

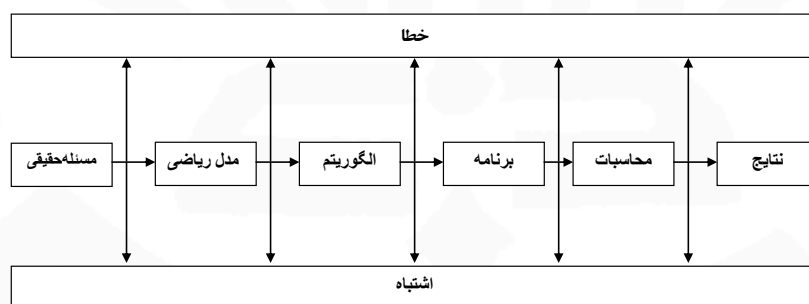
فصل ۲

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 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 یا 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 را نیز نشان می‌دهد.



شکل ۱.۲: فرایند تولید جواب عددی (تقریبی)

تذکر ۱.۲ در مدل‌سازی مسائل باید از بروز خطاهای ذاتی تا حد ممکن جلوگیری کرد در حالی که پرهیز از خطاهای محاسباتی و کنترل آن‌ها، از وظایف متخصص آنالیز عددی است. تفاوت بین جواب واقعی و تقریبی از اشتباهات و خطاها ناشی می‌شود. اشتباه را می‌توان برطرف کرد ولی خطا ممکن است اجتناب‌ناپذیر باشد. به عنوان مثال قرار دادن 2323 به جای 2332 یک اشتباه است حال آن که عدد π بسط داده‌ی نامختوم دارد و در وسایل محاسباتی باید یک بسط داده‌ی مختوم برای آن در نظر گرفت که این موجب بروز خطا می‌شود.

مراحل مختلف این روندنما را در مثال بعد دنبال می‌کنیم.

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 در زمان t باشد. 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

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

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

$$\theta = 6^\circ \simeq 0.1047 \text{ rad} \rightarrow \sin \theta \simeq 0.1045, \quad \theta = 15^\circ \simeq 0.262 \text{ rad} \rightarrow \sin \theta \simeq 0.259,$$

می‌توان فرض کرد $\sin \theta \simeq \theta^{rad}$ و معادله دیفرانسیل را به صورت $\frac{d^2\theta}{dt^2} + \omega^2 \theta = 0$ 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 $T_L = 2\pi \sqrt{\frac{l}{g}}$. 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 یک الگوریتم برای این مدل ریاضی 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_N , T_L و اختلاف آن‌ها

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

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

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

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

(5) قرار دهید $T_N = T_k$ و T_L و $T_N - T_L$ را چاپ کنید.

حال برای الگوریتم ۱.۲ یک برنامه با استفاده از بسته نرم‌افزاری Mathematica می‌نویسیم. پس از اجرای برنامه ۱.۲ در محیط Mathematica محاسبات شروع شده و به ازای $\epsilon = 0.0001$ می‌توان نتایج گزارش شده در جدول ۱.۲ را تولید کرد. باید توجه داشت که برای $\theta_m < 15^\circ$ اختلاف T_N و T_L از 0.005 بیشتر نیست. Δ

با توجه به روندنمای ارائه شده در شکل ۱.۲ و مثال ۱.۲، خطاها از نظر منابع تولید به صورت زیر تقسیم‌بندی می‌شوند.

۱. خطای ذاتی

• خطای مدل (ناشی از صرف نظرها، چشم‌پوشی‌ها و ساده‌سازی‌ها مانند فرض $\sin \theta \simeq \theta^{rad}$)

• خطای داده‌های مدل (ناشی از آزمایشات و اندازه‌گیری‌ها مانند g, l)

۲. خطای محاسباتی

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

```

l = Input["طول آونگ را بر حسب متر وارد کنید"];
theta_m = Input["جایابی زاویه‌ای ماکزیمم را بر حسب درجه وارد کنید"];
epsilon = Input["مقدار مجاز خطا را وارد کنید"];
st = Sin[theta_m * Pi / 180];
g = 9.8; TL = 2 * pi * sqrt(l/g); k = 1; M = ((2^k - 1) / 2^k)^2;
CC = TL * M * st^k; T1 = TL; T2 = T1 + CC;
While[CC > epsilon, {
    k++;
    M = ((2^k - 1) / 2^k)^2;
    T1 = T2;
    CC = TL * M * st^k;
    T2 = T1 + CC;
}];
TN = T2;
Print["TL = ", TL, " TN = ", TN, " TN - TL = ", TN - TL];

```

برنامه ۱.۲: برنامه یافتن دوره تناوب حرکت یک آونگ ساده.

$l(cm)$	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 ۱.۲: دوره تناوب آونگ ساده با طول‌های مختلف

۲.۲ نمایش اعداد

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

قضیه ۱.۲ هر عدد حقیقی مثبت x نمایشی به صورت

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

دارد که در آن $m \in \mathbb{Z}$ و $a_i \in \{0, 1, 2, \dots, \beta - 1\}$. برای منحصر به فرد بودن این 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 $a_{-j} \neq \beta - 1$.

