

به نام خداوند بخشندهی مهربان

آنالیز عددی

حسین زارع



خطاها و دستگاه‌های نمایش اعداد

۱-۱ منابع خطا

در طی روند حل عددی یک مسئله، عوامل گوناگونی دست به دست هم می‌دهند تا جواب حاصل با جواب دقیق برابر نباشد. به‌طور کلی خطاهای موجود در جواب تقریبی یک مسئله از چهار منبع ناشی می‌شوند^۱. این منابع عبارتند از:

• مدل‌سازی مسئله

اولین مرحله از حل هر مسئله‌ی واقعی، مدل‌سازی آن مسئله است. از آن‌جا که دخالت دادن تمام عوامل مؤثر در مسئله، احتمالاً به یک مدل پیچیده و گاه غیر مفید منجر می‌شود، برخی از عوامل کم‌اهمیت‌تر نادیده گرفته می‌شوند. برای مثال در مدل‌سازی مسئله‌ی «تعیین زمان تناوب حرکت یک آونگ ساده» در مکانیک، از جرم نخ متصل به آونگ، اصطکاک بین نخ و نقطه‌ی آویز و نیز مقاومت هوا صرف‌نظر می‌شود. بدیهی است که این عوامل، تأثیراتی هر چند کوچک در جواب مسئله خواهند گذاشت. بنابراین با در نظر نگرفتن آن‌ها خطایی در جواب مسئله به‌وجود می‌آید. این خطا را خطای مدل می‌نامیم.

• داده‌های ورودی

پارامترهای معلوم یک مدل اغلب کمیت‌هایی هستند که از طریق اندازه‌گیری به‌دست می‌آیند و چون وسایل اندازه‌گیری صد در صد دقیق نیستند، اندازه‌گیری‌ها با خطا همراه هستند. البته هرچه وسایل اندازه‌گیری دقیق‌تر باشند، این خطا کمتر می‌شود.

^۱ برخی کتاب‌ها منابع خطا را به سه دسته و برخی دیگر به پنج دسته تقسیم می‌کنند. در این دسته‌بندی که منسوب به جان فون نویمان* و هرمان گلدشتین* است [کتاب هوسهولدر]، منابع خطا به چهار دسته تقسیم شده است.

• روند کردن اعداد


انجام محاسبات تقریبی بر روی اعدادی که دارای نامتناهی رقم هستند، جز با انتخاب تعدادی متناهی از ارقام آن‌ها ممکن نیست. علاوه بر این، همواره در وسایل محاسباتی به دلیل محدودیت حافظه، تعداد محدودی از ارقام یک عدد قابل ذخیره‌سازی و نمایش است. از آن‌جا که بسیاری از اعداد حقیقی مانند $1/3$ ، $\sqrt{2}$ ، e ، π و ... دارای نامتناهی رقم در بسط اعشاری یا دودویی خود هستند، می‌توان گفت که بیشتر اعداد به صورت تقریبی و در نتیجه توأم با خطا در نظر گرفته می‌شوند. به خطای ناشی شده در این موارد خطای نمایش اعداد می‌گوییم.

از طرف دیگر، گاهی در حین محاسبات به دلیل انجام یک عمل حسابی بر دو عدد حقیقی، حاصل عمل دارای نامتناهی رقم است که با توجه به آنچه که در بند قبل گفته شد، تنها تعدادی متناهی از ارقام آن قابل انتخاب است. خطای حاصل در این حالت را خطای تولید شده توسط عمل حسابی می‌نامیم. در یک ماشین محاسباتی، حتی نیازی نیست که حاصل یک عمل دارای نامتناهی رقم باشد، زیرا حتی با وجود متناهی بودن تعداد ارقام حاصل از یک عمل، اگر این تعداد رقم بیش از حافظه‌ی تخصیص داده شده برای نگهداری هر عدد باشد، نتیجه‌ی عمل با خطا در ماشین نمایش داده می‌شود.

در حالت کلی به خطای نمایش اعداد و خطای تولید شده توسط اعمال حسابی، خطاهای روند کردن می‌گوییم. خطاهای روند کردن آثار خود را در حین فرآیند محاسبات برجای می‌گذارند. زیرا اگر مقادیری که در یک محاسبه به کار می‌روند خود با خطایی همراه باشند، انجام عملیات حسابی بر آن‌ها موجب ایجاد خطاهای جدید می‌شود که خطای حاصل در این حالت را خطای منتشر شده توسط عمل حسابی می‌نامیم. انتشار خطا توسط توابع نیز صورت می‌گیرد که در ادامه‌ی فصل به‌طور کامل در مورد آن بحث خواهد شد.

• روش عددی

چنانچه حل تحلیلی یک مسئله‌ی ریاضی امکان‌پذیر یا از لحاظ محاسباتی مقرون به‌صرفه نباشد، از یک روش عددی برای تعیین جواب تقریبی مسئله استفاده می‌شود. برای این منظور ممکن است روش‌های عددی گوناگونی وجود داشته باشند. دقت تقریب حاصل از هر روش عددی به ویژگی‌های آن روش (از جمله پایداری و سرعت همگرایی) و مرحله‌ی توقف آن بستگی دارد. معمولاً در مباحث عددی، خطای روش نتیجه‌ی تقریب یک مسئله‌ی اولیه با یک مسئله‌ی ساده‌تر توسط گسسته‌سازی یا مختوم نمودن یک فرآیند نامتناهی است.

 **تذکر:** خطاها اغلب جزء اجتناب‌ناپذیری از حل عددی یک مسئله هستند و نباید آن‌ها را با اشتباهاتی که ممکن است در روند حل هر مسئله‌ای رخ دهند، یکی گرفت. از این رو همواره میان دو مفهوم خطا (*error*) و اشتباه (*mistake*) تمایز قائل می‌شویم و اشتباهات را جزء منابع خطا در نظر نمی‌گیریم.

خطاهای مدلسازی و داده‌های ورودی در آنالیز عددی مورد بررسی قرار نمی‌گیرند و افرادی که در علوم مختلف به تعیین مدل مسئله و اندازه‌گیری داده‌ها می‌پردازند، مسئول آن‌ها هستند. اما خطاهای ناشی از روند کردن اعداد و روش عددی به آنالیز عددی مربوط می‌شوند.

۱-۲ انواع خطا

بخش عمده‌ای از علم آنالیز عددی به یافتن تقریب مناسبی از جواب یک مسئله مربوط می‌شود. در این جا عبارت «تقریب مناسبی» باید به نحوی مطلوب تعریف شود. نخست تعریفی از تقریب یک عدد ارائه می‌دهیم.

تعریف ۱-۱:

فرض کنیم A و a دو عدد حقیقی باشند. a را یک تقریب A می‌گوییم هرگاه داشته باشیم $A - a \neq 0$. اگر $a < A$ آن‌گاه a را تقریب نقصانی A و چنانچه $a > A$ در آن صورت a را تقریب اضافی A می‌نامیم.

مثال ۱-۱: در مورد عدد π داریم:

$$3,14 < \pi < 3,15$$

بنابراین عدد $3,14$ یک تقریب نقصانی و عدد $3,15$ یک تقریب اضافی عدد π می‌باشند.

از تعریف فوق معلوم می‌شود که تقریب یک عدد حقیقی منحصر به فرد نیست. بنابراین مناسب بودن یک تقریب به چه معناست؟ بدیهی است که اگر ناچار باشیم به یافتن جواب تقریبی یک مسئله رضایت دهیم، خواهان تقریبی هستیم که تا حدود زیادی با جواب دقیق مسئله مطابقت داشته باشد. بنابراین بایستی معیار یا معیارهایی برای مقایسه‌ی دقت تقریب‌های یک عدد داشته باشیم. انواع خطا در حقیقت همان معیارهای سنجش دقت یک تقریب هستند که در ادامه معرفی می‌شوند.

تعریف ۱-۲: اگر a تقریبی از عدد حقیقی A باشد، آن‌گاه خطای مطلق a را با $e(a)$ نشان داده و به صورت زیر تعریف می‌کنیم.

$$e(a) = |A - a|.$$

برای مثال اگر $A = 10,12$ ، $a = 10,05$ و $a' = 10,15$ ، آن‌گاه خطاهای مطلق a و a' به عنوان تقریب‌هایی از A برابرند با:

$$e(a) = |10,12 - 10,05| = 0,07,$$

$$e(a') = |10,12 - 10,15| = 0,03.$$

واضح است که هر چه $e(a)$ کوچک‌تر باشد، a تقریب دقیق‌تری برای A خواهد بود.

با اینکه خطای مطلق، معیار ساده‌ای برای تعیین دقت یک تقریب است، اما میزان دقت دو یا چند تقریب مختلف را نمی‌توان به کمک آن مقایسه کرد؛ زیرا این خطا فقط اختلاف بین مقدار واقعی و مقدار تقریبی را مشخص می‌کند و به اینکه این مقادیر چقدر بزرگ یا چقدر کوچک باشد، توجهی ندارد. برای مثال فرض کنید که $A = 1000$ و $B = 10^10$. اگر $a = 1500$ و $b = 10^10 + 500$ را به ترتیب به عنوان تقریب‌هایی از A و B در نظر بگیریم، در این صورت خطای مطلق هر دو تقریب برابر با 500 است. اما آیا دقت a و b نیز یکسان است؟ واضح است که نه و دقت b به مراتب بیش‌تر از دقت a است. به عبارت دیگر عدد 500 در مقایسه با عدد 10^10 مقداری ناچیز می‌نماید

اما در مقایسه با عدد 1000 ، رقم قابل ملاحظه‌ای است. با این توضیحات به معیار مناسب‌تری برای تعیین دقت یک تقریب نیاز داریم. این معیار که نشان‌دهنده‌ی خطا در واحد کمیت است، خطای نسبی نامیده می‌شود.

تعریف ۱-۳: اگر a تقریبی از عدد حقیقی $A \neq 0$ باشد، آنگاه خطای نسبی a را با $\delta(a)$ نشان داده و به صورت زیر تعریف می‌کنیم:

$$\delta(a) = \frac{|A - a|}{|A|} = \frac{e(a)}{|A|}.$$

برای مثال خطاهای نسبی $a = 1500$ به عنوان تقریبی از $A = 1000$ و $b = 10^{10} + 500$ به عنوان تقریبی از $B = 10^{10}$ برابرند با:

$$\delta(a) = \frac{500}{1000} = \frac{1}{2},$$

$$\delta(b) = \frac{500}{10^{10}} = \frac{1}{2} \times 10^{-7}.$$

واضح است که هر چه $\delta(a)$ کوچک‌تر باشد، a تقریب دقیق‌تری برای A خواهد بود.

اگرچه کوچک بودن خطای نسبی یک تقریب نشان‌دهنده‌ی دقت بالای آن تقریب است، اما از بزرگ بودن اندازه‌ی آن نمی‌توان نادقیق بودن یک تقریب را نتیجه گرفت. به عبارت دیگر، گاهی ممکن است که یک تقریب a از عدد ناصفر A از دقت بسیار خوبی برخوردار باشد اما خطای نسبی آن بسیار زیاد باشد و آن هنگامی است که $|A|$ خود کوچک باشد. برای مثال فرض کنید $A = 10^{-11}$ و $a = 10^{-10}$. عدد a تا حد زیادی به A نزدیک و تقریب مناسبی از آن است. اکنون خطاهای مطلق و نسبی a را به عنوان تقریبی از A محاسبه می‌کنیم:

$$e(a) = |A - a| = |10^{-11} - 10^{-10}| = 9 \times 10^{-11},$$

$$\delta(a) = \frac{e(a)}{|A|} = \frac{9 \times 10^{-11}}{10^{-11}} = 9.$$

در این جا خطای نسبی a در مقایسه با مقادیر A و a عدد بسیار بزرگی است و به نظر می‌رسد که خطای مطلق، دقت a را بهتر مشخص می‌کند. اکنون معیاری را معرفی می‌کنیم که می‌تواند هم برای مقادیر خیلی بزرگ و هم برای مقادیر خیلی کوچک A مورد استفاده قرار گیرد.

تعریف ۱-۴: اگر a تقریبی از عدد حقیقی A باشد، آنگاه خطای ترکیبی a را با $\mu(a)$ نشان داده و به صورت زیر تعریف می‌کنیم:

$$\mu(a) = \frac{|A - a|}{1 + |A|} = \frac{e(a)}{1 + |A|}.$$

با توجه به تعریف فوق می‌بینیم که اگر $|A| \ll 1$ ، یعنی $|A|$ خیلی کوچک‌تر از یک باشد، آنگاه $\mu(a) \approx e(a)$ (یعنی خطای ترکیبی تقریباً برابر با خطای مطلق است) و چنانچه $|A| \gg 1$ ، یعنی $|A|$ خیلی بزرگ‌تر از یک باشد،

^۱ به این خطا، خطای نسبی - مطلق نیز می‌گویند.

آن‌گاه $\delta(a) \approx \mu(a)$. همچنین واضح است که $0 < \mu(a) \leq \min\{e(a), \delta(a)\}$.

