3} &= 0.666\dots = 0.\bar{6} = (0/2)_3, & \frac{2}{8} &= 0.25 = (0/3)_8, & \sqrt{2} &= 1.4142\dots, \\ 0/1 &= (0/00011)_2, & \frac{1}{4} &= 0.25 = (0/01)_2, & \pi &= 3.141592\dots \end{aligned}$$

△

تذکر ۵.۱ اگر چه فرض $a_m \neq 0$ برای یکنایی نمایش (۱.۱) در نظر گرفته شده است، برای منحصر به فرد بودن نمایش (۱.۱) به فرض‌های دیگری نیز نیاز است. به مثال زیر توجه کنید

$$\begin{aligned} 3/47999\dots &= 3/47\bar{9} = 3 \times 10^0 + 4 \times 10^{-1} + 7 \times 10^{-2} + 9 \times 10^{-3} + 9 \times 10^{-4} + \dots \\ &= 3/47 + \frac{9 \times 10^{-3}}{1 - 10^{-1}} = 3/47 + 0/01 = 3/48. \end{aligned}$$

یعنی برای اعداد $3/47\bar{9}$ و $3/48$ یک نمایش وجود دارد. اگر فرض کنیم عدد صحیح z چنان وجود داشته باشد که $0 = a_j = a_{j-1} = \dots$ (به عبارتی فرض کنیم بسط مختوم باشد) این مشکل برطرف می‌شود.

تذکر ۶.۱ هنگام کار با رایانه (ماشین حساب) اعداد را در مبنای 10 وارد کرده و انتظار داریم نتایج (خروجی) نیز در همین مبنا نمایش داده شود ولی این وسایل با مبنای دیگری (امروزه مبنای 2 و در قدیم مبنای دیگری مانند 16) کار می‌کنند. بنابراین مسئله تغییر مبنا مطرح می‌شود که ممکن است خطایی به دنبال داشته باشد که در اینجا از بررسی آن صرف نظر می‌کنیم.

۳.۱ نمایش اعداد در رایانه

برای نمایش اعداد در ماشین، ابتدا نمایشی به نام ممیز ثابت^۲ در نظر گرفته شد که در آن هر عدد حقیقی x به صورت زیر نمایش داده می‌شود

$$x = \pm(a_n a_{n-1} \dots a_1 a_0 / a_{-1} a_{-2} \dots a_{-m}) \beta,$$

که در آن m و n اعداد مشخص و ثابتی هستند. در اصل در این نمایش مکان ممیز مشخص و ثابت است. 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.

یک روش جدید و متداول برای نمایش اعداد در رایانه، نمایش ممیز (نقطه) شناور (سیار)^۳ است که از بدو پیدایش مورد توجه سازندگان سخت‌افزار رایانه قرار گرفت و تا حدودی به طور سلیقه‌ای با آن برخورد شد تا زمانی که استنادی توسط IEEE^۴ وضع شد.

۱.۳.۱ نمایش ۶۴-بیتی ممیز شناور

این نمایش پیش از این به دقت دو برابر (مضاعف)^۵ معروف بوده و متناظر با نوع double در زبان C است. در این نمایش برای نمایش هر عدد در مبنای 2 ، ابتدا یک ساختار به طول 64 بیت در نظر گرفته می‌شود. اولین بیت به بیت علامت

^۲ Fixed point

^۳ Floating point

^۴ IEEE standard 754-1985

^۵ Double precision

معروف است و با s نمایش داده شود و بلافاصله بعد از آن ۱۱ بیت برای مشخصه^۶ در نظر گرفته می شود و با c نمایش داده می شود و ۵۲ بیت باقی مانده به نام مانتیس منظور می شود که آن را با f نشان می دهند. اگرچه برای مانتیس ۵۲ بیت در نظر گرفته شده است ولی در واقع ساختار $1+f$ موجب می شود که ۵۳ بیت داشته باشیم که آن یک بیت اضافه به بیت پنهان معروف است و برای یکتایی نمایش لازم است. برای هر عدد نمایشی به صورت $(-1)^s \times 2^e \times (1+f)$ در نظر گرفته می شود که در آن $e = c - 1023$ (نما) است. محدودیت $2047 = 2^{11} - 1 = (110001)_2 \leq c \leq 2054$ موجب می شود که

$$-1023 \leq e \leq 1024.$$

از طرف دیگر محدودیت $(0.110001)_2 \leq f \leq (1.110001)_2$ باعث می گردد که $0 \leq f \leq 2$ بنابراین اگر کوچک ترین و بزرگ ترین عدد مثبت قابل نمایش را به ترتیب با mN (در محیط MATLAB) و MN (در محیط MATLAB) نشان دهیم آنگاه

$$mN = 2^{-1022} \times (1.000000)_2 \approx 2.22507 \times 10^{-308},$$

$$MN = 2^{1024} \times (1.000000)_2 \approx 1.79769 \times 10^{308}.$$

تذکره ۷.۱ باید توجه داشت که توان -1023 برای نمایش صفر و توان 1024 با مانتیس مثبت برای نمایش ∞ (در محیط MATLAB) و صور مبهم (NaN در محیط MATLAB) مورد استفاده قرار می گیرد.

تعریف ۱.۱ در نمایش اعداد ماشینی، کوچک ترین عدد مثبت ماشینی که اگر به ۱ اضافه شود عددی بزرگ تر از ۱ به دست می آید به اپسیلون ماشینی^۷ معروف است و با eps نمایش داده می شود.

چون در نمایش ۶۴-بیتی بعد از ۱ عدد $(1.000000)_2$ قرار می گیرد، پس

$$\text{eps} = (1.000000)_2 - 1 = (0.000000)_2 = 2^{-52} \approx 2.220446 \times 10^{-16}.$$

حال بزرگترین عدد صحیح مثبت M را تعیین می کنیم 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 $(1.000000)_2 \times 2^{52}$ قابل نمایش است. اما تعداد ارقام در مانتیس جهت نمایش $2^{52} + 1$ کافی نیست (۵۳ رقم در مانتیس لازم است). بنابراین $M = 2^{52} \approx 9.007199 \times 10^{15}$ و در نتیجه تمام اعداد صحیح ۱۵ رقمی و بسیاری از اعداد ۱۶ رقمی در این نمایش به طور دقیق قابل نمایش هستند (گفته می شود دقت در نمایش ۶۴-بیتی ۱۵ الی ۱۶ رقم است). البته بسیاری از توان های ۲ نیز به طور دقیق قابل نمایش هستند.