رابطه (۲.۲) به نمایش (بسط) عدد x در مبنای β معروف است. اگر $\beta = 10$ اختیار شود به (۲.۲) نمایش (بسط) ده‌دهی (اعشاری) گویند که در زندگی روزمره با آن سرو کار داریم. در حالتی که $\beta = 2$ در نظر گرفته شود، (۲.۲) به نمایش دودویی (باینری) عدد x معروف بوده و این مبنا، اساس کار رایانه و وسایل دیجیتال است.

تذکر ۲.۲ شرط دوم منحصر به فرد بودن نمایش، لازم است چه در غیر این صورت به عنوان مثال می‌توان نوشت

$$\begin{aligned} 3,47999\dots &= 3,47\bar{9} = 3 \times 10^0 + 4 \times 10^{-1} + 7 \times 10^{-2} + 9 \times 10^{-3} + 9 \times 10^{-4} + \dots \\ &= 3,47 + \frac{9 \times 10^{-3}}{1 - 10^{-1}} = 3,47 + 0,1 = 3,57. \end{aligned}$$

یعنی بسط ده‌دهی بیان شده را می‌توان برای دو عدد $3,47999\dots$ و $3,57$ در نظر گرفت.

قضیه ۲.۲ اگر بسط ده‌دهی یک عدد مختوم یا نامختوم متناوب باشد، آن عدد گویا است.

نتیجه ۱.۲.۲ بسط ده‌دهی یک عدد گنگ، نامختوم نامتناوب است.

مثال ۲.۲ به نمایش اعداد زیر توجه کنید.

$$\begin{aligned} \frac{2}{3} &= 0,666\dots = 0,6\bar{6} = (0,2)_3, & \frac{3}{8} &= 0,375 = (0,011)_2, & \frac{1}{7} &= 0,142857\dots = (0,1)_7, \\ \sqrt{2} &= 1,414213\dots, & \pi &= 3,141592\dots, & \exp(1) = e &= 2,718281\dots \end{aligned}$$

عدد $0,191919\dots$ گویا است ولی اعداد $0,122333444\dots$ و $0,10200300040000\dots$ گنگ هستند. Δ

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 مانند $2, 8$ و 16 کار می‌کنند. بنابراین مسئله تغییر مبنا مطرح می‌شود.

پرسش ۱.۲ مبنای 2 در مقایسه با مبنای 16 چه معایب و مزایایی دارد؟

مثال ۳.۲ در این مثال یک عدد در مبنای ۲ را به مبنای ۱۰ می‌بریم.

$$\begin{aligned}(1011001/111001)_2 &= 2^6 + 2^4 + 2^3 + 2^0 + 2^{-1} + 2^{-2} + 2^{-3} + 2^{-6} \\ &= 64 + 16 + 8 + 1 + 0,5 + 0,25 + 0,125 + 0,015625 = 89,890625\end{aligned}$$

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

مثال ۴.۲ در این مثال یک عدد صحیح مثبت را از مبنای ۱۰ به مبنای ۲ می‌بریم.

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

پس $(10111)_2 = 23$. بنابراین از تقسیم‌های متوالی بر ۲ استفاده می‌کنیم. \triangle

اگر عدد داده‌شده 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 $2x = b_1 + b_2 \times 2^{-1} + b_3 \times 2^{-2} + \dots$ و 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 و سایر رقم‌ها 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 در این مثال این روند را برای حالت‌هایی دنبال می‌کنیم که عدد نمایشی مختوم یا نامختوم متناوب دارد. با توجه به جدول ۲.۲ می‌توان نوشت

$$0,578125 = (0,100101)_2, \quad 0,1 = (0,00011)_2.$$

با توجه به نتایج این مثال، نمایش یک عدد گویا در هر مبنایی یا مختوم است یا نامختوم متناوب [۱]. \triangle

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.

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

۳.۲ نمایش اعداد در رایانه

برای نمایش اعداد در ماشین، ابتدا نمایشی به نام ممیز ثابت^۱ در نظر گرفته شد که در آن هر عدد حقیقی x به صورت زیر نمایش داده می‌شود

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

که در آن n و m 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)^۲ 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. آنچه در ادامه درباره این نمایش خواهد آمد بر اساس استاندارد IEEE^۳ وضع کرده است.

۱.۳.۲ نمایش ۶۴-بیتی ممیز شناور

این نمایش پیش از این به دقت دو برابر (مضاعف)^۴ معروف بوده و متناظر با نوع double در زبان C است. در این نمایش برای ذخیره هر عدد در مبنای ۲، ابتدا یک ساختار به طول ۶۴ بیت در نظر گرفته می‌شود. اولین بیت به بیت علامت معروف است و با s نمایش داده می‌شود و بلافاصله بعد از آن ۱۱ بیت برای مشخصه^۵ در نظر گرفته می‌شود و با c نمایش داده می‌شود و ۵۲ بیت باقی‌مانده نیز برای مانتیس کنار گذاشته می‌شود و آن را با f نشان می‌دهند. سپس برای هر عدد نمایشی به صورت $(1 + 0/f)2^e \times (-1)^s$ در نظر گرفته می‌شود که در آن $e = c - 1023$ می‌توان (نما) است. محدودیت $2047 - 1 = 2^{11} - 1 = 2047$ است. $0 \leq c \leq (10001)_2 = 2047$ موجب می‌شود که $-1023 \leq e \leq 1024$.