تاکنون سه معیار برای سنجش دقت یک تقریب معرفی کرده‌ایم. یک معیار دیگر تعداد ارقام بامعناى درست است که چون از مفهوم گرد کردن اعداد ناشی می‌شود، در بخش‌های بعد معرفی خواهد شد. در معیارهایی که تاکنون بیان داشته‌ایم فرض بر این بوده است که مقادیر دقیق و تقریبی (به ترتیب A و a) برای ما معلوم هستند. اما در بسیاری از موارد A که جواب دقیق مسئله است، مجهول می‌باشد و تنها تقریب‌هایی از آن را در دست داریم (مانند ریشه‌ی یک معادله). از این‌رو عملاً تعیین هیچ‌یک از خطاهای گفته شده برای سنجش دقت این تقریب‌ها امکان‌پذیر نیست. به همین دلیل در روش‌های عددی اغلب کران بالایی برای خطای مطلق یا نسبی یک تقریب به دست می‌آوریم. بدیهی است که اگر کران بالای خطایی بسیار بزرگتر از واقعیت موجود باشد، نمی‌توان آن را برای برآورد معقولی از خطا مورد استفاده قرار داد و باید سعی کنیم تا کران بالای کوچک‌تری برای خطا به دست آوریم.

تعریف ۱-۵: اگر a تقریبی از عدد حقیقی A باشد، آن‌گاه هر عدد نا کمتر از $e(a)$ را یک خطای مطلق حدی a می‌نامیم و آن را با e_a نشان می‌دهیم. بنابراین $e(a) \leq e_a$ و e_a منحصر به فرد نیست.

توجه داریم که اگر e_a یک خطای مطلق حدی a به عنوان تقریبی از A باشد، آن‌گاه:

$$e(a) \leq e_a \Rightarrow |A - a| \leq e_a \Rightarrow a - e_a \leq A \leq a + e_a$$

یعنی A در بازه‌ی $[a - e_a, a + e_a]$ قرار دارد. طول این بازه $2e_a$ است؛ بنابراین هرچه e_a کوچک‌تر باشد، حدود دقیق‌تری برای A حاصل می‌شود.

تذکر: بنابر قرارداد نامساوی $a - e_a \leq A \leq a + e_a$ را به اختصار به صورت زیر می‌نویسیم:

$$A = a \pm e_a.$$

مثال ۱-۲: با فرض آن‌که $3,141 < \pi < 3,142$ ، یک کران بالا برای خطای مطلق $3,14$ به عنوان تقریبی از π به دست آورید.

حل: داریم

$$3,141 < \pi < 3,142 \Rightarrow 3,141 - 3,14 < \pi - 3,14 < 3,142 - 3,14$$

$$\Rightarrow 0,001 < \pi - 3,14 < 0,002 \Rightarrow |\pi - 3,14| < 0,002$$

بنابراین $e(3,14) < 0,002$.

معمولاً در اکثر روش‌های آنالیز عددی حدود جواب، یعنی کران بالا و پایینی برای جواب، قابل محاسبه است و از آن‌جا می‌توان مانند مثال فوق کران بالایی برای $e(a)$ به دست آورد.

تعریف ۱-۶: اگر a تقریبی از عدد حقیقی A باشد، آن‌گاه هر عدد نا کمتر از $\delta(a)$ را یک خطای نسبی حدی a می‌نامیم و آن را با δ_a نشان می‌دهیم. بنابراین $\delta(a) \leq \delta_a$ و δ_a منحصر به فرد نیست.

قضیه‌ی زیر نشان می‌دهد که چگونه می‌توان با دانستن خطای مطلق حدی یک تقریب، کران بالایی برای خطای نسبی آن به دست آورد.

قضیه ۱-۱: اگر a تقریبی از A و e_a یک خطای مطلق حدی a باشد، آنگاه:

$$\delta(a) \leq \frac{e_a}{|a| - e_a}.$$

اثبات: از خواص قدرمطلق می‌دانیم که $|a| - |A| \leq |A - a|$ و با توجه به فرض قضیه داریم:

$$|A - a| \leq e_a$$

بنابراین، $|a| - |A| \leq e_a$ و در نتیجه:

$$\begin{aligned} |a| - e_a \leq |A| &\Rightarrow \frac{1}{|A|} \leq \frac{1}{|a| - e_a} \\ &\Rightarrow \frac{e(a)}{|A|} \leq \frac{e(a)}{|a| - e_a} \left(\leq \frac{e_a}{|a| - e_a} \right) \\ &\Rightarrow \delta(a) \leq \frac{e_a}{|a| - e_a}. \end{aligned}$$

مثال ۳-۱: با فرض آن‌که $\sqrt{10} < 3,17 < 3,1$ ، یک کران بالا برای خطای نسبی $3,15$ به‌عنوان تقریبی از $\sqrt{10}$ به‌دست آورید.

حل: داریم

$$\begin{aligned} 3,1 < \sqrt{10} < 3,17 &\Rightarrow 3,1 - 3,15 < \sqrt{10} - 3,15 < 3,17 - 3,15 \\ &\Rightarrow -0,05 < \sqrt{10} - 3,15 < 0,02 \Rightarrow |\sqrt{10} - 3,15| < 0,05 \end{aligned}$$

بنابراین $e(3,15) < 0,05$. اکنون با توجه به قضیه‌ی قبل داریم:

$$\delta(3,15) \leq \frac{0,05}{3,15 - 0,05} = \frac{0,05}{3,1} = \frac{1}{62}.$$

تذکره: اگر e_a در مقایسه با $|a|$ کوچک باشد، آنگاه $|a| - e_a \approx |a|$ و بنابراین می‌توان نوشت:

$$\delta(a) \lesssim \frac{e_a}{|a|}$$

(نماد \lesssim به معنی «تقریباً کوچک‌تر از» می‌باشد.)

معمولاً اگر e_a در مقایسه با $|a|$ ناچیز باشد $\delta(a)$ تقریباً مساوی $\frac{e_a}{|a|}$ گرفته می‌شود. یعنی:

$$\delta(a) \approx \frac{e_a}{|a|}$$

برای نمونه، در مثال قبل می‌توانیم بنویسیم:

$$\delta(3,15) \approx \frac{0,05}{3,15} = \frac{5}{315} = \frac{1}{63}.$$

۱-۳ دستگاه‌های نمایش اعداد

استفاده از دستگاه اعشاری (دهدی)، متداول‌ترین شیوهی نمایش اعداد در زندگی روزمره‌ی ماست. هنگامی که مقدار دقیق یک کمیت یا تقریبی از آن را می‌نویسیم از ده رقم ۰، ۱، ۲، ۳، ۴، ۵، ۶، ۷، ۸ و ۹ به همراه یک ممیز استفاده می‌کنیم. این در حالی است که در کامپیوترها معمولاً از دستگاه دودویی (binary) با تنها دو رقم ۰ و ۱ برای نمایش اعداد استفاده می‌شود. از این رو شیوهی نمایش اعداد در دستگاه‌های مختلف و چگونگی تبدیل آن‌ها از دستگاهی به دستگاه دیگر را مورد مطالعه قرار می‌دهیم.

تعریف ۱-۷: فرض کنید $\beta > 1$ عددی طبیعی و A یک عدد حقیقی مثبت باشد. در این صورت منظور از نمایش A در مبنای β نمایش آن به صورت

$$A = a_m \times \beta^m + a_{m-1} \times \beta^{m-1} + \dots$$

می‌باشد که در آن

$$a_m \neq 0; m, a_i \in \mathbb{Z}; 0 \leq a_i \leq \beta - 1; i = m, m-1, \dots$$

در بسط فوق چنانچه $\beta = 10$ بسط اعشاری (دهدی) A و چنانچه $\beta = 2$ بسط دودویی A حاصل می‌شود. معمولاً نمایش فوق را به صورت $A = (a_m a_{m-1} a_{m-2} \dots)_\beta$ می‌نویسند.

قضیه ۱-۲: (وجود و یکتایی نمایش اعداد طبیعی در مبنای $\beta > 1$) اگر A عددی صحیح و مثبت باشد، آن‌گاه دارای نمایش منحصر به فرد در مبنای $\beta > 1$ است.

$$A = a_m \beta^m + a_{m-1} \beta^{m-1} + \dots + a_1 \beta + a_0 \quad (*)$$

***اثبات:** ابتدا ثابت می‌کنیم A دارای چنین نمایشی است و سپس نشان می‌دهیم این نمایش منحصر به فرد است. اگر $A < \beta$ ، آن‌گاه (*) به ازای $a_0 = A$ و $m = 0$ برقرار است. یعنی هر عدد A کوچکتر از β نمایش خودش در مبنای β است. برای $A \geq \beta$ فرض می‌کنیم تمام اعداد صحیح مثبت کوچکتر از A دارای نمایشی در مبنای β باشند (فرض استقرای قوی). بنابر قضیه‌ی تقسیم داریم $A = \beta q + r$ که در آن $0 < q < A$ و $0 \leq r < \beta$. طبق فرض استقرا، q نمایشی در مبنای β دارد. با جایگذاری نمایش q در رابطه‌ی بالا، نمایشی در مبنای β برای A حاصل می‌شود.

اکنون نشان می‌دهیم که نمایش A در مبنای β منحصر به فرد است. اگر $A < \beta$ آن‌گاه $A = a_0$ و واضح است که A به جز این، نمایش دیگری ندارد. حال فرض کنید $A \geq \beta$ و نیز هر عدد کوچکتر از A دارای نمایش یکتایی در مبنای β باشد (فرض استقرا). هرگاه (*) برقرار باشد داریم:

$$A = (a_m \beta^{m-1} + a_{m-1} \beta^{m-2} + \dots + a_1) \times \beta + a_0$$

و چون $0 \leq r < \beta - 1$ پس a_0 باقیمانده و

$$(a_m \beta^{m-1} + a_{m-1} \beta^{m-2} + \dots + a_1)$$

خارج قسمت تقسیم A بر β هستند. علاوه بر این، خارج قسمت بزرگتر از صفر است زیرا $A \geq \beta$. حال اگر A

دارای نمایش دیگری مانند $A = c_t \beta^t + a_{t-1} \beta^{t-1} + \dots + c_1 \beta + c_0$ باشد، آن گاه همانند آنچه که در مورد (*) بحث شد، c_0 باقیمانده و

$$(c_t \beta^{t-1} + c_{t-1} \beta^{t-2} + \dots + c_1)$$

خارج قسمت تقسیم A بر β هستند. اما در قضیه‌ی تقسیم، باقیمانده و خارج قسمت یکتا هستند. پس $a_0 = c_0$ و $a_m \beta^{m-1} + \dots + a_1 = c_t \beta^{t-1} + \dots + c_1$. چون طرفین تساوی اخیر، هر دو نمایش‌های یک عدد طبیعی کوچکتر از A در مبنای β هستند، بنابر فرض استقرا، این دو نمایش یکسان هستند. چون $a_0 = c_0$ پس $m = t$. بنابراین برای هر $0 \leq i \leq m$ داریم $a_i = c_i$. در نتیجه نمایش A در مبنای β منحصر به فرد است.

با توجه به مطالب فوق می‌توان روند تکراری زیر را برای نمایش عدد طبیعی A در مبنای β ارائه نمود:

۱- قرار می‌دهیم $i = 0$ و $q = \lfloor \frac{A}{\beta} \rfloor$. (q خارج قسمت تقسیم A بر β است.)

۲- تا زمانی که $q \neq 0$ ، اعمال زیر را به ترتیب انجام می‌دهیم:

$$a_i = A - \beta q$$

$$A \leftarrow q$$

$$q \leftarrow \lfloor A/\beta \rfloor$$

$$i \leftarrow i + 1$$

بدین ترتیب برای A نمایشی مانند $A = (a_m a_{m-1} \dots a_1 a_0)_\beta$ به دست می‌آید. روش فوق برای نمایش عدد طبیعی A در مبنای β ، روش تقسیم‌های متوالی نام دارد.

مثال ۴-۱: نمایش دودویی عدد ۲۲ را به دست آورید.

باقیمانده خارج قسمت

حل:

| i | A | $q = \lfloor \frac{A}{\beta} \rfloor$ | $A - \beta q$ |
|-----|-----|---------------------------------------|---------------|
| ۰ | ۲۲ | ۱۱ | $0 = a_0$ |
| ۱ | ۱۱ | ۵ | $1 = a_1$ |
| ۲ | ۵ | ۲ | $1 = a_2$ |
| ۳ | ۲ | ۱ | $0 = a_3$ |
| ۴ | ۱ | ۰ | $1 = a_4$ |

→ به خارج قسمت ۰ رسیدیم.

بنابراین: $22 = (10110)_2$.