تذکر ۸.۱ مجموعه اعداد با ممیز شناور را با \mathbb{FP} نمایش می‌دهیم. محور اعداد با ممیز شناور بر خلاف محور اعداد حقیقی، متناهی و گسسته است و بر خلاف آنچه به نظر می‌رسد، نقاط روی این محور هم‌فاصله نیستند. اگر در هنگام انجام محاسبات، عددی در فاصله‌ی $(-mN, mN)$ تولید شود پیام پاریز^۸ و اگر عدد تولید شده از MN بزرگ‌تر یا از $-MN$ کوچک‌تر شود پیام سرریز^۹ صادر می‌شود. از توان‌های 10^{23} و 10^{24} برای نمایش پیام‌های پاریز و سرریز نیز استفاده می‌شود.

تذکر ۹.۱ بسیاری از نرم‌افزارها مانند MATLAB از محدودیت‌های سخت‌افزار پیروی می‌کنند و بعضی از نرم‌افزارها مانند Mathematica، محدودیت‌های سخت‌افزار را از طریق برنامه‌های نرم‌افزاری برطرف می‌کنند و به اصطلاح دقت را بالا می‌برند که ممکن است با کاهش سرعت انجام محاسبات همراه باشد.

۲.۳.۱ اعداد ماشینی

به منظور سادگی نوشتار، برای هر عدد ماشینی^{۱۰} (یا عناصر مجموعه اعداد \mathbb{FP})، نمایشی به صورت

$$\text{نما (توان)} \leftarrow \pm \underbrace{d_1 \dots d_k}_{\text{مانتیس}} \times 10^n$$

در نظر گرفته می‌شود که در آن n به زیرمجموعه‌ای از \mathbb{Z} مانند $[L, U]$ تعلق دارد و به ازای $i = 1, 2, \dots, k$ داریم $d_i \in \{0, 1, \dots, 9\}$ و برای یکتایی نمایش فرض می‌شود $d_1 \neq 0$. به این نمایش، نمایش ممیز شناور ده‌دهی نرمال شده گفته می‌شود و چنین اعدادی را، اعداد ماشینی ده‌دهی k -رقمی نیز می‌نامند. ایده چنین نمایشی، از نمایش علمی نرمال شده اعداد ناشی شده است. در نمایش علمی نرمال شده، هر عدد حقیقی مخالف صفر x را می‌توان به صورت

$$\text{نما (توان)} \leftarrow \pm \underbrace{d_1 d_2 \dots d_k d_{k+1} d_{k+2} \dots}_{\text{مانتیس}} \times 10^n$$

نمایش داد که در آن $n \in \mathbb{Z}$ و به ازای $i = 1, 2, \dots, k$ داریم $d_i \in \{0, 1, \dots, 9\}$ و $d_1 \neq 0$. واضح است که برای نمایش x به صورت ممیز شناور ده‌دهی نرمال شده در یک ماشین $(k$ -رقمی)، باید k رقم از مانتیس آن را حفظ کرده و بقیه را کنار گذاشت، که برای این کار روش‌های زیر موجود است

۱. روش قطع کردن (برش)^{۱۱}

۲. روش گرد کردن معمولی^{۱۲}

۳. روش گرد کردن به زوج^{۱۳}

^۸ Underflow

^۹ Overflow

^{۱۰} Machine number

^{۱۱} Chopping

^{۱۲} Rounding

^{۱۳} Rounding to even

در روش قطع کردن، k رقم از ماتیس حفظ و بقیه کنار گذاشته می‌شود در حالی که در روش گرد کردن معمولی، ابتدا روش قطع کردن اعمال شده، سپس اگر $d_{k+1} \geq 5$ یک واحد به d_k 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 d_k اضافه می‌گردد

$$d_{k+1} > 5 \bullet$$

$$d_{k+1} = 5 \text{ و رقم مخالف صفری در سمت راست } d_{k+1} \text{ مشاهده شود}$$

$$d_{k+1} = 5 \text{ و رقم مخالف صفری در سمت راست } d_{k+1} \text{ مشاهده نشود و } d_k \text{ فرد باشد.}$$

قضیه ۳.۱ اگر $x = \circ/d_1 d_2 d_3 \dots d_k d_{k+1} \dots \times 10^n$ یک عدد حقیقی ناصفر در نمایش علمی نرمال شده باشد، $\tilde{fl}(x) = \circ/d_1 d_2 d_3 \dots d_k \times 10^n$ عدد ماشینی k -رقمی متناظر با x است که از روش قطع کردن به دست می‌آید و

$$|x - \tilde{fl}(x)| \leq 1 \times 10^{n-k}.$$

هم‌چنین $fl(x) = \circ/\delta_1 \delta_2 \dots \delta_m \times 10^m$ ($m = n + 1$ یا $m = n$) عدد ماشینی k رقمی متناظر با x است که از روش گرد کردن معمولی به دست می‌آید و خواهیم داشت

$$|x - fl(x)| \leq 5 \times 10^{n-k-1} = \circ/5 \times 10^{n-k}.$$

پرسش ۱.۱ تفاوت 47km با 47000m یا تفاوت $3/7$ با $3/70$ یا $3/700$ در چیست؟

تعریف ۲.۱ منظور از ارقام بامعنا 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 است (تمام ارقام ماتیس در نمایش علمی نرمال شده).

مثال ۳.۱ عدد $\circ/0007045000$ حداقل ۴ رقم بامعنا و حداکثر ۷ رقم بامعنا دارد. \triangle

قرارداد ۱.۱ rD یعنی r رقم اعشار (تنظیم ماشین حساب روی Mode fix r) و rS یعنی r رقم بامعنا (تنظیم ماشین حساب روی Mode sci r) ^{۱۴}.