^۱ Fixed point

^۲ Floating point

^۳ IEEE standard 754-1985

^۴ Double precision

^۵ Characteristic

طرف دیگر محدودیت $(\underbrace{0 \dots 0}_{{52 \text{ times}}})_2 \leq f \leq (\underbrace{0 \dots 0}_{{52 \text{ times}}})_2 < 2$ باعث می‌گردد که $0 \leq f \leq (\underbrace{0 \dots 0}_{{52 \text{ times}}})_2 < 2$. بنابراین اگر کوچک‌ترین و بزرگ‌ترین عدد مثبت قابل نمایش را به ترتیب با $+mN$ و $+MN$ نشان دهیم آن‌گاه

$$+mN = 2^{-1022} \times 1/0 \simeq 2/225074 \times 10^{-208}, \quad +MN = 2^{1024} \times 1/0 \simeq 1/797693 \times 10^{208}.$$

همچنین برای نمایش صفر و بی‌نهایت می‌توان از قراردادهای زیر استفاده کرد

$$+0 = (\underbrace{1 \dots 1}_{{52 \text{ times}}})_2 \times 2^{-1023}, \quad +\infty = (\underbrace{1 \dots 0}_{{51 \text{ times}}})_2 \times 2^{1024}.$$

تمرین ۲.۲ نمایش $-mN$ ، $-MN$ ، -0 و $-\infty$ را بنویسید.

تعریف ۱.۲ در نمایش اعداد ماشینی، کوچک‌ترین عدد مثبت ماشینی که اگر به ۱ اضافه شود عددی بزرگ‌تر از ۱ به دست می‌آید به اپسیلون ماشین^۶ معروف است و با eps نمایش داده می‌شود. نصف eps به رند واحد^۷ یا واحد گرد کردن^۸ یا دقت ماشین^۹ معروف است.

چون در نمایش ۶۴-بیتی بعد از ۱ عدد $(\underbrace{1 \dots 0}_{{51 \text{ times}}})_2$ قرار می‌گیرد، پس

$$\text{eps} = (\underbrace{1 \dots 0}_{{51 \text{ times}}})_2 - 1 = (\underbrace{0 \dots 0}_{{51 \text{ times}}})_2 = 2^{-52} \simeq 2/220446 \times 10^{-16}.$$

در ادامه بزرگ‌ترین عدد صحیح مثبت M را تعیین می‌کنیم که هر عدد صحیح x با شرط $0 < x \leq M$ را بتوان به طور دقیق نمایش داد. با توجه به موارد بیان‌شده، تمام اعداد صحیح نامنفی که بزرگ‌تر از $2^{52} \times 2 = 2^{53} - 1 = (\underbrace{1 \dots 1}_{{52 \text{ times}}})_2$ نباشند به طور دقیق قابل نمایش هستند و به علاوه 2^{52} نیز به صورت $(1/0)_2 \times 2^{52}$ قابل نمایش است. اما تعداد ارقام در مانتیس جهت نمایش $2^{52} + 1$ کافی نیست (۵۳ رقم در مانتیس لازم است). بنابراین $10^{15} \times 0.7199 \times 10^7 \simeq 9.0 \times 10^{22} = M$ و در نتیجه تمام اعداد صحیح ۱۵ رقمی و بسیاری از اعداد ۱۶ رقمی در این نمایش به طور دقیق قابل نمایش هستند به بیان دیگر نمایش ۶۴-بیتی ۱۵ الی ۱۶ رقم دقت دارد. بنابراین اعدادی که بیشتر از ۱۶ رقم داشته باشند به طور دقیق قابل نمایش نیستند.

تذکر ۳.۲ مجموعه اعداد با ممیز شناور را با \mathbb{F} 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 فاصله $(-0, +0)$ تولید شود پیام پاریز^{۱۰} و اگر عدد تولیدشده از $+MN$ بزرگ‌تر یا از $-MN$ کوچک‌تر باشد پیام سرریز^{۱۱} صادر می‌شود. برنامه `uoflow.nb` را ببینید.

تذکر ۴.۲ بعضی از نرم‌افزارها مانند Mathematica، محدودیت‌های سخت‌افزار را از طریق برنامه‌های نرم‌افزاری برطرف می‌کنند و به اصطلاح دقت را بالا می‌برند که ممکن است تا حدودی سرعت محاسبات کاهش یابد.

Machine epsilon^۶Unit round^۷Roundoff unit^۸Machine precision^۹Underflow^{۱۰}Overflow^{۱۱}

تمرین ۳.۲ در نمایش ۳۲-بیتی (دقت معمولی^{۱۲} متناظر با نوع float در زبان C) هشت بیت برای مشخصه منظور می‌شود. تمام کمیت‌های معرفی‌شده در این بخش را برای این نمایش تعیین کنید.