قضیه ۳-۱: هر عدد $0 < r < 1$ دارای نمایش منحصر به فرد در مبنای $\beta > 1$ است،

$$r = r_1 \beta^{-1} + r_2 \beta^{-2} + \dots + r_m \beta^{-m} + \dots \quad (**),$$

که در آن:

بی‌نهایت بار $1 - \beta$ ؛ $i = 1, 2, \dots$ ؛ $0 \leq r_i \leq \beta - 1$ ؛

*اثبات: برای آن که ثابت کنیم r دارای چنین نمایشی است کافی است r_i ها را معرفی کنیم. برای معرفی r_1 ابتدا طرفین (***) را در β ضرب می‌کنیم. خواهیم داشت:

$$\beta r = r_1 + r_2 \beta^{-1} + r_3 \beta^{-2} + \dots + r_m \beta^{-(m-1)} + \dots$$

$$. h = r_2 \beta^{-1} + r_3 \beta^{-2} + \dots + r_m \beta^{-(m-1)} + \dots$$

از آنجا که $0 \leq r_i \leq \beta - 1$ و بی‌نهایت بار $r_i \neq \beta - 1$ داریم،

$$0 \leq h < (\beta - 1)(\beta^{-1} + \beta^{-2} + \dots) = (\beta - 1) \left(\frac{1}{\beta - 1} \right) = 1$$

یعنی h عددی بزرگتر مساوی صفر و اکیداً کوچکتر از ۱ است. چون $\beta r = r_1 + h$ و $0 \leq h < 1$ در نتیجه $r_1 = [\beta r]$. برای معرفی r_2 می‌نویسیم:

$$\beta r - r_1 = r_2 \beta^{-1} + r_3 \beta^{-2} + \dots + r_m \beta^{-(m-1)} + \dots$$

و سپس مشابه روشی که در مرحله‌ی قبل برای یافتن r_1 به کار گرفتیم، در این مرحله r_2 را بدست می‌آوریم:

$$r_2 = [\beta r - r_1]$$

اکنون اگر مجدداً $\beta r - r_1$ را بنامیم، یعنی $r \leftarrow \beta r - r_1$ خواهیم داشت $r_2 = [\beta r]$ و می‌توانیم سایر r_i ها را به روشی مشابه به دست آوریم. واضح است که اگر در مرحله‌ای داشته باشیم $r = 0$ ، آن‌گاه نمایش r در مبنای β مختوم (دارای تعداد متناهی رقم) است. چنانچه $r \neq 0$ آن‌گاه نمایش r در مبنای β نامختوم (بی‌پایان) است و می‌توانیم روش فوق را تا هر جا که مایل باشیم تکرار کنیم.

با توجه به مطالب فوق، روند تکراری زیر را برای یافتن r_i ها داریم.

$$1 - r_1 = [\beta r]$$

۲- برای $i \geq 2$ تا رسیدن به نتیجه‌ی مطلوب، اعمال زیر را انجام می‌دهیم:

$$r \leftarrow \beta r - r_{i-1}$$

$$r_i = [\beta r]$$

۳- r_i ها به ازای $i = 1, 2, \dots$ نمایش عدد $0 < r < 1$ در مبنای β هستند.

تذکر: شرط بی‌نهایت بار $r_i \neq \beta - 1$ برای یکتایی بسط لازم است. زیرا مثلاً در مبنای $10 = \beta$ با بی‌نهایت

بار $r_i = 9$ با توجه به فرمول حد مجموع یک سری هندسی داریم:

$$0.9999 \dots = 0.9 + 0.09 + 0.009 + \dots = \frac{0.9}{1 - 0.1} = 1$$

و اگر شرط مذکور را برداریم برای عدد ۱، دو بسط اعشاری خواهیم داشت. همچنین:

$$0.24999 \dots = 0.25 \quad , \quad 0.29999 \dots = 0.3$$

و به‌طور کلی، هر عدد اعشاری مختوم (با تعداد متناهی رقم) دارای دو بسط اعشاری خواهد بود.

مثال ۱-۵: نمایش دودویی عدد 0.4 را به دست آورید.

حل:

| i | r | $2r$ | $r_i = [2r]$ | $2r - r_i$ |
|-----|-------|-------|--------------|------------|
| ۱ | 0.4 | 0.8 | ۰ | 0.8 |
| ۲ | 0.8 | 1.6 | ۱ | 0.6 |
| ۳ | 0.6 | 1.2 | ۱ | 0.2 |
| ۴ | 0.2 | 0.4 | ۰ | 0.4 |
| ۵ | 0.4 | | | |

مجدداً به $r = 0.4$ رسیدیم. پس نمایش دودویی 0.4 متناوب است. اعداد ستون r_i ، به ترتیب از بالا به پایین بیانگر نمایش دودویی عدد 0.4 هستند، یعنی $0.4 = (0.0110)_2$.

نتیجه: فرض کنید A یک عدد حقیقی مثبت باشد. می توان A را به صورت مجموع جزء صحیح و جزء کسری آن نوشت. یعنی $A = [A] + (A)$ که در آن (A) نشان دهندهی جزء کسری A است و برای آن داریم:

$$0 \leq (A) < 1$$

اکنون با توجه به دو قضیهی قبل $[A]$ و (A) هر یک دارای نمایش منحصر به فردی در مبنای β هستند. بنابراین با مجموع این دو نمایش، نمایش منحصر به فردی برای A حاصل می شود. لذا هر عدد حقیقی مثبت دارای یک نمایش منحصر به فرد در مبنای $\beta > 1$ است.

مثال ۱-۶: نمایش دودویی عدد $22/4$ را به دست آورید.

حل: برای نمایش دودویی عدد $22/4$ ، با توجه به نتیجهی قبل ابتدا باید آن را به صورت مجموع جزء صحیح و جزء کسری اش بنویسیم و سپس نمایش دو جزء آن را با هم جمع کنیم. با توجه به دو مثال قبل داریم:

$$22/4 = (10110.0110)_2$$

تذکره: در مبنای ۱۶ از حروف بزرگ الفبای لاتین به صورت زیر برای نمایش ارقام پس از ۹ استفاده می شود:

$$\begin{aligned} A = 10, & \quad B = 11, & \quad C = 12, \\ D = 13, & \quad E = 14, & \quad F = 15. \end{aligned}$$

همچنین در مبنای ۲۶، از حروف انگلیسی کوچک $a - z$ به جای ارقام ۰ تا ۲۵ استفاده می شود.

قضیه ۱-۴: اگر نمایش عدد A در مبنای β مختوم یا نامختوم و متناوب باشد، A یک عدد گویاست.

اثبات: ابتدا فرض می کنیم نمایش A مختوم باشد. مثلاً $A = (a_m a_{m-1} \dots a_1 a_0 / r_1 r_2 \dots r_n)_\beta$. بنابراین:

$$A = a_m \beta^m + a_{m-1} \beta^{m-1} + \dots + a_1 \beta + a_0 + r_1 \beta^{-1} + r_2 \beta^{-2} + \dots + r_n \beta^{-n}$$

طرفین تساوی را در β^n ضرب می‌کنیم.

$$\beta^n A = \beta^n (a_m \beta^m + a_{m-1} \beta^{m-1} + \dots + a_1 \beta + a_0) + (r_1 \beta^{n-1} + r_2 \beta^{n-2} + \dots + r_n)$$

در نتیجه:

$$A = \frac{\beta^n (a_m \beta^m + a_{m-1} \beta^{m-1} + \dots + a_1 \beta + a_0) + (r_1 \beta^{n-1} + r_2 \beta^{n-2} + \dots + r_n)}{\beta^n}$$

اکنون از اینکه صورت و مخرج کسر عددهای صحیح هستند نتیجه می‌شود A عددی گویاست.

حال فرض می‌کنیم که نمایش A نامختوم و متناوب باشد. مثلاً:

$$A = (a_m a_{m-1} \dots a_1 a_0 / r_1 r_2 \dots r_n \overline{c_1 c_2 \dots c_k})_\beta$$

(توجه کنید که اگر m یا n صفر باشند، آن‌گاه $a_m \dots a_0$ یا $r_1 \dots r_n$ وجود نخواهد داشت.) از تعریف A داریم:

$$A = (a_m a_{m-1} \dots a_1 a_0 / r_1 r_2 \dots r_n)_\beta + (\underbrace{\circ / \circ \circ \dots \circ}_{n \text{ مرتبه}} \overline{c_1 c_2 \dots c_k})_\beta$$

با توجه به قسمت قبل، نخستین عبارت سمت راست تساوی یک عدد گویاست. در مورد عبارت دوم داریم:

$$\begin{aligned} (\underbrace{\circ / \circ \circ \dots \circ}_{n \text{ مرتبه}} \overline{c_1 c_2 \dots c_k})_\beta &= \beta^{-n} (\circ / \overline{c_1 c_2 \dots c_k})_\beta \\ &= \beta^{-n} (c_1 c_2 \dots c_k)_\beta (\beta^{-k} + \beta^{-2k} + \dots) \\ &= \beta^{-n} (c_1 c_2 \dots c_k)_\beta \left(\frac{\beta^{-k}}{1 - \beta^{-k}} \right) \\ &= \beta^{-n} (c_1 c_2 \dots c_k)_\beta \left(\frac{1}{\beta^k - 1} \right) \\ &= \frac{(c_1 c_2 \dots c_k)_\beta}{\beta^n} \left(\frac{1}{\beta^k - 1} \right) \end{aligned}$$

که این هم یک عدد گویاست. بنابراین A مجموع دو عدد گویا و در نتیجه، خود عددی گویاست.

نتیجه: اگر A یک عدد گنگ باشد، نمایش A در مبنای β نامختوم و نامتناوب است.

تبدیل مبنای اعداد:

نحوه‌ی تبدیل یک عدد از مبنای 10 به مبنای β را قبلاً در متن درس شرح داده‌ایم. برای تبدیل یک عدد از مبنای β به مبنای 10 چهار حالت در نظر می‌گیریم:

۱. برای تبدیل عدد $A = (a_m a_{m-1} \dots a_1 a_0)_\beta$ به مبنای 10 از رابطه‌ی زیر استفاده می‌کنیم:

$$A = a_m \beta^m + a_{m-1} \beta^{m-1} + \dots + a_1 \beta + a_0.$$

۲. برای تبدیل عدد $A = (\circ / r_1 r_2 \dots r_n)_\beta$ به مبنای 10 از رابطه‌ی زیر استفاده می‌کنیم:

$$A = r_1 \beta^{-1} + r_2 \beta^{-2} + \dots + r_n \beta^{-n}.$$

۳. برای تبدیل عدد $A = (\circ / r_1 r_2 \dots r_n \overline{c_1 c_2 \dots c_k})_\beta$ به مبنای 10 از رابطه‌ی زیر استفاده می‌کنیم:

$$A = \frac{(r_1 r_2 \dots r_n \overline{c_1 c_2 \dots c_k})_\beta - (r_1 r_2 \dots r_n)_\beta}{(\underbrace{z z \dots z}_k \text{ مرتبه } \underbrace{\circ \circ \dots \circ}_n \text{ مرتبه } \circ)_\beta}$$

که در آن منظور از z رقم $\beta - 1$ است. (این رابطه در ادامه ثابت می‌شود).

۴. برای تبدیل عدد $A = (a_m a_{m-1} \dots a_1 a_0 / r_1 r_2 \dots r_n \overline{c_1 c_2 \dots c_k})_\beta$ به مبنای 10 ابتدا بخش‌های صحیح و کسری A را جداگانه به مبنای 10 می‌بریم و سپس آن‌ها را باهم جمع می‌کنیم.

نکته: به کمک نکات قبل می‌توان کسر متعارفی مولد یک عدد اعشاری را نیز تعیین نمود.

اثبات رابطه‌ی قسمت ۳:

$$\begin{aligned} A &= (\circ / r_1 r_2 \dots r_n \overline{c_1 c_2 \dots c_k})_\beta = (\circ / r_1 r_2 \dots r_n)_\beta + (\circ / \underbrace{\circ \circ \dots \circ}_n \text{ مرتبه } \overline{c_1 c_2 \dots c_k})_\beta \\ \Rightarrow A &= \frac{(r_1 r_2 \dots r_n)_\beta}{\beta^n} + \frac{(c_1 c_2 \dots c_k)_\beta}{\beta^n} \left(\frac{1}{\beta^k - 1} \right) = \frac{(\beta^k - 1)(r_1 r_2 \dots r_n)_\beta + (c_1 c_2 \dots c_k)_\beta}{\beta^n (\beta^k - 1)} \\ &= \frac{(\beta^k)(r_1 r_2 \dots r_n)_\beta + (c_1 c_2 \dots c_k)_\beta - (r_1 r_2 \dots r_n)_\beta}{\beta^n (\beta^k - 1)} = \frac{(r_1 r_2 \dots r_n c_1 c_2 \dots c_k)_\beta - (r_1 r_2 \dots r_n)_\beta}{\beta^n (\beta^k - 1)} \end{aligned}$$

از طرفی:

$$\begin{aligned} \beta^n (\beta^k - 1) &= \beta^n (\beta - 1) (\beta^{k-1} + \beta^{k-2} \dots + \beta + 1) \\ &= (\beta - 1) \beta^{n+k-1} + (\beta - 1) \beta^{n+k-2} + \dots + (\beta - 1) \beta^n \\ &= \underbrace{((\beta - 1) \dots (\beta - 1))}_k \text{ مرتبه } \underbrace{\circ \circ \dots \circ}_n \text{ مرتبه } \circ)_\beta \end{aligned}$$

در نتیجه با قرار دادن $z := \beta - 1$ خواهیم داشت:

$$A = \frac{(r_1 r_2 \dots r_n c_1 c_2 \dots c_k)_\beta - (r_1 r_2 \dots r_n)_\beta}{(\underbrace{z z \dots z}_k \text{ مرتبه } \underbrace{\circ \circ \dots \circ}_n \text{ مرتبه } \circ)_\beta}$$

مثال ۷-۱: کسر متعارفی معادل با عدد $1/23\bar{4}$ را تعیین کنید.