مثال ۴.۱ در جدول ۱.۱، چند عدد و مقادیر تقریبی متناظر با آن‌ها با روش‌های قطع کردن (\hat{x}) و گرد کردن معمولی (\tilde{x}) با دقت $3S$ و $3D$ داده شده است. \triangle

تذکر ۱۰.۱ از این به بعد، گرد کردن معمولی را به کار می‌بریم.

x	۲۸,۶۴۲۴	۰,۰۰۵۷۶۷۱	۴,۹۸۵۰	-۲۱۷۵,۳۴۵۱۲
$\tilde{x}(3S)$	۲۸,۶	۰,۰۰۵۷۷	۴,۹۹	-۲۱۸۰
$\tilde{x}(3D)$	۲۸,۶۴۲	۰,۰۰۰۶	۴,۹۸۵	-۲۱۷۵,۳۴۵
$\hat{x}(3S)$	۲۸,۶	۰,۰۰۵۷۶	۴,۹۸	-۲۱۷۰
$\hat{x}(3D)$	۲۸,۶۴۲	۰,۰۰۰۵	۴,۹۸۵	-۲۱۷۵,۳۴۵

جدول ۱.۱: مثال‌هایی از گرد کردن و قطع کردن معمولی

پرسش ۲.۱ در نظر بگیرید \hat{x} تقریبی از x باشد. آیا هر چه تعداد ارقام با معنای \hat{x} بیشتر باشد، می‌توان گفت \hat{x} تقریب بهتری (به معنای خطای کمتر) است؟

تعریف ۳.۱ فرض کنید $\hat{x} = a_m \times 10^m + \dots + a_1 \times 10^1 + a_0 \times 10^0 + a_{-1} \times 10^{-1} + \dots + a_{-k} \times 10^{-k}$ که در آن $a_m \neq 0$ تقریبی از عدد حقیقی مثبت x باشد. منظور از تعداد ارقام با معنای درست \hat{x} نسبت به x بزرگ‌ترین n می‌باشد ($n \in \mathbb{N}_0$) است که در نابرابری‌های زیر صدق کند.

$$n \leq m + k + 1, \quad |x - \hat{x}| \leq 5 \times 10^{m-n}$$

تذکر ۱۱.۱ m در تعریف اخیر، نشان دهنده با ارزش‌ترین مکان عدد است.

مثال ۵.۱ از اعداد $\hat{x} = ۹۹/۹۶$ و $\tilde{x} = ۱۰۰/۷$ کدام یک تقریب بهتری برای $x = ۱۰۰$ است؟

$$\hat{x} = ۹۹/۹۶ \rightarrow m = ۱, \quad |x - \hat{x}| = ۰,۰۰۴ \leq 5 \times 10^{1-n} \rightarrow n = ۳$$

$$\tilde{x} = ۱۰۰/۷ \rightarrow m = ۲, \quad |x - \tilde{x}| = ۰,۷ \leq 5 \times 10^{2-n} \rightarrow n = ۲$$

△

مثال ۶.۱ تعداد ارقام با معنای درست $\hat{x} = ۱۰۰/۳۱$ را نسبت به $x = ۱۰۰/۳۱۰۴$ مشخص کنید.

$$\hat{x} = ۱۰۰/۳۱ \rightarrow m = ۲, \quad |x - \hat{x}| = ۰,۰۰۰۰۴ \leq 5 \times 10^{2-n} \rightarrow n = ۶$$

△

و چون \hat{x} فقط پنج رقم با معنا دارد پس $n = ۵$.

تذکر ۱۲.۱ در نظر بگیرید \hat{x} تقریبی از x باشد. هر چه تعداد ارقام با معنای درست \hat{x} نسبت به x بیشتر باشد، \hat{x} تقریب بهتری خواهد بود.

۴.۱ انواع خطا

تعریف ۴.۱ اگر \hat{x} تقریبی از x باشد $\Delta x = |x - \hat{x}|$ را خطای مطلق \hat{x} نسبت به x نامند. Δx منحصر به فرد است و در عمل بیشتر مواقع قابل تعیین نیست و به جای آن از هر عدد b_x استفاده می‌شود که کمتر از Δx نباشد. b_x منحصر

به فرد نیست و به آن کران خطای مطلق گویند. بنابراین $\Delta x \leq b_x$ و در نتیجه $\hat{x} - b_x \leq x \leq \hat{x} + b_x$. بعضی مواقع از نمایش $x = \hat{x} \pm b_x$ استفاده می‌شود.

مثال ۷.۱ اگر عدد $\hat{x} = ۱,۷۳۲$ را به عنوان تقریبی از $x = \sqrt{۳}$ در نظر بگیریم، آنگاه

$$\Delta x = |\sqrt{۳} - ۱,۷۳۲| = ۱,۷۳۲۰۵۰۸۰۷۵\dots - ۱,۷۳۲ = ۰,۰۰۰۰۵۰۸۰۷۵\dots,$$

از طرفی می‌دانیم $۱,۷۳۲۰ < \sqrt{۳} < ۱,۷۳۲۱$ بنابراین $۰ < \sqrt{۳} - ۱,۷۳۲ < ۰,۰۰۰۱$ پس $b_x = ۰,۰۰۰۱$ که معیاری برای نزدیکی $۱,۷۳۲$ به $\sqrt{۳}$ است. Δ

پرسش ۳.۱ آیا خطای مطلق معیار مناسبی برای مقایسه خطاها است؟

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

تعریف ۵.۱ اگر $\hat{x} \neq ۰$ تقریبی از x باشد $\delta x = \frac{|x - \hat{x}|}{|x|} = \frac{\Delta x}{|x|}$ 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. $۱۰۰ \times \delta x$ به درصد خطا معروف است.

قضیه ۴.۱ اگر \hat{x} تقریبی از x باشد و b_x یک کران خطای مطلق برای این تقریب باشد آنگاه

$$\delta x \leq \frac{b_x}{|\hat{x}| - b_x}.$$

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

$$\delta x \leq \frac{b_x}{|\hat{x}|}.$$

مثال ۸.۱ اگر عدد $\hat{x} = ۱,۷۳۲$ را به عنوان تقریبی از $x = \sqrt{۳}$ در نظر بگیریم، آنگاه

$$\delta x = \frac{|\sqrt{۳} - ۱,۷۳۲|}{\sqrt{۳}} = \frac{۱,۷۳۲۰۵۰۸۰۷۵\dots - ۱,۷۳۲}{۱,۷۳۲۰۵۰۸۰۷۵\dots} = \frac{۰,۰۰۰۰۵۰۸۰۷۵\dots}{۱,۷۳۲۰۵۰۸۰۷۵\dots}.$$

پس

$$\delta x = ۰,۰۰۰۰۲۹۳۳۳۷\dots < ۰,۰۰۰۰۰۳.$$

اما با توجه به قضیه ۴.۱ و $b_x = ۰,۰۰۰۱$ می‌توان نوشت

$$\delta x \leq \frac{۰,۰۰۰۱}{۱,۷۳۲ - ۰,۰۰۰۱} = \frac{۰,۰۰۰۱}{۱,۷۳۱۹} = ۰,۰۰۰۰۵۷۷۴۰۰۵\dots < ۰,۰۰۰۰۰۶$$

$$\delta x \leq \frac{۰,۰۰۰۱}{۱,۷۳۲} = ۰,۰۰۰۰۰۵۷۷۳۶۷۲\dots < ۰,۰۰۰۰۰۶.$$

Δ

قضیه ۵.۱ اگر \hat{x} تقریبی از x با n رقم بامعنای درست باشد و $\hat{y} = 10^t \times \hat{x}$ و $y = 10^t \times x$ که در آن t عددی صحیح است آنگاه \hat{y} نیز تقریبی از y با n رقم بامعنای درست است و خطای نسبی \hat{y} و \hat{x} برابر است.