۲.۳.۲ اعداد ماشینی

به اعداد \mathbb{F} اعداد ماشینی نیز گفته می‌شود و برای سادگی، نمایشی به صورت $\pm a \times 10^b$ برای آن‌ها در نظر گرفته می‌شود. به عبارتی یک عدد ماشینی با زوج مرتب (a, b) نظیر می‌شود. در اینجا b توان (نما) است و به زیرمجموعه‌ای از \mathbb{Z} تعلق دارد. به بیان دقیق‌تر $L \leq b \leq U$ که در آن L و U با توجه به محدودیت‌های سخت‌افزاری و نرم‌افزاری تعیین می‌شوند. نمایش فرض می‌شود $a = \circ/d_1 \dots d_k$ بیانگر مانتیس (جزء کسری) است و به ازای $i = 1, \dots, k$ داریم $d_i \in \{0, 1, \dots, 9\}$ و برای یکتایی نمایش $d_1 \neq 0$. بنابراین $0 < a < 1$. این نمایش، به نمایش ممیز شناور ده‌دهی نرمال‌شده نیز معروف است و چنین اعدادی را، اعداد ماشینی ده‌دهی k -رقمی می‌نامند. ایده چنین نمایشی، از نمایش علمی نرمال‌شده اعداد ناشی شده است. در نمایش علمی نرمال‌شده، هر عدد حقیقی مخالف صفر x را می‌توان به صورت $x = \pm a \times 10^b$ نمایش داد که در آن $b \in \mathbb{Z}$ و $a = \circ/d_1 \dots d_k d_{k+1} \dots$ و $d_i \in \{0, 1, \dots, 9\}$ داریم $i = 1, 2, \dots$ که در آن به ازای k رقم از مانتیس آن را حفظ کرده و رقم‌های اضافی را به کمک یکی از روش‌های زیر کنار گذاشت

روش قطع کردن (برش)^{۱۳}، روش گرد کردن معمولی^{۱۴}، روش گرد کردن به زوج^{۱۵}.

در روش قطع کردن، k رقم از مانتیس 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

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

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

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

قضیه ۳.۲ اگر $x = \circ/d_1 \dots d_k d_{k+1} \dots \times 10^n$ یک عدد حقیقی ناصفر در نمایش علمی نرمال‌شده باشد آن‌گاه $c(x) = \circ/d_1 \dots d_k \times 10^n$ و $r(x) = \circ/d_1 \dots d_k \times 10^m$ ($m = n$ یا $m = n + 1$) اعداد ماشینی k -رقمی متناظر با x هستند که به ترتیب از روش قطع کردن و گرد کردن معمولی به دست می‌آیند و خواهیم داشت

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

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

^{۱۲} Single precision

^{۱۳} Chopping

^{۱۴} Rounding

^{۱۵} Rounding to even

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 ۳.۲: مثال‌هایی از گرد کردن و قطع کردن معمولی

تعریف ۲.۲ 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 مانع در نمایش علمی نرمال شده).

مثال ۶.۲ عدد $۰,۰۰۰۷۰۴۵۰۰۰$ حداقل ۴ رقم بامعنا و حداکثر ۷ رقم بامعنا دارد. \triangle

قرارداد ۱.۲ یعنی rD رقم اعشار و rS یعنی r رقم بامعنا^{۱۶}.

مثال ۷.۲ در جدول ۳.۲، چند عدد و مقدارهای 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 ۵.۲ از این به بعد، همانند بسیاری از وسایل محاسباتی از گرد کردن معمولی استفاده می‌کنیم.

پرسش ۳.۲ در نظر بگیرید \hat{x} تقریبی از x باشد. آیا هر چه تعداد ارقام بامعنای \hat{x} بیشتر باشد، \hat{x} تقریب بهتری است؟

تعریف ۳.۲ فرض کنید $\hat{x} = d_m \times 10^m + \dots + d_1 \times 10^1 + d_0 \times 10^0 + d_{-1} \times 10^{-1} + \dots + d_{-k} \times 10^{-k}$ که در آن $d_m \neq 0$ ، تقریبی برای $x > 0$ باشد (در اینجا m نشان‌دهنده با ارزش‌ترین مکان است). تعداد ارقام بامعنای درست \hat{x} نسبت به x ، بزرگ‌ترین عدد طبیعی n است که در نابرابری‌های زیر صدق کند.

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

به وضوح هر چه تعداد ارقام بامعنای درست \hat{x} نسبت به x بیشتر باشد، \hat{x} تقریب بهتری خواهد بود.