حل:

$$1/23\bar{4} = 1 + \frac{234 - 23}{900} = 1 + \frac{211}{900} = \frac{1111}{900}$$

روش کلی برای تبدیل مبنای اعداد، استفاده از مبنای 10 به عنوان مبنای واسطه است. به عبارت دیگر برای تبدیل عدد X از مبنای β_1 به مبنای β_2 ابتدا عدد X را از مبنای β_1 به مبنای 10 و سپس از مبنای 10 به مبنای β_2 می‌بریم. البته در حالتی که β_1 یا β_2 توانی از دیگری باشد یا هر دو توانی از یک عدد دیگر باشند، راهکارهای سریع‌تری وجود دارد که به آن‌ها نیز اشاره خواهیم کرد.

مثال ۸-۱: عدد $(11001/1011)_7$ را به مبنای ۳ برید.

حل: ابتدا A را به مبنای ۱۰ می‌بریم. داریم:

$$(11001)_2 = 1 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 = 16 + 8 + 1 = 25,$$

$$(0,1011)_2 = \frac{(1011)_2 - (10)_2}{(1100)_2} = \frac{9}{12} = 0,75.$$

بنابراین $A = 25,75$. قرار می‌دهیم $r = (A) = 0,75$ ، ابتدا A_1 را به مبنای ۳ می‌بریم:

باقیمانده خارج قسمت

| i | A_1 | $q = \left[\frac{A_1}{3}\right]$ | $A_1 - 3q$ |
|-----|-------|----------------------------------|------------|
| ۰ | ۲۵ | ۸ | $1 = a_0$ |
| ۱ | ۸ | ۲ | $2 = a_1$ |
| ۲ | ۲ | ۰ | $2 = a_2$ |

بنابراین: $(221)_3 = 25$. حال r را به مبنای ۳ می‌بریم:

| i | r | $3r$ | $r_i = [3r]$ | $3r - r_i$ |
|-----|------|------|--------------|------------|
| ۱ | ۰,۷۵ | ۲,۲۵ | ۲ | ۰,۲۵ |
| ۲ | ۰,۲۵ | ۰,۷۵ | ۰ | ۰,۷۵ |
| ۳ | ۰,۷۵ | | | |

مجدداً به $r = 0,75$ رسیدیم. پس نمایش عدد $0,75$ در مبنای سه متناوب است. اعداد ستون r_i به ترتیب از بالا به پایین بیانگر نمایش عدد $0,75$ در مبنای ۳ هستند، یعنی $(0,20)_3 = 0,75$. اکنون با توجه به دو قسمت قبل داریم:

$$A = (11001,1011)_2 = 25,75 = (221,20)_3$$

مثال ۹-۱: نمایش در مبنای هشت عدد $A = (6F9)_{16}$ را تعیین کنید.

حل: داریم $2^4 = 16 = \beta_1$ و $2^3 = 8 = \beta_2$. حال ابتدا معادل دودویی هر یک از اعداد ۶، F و ۹ را با در نظر گرفتن چهار مکان می‌نویسیم. با این کار معادل دودویی عدد A مشخص می‌شود:

$$A = (6F9)_{16} = \left(\underbrace{0110}_6 \underbrace{1111}_F \underbrace{1001}_9 \right)_2$$

اکنون در نمایش دودویی A از سمت راست سه رقم سه رقم جدا می‌کنیم و معادل آن را در مبنای ۸ قرار می‌دهیم:

$$A = (6F9)_{16} = \left(\underbrace{0110}_6 \underbrace{1111}_F \underbrace{1001}_9 \right)_2 = \left(\underbrace{011}_3 \underbrace{011}_3 \underbrace{111}_7 \underbrace{001}_1 \right)_2 = (3371)_8.$$

مثال ۱۰-۱: نمایش در مبنای چهار عدد $A = (564,7)_8$ را تعیین کنید.

حل: داریم $\beta_1 = 8 = 2^3$ و $\beta_2 = 4 = 2^2$. پس ابتدا معادل دودویی هر یک از اعداد ۵، ۶، ۴ و ۷ را با در نظر گرفتن سه مکان می‌نویسیم. با این کار معادل دودویی عدد A مشخص می‌شود:

$$A = (5647)_8 = \left(\underbrace{101}_5 \ \underbrace{110}_6 \ \underbrace{100}_4 / \underbrace{111}_7 \right)_2$$

حال در قسمت صحیح عدد، از سمت راست و در قسمت کسری از سمت چپ، دو رقم دو رقم جدا می‌کنیم و معادل آن را در مبنای ۴ قرار می‌دهیم. توجه کنید که در قسمت صحیح (کسری)، در سمت چپ (راست) یک رقم اضافه می‌آید که با قرار دادن یک صفر در سمت چپ (راست) آن تعداد مکان‌های آن را به دو می‌رسانیم.

$$(5647)_8 = \left(\underbrace{101}_5 \ \underbrace{110}_6 \ \underbrace{100}_4 / \underbrace{111}_7 \right)_2 = \left(\underbrace{01}_1 \ \underbrace{01}_1 \ \underbrace{11}_3 \ \underbrace{01}_1 \ \underbrace{00}_2 / \underbrace{11}_3 \ \underbrace{10}_2 \right)_2 = (11310/32)_4$$

۱-۴ نمایش ماشینی اعداد

در کامپیوترها مقدار حافظه‌ای که برای ذخیره‌ی هر عدد اختصاص داده می‌شود (طول کلمه) محدود است. به همین دلیل همواره تعداد محدودی از اعداد حقیقی قابل ذخیره در کامپیوتر هستند. به‌طور کلی برای ذخیره‌ی اعداد در کامپیوترها از یکی از دو سیستم نمایش ممیز ثابت و نمایش ممیز شناور استفاده می‌شود. البته سیستم نمایش ممیز ثابت بیشتر در کامپیوترهای اولیه مورد استفاده قرار می‌گرفت و امروزه در کامپیوترها از سیستم نمایش ممیز شناور استفاده می‌شود.

سیستم نمایش ممیز ثابت:

در این سیستم نمایش، اگر طول کلمه‌ی کامپیوتر ℓ باشد، n مکان آن برای نگهداری ارقام قبل از ممیز، m مکان آن برای نگهداری ارقام بعد از ممیز و ۱ مکان برای علامت (منفی یا مثبت) در نظر گرفته می‌شود که

$$\ell = m + n + 1$$

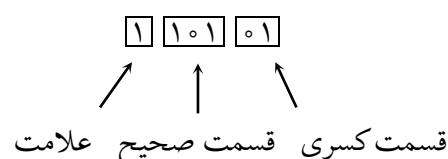
به عبارت دیگر مجموعه‌ی تمام اعداد در این سیستم نمایش به‌صورت زیر است:

$$\mathbb{F}_\beta^\circ = \left\{ \pm (a_{n-1} a_{n-2} \dots a_1 a_0 / b_1 b_2 \dots b_m)_\beta \mid a_i, b_j \in \{0, 1, \dots, \beta - 1\} \right\}$$

توجه کنید که همواره، مکان علامت با یکی از ارقام ۰ و ۱ پر می‌شود که رقم ۰ برای علامت مثبت و رقم ۱ برای علامت منفی به‌کار می‌رود. مثلاً اگر در یک سیستم نمایش ممیز ثابت داشته باشیم:

$$\beta = 2, \quad m = 2, \quad n = 3,$$

آن‌گاه عدد $(101/01)_2 = -5/25$ در این سیستم به‌صورت زیر نمایش داده می‌شود:



توجه کنید که در همین سیستم، مثلاً عدد ۱۰ قابل نمایش نیست زیرا $۱۰ = (۱۰۱۰)_۲$ که نیاز به چهار مکان در بخش صحیح دارد. بنابراین سؤالی که در این جا طبیعی است پرسیده شود این است که اساساً چند عدد در این سیستم قابل نمایش هستند و بزرگترین و کوچکترین آن‌ها کدامند؟ برای پاسخ به این سؤال ابتدا قضیه‌ی زیر را بیان می‌کنیم.

قضیه ۱-۵: در یک سیستم نمایش ممیز ثابت با مبنای β که در آن، n مکان برای نگهداری ارقام قبل از ممیز، m مکان برای نگهداری ارقام بعد از ممیز و ۱ مکان برای علامت (منفی یا مثبت) در نظر گرفته می‌شود:

(الف) تعداد اعداد حقیقی قابل نمایش عبارت است از $۲\beta^{m+n} - ۱$ ؛

(ب) بزرگترین عدد مثبت قابل نمایش برابر است با $\beta^n - \beta^{-m}$ ؛

(ج) کوچکترین عدد مثبت قابل نمایش برابر است با β^{-m} .

اثبات: الف) تعداد حالت‌های پر کردن مکان علامت برابر است با ۲ و هر یک از دیگر مکان‌ها را نیز می‌توان با یکی از β رقم ۰ تا $\beta - ۱$ پر کرد که چون تعداد این مکان‌ها $m + n$ است، طبق اصل ضرب این کار به β^{m+n} حالت امکان‌پذیر است. از طرفی ۰ و -۰ هر دو نمایش عدد صفر هستند. در نتیجه تعداد تمام اعداد قابل نمایش در این سیستم عبارت است از $۲\beta^{m+n} - ۱$.

(ب) فرض می‌کنیم x_{max} بزرگترین عدد مثبت قابل نمایش در این سیستم باشد و قرار می‌دهیم $\beta - ۱ := z$. لذا:

$$x_{max} = +(\underbrace{zz \dots z}_n / \underbrace{zz \dots z}_m)_\beta = \sum_{i=-m}^{n-1} z\beta^i = \sum_{i=-m}^{n-1} (\beta - 1)\beta^i = \sum_{i=-m}^{n-1} (\beta^{i+1} - \beta^i) = \beta^n - \beta^{-m}.$$

(ج) فرض می‌کنیم x_{min} کوچکترین عدد مثبت قابل نمایش در این سیستم باشد. بنابراین:

$$x_{min} = +(\underbrace{00 \dots 0}_n / \underbrace{00 \dots 0}_{m-1} 1)_\beta = \beta^{-m}.$$

نکته: توزیع اعداد در سیستم نمایش ممیز ثابت یکنواخت است. به این معنی که تمام اعداد با فواصل یکسان در بازه $[-x_{max}, x_{max}]$ قرار دارند. فاصله‌ی هر دو عدد متوالی برابر β^{-m} است.

حال سؤالی را که قبلاً مطرح کردیم در مثال زیر پاسخ می‌دهیم.

مثال ۱-۱۱: در سیستم نمایش ممیز ثابت با

$$\beta = ۲, \quad m = ۲, \quad n = ۳,$$

تعداد اعداد حقیقی قابل نمایش و نیز بزرگترین و کوچکترین اعداد مثبت قابل نمایش را مشخص کنید.

حل: تعداد اعداد حقیقی قابل نمایش در این سیستم برابر است با $۶۳ = ۲ \times ۲^{۲+۳} - ۱$. همچنین بزرگترین و کوچکترین اعداد مثبت قابل نمایش در این سیستم عبارتند از:

$$x_{max} = ۲^۳ - ۲^{-۲} = ۸ - ۰٫۲۵ = ۷٫۷۵, \quad x_{min} = ۲^{-۲} = ۰٫۲۵.$$

ضعف اصلی سیستم نمایش ممیز ثابت در این است که بازه‌ی $[-\beta^n, \beta^n] \approx [-x_{max}, x_{max}]$ محدوده‌ی چندان

بزرگی از اعداد حقیقی نیست مگر آن که n بسیار بزرگ انتخاب شود که حافظه‌ی محدود کامپیوتر این اجازه را نمی‌دهد. این ضعف از آنجا ناشی می‌شود که توزیع اعداد در این سیستم یکنواخت است. به عبارت دیگر حساسیت این سیستم روی اعداد کوچک و بزرگ یکی است. برای رفع این مشکل از سیستم ممیز شناور استفاده می‌شود که در ادامه به معرفی آن می‌پردازیم.