قضیه ۶.۱ اگر \hat{x} گردشده x تا n رقم بامعنا باشد آنگاه \hat{x} دارای n رقم بامعنای درست است.

قضیه ۷.۱ ارتباط دقت (خطای نسبی) با تعداد ارقام بامعنای درست

اگر \hat{x} دارای n رقم بامعنای درست باشد، آنگاه $\delta x < 5 \times 10^{-n}$ به شرط آن که ارقام بامعنای درست \hat{x} از یک رقم ۱ و $n-1$ رقم صفر در جلوی آن تشکیل نشده باشد. برعکس اگر $\delta x \leq 5 \times 10^{-n-1} = 0.5 \times 10^{-n}$ آنگاه \hat{x} دارای n رقم بامعنای درست است.

مثال ۹.۱ تقریبی از $\sqrt{3}$ ارایه دهید که خطای نسبی آن از 10^{-4} کمتر باشد. بنابر قضیه ۷.۱ اگر \hat{x} تقریبی از $\sqrt{3}$ باشد که ۵ رقم بامعنای درست داشته باشد آنگاه $10^{-4} < \delta x < 5 \times 10^{-5}$. از این رو، با توجه به قضیه ۶.۱ کافی است \hat{x} گرد شده $\sqrt{3}$ تا ۵ رقم بامعنا باشد، یعنی $\hat{x} = 1.7321$. \triangle

۵.۱ خطای محاسبات (فرمول)

فرض کنید $z = f(x_1, \dots, x_n)$ تابعی باشد که می‌خواهیم آن را در نقطه (x_1, \dots, x_n) ارزیابی کنیم. بدون کاستن از کلیت فرض کنید $\Delta x_i = x_i - \hat{x}_i$ که در آن \hat{x}_i مقدار تقریبی x_i است. پس می‌توان نوشت

$$z = f(x_1, \dots, x_n) = f(\hat{x}_1 + \Delta x_1, \dots, \hat{x}_n + \Delta x_n).$$

بنابر بسط تیلور توابع n متغیره داریم

$$z = f(\hat{x}_1, \dots, \hat{x}_n) + \left(\Delta x_1 \frac{\partial f}{\partial x_1} + \dots + \Delta x_n \frac{\partial f}{\partial x_n} \right) (\hat{x}_1, \dots, \hat{x}_n) + R,$$

که در آن R جمله خطا بوده و شامل حاصل‌ضرب‌ها و توان‌های Δx_i ها است و چون Δx_i ها کوچک هستند از R چشم‌پوشی کرده، خواهیم داشت

$$f(x_1, \dots, x_n) - f(\hat{x}_1, \dots, \hat{x}_n) \simeq \left(\Delta x_1 \frac{\partial f}{\partial x_1} + \dots + \Delta x_n \frac{\partial f}{\partial x_n} \right) (\hat{x}_1, \dots, \hat{x}_n),$$

و اگر $\hat{z} = f(\hat{x}_1, \dots, \hat{x}_n)$ آنگاه

$$\Delta z = |z - \hat{z}| \simeq \left| \left(\Delta x_1 \frac{\partial f}{\partial x_1} + \dots + \Delta x_n \frac{\partial f}{\partial x_n} \right) (\hat{x}_1, \dots, \hat{x}_n) \right|,$$

و بلافاصله داریم

$$\delta z = \frac{\Delta z}{|z|} \simeq \left| \frac{\left(\Delta x_1 \frac{\partial f}{\partial x_1} + \dots + \Delta x_n \frac{\partial f}{\partial x_n} \right) (\hat{x}_1, \dots, \hat{x}_n)}{f(\hat{x}_1, \dots, \hat{x}_n)} \right|.$$

مثال ۱۰.۱ (مستقیم) یک استوانه به شعاع قاعده $\frac{4}{\pi}$ و ارتفاع $\sqrt{2}$ را در نظر بگیرید. اگر شعاع و ارتفاع استوانه و عدد π را با دقت $4D$ وارد محاسبات کنیم، حجم این استوانه با چه خطایی به دست می آید؟ می دانیم حجم یک استوانه از قاعده $\pi r^2 h$ تعیین می شود که در آن شعاع قاعده h و ارتفاع استوانه است. اگر تعریف کنیم $z = V(p, r, h) = \pi r^2 h$ به ازای $p = \pi$ ، $r = \frac{4}{\pi}$ و $h = \sqrt{2}$ ارزیابی شود. تمام محاسبات را با دقت $4D$ دنبال می کنیم. بنابراین

$$p = \pi \quad \rightarrow \quad \hat{p} = 3,1416$$

$$r = \frac{4}{\pi} \quad \rightarrow \quad \hat{r} = 1,27323$$

$$h = \sqrt{2} \quad \rightarrow \quad \hat{h} = 1,4142$$

و داریم $\Delta p, \Delta r, \Delta h \leq 0,5 \times 10^{-4}$. بنابراین

$$\hat{z} = V(\hat{p}, \hat{r}, \hat{h}) = 3,1416 \times 1,27323^2 \times 1,4142 = 7,8980.$$

از طرفی

$$\Delta z \simeq \Delta p \hat{r}^2 \hat{h} + 2 \Delta r \hat{p} \hat{r} \hat{h} + \Delta h \hat{p} \hat{r}^2 \leq (\hat{r}^2 \hat{h} + 2 \hat{p} \hat{r} \hat{h} + \hat{p} \hat{r}^2) \times 0,5 \times 10^{-4} < 10^{-3}.$$

در نتیجه حجم استوانه برابر است با $7,898 \pm 0,001$. به کمک یک ماشین حساب 10 رقمی به دست می آوریم $z = 7,898458555$ و بنابراین خطای واقعی عبارت است از

$$\Delta z = |z - \hat{z}| = 4,585554 \times 10^{-4} < 0,001.$$

△

مثال ۱۱.۱ (معکوس) 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 دقت $2D$ حساب شود؟ اگر فرض کنیم