مثال ۸.۲ از اعداد $\hat{x} = ۹۹,۹۶$ و $\tilde{x} = ۱۰۰,۷$ کدام یک تقریب بهتری برای $x = ۱۰۰$ است؟

$$\hat{x} = ۹۹,۹۶ \rightarrow m = 1, \quad |x - \hat{x}| = ۰,۰۴ \leq 5 \times 10^{1-n} \rightarrow n = 3$$

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

بنابراین \hat{x} تقریب بهتری برای x است زیرا تعداد ارقام بامعنای درست بیشتری دارد. \triangle

مثال ۹.۲ تعداد ارقام بامعنای درست $\hat{x} = ۱۰۰,۳۱$ را نسبت به $x = ۱۰۰,۳۱۰۴$ مشخص کنید.

$$\hat{x} = ۱۰۰,۳۱ \rightarrow m = 2, \quad |x - \hat{x}| = ۰,۰۰۰۴ \leq 5 \times 10^{2-n} \rightarrow n = 6$$

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

^{۱۶} D حرف اول کلمه Decimal و S حرف اول کلمه Significant

۴.۲ انواع خطا

تعریف ۴.۲ اگر \hat{x} تقریبی از x باشد $\Delta x = |x - \hat{x}|$ خطای مطلق \hat{x} نسبت به x نامیده می‌شود. Δx یکتا بوده و در عمل بیشتر مواقع قابل تعیین نیست و به جای آن از هر عدد b_x استفاده می‌شود که از Δx کمتر نباشد. b_x یکتا نیست و به آن کران خطای مطلق گویند. بنابراین $\Delta x \leq b_x$ و در نتیجه $\hat{x} - b_x \leq x \leq \hat{x} + b_x$ و یا می‌نویسیم $x = \hat{x} \pm b_x$.

مثال ۱۰.۲ برای $\hat{x} = ۱/۷۳۲$ به عنوان تقریبی از $x = \sqrt{۳}$ می‌توان نوشت

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

از طرفی می‌دانیم $۱/۷۳۲۰ < \sqrt{۳} < ۱/۷۳۲۱$ بنابراین $۱/۷۳۲۰ < \sqrt{۳} - ۱/۷۳۲ < ۰/۰۰۰۰۱$ پس $b_x = ۰/۰۰۰۰۱$ که معیاری برای نزدیکی $۱/۷۳۲$ به $\sqrt{۳}$ است. \triangle

پرسش ۴.۲ 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 ۵.۲ اگر $\hat{x} \neq ۰$ تقریبی از $x \neq ۰$ باشد $\delta x = \frac{|x - \hat{x}|}{|x|} = \frac{\Delta x}{|x|}$ خطای نسبی \hat{x} نسبت به x نامیده می‌شود و یکتا بوده و بیشتر مواقع قابل تعیین نیست و از کران خطای نسبی استفاده می‌شود. $۱۰۰ \times \delta x$ به درصد خطا معروف است.

قضیه ۴.۲ 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 می‌توان نوشت $\delta x \leq \frac{b_x}{|\hat{x}|}$.

□

برهان. بنابر تعریف b_x و با استفاده از خواص نابرابری‌ها واضح است.

مثال ۱۱.۲ برای $\hat{x} = ۱/۷۳۲$ به عنوان تقریبی از $x = \sqrt{۳}$ می‌توان نوشت

$$\delta x = \frac{\left| \sqrt{۳} - ۱/۷۳۲ \right|}{\sqrt{۳}} = \frac{۱/۷۳۲۰۵۰۸۰۷۵\dots - ۱/۷۳۲}{۱/۷۳۲۰۵۰۸۰۷۵\dots} = \frac{۰/۰۰۰۰۵۰۸۰۷۵\dots}{۱/۷۳۲۰۵۰۸۰۷۵\dots}$$

پس $\delta x = ۰/۰۰۰۰۰۲۹۳۳۳۷\dots < ۰/۰۰۰۰۰۳$ و با توجه به قضیه ۴.۲ و $\Delta x = ۰/۰۰۰۰۱$ خواهیم داشت

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

و یا

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

△

قضیه ۵.۲ 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 بوده و خطای نسبی \hat{y} و \hat{x} برابر است.

قضیه ۶.۲ اگر \hat{x} گردشده x تا n رقم بامعنا باشد آن‌گاه \hat{x} دارای n رقم بامعنای درست است.

ارتباط خطای نسبی با تعداد ارقام بامعنای درست در قضیه بعدی مشخص می‌شود.

قضیه ۷.۲ 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} یک رقم ۱ و $n-1$ رقم صفر جلوی آن نباشد. اگر $\delta x \leq 0.5 \times 10^{-n}$ آن‌گاه \hat{x} حداقل n رقم بامعنای درست دارد.

مثال ۱۲.۲ تقریبی از $\sqrt{3} = 1.73205000\dots$ ارائه دهید که خطای نسبی آن از 10^{-4} کمتر باشد.

بنابر قضیه ۷.۲ اگر \hat{x} تقریبی از $\sqrt{3}$ باشد که ۵ رقم بامعنای درست داشته باشد آن‌گاه $10^{-4} < 5 \times 10^{-5} < \delta x$. از

این‌رو، با توجه به قضیه ۶.۲ کافی است \hat{x} گردشده $\sqrt{3}$ تا ۵ رقم بامعنا باشد، یعنی $\hat{x} = 1.7321$ \triangle

۵.۲ خطای محاسبات

فرض کنید $z = f(x_1, \dots, x_n)$ تابعی باشد که می‌خواهیم آن را در نقطه (x_1, \dots, x_n) ارزیابی کنیم و $x_i - \hat{x}_i = \Delta x_i$ 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 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

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

تذکر ۶.۲ اگر می‌خواهیم جوابی از یک مسئله را با دقت rD به دست آوریم، باید محاسبات میانی را با دقت $(r+1)D$ و حتی بیشتر انجام دهیم و نتیجه نهایی را با دقت rD گرد کنیم. اگر محاسبات میانی با دقت rD انجام شود، به دلیل اثرات خطای گرد کردن اعتباری به رقم (ارقام) سمت راست جواب نیست. به طور مشابه، هنگام کار با دقت rS باید ملاحظاتی را در نظر گرفت.