سیستم نمایش ممیز شناور:

تعریف ۸-۱: نمایش ممیز شناور یک عدد حقیقی دلخواه مانند A در حالت کلی (مبنای β) به صورت زیر است:

$$A = \pm a \times \beta^e,$$

در نمایش فوق، a عددی بین صفر و یک در مبنای β است که به آن مانتیس می‌گویند و e عددی صحیح در مبنای β است که به آن نما گفته می‌شود.^۱

با توجه به تعریف فوق، نمایش ممیز شناور هر عدد حقیقی از چهار پارامتر علامت، مانتیس، مبنای β و نما تشکیل می‌شود. هم‌چنین، یک عدد حقیقی ممکن است نمایش‌های ممیز شناور متفاوتی داشته باشد. مثلاً هر یک از اعداد $\beta^3 \times (01001)_\beta$ ، $\beta^2 \times (0101)_\beta$ و $\beta^1 \times (01)_\beta$ نمایش ممیز شناور عدد ۱ در مبنای β هستند. در حقیقت در نمایش ممیز شناور می‌توان ممیز را در طول ارقام مانتیس جابه‌جا نمود مشروط بر اینکه نما به‌طور مناسب تغییر کند. اگر ممیز یک رقم به سمت راست حرکت کند از نما یک واحد کم می‌شود. چنانچه شرط کنیم که اولین رقم پس از ممیز در مانتیس همواره مخالف صفر باشد (مگر در حالتی که مانتیس صفر باشد)، آن‌گاه هر عدد حقیقی، نمایش ممیز شناور یکتایی خواهد داشت که به آن نمایش ممیز شناور نرمال گفته می‌شود. بنابراین در نمایش ممیز شناور نرمال، حداقل مقدار مانتیس برابر با $\beta^{-1} \times (01)_\beta$ است.

در کامپیوترها به دلیل محدودیت حافظه، تعداد ارقام مانتیس و نما محدود انتخاب می‌شوند. بنابراین در یک سیستم ممیز شناور نرمال، اگر فرض کنیم که طول کلمه‌ی کامپیوتر l باشد، آن‌گاه n مکان برای نگهداری ارقام نما، m مکان برای نگهداری ارقام مانتیس، ۱ مکان برای علامت مانتیس و ۱ مکان برای علامت نما در نظر گرفته می‌شود که $l = m + n + 2$. به عبارت دیگر، مجموعه‌ی تمام اعداد قابل نمایش در این سیستم به صورت زیر است:

$$\mathbb{F} = \left\{ \pm (0b_1b_2 \dots b_m)_\beta \times \beta^{\pm(e_1e_2 \dots e_n)_\beta} \mid b_i, e_j \in \{0, 1, \dots, \beta - 1\}, b_1 \neq 0 \right\} \cup \{0\}.$$

در این سیستم نمایش نیز مکان‌های علامت با یکی از دو رقم ۰ و ۱ پر می‌شوند که رقم ۰ برای علامت مثبت و رقم ۱ برای علامت منفی به‌کار می‌رود. مثلاً اگر در یک سیستم نمایش ممیز شناور نرمال داشته باشیم:

$$\beta = 2, \quad m = 3, \quad n = 2,$$

آن‌گاه عدد $-0.375_{10} = -0.110_2 \times 2^{-(01)_2} = -(0110)_2 \times 2^{-(01)_2}$ در این سیستم به صورت زیر نمایش می‌شود:

$$\boxed{1} \boxed{110} \boxed{1} \boxed{01}$$

در نمایش فوق اولین رقم از چپ علامت مانتیس، سه رقم بعدی ارقام مانتیس (آن‌هایی که در سمت راست ممیز

^۱ چنانچه $1 \leq a < \beta$ نمایش حاصل را نمایش علمی A می‌گویند.

در نمایش دودویی قرار دارند)، رقم بعدی علامت نما و دو رقم آخر ارقام نما را مشخص می‌کنند.

قضیه ۱-۶: در یک سیستم نمایش ممیز شناور نرمال با مبنای β که در آن، n مکان برای نگهداری ارقام نما، m مکان برای نگهداری ارقام مانتیس، و 2 مکان برای علامت‌های مانتیس و نما (منفی یا مثبت) در نظر گرفته می‌شود:

الف) تعداد اعداد حقیقی قابل نمایش عبارت است از $1 + (2\beta^n - 1)\beta^{m-1}(\beta - 1)$ ؛

ب) بزرگترین عدد مثبت قابل نمایش برابر است با $(1 - \beta^{-m})\beta^{\beta^n - 1}$ ؛

ج) کوچکترین عدد مثبت قابل نمایش برابر است با $\beta^{-\beta^n}$.

اثبات: الف) ابتدا تعداد مانتیس‌های ممکن را می‌شماریم. علامت مانتیس می‌تواند مثبت یا منفی باشد (۲-حالت). از آنجا که هر عدد ناصفر قابل نمایش در این سیستم، نرمال شده است، اولین رقم مانتیس نمی‌تواند صفر باشد. پس اولین رقم مانتیس یکی از ارقام 1 تا $\beta - 1$ است ($\beta - 1$ حالت). سایر ارقام مانتیس می‌توانند هر یک از رقم‌های 0 تا $\beta - 1$ را اختیار کنند (در مجموع β^{m-1} حالت). بنابراین طبق اصل ضرب، تعداد مانتیس‌های ممکن برابر با $2(\beta - 1)\beta^{m-1}$ است. حال تعداد نماهای ممکن را می‌شماریم. علامت نما می‌تواند مثبت یا منفی باشد (۲-حالت) و هر یک از ارقام نما می‌توانند یکی از β رقم 0 تا $\beta - 1$ را اختیار کنند (در مجموع β^n حالت). از طرفی در نما $0+$ و $0-$ هر دو نمایش عدد صفر هستند. در نتیجه تعداد نماهای ممکن برابر با $2\beta^n - 1$ است. لذا تعداد تمام اعداد ناصفر قابل نمایش در این سیستم عبارت است از $(2\beta^n - 1)\beta^{m-1}(\beta - 1)$. اکنون با افزودن عدد صفر به مجموعه‌ی اعداد فوق، نتیجه‌ی مطلوب حاصل می‌شود.

ب) فرض می‌کنیم x_{max} بزرگترین عدد مثبت قابل نمایش در این سیستم باشد و قرار می‌دهیم $z = \beta - 1$. لذا:

$$x_{max} = +(\underbrace{z z \dots z}_m)_\beta \times \beta^{\underbrace{z z \dots z}_n} = \left(\sum_{i=-m}^{-1} z\beta^i \right) \beta^{\left(\sum_{i=0}^{n-1} z\beta^i \right)} = (1 - \beta^{-m})\beta^{\beta^n - 1}.$$

ج) فرض می‌کنیم x_{min} کوچکترین عدد مثبت قابل نمایش در این سیستم باشد. بنابراین با توجه به این‌که اولین رقم مانتیس نمی‌تواند صفر باشد داریم:

$$x_{min} = +(\underbrace{0 1 \dots 0}_{m-1})_\beta \times \beta^{\underbrace{-(z z \dots z)}_n} = \beta^{-1} \times \beta^{-(\beta^n - 1)} = \beta^{-\beta^n}.$$

نکته: توزیع اعداد در سیستم نمایش ممیز شناور نرمال یکنواخت نیست. در حقیقت اگر $x \in \mathbb{F}$ یک عدد مثبت و دارای نمایش $(0.b_1 b_2 \dots b_m)_\beta \times \beta^e$ باشد، عدد بزرگ‌تر بلافاصله بعد از آن در این سیستم نمایش که با x_+ نشان داده می‌شود، دارای نمایش $(0.b_1 b_2 \dots b_m + \underbrace{0 1 \dots 0}_{m-1})_\beta \times \beta^e$ است. بنابراین فاصله‌ی بین دو عدد

قابل نمایش متوالی x و x_+ برابر است با $\beta^{e-m} = (\underbrace{0 1 \dots 0}_{m-1})_\beta \times \beta^e$. این فاصله به‌وضوح به نما وابسته

است، یعنی با افزایش نما، فاصله‌ی بین اعداد بیشتر خواهد شد.

مثال ۱-۱۲: در سیستم نمایش ممیز شناور نرمال با

$$\beta = 2, \quad m = 3, \quad n = 1,$$

(الف) تعداد اعداد حقیقی قابل نمایش و نیز بزرگترین و کوچکترین اعداد مثبت قابل نمایش را مشخص کنید.

(ب) اعداد مثبت قابل نمایش در این سیستم را لیست کرده و بر روی محور اعداد مشخص کنید.

حل: (الف) تعداد اعداد حقیقی قابل نمایش در این سیستم برابر است با:

$$2 \times (2 - 1) \times 2^{3-1} \times (2 \times 2^1 - 1) + 1 = 25,$$

همچنین بزرگترین و کوچکترین اعداد مثبت قابل نمایش در این سیستم عبارتند از:

$$x_{max} = (1 - 2^{-3}) \times 2 = 2 - 2^{-2} = \frac{7}{4}, \quad x_{min} = 2^{-2} = \frac{1}{4}.$$

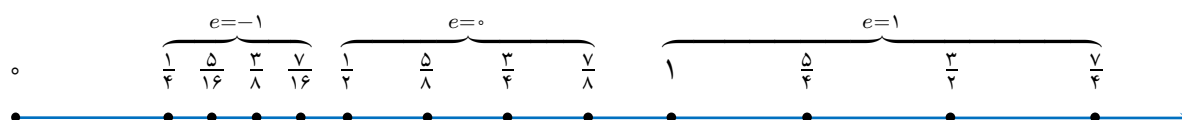
(ب) هر عدد مثبت قابل نمایش در این سیستم به شکل $x = (0.b_1 b_2 b_3)_2 \times 2^{\pm e}$ است که در آن:

$$e, b_i \in \{0, 1\}, b_1 \neq 0.$$

بنابراین اعداد مثبت قابل نمایش در این سیستم عبارتند از:

| | | |
|--------------------------------------|----------------------------------|----------------------------------|
| $0.100 \times 2^{-1} = \frac{1}{4}$ | $0.100 \times 2^0 = \frac{1}{2}$ | $0.100 \times 2^1 = 1$ |
| $0.101 \times 2^{-1} = \frac{5}{16}$ | $0.101 \times 2^0 = \frac{5}{8}$ | $0.101 \times 2^1 = \frac{5}{4}$ |
| $0.110 \times 2^{-1} = \frac{3}{8}$ | $0.110 \times 2^0 = \frac{3}{4}$ | $0.110 \times 2^1 = \frac{3}{2}$ |
| $0.111 \times 2^{-1} = \frac{7}{16}$ | $0.111 \times 2^0 = \frac{7}{8}$ | $0.111 \times 2^1 = \frac{7}{4}$ |

نمایش این اعداد بر روی محور، به صورت زیر است:



همان طور که در شکل فوق ملاحظه می شود، فاصله‌ی بین اعداد با نمای بزرگتر، بیشتر از فاصله‌ی بین اعداد با نمای کوچکتر است و اعداد با نمای یکسان، به طور متساوی الفاصله از یکدیگر قرار دارند.

گاهی به جای آن که در نمایش ممیز شناور نرمال تعداد ارقام نما مشخص شود، محدوده‌ای مانند $L \leq e \leq U$ برای نما در نظر گرفته می شود که در آن $L < 0$ و $U > 0$ دو عدد صحیح هستند. یعنی در این حالت مجموعه‌ی تمام اعداد ممیز شناور نرمال به صورت زیر است:

$$\mathbb{F} = \left\{ \pm (0.b_1 b_2 \dots b_m)_\beta \times \beta^e \mid b_i \in \{0, 1, \dots, \beta - 1\}, b_1 \neq 0, L \leq e \leq U \right\} \cup \{0\}.$$

این سیستم نمایش را اغلب به همراه پارامترهای مشخص کننده‌ی آن و به صورت $\mathbb{F}(\beta, m, L, U)$ نشان می دهند.

ویژگی‌های سیستم نمایش ممیز شناور نرمال در این حالت به صورت زیر بیان می‌شوند:

الف) تعداد اعداد حقیقی قابل نمایش عبارت است از $1 + (\beta - 1)\beta^{m-1}(U - L + 1)$ ؛

ب) بزرگترین عدد مثبت قابل نمایش برابر است با $(1 - \beta^{-m})\beta^U$ ؛

ج) کوچکترین عدد مثبت قابل نمایش برابر است با β^{L-1} .

د) اگر $x \in \mathbb{F}$ یک عدد مثبت و دارای نمایش $\beta^e \times (b_m \dots b_1 b_0)_\beta$ باشد، آنگاه فاصله‌ی بین x و x_+ برابر است با β^{e-m} . بنابراین توزیع اعداد در این سیستم نمایش یکنواخت نیست. کمترین فاصله‌ی ممکن بین دو عدد قابل نمایش به ازای $e = L$ و بیشترین فاصله‌ی ممکن به ازای $e = U$ رخ می‌دهد.