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

آنگاه

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

و در نتیجه با فرض $x = ۲/۲$ و $y = ۰/۳$ در نظر می‌گیریم داریم

$$\Delta z \simeq \left| \Delta x \frac{\partial f(۲/۲, ۰/۳)}{\partial x} + \Delta y \frac{\partial f(۲/۲, ۰/۳)}{\partial y} \right| \simeq |۴۸/۹ \Delta x + ۴۷/۹ \Delta y|,$$

و برای برقراری نابرابری $\Delta z \leq ۰/۵ \times ۱۰^{-۲}$ باید داشته باشیم

$$|۴۸/۹ \Delta x + ۴۷/۹ \Delta y| \leq ۰/۵ \times ۱۰^{-۲}.$$

یک جواب نامعادله اخیر عبارت است از

$$\Delta x \leq ۰/۵ \times ۱۰^{-۴}, \quad \Delta y \leq ۰/۵ \times ۱۰^{-۴}.$$

بنابراین برای رسیدن به نتیجه مطلوب، x و y را می‌توان با دقت $۴D$ در نظر گرفت.

مثال ۱۲.۱ اگر در محاسبه $z = ab^2c^3$ ، خطای نسبی تقریب‌های a ، b و c حداکثر $۰/۰۱$ باشد، بیشترین خطای نسبی قابل انتظار برای z چقدر است؟ با فرض $z = f(a, b, c) = ab^2c^3$ داریم

$$\delta z \simeq \frac{\left| \Delta a \frac{\partial f}{\partial a} + \Delta b \frac{\partial f}{\partial b} + \Delta c \frac{\partial f}{\partial c} \right|}{|ab^2c^3|} = \frac{|b^2c^3 \Delta a + 2abc^3 \Delta b + 3ab^2c^2 \Delta c|}{|ab^2c^3|} \leq \delta a + 2\delta b + 3\delta c = ۰/۰۶.$$

△

۱.۵.۱ خطای اعمال ریاضی

هنگام کار با اعداد ماشینی (ممیز شناور) دقت شود که بعضی از اصول میدان \mathbb{R} نظیر شرکت پذیری، منحصر به فرد بودن عضو خنثی و غیره برقرار نیست. به طور کلی، عبارت‌هایی که از نظر ریاضی معادل هستند ممکن است از نظر محاسباتی معادل نباشند. علت اصلی بروز این مشکلات، خطای گرد کردن است.

تعریف ۶.۱ فرض کنید A و B دو عدد حقیقی، a و b تقریب‌هایی از آن‌ها و \otimes بیان‌گر یک عمل دوتایی باشد. متناظر با $A \otimes B$ در ماشین عمل $a \otimes^* b$ انجام می‌شود و داریم

$$|A \otimes B - a \otimes^* b| = |(A \otimes B - a \otimes b) + (a \otimes b - a \otimes^* b)| \leq \underbrace{|A \otimes B - a \otimes b|}_{\text{خطای تولید شده}} + \underbrace{|a \otimes b - a \otimes^* b|}_{\text{خطای منتشر شده}}$$

قضیه ۸.۱ اگر \hat{x} و \hat{y} تقریب‌هایی از x و y بوده و همه این اعداد مثبت باشند، آنگاه

$$\delta(x+y) \leq \max\{\delta x, \delta y\} \quad \delta(x \pm y) \leq \frac{x}{|x \pm y|} \delta x + \frac{y}{|x \pm y|} \delta y \quad \Delta(x \pm y) \leq \Delta x + \Delta y . ۱$$

$$.۲ \quad \delta(xy) \leq \delta x + \delta y \quad , \quad \Delta(xy) \leq x\Delta y + y\Delta x$$

$$.۳ \quad \delta\left(\frac{x}{y}\right) \leq \delta x + \delta y \quad , \quad \Delta\left(\frac{x}{y}\right) \leq \frac{y\Delta x + x\Delta y}{y^2}$$

مثال ۱۳.۱ اگر \hat{x} و \hat{y} هر یک n رقم بامعنای درست داشته باشند، حداقل تعداد ارقام بامعنای درست $\hat{x} + \hat{y}$ و $\hat{x}\hat{y}$ را تعیین کنید. بنابر قضیه ۷.۱، $\delta x < 5 \times 10^{-n}$ و $\delta y < 5 \times 10^{-n}$. هم چنین داریم

$$\delta(x + y) \leq \max\{\delta x, \delta y\} < 5 \times 10^{-n} = 5 \times 10^{-(n-1)-1}.$$

پس $\hat{x} + \hat{y}$ حداقل $n - 1$ رقم بامعنای درست دارد. به علاوه می توان نوشت

$$\delta(\hat{x}\hat{y}) \leq \delta x + \delta y < 5 \times 10^{-n} + 5 \times 10^{-n} = 10 \times 10^{-n} = 10^{-(n-2)-1} < 5 \times 10^{-(n-2)-1}.$$

بنابراین $\hat{x}\hat{y}$ دست کم $n - 2$ رقم بامعنای درست دارد. \triangle

تذکر ۱۳.۱ با توجه به خطای نسبی $\delta(x - y) \simeq \frac{\Delta(x-y)}{|x-y|}$ واضح است که اگر \hat{x} و \hat{y} دو عدد هم علامت نزدیک به هم باشند $|x - y|$ کوچک و در نتیجه $\delta(x - y)$ بزرگ خواهد شد و در نتیجه تعداد ارقام بامعنای $\hat{x} - \hat{y}$ کم خواهد بود. بنابراین در عمل بهتر است از تفاضل دو عدد هم علامت نزدیک به هم جلوگیری شود (تفاضل دو عدد هم علامت نزدیک به هم موجب از بین رفتن ارقام بامعنا^{۱۵} می شود مانند $1/42 - 1/41 = 0/01$). اگر تفاضل اجتناب ناپذیر است باید عمل با دقت دو برابر (یا بیشتر) انجام شود.