مثال ۱۳.۲ (مستقیم) یک استوانه به شعاع قاعده $\frac{4}{3}$ 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 می‌آید؟ می‌دانیم حجم یک استوانه از قاعده $\pi r^2 h$ تعیین می‌شود که در آن شعاع قاعده h و ارتفاع استوانه است. اگر تعریف کنیم $z = V(p, r, h) = \pi r^2 h$ به ازای $p = \pi$ ، $r = \frac{4}{3}$ و $h = \sqrt{2}$ ارزیابی شود. محاسبات میانی را با دقت $4D$ دنبال می‌کنیم (ماشین حساب را به صورت MODE FIX 4 تنظیم می‌کنیم). بنابراین

$$p = \pi \rightarrow \hat{p} = 3,1416, \quad r = \frac{4}{3} \rightarrow \hat{r} = 1,3333, \quad h = \sqrt{2} \rightarrow \hat{h} = 1,4142.$$

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

$$\hat{z} = V(\hat{p}, \hat{r}, \hat{h}) = 3,1416 \times 1,3333^2 \times 1,4142 = 7,8980.$$

از طرفی

$$\Delta z \simeq \Delta p \hat{r}^2 \hat{h} + 2 \Delta r \hat{p} \hat{r} \hat{h} + \Delta h \hat{p} \hat{r}^2 \leq (\hat{r}^2 \hat{h} + 2 \hat{p} \hat{r} \hat{h} + \hat{p} \hat{r}^2) \times 0,5 \times 10^{-4} \simeq 9,9731 \times 10^{-4}.$$

در نتیجه حجم استوانه با دقت $3D$ برابر است با $7,898$ و حداکثر خطای مرکب شده عبارت است از

$$\Delta z + 0,5 \times 10^{-4} \lesssim 0,0010.$$

به کمک یک ماشین حساب 10 رقمی نتیجه $z = 7,898458555$ به دست می‌آید و بنابراین

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

△

تمرین ۴.۲ (معکوس) یک استوانه به شعاع قاعده $\frac{4}{3}$ 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؟

مثال ۱۴.۲ (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 $y = ۰/۳$ و (x, y) را با دقت $۱S$ یا $۲S$ 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 این نامعادله 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، می‌توان x و y را با دقت $۴D$ در نظر گرفت. \triangle

۱.۵.۲ خطای اعمال ریاضی

هنگام کار با اعداد ماشینی (ممیز شناور) باید توجه داشت که بعضی از اصول میدان \mathbb{R} نظیر شرکت‌پذیری، یکتایی 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 ۶.۲ فرض کنید x و y دو عدد حقیقی، \hat{x} و \hat{y} تقریب‌هایی از آن‌ها و \otimes بیانگر یک عمل دوتایی باشد. متناظر با $x \otimes y$ در ماشین عمل $\hat{x} \otimes \hat{y}$ انجام می‌شود و می‌توان نوشت

$$|x \otimes y - \hat{x} \otimes \hat{y}| = |(x \otimes y - \hat{x} \otimes \hat{y}) + (\hat{x} \otimes \hat{y} - \hat{x} \otimes \hat{y})| \leq |x \otimes y - \hat{x} \otimes \hat{y}| + |\hat{x} \otimes \hat{y} - \hat{x} \otimes \hat{y}|.$$

بنابراین خطا دو بخش دارد: خطای منتشرشده یا انباشته‌شده (جمله اول) و خطای تولیدشده یا گرد کردن (جمله دوم).

قضیه ۸.۲ اگر \hat{x} و \hat{y} تقریب‌هایی از x و y بوده و همه این اعداد مثبت باشند، آنگاه

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

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

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

□

برهان. به عنوان تمرین.

مثال ۱۵.۲ \hat{x} و \hat{y} هر یک n رقم بامعنای درست دارند. حداقل تعداد ارقام بامعنای درست $\hat{x} + \hat{y}$ و $\hat{x}\hat{y}$ را تعیین کنید. بنابر قضیه ۷.۲، $\delta x < ۵ \times ۱۰^{-n}$ و $\delta y < ۵ \times ۱۰^{-n}$. هم‌چنین بنابر قضیه ۸.۲ می‌توان نوشت

$$\delta(x + y) \leq \max\{\delta x, \delta y\} \leq ۵ \times ۱۰^{-n} = ۰/۵ \times ۱۰^{-(n-1)}.$$

پس $\hat{x} + \hat{y}$ حداقل $n - ۱$ رقم بامعنای درست دارد. به علاوه می‌توان نوشت

$$\delta(xy) \leq \delta x + \delta y < ۵ \times ۱۰^{-n} + ۵ \times ۱۰^{-n} = ۱۰ \times ۱۰^{-n} = ۱۰^{-(n-2)-1} < ۰/۵ \times ۱۰^{-(n-2)}.$$

△

بنابراین \widehat{xy} دست کم $n - 2$ رقم بامعنای درست دارد.