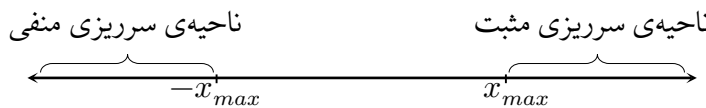
اعداد ماشینی و غیرماشینی:

تعریف ۹-۱: در یک ماشین محاسباتی که از یک سیستم نمایش اعداد استفاده می‌کند، تنها تعدادی متناهی از اعداد حقیقی را می‌توان نمایش داد. به اعداد قابل نمایش در ماشین، اعداد ماشینی و به اعداد غیر قابل نمایش در ماشین، اعداد غیرماشینی می‌گویند. مجموعه‌ی اعداد ماشینی را با A و مجموعه‌ی اعداد غیرماشینی را با A^c نشان می‌دهند.

سرریزی (overflow):

تعریف ۱۰-۱: اگر در یک سیستم نمایش برای عدد حقیقی x داشته باشیم $|x| > x_{max}$ ، آنگاه نمایش x در این سیستم امکان‌پذیر نیست. به چنین حالتی اصطلاحاً پدیده‌ی سرریزی گفته می‌شود.

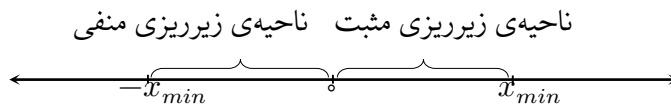
با رخ دادن پدیده‌ی سرریزی در یک ماشین، انجام عملیات توسط ماشین با دادن یک پیغام خطا متوقف می‌شود. پدیده‌ی سرریزی در سیستم نمایش ممیز شناور به علت محدود بودن نما از بالا رخ می‌دهد. ناحیه‌ی سرریزی را می‌توان به شکل زیر روی محور اعداد حقیقی مشخص کرد:



زیرریزی (under flow):

تعریف ۱۱-۱: اگر در یک سیستم نمایش برای عدد حقیقی x داشته باشیم $0 < |x| < x_{min}$ ، آنگاه نمایش x در این سیستم امکان‌پذیر نیست. به چنین حالتی اصطلاحاً پدیده‌ی زیرریزی گفته می‌شود.

با رخ دادن پدیده‌ی زیرریزی در یک ماشین، مقدار x برابر با صفر در نظر گرفته می‌شود. پدیده‌ی زیرریزی در سیستم نمایش ممیز شناور به علت محدود بودن نما از پایین رخ می‌دهد. ناحیه‌ی زیرریزی را می‌توان به شکل زیر روی محور اعداد حقیقی مشخص کرد. به ناحیه‌ی زیرریزی اصطلاحاً «حفره در صفر» نیز می‌گویند.



استاندارد IEEE برای نمایش اعداد ممیز شناور:

به لحاظ تاریخی، کامپیوترهای گوناگون انتخاب‌های متفاوتی در تعیین مبنا، کران‌های نما و ارقام مانتیس داشته‌اند. این تفاوت‌ها منجر به انتقال ناپذیری نرم‌افزارها روی ماشین‌های متفاوت می‌گردید تا این که در سال ۱۹۸۵ با تلاش‌های گروهی از ریاضیدانان، دانشمندان علوم کامپیوتر و شرکت‌های تولید سخت‌افزار به سرپرستی ویلیام کاهان* از دانشگاه کالیفرنیا، استاندارد برای نمایش اعداد ممیز شناور تحت عنوان IEEE۷۵۴ به سازندگان سخت‌افزارها عرضه شد. هم‌اکنون در بیشتر کامپیوترها از این استاندارد استفاده می‌شود. استاندارد IEEE، چند قالب کلی با دقت‌های مختلف از جمله دقت معمولی، دقت مضاعف و دقت‌های معمولی و مضاعف توسعه یافته برای نمایش اعداد ارائه می‌کند. در این جا به منظور آشنایی بیشتر با شیوهی نمایش اعداد در این استاندارد، نحوهی نمایش در دقت معمولی و مضاعف را شرح می‌دهیم.

استاندارد IEEE برای دقت معمولی، سیستم ممیز شناور نرمال $\mathbb{F}(2, 24, -126, 127)$ و برای دقت مضاعف، سیستم ممیز شناور نرمال $\mathbb{F}(2, 53, -1022, 1023)$ را پیشنهاد می‌کند. نحوهی ذخیره‌سازی اعداد در این استاندارد کمی با شیوهی عادی متفاوت است که این تفاوت به منظور استفاده هرچه بهینه‌تر از حافظه‌ی ماشین است. مطابق این استاندارد، در دقت معمولی از ۳۲ بیت و در دقت مضاعف از ۶۴ بیت برای نمایش یک عدد استفاده می‌شود. هر نمایش از سه بخش تشکیل می‌شود که عبارتند از علامت (s)، نمای تعدیل یافته (c) و قسمت کسری مانتیس نرمال شده (f). این سه بخش با استفاده از روابط زیر مشخص می‌شوند و به صورت $s|c|f$ در کنار هم قرار می‌گیرند:

$$x = \pm(1.f)_2 \times 2^e = (-1)^s(1.f)_2 \times 2^{c-127} \quad \text{دقت معمولی:}$$

$$x = \pm(1.f)_2 \times 2^e = (-1)^s(1.f)_2 \times 2^{c-1023} \quad \text{دقت مضاعف:}$$

همان‌طور که ملاحظه می‌کنید شکل کلی قالب‌های ذکر شده در دقت‌های معمولی و مضاعف، کمی شبیه به نمایش ممیز شناور نرمال است. فقط باید توجه داشت که در استاندارد IEEE مانتیس x به صورت $(1.f)_2$ نرمال شده است و تنها قسمتی از مانتیس که با f نشان داده شده است نمایش داده می‌شود. در واقع، چون اولین بیت مانتیس نرمال شده همواره برابر با ۱ است نیازی به ذخیره‌سازی آن نیست. در عوض، این بیت برای نمایش نما مورد استفاده قرار می‌گیرد. همچنین برای آنکه در ذخیره‌سازی نمای واقعی (e) نیازی به بیت علامت نباشد، این نما به صورت اریبی تعدیل و آنگاه ذخیره می‌شود. به عبارت دیگر یک مقدار ثابت به نام اُریب به نما (e) اضافه می‌شود تا یک عدد بدون علامت (c) در حافظه ذخیره شود. به سادگی می‌توان دید که برای دقت ساده اریب برابر با ۱۲۷ و برای دقت مضاعف برابر با ۱۰۲۳ است.

مثال ۱-۱۳: عدد $x = -45/75$ را در نظر بگیرید. می‌خواهیم این عدد را در استاندارد IEEE با دقت معمولی نمایش دهیم. برای این منظور، ابتدا نمایش دودویی آن را می‌یابیم. داریم $x = -(101101/11)$. حال این عدد باید به فرم $(1.f)_2 \times 2^e$ تبدیل شود. به سادگی می‌توان دید که $x = -(1/0110111) \times 2^5$.

از آن را قطع می‌کنیم.

۳. اگر اولین رقم ناخواسته مساوی با ۵ باشد و بعد از آن رقم مخالف صفری وجود داشته باشد، به رقم قبل از ۵ یک واحد اضافه می‌کنیم و رقم ۵ و ارقام پس از آن را قطع می‌کنیم.

۴. اگر اولین رقم ناخواسته مساوی با ۵ باشد و بعد از آن رقم مخالف صفری وجود نداشته باشد، آن‌گاه در صورتی که رقم قبل از ۵ فرد باشد یک واحد به آن اضافه می‌کنیم و در صورتی که رقم قبل از ۵ زوج باشد آن را تغییر نمی‌دهیم و رقم ۵ و ارقام پس از آن را قطع می‌کنیم.^۱

مثال ۱-۱۷: گرد شده‌ی عدد $e = ۲,۷۱۸۲ \dots$ تا یک رقم اعشار عبارت است از: $e = ۲,۷$ (۱D).

مثال ۱-۱۸: گرد شده‌ی عدد $e = ۲,۷۱۸۲ \dots$ تا دو رقم اعشار عبارت است از: $e = ۲,۷۲$ (۲D).

مثال ۱-۱۹: گرد شده‌ی عدد $w = ۰,۷۵۰۲$ تا یک رقم اعشار عبارت است از: $w = ۰,۸$ (۱D).

مثال ۱-۲۰: گرد شده‌ی عدد $x = ۱,۴۳۵$ تا دو رقم اعشار عبارت است از: $x = ۱,۴۴$ (۲D).

مثال ۱-۲۱: گرد شده‌ی عدد $y = ۰,۸۴۵$ تا دو رقم اعشار عبارت است از: $y = ۰,۸۴$ (۲D).

مثال ۱-۲۲: گرد شده‌ی عدد $z = ۳,۹۹۵$ تا دو رقم اعشار عبارت است از: $z = ۴,۰۰$ (۲D).

مقایسه‌ی خطای قطع کردن با خطای گرد کردن

فرض کنید a گرد شده‌ی عدد $\dots ۰,۶۶۶ = \frac{۲}{۳}$ و b قطع شده‌ی $\frac{۲}{۳}$ تا دو رقم اعشار باشند. داریم:

$$a = ۰,۶۷ \Rightarrow \left| \frac{۲}{۳} - a \right| = \frac{۱}{۳۰۰},$$

$$b = ۰,۶۶ \Rightarrow \left| \frac{۲}{۳} - b \right| = \frac{۲}{۳۰۰}.$$

یعنی فاصله‌ی b تا $\frac{۲}{۳}$ دو برابر فاصله‌ی a تا $\frac{۲}{۳}$ است. به عبارت دیگر، خطای قطع کردن تا دو برابر خطای گرد کردن است. به همین دلیل در عمل بیشتر از گرد کردن استفاده می‌شود (هرچند که قطع کردن ساده‌تر است).

قضیه ۱-۷: فرض کنید a و b به ترتیب گرد شده و قطع شده‌ی عدد A تا n رقم اعشار باشند. در این صورت:

$$e(a) \leq ۵ \times ۱۰^{-(n+1)} = \frac{۱}{۴} \times ۱۰^{-n},$$

$$e(b) \leq ۱۰^{-n}.$$

اثبات: فرض کنید که $r = ۰, d_1 d_2 \dots d_n d_{n+1} \dots$ نمایش قسمت کسری A باشد و a را گرد شده‌ی A تا n رقم اعشار در نظر بگیرید. اگر d_{n+1} کوچک‌تر از ۵ باشد، آن‌گاه:

$$e(a) = |A - a| = ۰, \underbrace{۰ \dots ۰}_{n \text{ مرتبه}} d_{n+1} \dots \leq ۵ \times ۱۰^{-(n+1)}$$

^۱ وضع این قاعده که منسوب به گاوس* است، به این دلیل است که در یک مسئله با محاسبات زیاد، احتمال وقوع اعداد اعشاری با رقم زوج قبل از ۵ و احتمال وقوع اعداد اعشاری با رقم فرد قبل از ۵ با هم برابرند و از این رو میانگین خطای ناشی از گرد کردن متناظر با آن‌ها در محاسبات صفر است.

اگر d_{n+1} بزرگتر از ۵ باشد به رقم d_n یک واحد اضافه می‌شود. بنابراین:

$$e(a) = |A - a| = 1 \times 10^{-n} - \underbrace{0, \dots, 0}_{n \text{ مرتبه}} d_{n+1} \dots \leq 10^{-n} - 5 \times 10^{-(n+1)} = 5 \times 10^{-(n+1)}.$$

اگر $d_{n+1} = 5$ ، آنگاه با توجه به زوج یا فرد بودن d_n یکی از دو حالت قبل رخ می‌دهد و باز هم حکم برقرار است (در این حالت علامت‌های \leq به $=$ تبدیل می‌شوند).

برای اثبات قسمت دوم قضیه، مجدداً فرض کنید که $r = 0, d_1 d_2 \dots d_n d_{n+1} \dots$ نمایش قسمت کسری A باشد و این بار b را قطع شده‌ی A تا n رقم اعشار در نظر بگیرید. در این صورت:

$$e(b) = |A - b| = \underbrace{0, \dots, 0}_{n \text{ مرتبه}} d_{n+1} \dots \leq 10 \times 10^{-(n+1)} = 10^{-n}.$$

نمایش علمی اعداد

گاهی اندازه‌ی یک کمیت ممکن است با مقادیر خیلی بزرگ یا خیلی کوچک بیان شود. مثلاً فاصله‌ی زمین تا خورشید در حدود $150,000,000,000$ متر (یکصد و پنجاه میلیارد متر) است، یا مثلاً برای نوشتن جرم یک الکترون برحسب گرم باید بعد از ممیز ۲۷ صفر قرار داد و پس از آن عدد 9109 را نوشت. بدیهی است که نوشتن چنین عددهایی با صفرهای زیاد، علاوه بر آن که خواندن و نوشتن را مشکل می‌کنند، احتمال اشتباه را هم زیاد می‌کنند. به همین دلیل از نمایشی موسوم به نمایش علمی برای اعداد استفاده می‌شود. فرض کنید A عددی مخالف صفر باشد. واضح است که A را همواره می‌توان به صورت $A = \pm a \times 10^e$ نوشت که در آن e عددی صحیح است و

$$1 \leq a < 10.$$

نمایش عدد A به صورت فوق را نمایش علمی A می‌نامند. در این نمایش a ماننسیس و b نمای عدد A نامیده می‌شود. بنابراین در نمایش علمی، حداقل مقدار ماننسیس برابر با ۱ است.^۱

مثال ۱-۲۳: نمایش علمی فاصله‌ی زمین تا خورشید، عبارت است از: $d = 1,5 \times 10^{11}$ (m).