مثال ۱۴.۱ جهت اجتناب از تفاضل در محاسبات می توان از اتحادها کمک گرفت. به عنوان مثال

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

\triangle

مثال ۱۵.۱ با فرض $g(x) = x \left(\sqrt{1 + \frac{1}{x}} - 1 \right)$ مقدار $g(10^9)$ را با دقت یک ماشین حساب 10^0 رقمی به دست آورید. تمام محاسبات را با دقت $9D$ انجام می دهیم. پس

$$1 + \frac{1}{x} = 1.0000000001, \quad \sqrt{1 + \frac{1}{x}} = 1.0000000005.$$

بنابراین با دقت $9D$ ، ماشین حساب نتیجه زیر را به دست می دهد

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

به این ترتیب با یک ماشین حساب ۱۰ رقمی $g(10^9) = 0.0000000000$ که نتیجه‌ای نادرست است. در حالی که اگر از یک رایانه استفاده شود، نتیجه درست $g(10^9) = 0.3333333333$ است. دلیل این خطای فاحش، تفاضل دو عدد نزدیک به هم در محاسبات است که منجر به از بین رفتن ارقام بامعنا می‌شود. برای به دست آوردن تقریب بهتر، روش محاسبه را به صورت زیر تغییر می‌دهیم

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

در این صورت با توجه به

$$\left(1 + \frac{1}{x}\right)^2 = 1.0000000002, \quad \sqrt{\left(1 + \frac{1}{x}\right)^2} = 1.0000000000,$$

با دقت ۹D نتیجه زیر به دست می‌آید

$$g(10^9) = \frac{1}{1.0000000000 + 1.0000000000 + 1} = 0.3333333333.$$

△

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

تذکر ۱۴.۱ با توجه به قضیه ۸.۱، در محاسبات باید از ضرب اعداد بزرگ در اعداد تقریبی (تقسیم اعداد تقریبی به اعداد کوچک) پرهیز کرد.

مثال ۱۶.۱ در محاسبه $x = 10000\pi$ داریم

$$\pi = 3/14 \quad \rightarrow \quad \hat{x} = 31400,$$

$$\pi = 3/142 \quad \rightarrow \quad \hat{x} = 31420,$$

$$\pi = 3/1416 \quad \rightarrow \quad \hat{x} = 31416.$$

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

△

تذکر ۱۵.۱ 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 ۱۷.۱ به جای عبارت $ax^3 + bx^2 + cx + d$ از عبارت $((ax + b)x + c)x + d$ استفاده شود.

۲.۵.۱ تقریب توابع یک متغیره

قضیه ۹.۱ (تیلور با باقیمانده لاگرانژ) فرض کنید $f \in C^n[a, b]$ (یعنی تابع f و مشتقات تا مرتبه n آن روی بازه $[a, b]$ پیوسته هستند) 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 $y(x)$ 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

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

در این جا p_n چند جمله‌ای تیلور مرتبه n م f حول x_0 و $R_n(x, x_0)$ جمله باقی‌مانده (یا خطای برش^{۱۶}) متناظر با p_n نامیده می‌شود. اگر $n \rightarrow \infty$ آنگاه p_n به یک سری بی‌پایان تبدیل می‌شود که به آن سری تیلور f حول نقطه x_0 گویند. در این حالت، شرط بی‌نهایت بار مشتق‌پذیر بودن f در x_0 لازم است.

تذکر ۱۶.۱ در قضیه قبل اگر $x_0 = 0$ آنگاه واژه تیلور به مک‌لورن تبدیل می‌شود.

تذکر ۱۷.۱ مقدار خطا در استفاده از p_n به جای f را نشان می‌دهد. در عمل با یافتن کرانی برای جمله باقی‌مانده، در واقع برای خطای تقریب f با p_n کرانی پیدا می‌کنیم.

تذکر ۱۸.۱ ویژگی مهم چند جمله‌ای تیلور مرتبه n م آن است که p_n و مشتقات تا مرتبه n م آن با f و مشتقات تا مرتبه n م آن در نقطه x_0 برابر هستند.

تذکر ۱۹.۱ شکل دیگر (کاربردی) قضیه تیلور به صورت زیر است

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

و یا

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

که در آن $h = x - x_0$.

مثال ۱۸.۱ فرض کنید $f(x) = 2 + 4x - x^2$. آنگاه به وضوح برای $n \geq 3$ داریم

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

و برای $n = 2$ خواهیم داشت

$$p_2(x) = 2 + 4x, \quad R_2(x, 0) = -x^2,$$

و برای $n = 1$ می‌توان نوشت

$$p_1(x) = 2 + 4x, \quad R_1(x, 0) = -3x^2 y(x),$$

که در آن $y(x)$ نقطه‌ای بین 0 و x است. \triangle

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 ۱۹.۱ مطلوب است محاسبه مقدار $e^{\frac{\pi}{10}}$ با دقت $2D$.

روش اول- به کمک قضیه تیلور داریم

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

که در آن $0 < y(x) < x$ در نتیجه

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

که در آن $0 < y < \frac{\pi}{10} < 1 < \frac{\pi}{10} < 1$. چون e^x تابعی صعودی است پس $3 < e^1 < e^y < e^0 = 1$. از طرفی $0.315 < \frac{\pi}{10}$ (دقت $3D$ منظور می‌شود) و بنابراین

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

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

$$e^{\frac{\pi}{10}} \simeq 1 + 0.315 + \frac{(0.315)^2}{2!} + \frac{(0.315)^3}{3!} = 1.370,$$

و با دقت $2D$ داریم $e^{\frac{\pi}{10}} \simeq 1.37$ که با جواب ماشین حساب یعنی 1.369107771 کمتر از 0.001 اختلاف دارد. روش دوم- به کمک سری تیلور داریم

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

در نتیجه

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

با توجه به $\frac{\pi}{10} < 0,315$ (دقت ۳D منظور می شود) از $0,5 \times 10^{-2} < \left|\frac{(0,315)^n}{n!}\right|$ نتیجه می شود $n \geq 4$. پس

$$\begin{aligned} e^{\frac{\pi}{10}} &\approx 1 + 0,315 + \frac{(0,315)^2}{2!} + \frac{(0,315)^3}{3!} + \frac{(0,315)^4}{4!} \\ &= 1,315 + 0,050 + 0,005 + 0,000 = 1,370. \end{aligned}$$

△

تذکر ۲۰.۱ بعضی مواقع ممکن است تعداد جملاتی که از روش دوم به دست می آید کافی نباشد و بهتر است با جملات بیشتر هم مقایسه کرد.

مثال ۲۰.۱ (همگرایی سریع) می خواهیم تابع $\cos x$ را به ازای مقادیر $|x| < \frac{\pi}{4}$ با دقت ۵D 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 می دهد $n \geq 6$.