تذکر ۷.۲ با توجه به خطای نسبی $\delta(x-y) \simeq \frac{\Delta(x-y)}{|x-y|}$ واضح است که اگر x و y دو عدد هم‌علامت و به هم نزدیک باشند $|x-y|$ کوچک و در نتیجه $\delta(x-y)$ بزرگ خواهد شد و در نتیجه تعداد ارقام بامعنای درست $x-y$ کم خواهد بود. بنابراین در عمل بهتر است از تفاضل دو عدد هم‌علامت نزدیک به هم جلوگیری شود (به اصطلاح تفاضل دو عدد نزدیک به هم موجب از بین رفتن ارقام بامعنا می‌شود مانند $1/41 - 1/42 = 0/01$). اگر تفاضل اجتناب‌پذیر نیست باید دقت را بالا برد (دقت دو برابر یا بیشتر). البته به کمک قضیه ۸.۲، می‌توان پدیده از بین رفتن ارقام بامعنا را بهتر تفسیر کرد.

مثال ۱۶.۲ جهت اجتناب از تفاضل در محاسبات می‌توان از اتحادها کمک گرفت. به عنوان مثال

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

△

مثال ۱۷.۲ در صورتی که $g(x) = x \left(\sqrt{1 + \frac{1}{x}} - 1 \right)$ مقدار $g(10^9)$ را با دقت $8D$ به دست آورید. همان‌گونه که بیان شد، چون دقت $8D$ مورد نظر است، محاسبات میانی را با دقت $9D$ انجام می‌دهیم

$$1 + \frac{1}{x} = 1/0000000001, \quad \sqrt{1 + \frac{1}{x}} = 1/0000000000.$$

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

$$g(10^9) = 10^9(1/0000000000 - 1) = 0/0000000000.$$

به این ترتیب با یک ماشین حساب 10 رقمی $g(10^9) = 0/0000000000$ که نتیجه‌ای نادرست است. در حالی که اگر از یک رایانه استفاده شود، نتیجه درست $g(10^9) = 0/2222222222$ است. 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

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

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 $9D$ نتیجه زیر به دست می‌آید

$$g(10^9) = \frac{1}{1/0000000000 + 1/0000000000 + 1} = 0/2222222222,$$

و بنابراین با دقت $8D$ داریم $g(10^9) = 0,333333333$.

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 قائل شویم.

تذکر ۸.۲ با توجه به قضیه ۸.۲، در محاسبات باید از ضرب اعداد بزرگ در اعداد تقریبی (تقسیم اعداد تقریبی به اعداد کوچک) پرهیز کرد.

مثال ۱۸.۲ برای به دست آوردن تقریبی از $x = 10000 \times \pi$ می‌توان نوشت

$$\pi \simeq 3/14 \rightarrow x \simeq 31400, \quad \pi \simeq 3/142 \rightarrow x \simeq 31420, \quad \pi \simeq 3/1416 \rightarrow x \simeq 31416.$$

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

تذکر ۹.۲ 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 ۱۹.۲ به جای عبارت $ax^3 + bx^2 + cx + d$ از عبارت $((ax + b)x + c)x + d$ استفاده شود.

۲.۵.۲ تقریب توابع یک متغیره

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

روش اول- به کمک قضیه تیلور داریم

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

که در آن $0 < \xi(x) < x$ در نتیجه

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

که در آن $0 < \xi < \frac{\pi}{10} < 1$ چون e^x تابعی صعودی است پس $1 < e^1 < e^\xi < e^0 = 1$. از طرفی $0,315 < \frac{\pi}{10}$ (دقت $3D$ منظور می‌شود) و بنابراین

$$\text{خطا} \leq \left| \frac{(\frac{\pi}{10})^{n+1}}{(n+1)!} e^{\xi} \right| < \frac{3 \times (0,315)^{n+1}}{(n+1)!}.$$

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

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

با توجه به $\frac{\pi}{10} < 0.315$ (دقت $3D$ منظور می‌شود) از $0.5 \times 10^{-2} < \left| \frac{(0.315)^n}{n!} \right|$ نتیجه می‌شود $n \geq 4$. پس

$$e^{\frac{\pi}{10}} \simeq 1 + 0.315 + \frac{(0.315)^2}{2!} + \frac{(0.315)^3}{3!} + \frac{(0.315)^4}{4!} = 1.315 + 0.050 + 0.005 + 0.000 = 1.370.$$

△

مثال ۲۱.۲ (همگرایی سریع) می‌خواهیم تابع $\cos x$ را به ازای مقادیر $|x| < \frac{\pi}{4}$ با دقت $5D$ ارزیابی کنیم. با توجه به سری تیلور

$$\cos x = \sum_{i=0}^{\infty} (-1)^i \frac{x^{2i}}{(2i)!} = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \dots,$$

و این که $\left| \frac{x^{2n}}{(2n)!} \right| < \frac{(\frac{\pi}{4})^{2n}}{(2n)!} < \frac{1}{6^{2n}} < 0.5 \times 10^{-5}$ باید داشته باشیم که نتیجه می‌دهد $n \geq 6$. △

مثال ۲۲.۲ (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 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 $5D$ ، 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 ۲۳.۲ تابع $x \geq 0$ $Si(x) = \int_0^x \frac{\sin t}{t} dt$ را 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 می توان 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 $10^{-6} \times 0.5 < \frac{1}{(2n+2) \times (2n+2)!}$ نتیجه می شود $n \geq 4$ 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 $6D$ داریم $Si(1) \approx 0.946083$.