مثال ۱-۲۴: نمایش علمی جرم الکترون عبارت است از: $m = 9,109 \times 10^{-28}$ (gr).

مثال ۱-۲۵: نمایش علمی عدد $0,0000341$ عبارت است از: $3,41 \times 10^{-5}$.

مثال ۱-۲۶: نمایش علمی عدد 2050 عبارت است از: $2,05 \times 10^3$.

مثال ۱-۲۷: نمایش علمی عدد 10000000 عبارت است از: 1×10^7 .

^۱ اگر به جای عدد 10 مبنای β در نظر گرفته شود، آنگاه نمایش علمی A در مبنای β حاصل می‌شود. در این صورت e عددی صحیح در مبنای β است و محدوده‌ی ماننسیس عبارت است از $1 \leq a < \beta$.

ارقام بامعنا

تعریف ۱-۱۲: ارقام بامعناى یک عدد شامل ممیز عبارتند از اولین رقم غیر صفر سمت چپ و ارقام صفر و غیرصفر پس از آن. ارقام بامعناى یک عدد صحیح بسته به این که در نمایش علمی آن، ماننیتس شامل چند رقم بامعنا باشد، برابر با تعداد ارقام بامعناى ماننیتس تعریف می شود.

مثال ۱-۲۸: اعداد زیر همگی شامل ۵ رقم بامعنا هستند:

$$۱/۵۰۳۵, \quad ۱/۵۰۳۰, \quad ۰/۰۱۵۰۳۵, \quad ۰/۰۰۱۵۰۳۰, \quad ۱۵۰/۳۰, \quad ۱۵۰۳/۰.$$

با توجه به تعریف فوق، اگر یک عدد صحیح مانند A شامل هیچ صفری در سمت راست خود نباشد، آنگاه تعداد ارقام بامعناى آن دقیقاً با تعداد ارقام بامعناى ماننیتس خود در نمایش علمی اش برابر است. اما چنانچه حداقل یک صفر، ارقام سمت راست A را تشکیل دهند، آنگاه این صفرها ممکن است بامعنا یا بی معنا باشند.

مثال ۱-۲۹: عدد $A = ۱۵۰۳۰۰۰$ را در نظر بگیرید. این عدد دست کم چهار رقم بامعنا دارد. اما بدون هیچ اطلاع دیگری نمی توان گفت تعداد ارقام بامعنا بیشتر هستند یا خیر.

اگر نمایش علمی A به صورت $۱/۵۰۳۰ \times ۱۰^۶$ در نظر گرفته شود، آنگاه تعداد ارقام بامعناى A برابر با ۵ است و در این صورت دو صفر سمت راست A بی معنا هستند و اگر نمایش علمی A به صورت $۱/۵۰۳۰۰۰ \times ۱۰^۶$ در نظر گرفته شود، آنگاه تعداد ارقام بامعناى A برابر با هفت است و تمام رقم های A بامعنا هستند.

توجه داشته باشید که همواره در قطع کردن یا گرد کردن یک عدد داده شده تا n رقم بامعنا (nS)^۱ عدد اصلی با عددی شامل n رقم بامعنا جایگزین می شود.

مثال ۱-۳۰: گردشده‌ی عدد $u = ۲۵/۸$ تا دو رقم بامعنا عبارت است از: $u = ۲۶$ ($۲S$).

مثال ۱-۳۱: گردشده‌ی عدد $v = ۰/۷۵$ تا یک رقم بامعنا عبارت است از: $v = ۰/۸$ ($۱S$).

مثال ۱-۳۲: گردشده‌ی عدد $w = ۴/۰۳۴$ تا سه رقم بامعنا عبارت است از: $w = ۴/۰۳$ ($۳S$).

مثال ۱-۳۳: گردشده‌ی عدد $x = ۰/۳۰۵$ تا دو رقم بامعنا عبارت است از: $x = ۰/۳۰$ ($۲S$).

مثال ۱-۳۴: گردشده‌ی عدد $y = ۲۰۸$ را تا دو رقم بامعنا به دست آورید.

حل: این یک مثال آموزشی جالب است، زیرا نمی توانیم ۲۱ را به عنوان جواب بنویسیم. با این که ۲۱ دارای دو رقم بامعناست ولی تقریب مطلوب y نیست. برای این منظور ابتدا نمایش علمی y را می نویسیم:

$$y = ۲/۰۸ \times ۱۰^۲.$$

حال ماننیتس را تا دو رقم بامعنا گرد می کنیم. جواب $۲/۱ \times ۱۰^۲$ است.

مثال ۱-۳۵: گردشده‌ی عدد $z = ۲۱۶۹۹۸۰$ را تا ۵ رقم بامعنا به دست آورید.

^۱ S حرف اول واژه‌ی Significant به معنای «بامعنا» است. از این رو نمایش یک عدد تا n رقم بامعنا با نماد (nS) بیان می شود.

حل: ابتدا نمایش علمی z را می‌نویسیم:

$$z = 2,16998 \times 10^6.$$

حال مانتیس را تا پنج رقم بامعنا گرد می‌کنیم. بنابراین جواب برابر است با $2,1700 \times 10^6$.

ارقام بامعناى درست یک تقریب

همان‌گونه که قبلاً گفته شد یکی از راه‌های سنجش دقت یک تقریب، تعداد ارقام بامعناى درست آن است. با توجه به این معیار، هرچه تعداد ارقام بامعناى درست یک تقریب بیشتر باشد، آن تقریب دقیق‌تر است. در این جا باید توجه داشت که دقت یک تقریب هیچ ارتباط معنی‌داری با تعداد ارقام بامعناى آن ندارد. مثلاً عدد $x = 3,7182$ تقریبی از عدد e است که ۵ رقم بامعنا دارد، اما تقریب خوبی از e نیست و مثلاً عدد ۳ (تنها با یک رقم بامعنا) تقریب بهتری است. بنابراین برای آن‌که بهتر بتوانیم از روی ارقام یک تقریب به دقت آن پی ببریم، مفهوم تعداد ارقام بامعناى درست یک تقریب را تعریف می‌کنیم. این مفهوم را می‌توان به‌طور نادقیق این‌گونه تعریف کرد که تقریب a از A دارای n رقم بامعناى درست است اگر a و A هر دو به یک عدد سوم با n رقم بامعنا گرد شوند. این تعریف معمولاً کارساز و از نظر شهودی درست است، اما مثال زیر را در نظر بگیرید:

$$A = 0,9949, \quad a = 0,9951,$$

با توجه به تعریف بالا a نمی‌تواند دو رقم با معناى درست داشته باشد. زیرا اگر هر دوی a و A را تا دو رقم با معنا گرد کنیم، A به عدد $0,99$ گرد می‌شود در حالی که a به عدد $1,0$ گرد می‌شود. اما در همین مثال، a دارای یک رقم با معناى درست و نیز سه رقم بامعناى درست است! تعریف دقیق زیر، این ابهام را برطرف می‌کند.

تعریف ۱-۱۳: اگر $a \neq 0$ تقریبی از A با نمایش علمی $a = a_m/a_{m-1}a_{m-2} \dots \times 10^m$ و d تعداد ارقام بامعناى a باشد، آن‌گاه بزرگ‌ترین عدد صحیح نامنفی n که $n \leq d$ و

$$|A - a| \leq 5 \times 10^{m-n}$$

تعداد ارقام بامعناى درست a نامیده می‌شود.

مثال ۱-۳۶: اگر $A = 0,9949$ آن‌گاه تعداد ارقام بامعناى درست $a = 0,9951$ به‌عنوان تقریبی از A چقدر است؟

حل: ابتدا نمایش علمی a را برای تعیین m می‌نویسیم. با توجه به تعریف، m برابر با توان 10 در نمایش علمی a است:

$$a = 9,951 \times 10^{-1} \implies m = -1$$

حال باید بزرگ‌ترین n را بیابیم که $n \leq 4$ و $|A - a| = 2 \times 10^{-4} \leq 5 \times 10^{-1-n}$ داریم:

$$2 \times 10^{-4} \leq 5 \times 10^{-1-n} \implies 2 \leq 5 \times 10^{4-1-n} \implies n_{max} = 3,$$

لذا تعداد ارقام بامعناى درست a برابر است با ۳.

مثال ۱-۳۷: ثابت اویلر، ... $\gamma = 0,577215664$ را در نظر بگیرید^۱. فرض کنید که به شما اعداد

$$a = 0,577215494, \quad a' = 0,577215834$$

برای تقریب γ پیشنهاد می‌شود. تقریب مناسب‌تر را انتخاب کنید.

حل: واضح است که در این جا نمی‌توان خطای مطلق یا نسبی یا ترکیبی را حساب کرد؛ زیرا مقدار دقیق γ مجهول است. هم‌چنین توجه دارید که عددهای داده شده تا ۶ رقم بامعنا با γ یکسان هستند و نیز اختلاف رقم‌های هفتم به بعد آن‌ها یعنی ۴۹۴ و ۸۳۴ با ارقام متناظرشان در γ یعنی ۶۶۴ برابرند. تا این جا a و a' وضعیت یکسانی دارند. حال بیایید γ ، a و a' را تا شش رقم بامعنا، یعنی شماره‌ی آخرین رقم بامعنا را قطعی در دو تقریب داده شده، گرد کنیم:

$$\gamma = 0,577216 \text{ (6S)},$$

$$a = 0,577215 \text{ (6S)},$$

$$a' = 0,577216 \text{ (6S)}.$$

ملاحظه می‌شود که گرد شده‌ی a' و γ تا شش رقم بامعنا با هم یکسان شدند ولی برای a چنین اتفاقی نیفتاد. این یعنی تعداد ارقام بامعنا در a' بیشتر از تعداد ارقام بامعنا در a است. در نتیجه بهتر است a' را انتخاب کنید. این انتخاب از نظر ریاضی هم کاملاً درست است، زیرا سرانجام پس از ارقام ۶۶۴ در نمایش γ رقم مخالف صفری یافت می‌شود (چون اگر هرگز رقم مخالف صفری یافت نشود آن‌گاه نتیجه می‌شود که γ گویاست!). وجود همان رقم مخالف صفر، نزدیکی بیشتر a' به γ را ثابت می‌کند.

مثال ۱-۳۸: فرض کنید $A = 8,000$ ، $a = 7,997$ و $a' = 8,08$. تعداد ارقام بامعنا در a و a' را تعیین کنید.

حل: همان‌طور که مشاهده می‌شود که a' دو رقم بامعنا مساوی با ارقام A دارد (با حفظ ارزش هر رقم) اما هیچ‌یک از رقم‌های a مساوی با ارقام A نیست. آیا می‌توان گفت که ارقام درست a' بیشتر از ارقام درست a است؟ خواهیم دید که خیر. در این جا $e(a) = 0,003$ و $e(a') = 0,08$ و در واقع a باید تعداد ارقام درست بیشتری داشته باشد! اگر a را تا سه رقم بامعنا گرد کنیم، عدد A حاصل می‌شود. از این رو، a سه رقم بامعنا درست دارد. اگر a' را تا رقم یکان گرد کنیم، A حاصل می‌شود (توجه کنید که حتی گرد شده‌ی a' تا یک رقم اعشار، به عدد $8,1$ منجر می‌شود که مساوی با A نیست). یعنی a' تنها یک رقم بامعنا درست دارد. اکنون تعداد ارقام بامعنا در a و a' را با استفاده از تعریف نیز محاسبه می‌کنیم. چون $a = 7,997 \times 10^0$ و $a' = 8,08 \times 10^0$ بنابراین مقدار m برای هر دو صفر است. برای a داریم $e(a) = 0,003$ ، بنابراین باید بزرگ‌ترین n را بیابیم که در نامساوی زیر صدق کند:

^۱ این عدد در واقع برابر است با $\lim_{n \rightarrow \infty} (1 + \frac{1}{2} + \dots + \frac{1}{n} - \ln n)$ و تاکنون گنگ یا گویا بودن آن معلوم نشده است.

$$0.003 \leq 5 \times 10^{-n}$$

بدیهی است که بزرگترین n برابر است با ۳. یعنی، a دارای ۳ رقم بامعناى درست است. برای a' نیز داریم:

$$e(a') = 0.08 = 8 \times 10^{-2} < 5 \times 10^{-1}$$

یعنی a' تنها یک رقم بامعناى درست دارد.

مثال ۱-۳۹: اگر $A = 100$ ، $a = 99.98$ و $b = 100.6$. تعداد ارقام بامعناى درست a و b را تعیین کنید.