مثال ۲۱.۱ (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 یا

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

برای ارزیابی $\ln(1+x)$ با دقت ۵D، باید داشته باشیم

$$\left| \frac{x^n}{n} \right| < 0,5 \times 10^{-5},$$

△

که برای $x = 0,99$ نتیجه می‌دهد $n \geq 582$.

مثال ۲۲.۱ تابع $Si(x) = \int_0^x \frac{\sin t}{t} dt$ برای $x \geq 0$ را در نظر بگیرید. دست کم چند جمله از بسط مک‌لورن تابع $f(t) = \sin t$ لازم است تا $Si(1)$ با دقت ϵD مشخص شود؟
با اعمال قضیه تیلور برای تابع f داریم

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

و یا

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

در نتیجه

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

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

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

حال می‌توان نوشت

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

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

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

و یا

$$Si(1) \approx 1 - 0,0555556 + 0,0016667 - 0,0000283 + 0,0000003 = 0,9460831.$$

△

پس با دقت ϵD داریم $Si(1) \approx 0,946083$.

تمرین

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

۲. مطلوب است تعیین تقریبی از عدد π با دقت $3D$ به کمک بسط مکلورن تابع $f(x) = \tan^{-1} x$.

۳. حداقل چند جمله از بسط مکلورن تابع $f(x) = e^{-\frac{\pi}{3}}$ نیاز است تا $f(\ln(2))$ با دقت $3D$ به دست آید؟

۴. حداکثر خطای نسبی محاسبه حجم یک مخروط از فرمول $V = \frac{\pi}{3} r^2 h$ چقدر است، هرگاه تقریبی از عدد π با خطای نسبی حداکثر 0.003 داشته باشیم و کمیت‌های r و h را بتوان با حداکثر خطای نسبی 0.002 اندازه‌گیری کرد.

۵. برای محاسبه عبارت $e^{1-\ln(2/\sqrt{3})}$ در یک ماشین حساب با سه رقم بامعنا کدام گزینه مناسب‌تر است؟

(الف) $e^{1-\ln(2/\sqrt{3})}$ (ب) $\frac{e}{\sqrt{3}}$ (ج) $e \times e^{-\ln(2/\sqrt{3})}$ (د) $\frac{e}{e^{\ln(2/\sqrt{3})}}$

۶. برای همه گزینه‌ها a_0 تقریبی از عدد π با دقت $3S$ است. در کدام مورد خطای مطلق در محاسبه a_1 کمتر است؟

(الف) $a_n = 1 + 2a_{n-1}$ (ب) $a_n = 1 + na_{n-1}$ (ج) $a_n = 1 + a_{n-1}$ (د) $a_n = 1 + \frac{1}{n}a_{n-1}$

۷. برای محاسبه عبارت $\frac{1-\cos(x)}{x}$ در کامپیوتر، وقتی x عدد مثبت کوچکی است، کدام گزینه مناسب‌تر است؟

(الف) $\frac{1-\cos(x)}{x}$ (ب) $\frac{2 \sin^2(\frac{x}{2})}{x}$ (ج) $\frac{\sin^2(x)}{x(1+\cos(x))}$ (د) $\frac{1}{x} - \frac{\cos(x)}{x}$

۸. برای محاسبه $\sqrt{x^2 + \frac{1}{x^2}} - \sqrt{x^2 - \frac{1}{x^2}}$ برای x ‌های خیلی بزرگ در ماشین کدام گزینه مناسب‌تر است؟

(الف) $\frac{1}{x^2}$ (ب) $\frac{2}{x(\sqrt{x^2+1} + \sqrt{x^2-1})}$ (ج) $\sqrt{x^2 + \frac{1}{x^2}} - \sqrt{x^2 - \frac{1}{x^2}}$ (د) $\frac{1}{x} (\sqrt{x^4 + 1} - \sqrt{x^4 - 1})$

۹. در یک ماشین با دقت $4S$ و گرد کردن، کوچکترین عدد طبیعی x که در تساوی $x + 1 = x$ صدق می‌کند، کدام

گزینه است؟ (الف) $10^4 \times 0.5000$ (ب) $10^5 \times 0.1000$ (ج) $10^5 \times 0.5000$ (د) $10^6 \times 0.1000$

۱۰. برای زوایای کمتر از 6 درجه، $\sin(x)$ را با مقدار x بر حسب رادیان تقریب می‌زنند. کران بالای خطای مطلق این

کار چقدر است؟

۱۱. حداقل چند جمله از بسط مکلورن تابع $f(x) = \sin(\frac{x}{\pi})$ لازم است تا $\sin(0.1)$ با دقت $4D$ مشخص شود؟

۱۲. درصد خطای محاسبه عدد π و اندازه‌گیری شعاع دایره به ترتیب در هر یک از گزینه‌های زیر داده شده است. با

انتخاب کدام گزینه می‌توان مساحت دایره را با حداکثر خطای نسبی 10^{-2} مشخص کرد؟

(الف) 0.5% و 0.25% (ب) 0.6% و 0.3% (ج) 0.5% و 0.35% (د) 0.55% و 0.25%

۱۳. می‌خواهیم کمیت‌های U و V را از فرمول‌های $U = x^2 \div y$ و $V = xy^2$ محاسبه کنیم. درصد خطای x و y به

ترتیب در هر یک از گزینه‌های زیر داده شده است. با انتخاب کدام گزینه می‌توان U و V را با حداکثر خطای

نسبی 10^{-2} مشخص کرد؟

الف) $۰/۲\%$ و $۰/۳۵\%$ (ب) $۰/۲۵\%$ و $۰/۳\%$ (ج) $۰/۳\%$ و $۰/۲۵\%$ (د) $۰/۲\%$ و $۰/۳۵\%$

۱۴. کدام یک از گزینه‌های زیر در حساب ممیز شناور درست است؟

الف) عمل جمع دارای خاصیت شرکت پذیری است. (ب) عضو خنثای عمل جمع یکتا است.
ج) عمل جمع خاصیت جابجایی دارد. (د) جمع دو عدد مثبت نزدیک به هم دارای خطای زیادی است.