۶.۲ تمرین ها

۱. حاصل $(\sqrt[n]{z})_\beta$ که در آن $z = \beta - 1$ را به دست آورید.۲. فاصله بین اعداد در دستگاه ممیز ثابت $(a_1 a_2 a_3 a_4 / b_1 b_2 b_3 b_4)_r$ را به دست آورید.۳. آیا $\frac{1}{\beta+1}$ می تواند ماتیس یک عدد نرمال شده در مبنای β باشد؟ چرا؟

۴. بزرگ‌ترین عدد p را به دست آورید که نمایش آن در دستگاه نمایش اعداد نرمال شده در مبنای β با m رقم ماتنیس و $l \leq e \leq u$ برابر خود آن عدد باشد.

۵. فرض کنید x_+ و x_- به ترتیب اعداد ماشینی بلافاصله بعد و قبل از عدد غیرماشینی x باشند. کمترین مقدار $x_+ - x_-$ در دستگاه اعداد نرمال $\pm (\circ/b_1 \cdots b_m)_\beta \times \beta^e$ که در آن $l \leq e \leq u$ را محاسبه کنید.

۶. اگر نمایش دو عدد x و $x+y$ در دستگاه نمایش اعداد نرمال در مبنای β با m رقم ماتنیس و $l \leq e \leq u$ یکسان باشد، آن‌گاه چه رابطه‌ای بین x و y برقرار است.

۷. در یک دستگاه ممیز شناور نرمال برای اعداد حقیقی در مبنای ۲ با t رقم ماتنیس و روش گرد کردن را در نظر بگیرید. اگر x و y هر دو در این دستگاه قابل نمایش باشند، چه شرطی برقرار باشد تا نمایش $x + 2^{-t}y$ همان نمایش x باشد؟

۸. نشان دهید تساوی زیر برقرار است

$$(\circ/b_1 b_2 \cdots b_m \overline{c_1 c_2 \cdots c_k})_\beta = \frac{(b_1 b_2 \cdots b_m c_1 c_2 \cdots c_k)_\beta - (b_1 b_2 \cdots b_m)_\beta}{\underbrace{(z z \cdots z)_\beta}_k \underbrace{(z \circ \cdots \circ)_\beta}_m}, \quad z = \beta - 1$$

۹. مقدار محاسبه شده برای

$$\frac{f(a + 2^{-p}h) - f(a - 2^{-p}h)}{h 2^{1-p}}$$

و $p \geq 1$ در یک ماشین با روند عدد یک برابر 2^{1-t} به ازای چه مقادیری از $|h|$ برابر \circ است؟

۱۰. درستی عبارات $\frac{|r(x) - x|}{|x|} < \frac{1}{4} \beta^{1-k}$ و $\frac{|c(x) - x|}{|x|} < \beta^{1-k}$ را بررسی کنید.

۱۱. برای محاسبه مقدار $1 - \cos(x)$ در نزدیکی \circ از چه عبارتی استفاده می‌کنید؟ توضیح دهید.

$$\frac{\sin^2(x)}{1 + \cos(x)} \quad (۴) \quad x^3 - \frac{x^2}{3} \quad (۳) \quad \frac{x}{2} - \frac{x^2}{3} \quad (۲) \quad \circ \quad (۱)$$

۱۲. برای محاسبه $\sum_{k=1}^{2n} (-1)^k x^k$ با یک وسیله محاسباتی، زمانی که $\circ < x < 1$ و n خیلی بزرگ باشد، از چه عبارتی استفاده کنیم مناسب تر است؟

$$\frac{1}{(\sqrt{2} + 1)^6} \quad (۱) \quad 99 - 7\circ\sqrt{2} \quad (۲) \quad \frac{1}{99 + 7\circ\sqrt{2}} \quad (۳) \quad (\sqrt{2} - 1)^6 \quad (۴)$$

۱۴. برای تقریب تابع $\cos(x)$ به ازای $|x| \leq \circ/1$ با استفاده از بسط تیلور حول نقطه \circ به منظور دستیابی به خطای برشی حدود 10^{-5} به چه تعداد از جملات سری تیلور نیاز داریم؟

۱۵. \hat{x} و \hat{y} به ترتیب تقریبی از x و y هستند. اگر این اعداد مثبت باشند، در چه صورت خطای مطلق \widehat{xy} به عنوان تقریبی از xy در حد خطای مطلق $\widehat{x} + \widehat{y}$ است؟

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