۱۵. حداکثر خطای نسبی محاسبه حجم یک کره از فرمول $V = \frac{4}{3}\pi r^3$ چقدر است، هرگاه تقریبی از اعداد π و $\frac{1}{3}$ با خطای نسبی حداکثر $۰/۰۰۲$ داشته باشیم و شعاع را بتوان با حداکثر خطای نسبی $۰/۰۰۲$ اندازه‌گیری کرد؟

الف) $۰/۰۰۱$ (ب) $۰/۰۱۰$ (ج) $۰/۰۲۰$ (د) $۰/۰۰۲$

۱۶. کدام یک از گزینه‌های زیر در حساب ممیز شناور درست نیست؟

الف) در محاسبه تقریبی $\frac{\sqrt{2}}{100000}$ قطعا خطای زیادی داریم.
ب) در محاسبه تقریبی $100000\sqrt{2}$ قطعا خطای زیادی داریم.
ج) در تفاضل دو عدد منفی نزدیک به هم، انتشار خطا رخ می‌دهد.
د) در تفاضل دو عدد مثبت نزدیک به هم، انتشار خطا رخ می‌دهد.

۱۷. در محاسبه عبارت $\frac{e^{x^2} - \cos x}{x^2}$ کدام گزینه صحیح است؟

الف) به ازای تمام مقادیر x انتشار خطا رخ نمی‌دهد. (ب) برای مقادیر بزرگ $|x|$ انتشار خطا رخ می‌دهد.
ج) برای مقادیر کوچک $|x|$ انتشار خطا رخ می‌دهد. (د) به ازای همه مقادیر $|x|$ انتشار خطا رخ می‌دهد.

۱۸. در یک دستگاه ممیز شناور که در آن اعداد به صورت $10^n \times d_1 d_2 d_3 d_4 \pm 0$ با $d_1 \neq 0$ و $8 \leq n \leq 7 -$ و $9 \leq d_i \leq 0$ نمایش داده می‌شوند فاصله بین عدد 10000 و اولین عدد قابل نمایش بزرگتر از 10000 کدام است؟

الف) 10 (ب) $۰/۱$ (ج) $۰/۰۱$ (د) $۰/۰۰۵$

۱۹. اگر مقدار $\ln(1+x)$ را با استفاده از دو جمله اول بسط مکلورن آن تقریب بزنیم، مقدار تقریبی و حداکثر خطای محاسبه $\ln(1/1)$ به ترتیب کدام است؟

الف) $(۰/۰۰۱, ۰/۰۰۵)$ (ب) $(۰/۱, ۰/۰۰۵)$ (ج) $(۰/۰۰۵, ۰/۰۰۵)$ (د) $(۰/۱, ۰/۰۰۵)$

۲۰. در محاسبه دوره تناوب آونگ ساده $T = 2\pi\sqrt{\frac{l}{g}}$ ، درباره اثر خطای نسبی مقادیر π, l, g در خطای نسبی T چه می‌توان گفت؟

الف) اثر خطای نسبی π بیشتر است. (ب) اثر خطای نسبی l بیشتر است.
ج) اثر خطای نسبی g بیشتر است. (د) π, l, g اثر یکسانی دارند.

۲۱. در یک دستگاه ممیز شناور که در آن اعداد به صورت $10^n \times d_1 d_2 d_3 d_4 \pm 0$ با $d_1 \neq 0$ و $9 \leq d_i \leq 0$ و $8 \leq n \leq 7 -$ با گرد کردن نمایش داده می‌شوند، بزرگترین عدد x که در معادله $500 + x = 500$ صدق می‌کند، کدام است؟

الف) $۰/۰۴۹۹۹$ (ب) $۰/۰۹۹۹۹$ (ج) $۰/۰۵۰۰۰$ (د) $۰/۱۰۰۰$

۲۲. اگر مقدار دقیق و تقریبی در توانی از ده ضرب شوند، آنگاه

- (الف) خطای نسبی به توانی از ده تقسیم می‌شود.
 (ب) خطای مطلق به توانی از ده تقسیم می‌شود.
 (ج) خطای نسبی تغییری نمی‌کند.
 (د) خطای مطلق تغییری نمی‌کند.

۲۳. گزینه مناسب برای محاسبه $T = \frac{1 - \cos x}{\sin x}$ به ازای مقادیر x نزدیک صفر کدام است؟

- (الف) $\frac{x}{2}$ (ب) $\frac{\pi}{2}$ (ج) $\frac{1}{\sin x} - \cot x$ (د)

۲۴. در محاسبه تقریبی $\ln(1/0.01)$ با سه جمله از بسط تیلور، حداکثر خطا چقدر است؟

- (الف) $\frac{10^{-6}}{3}$ (ب) $\frac{10^{-9}}{3}$ (ج) $\frac{10^{-4}}{3}$ (د) $\frac{10^{-6}}{2}$

۲۵. در محاسبه حجم $(V = \frac{4}{3}\pi r^3)$ یک کره به شعاع r ، درباره اثر خطای نسبی داده‌ها در خطای نسبی V چه می‌توان گفت؟

- (الف) اثر خطای نسبی r بیشتر است.
 (ب) اثر خطای نسبی π بیشتر است.
 (ج) اثر خطای نسبی $\frac{4}{3}$ بیشتر است.
 (د) هر سه اثر خطای نسبی یکسانی دارند.

۲۶. گزینه مناسب برای محاسبه $T = \frac{\cos(2x) - 1 + 2x^2}{x^4}$ به ازای مقادیر x نزدیک صفر کدام است؟

- (الف) $\frac{3}{4}$ (ب) $\frac{2}{3}$ (ج) $\frac{2x}{3}$ (د) $\frac{3x}{4}$

۲۷. برای محاسبه $f(3)$ با خطایی کمتر از 0.0007 چند جمله از بسط تیلور تابع f حول نقطه $x_0 = 2/5$ لازم است (می‌دانیم $|f^{(n)}(x)| \leq \frac{1}{n}$)؟

- (الف) دو جمله (ب) سه جمله (ج) پنج جمله (د) چهار جمله

۲۸. در یک دستگاه ممیز شناور نرمال شده که هر عدد حقیقی به صورت $\pm 0.d_1d_2d_3d_4 \times 10^e$ با $d_1 \neq 0$ و $0 \leq d_i \leq 9$ برای $i = 1, 2, 3, 4$ و $-6 \leq e \leq 7$ ، نمایش داده می‌شود، اعداد قابل نمایش قبل و بعد از 10000 کدام هستند؟

- (الف) $9999, 10001$ (ب) $9990, 10010$ (ج) $9999, 10010$ (د) $9990, 10001$

۲۹. در دستگاه ممیز شناور نرمال شده سوال قبل، تعداد اعداد قابل نمایش در بازه $[1, 10]$ کدام است؟

- (الف) 9001 (ب) 9000 (ج) 8999 (د) 9002