حل: در مورد a داریم:

$$a = 9.998 \times 10^1 \implies m = 1,$$

$$e(a) = |A - a| = |100 - 99.98| = 0.02.$$

بنابراین باید بزرگترین n را بیابیم که $0.02 < 5 \times 10^{1-n}$.

بدیهی است که $n = 3$ جواب است. یعنی، a دارای ۳ رقم بامعناى درست است. در مورد b داریم:

$$b = 1.006 \times 10^2 \implies m = 2,$$

$$e(b) = |A - b| = |100 - 100.6| = 0.6.$$

بنابراین باید بزرگترین n را بیابیم که $0.6 < 5 \times 10^{2-n}$. در این جا $n = 2$ جواب است. یعنی، b تنها ۲ رقم بامعناى درست دارد.

قضیه ۱-۸: اگر a تقریبی از A با n رقم بامعناى درست باشد و $b = 10^k \times a$ و $B = 10^k \times A$ ، که در آن k عددی صحیح است، آنگاه b تقریبی از B با n رقم بامعناى درست است و خطای نسبی a و b یکسان هستند.

اثبات: فرض کنید $a = a_m/a_{m-1}a_{m-2} \dots \times 10^m$ نمایش علمی a باشد. در این صورت

$$b = 10^k \times a = a_m/a_{m-1}a_{m-2} \dots \times 10^{m+k},$$

یعنی نمای مربوط به b در نمایش علمی برابر با $m+k$ است. چون a دارای n رقم بامعناى درست است پس n بزرگترین عدد صحیحی است که در رابطه $|A - a| \leq 5 \times 10^{m-n}$ صدق می‌کند. بنابراین:

$$|B - b| = 10^k |A - a| \leq 5 \times 10^{(m+k)-n},$$

و n بزرگترین عدد صحیحی است که در این نامساوی صدق می‌کند. از این رو b نیز دارای n رقم بامعناى درست است. قسمت دوم حکم به سادگی از تعریف خطای نسبی b نتیجه می‌شود.

قضیه ۱-۹: اگر a گرد شده عدد مثبت A تا n رقم بامعنا باشد، آنگاه a دارای n رقم بامعناى درست است.

اثبات: با توجه به قضیه‌ی قبل می‌توانیم فرض کنیم که:

$$a = 0.c_1c_2c_3 \dots c_n, \quad c_1 \neq 0$$

(اگر a چنین نباشد، با انتخاب عدد صحیح و مناسب k می‌توان $10^k \times a$ را به این شکل در آورد). بنابراین در

مورد این a داریم $m = -1$.

از طرف دیگر می‌دانیم که اگر a گرد شده‌ی A تا n رقم اعشار باشد، آنگاه $|A - a| \leq 5 \times 10^{-(n+1)}$. اکنون با استفاده از این موضوع داریم:

$$|A - a| \leq 5 \times 10^{-(n+1)} = 5 \times 10^{-n-1},$$

که نتیجه می‌دهد a دارای n رقم بامعناى درست است. (توجه کنید که ممکن است داشته باشیم

$$|A - a| \leq 5 \times 10^{-1-n'}$$

که در آن $n' > n$. اما چون a بیش از n رقم بامعنا ندارد، کمترین n و همان n می‌شود.)

قضیه ۱-۱۰: اگر $a > 0$ تقریبی از A و دارای n رقم بامعناى درست باشد، آنگاه خطای نسبی a از 5×10^{-n} کمتر است، به شرط آن‌که رقم‌های درست a شامل یک رقم یک و $n-1$ رقم صفر در سمت راست آن نباشد.

اثبات: با توجه به قضایای قبل می‌توان فرض کرد که:

$$a = b_n \dots b_2 b_1 / c_1 c_2 c_3 \dots = b_n / b_{n-1} \dots b_2 b_1 c_1 c_2 c_3 \dots \times 10^{n-1}$$

که در آن $b_n \dots b_2 b_1$ عدد حاصل از n رقم بامعناى درست است. بنابراین نمای مربوط به این a برابر با $n-1$ است و با توجه به تعریف تعداد ارقام بامعناى درست هر تقریب داریم:

$$|A - a| \leq 5 \times 10^{(n-1)-n} = 0.5. \quad (1)$$

با توجه به این‌که $|a| - |A| \leq |A - a|$ ، از نامساوی فوق نتیجه می‌شود:

$$|A| \geq |a| - 0.5 = a - 0.5,$$

اما بنابر فرض قضیه $b_n \dots b_2 b_1 \neq 10^{n-1}$ ، پس $a > b_n \dots b_2 b_1$. لذا $a \geq 10^{n-1} + 1$ و در نتیجه:

$$|A| \geq a - 0.5 \geq 10^{n-1} + 1 - 0.5 = 10^{n-1} + 0.5 > 10^{n-1}. \quad (2)$$

$$(1), (2) \implies \delta(a) = \frac{|A - a|}{|A|} < \frac{0.5}{10^{n-1}} = 5 \times 10^{-n}.$$

در عمل به شرط انتهای این قضیه توجه نمی‌شود، زیرا به‌ندرت ممکن است تقریبی دارای این ویژگی باشد.

نتیجه: احکام دو قضیه‌ی قبل را می‌توان به صورت زیر خلاصه کرد:

$$a > 0 \implies \delta(a) < 5 \times 10^{-n} \implies a \text{ دارای } n \text{ رقم بامعناى درست است} \implies a \text{ گرد شده‌ی } A \text{ با } n \text{ رقم بامعنا}$$

مثال ۱-۴۰: تقریبی از عدد π ارائه دهید که خطای نسبی آن کمتر از 10^{-3} باشد.

حل: با توجه به نامساوی $10^{-3} < 5 \times 10^{-4}$ ، تقریب a از π را چنان ارائه می‌دهیم که دارای چهار رقم بامعناى درست باشد. چرا که در آن صورت، بنابر نتیجه‌ی قبل $\delta(a) < 5 \times 10^{-4}$ و بنابراین $\delta(a) < 10^{-3}$. برای این منظور کافی است a را گرد شده‌ی π تا چهار رقم بامعنا اختیار کنیم، که از آن‌جا $a = 3.142$ به دست می‌آید.

مثال ۱-۴۱: تقریبی از عدد $\sqrt{10}$ ارائه دهید که خطای نسبی آن کمتر از 10^{-4} باشد.

حل:

$$5 \times 10^{-n} < 10^{-4} \Rightarrow 5 < 10^{n-4} \Rightarrow n_{min} = 5$$

بنابراین کافی است تقریبی از $\sqrt{10}$ با پنج رقم بامعناى درست ارائه دهیم. از این رو گردشده‌ی $\sqrt{10}$ را تا پنج رقم بامعناى نویسیم:

$$\sqrt{10} = 3,16227766016... = 3,1623 \text{ (5S)}.$$

قضیه ۱-۱۱: اگر $a > 0$ تقریبی از A باشد به طوری که $\delta(a) \leq 0,5 \times 10^{-n}$ ، آنگاه a حداقل n رقم بامعناى درست دارد.

اثبات: فرض کنید $a = a_m/a_{m-1}/a_{m-2} \dots \times 10^m$ نمایش علمی a باشد. با توجه به این که $a_i \leq 9$ داریم:

$$a < 9,99 \dots \times 10^m = 10 \times 10^m = 10^{m+1},$$

از طرفی می توان نوشت $\delta(a) \approx \frac{e(a)}{a}$ که این به همراه فرض قضیه نتیجه می دهد:

$$e(a) \approx a \times \delta(a) \leq 10^{m+1} \times 0,5 \times 10^{-n} = 5 \times 10^{m-n}.$$

چون ممکن است عددی بزرگتر از n نیز در نامساوی بالا صدق کند، نتیجه می گیریم که a حداقل n رقم بامعناى درست دارد.

در برخی از روش های عددی (مانند حل دستگاه های معادلات خطی) می توان خطای نسبی یک تقریب و در برخی دیگر از روش ها (مانند حل معادلات غیرخطی به روش نیوتن) می توان تعداد ارقام بامعناى درست یک تقریب را به دست آورد. در این موارد، کاربرد دو قضیه ی اخیر بدین صورت است که با داشتن یکی از دو اطلاع خطای نسبی یا تعداد ارقام بامعناى درست یک تقریب، می توان دیگری را نیز برآورد نمود.

روند کردن ماشینی:

فرآیند محاسبات و اعلام نتیجه ی محاسبه توسط ماشین، به کمک عددهای ماشینی صورت می گیرد. اما محدودیت در تعداد ارقام مانع سبب می شود که مجموعه ی اعداد ماشینی تنها زیرمجموعه ای متناهی از اعداد گویای مختوم متعلق به بازه ی $\Omega = [-x_{max}, x_{max}]$ باشند. بنابراین در ماشین های محاسباتی امکان نگهداری و نمایش بسیاری از عددهای حقیقی وجود ندارد. از این رو، عددهای غیر ماشینی با اعداد ماشینی تقریب زده می شوند و فرآیند محاسبات بر روی تقریب آنها صورت می گیرد.

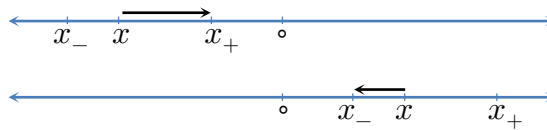
فرض کنید که حالت های سرریز یا زیرریز رخ ندهند و $x \in \Omega - A$ که A مجموعه ی اعداد ماشینی است. در این صورت می توان با تعریف نگاشت نمایش ساز $\text{fl}: \Omega \rightarrow A$ عدد x را با یک عدد ماشینی تقریب زد و به اصطلاح x را قابل نمایش نمود.^۱ عمل نگاشت fl را روند کردن می نامیم.

^۱ توجه کنید که دامنه ی نگاشت fl برابر با Ω است، یعنی نگاشت fl برای اعداد ماشینی نیز تعریف می شود. در حقیقت $x \in A \Leftrightarrow \text{fl}(x) = x$.

ضابطه‌ی نگاشت fl به شکل‌های مختلفی تعریف می‌شود، از جمله روند کردن به سمت صفر (بریدن یا قطع کردن)، روند کردن به نزدیک‌ترین عدد ماشینی (گرد کردن)، روند کردن به بالا (یا به سمت $+\infty$) و روند کردن به پایین (یا به سمت $-\infty$) که در ادامه نحوه‌ی تعریف هر یک از آن‌ها را بیان می‌کنیم. فرض کنید $x \in \Omega - A$ و x_+ و x_- به ترتیب نشان‌دهنده‌ی نزدیک‌ترین اعداد ماشینی بزرگ‌تر و کوچک‌تر از x باشند. در این صورت:

۱- در روند کردن به سمت صفر (قطع کردن)، x با اولین عددی که بر سر راه حرکت آن به سمت صفر قرار دارد تقریب زده می‌شود. به عبارت دیگر، اگر fl_c بیان‌گر این نگاشت باشد داریم:

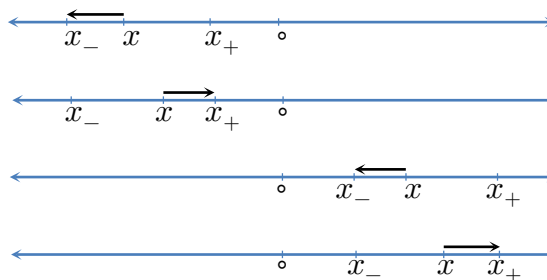
$$\text{fl}_c(x) = \begin{cases} x_+; & x < 0, \\ x_-; & x > 0. \end{cases}$$



۲- در گرد کردن، x با نزدیک‌ترین عدد ماشینی مجاور خود تقریب زده می‌شود. اگر این نگاشت را با fl_r نشان دهیم، داریم:

$$\text{fl}_r(x) = \begin{cases} x_-; & x < x_*, \\ x_+; & x > x_*, \\ x_+ \text{ OR } x_-; & x = x_*. \end{cases}$$

که در آن $x_* = (x_- + x_+)/2$. همان‌طور که ملاحظه می‌کنید در این روش چنانچه $x = x_*$ ، یعنی x دقیقاً بین x_- و x_+ واقع شده باشد، می‌توان به دلخواه گرد کردن را به بالا یا به پایین انجام داد. اما معمولاً در این حالت از قرارداد «گرد کردن به زوج» استفاده می‌شود. به این ترتیب که هر یک از اعداد x_- و x_+ که آخرین رقم سمت راست مانتیس آن زوج باشد، به عنوان تقریب x انتخاب می‌شود. می‌توان دید که با این قرارداد، نیمی از اعداد غیر ماشینی به بالا و نیم دیگر به پایین گرد می‌شوند. به همین دلیل، این قرارداد نتیجه‌ای مانند قرارداد گاوس در گرد کردن عادی دارد. توجه کنید که معمولاً $\text{fl}_r(x)$ را می‌توان مانند گرد کردن معمولی، با توجه به ارقام ناخواسته در مانتیس (اگر مانتیس m رقمی باشد، رقم $m+1$ ام به بعد)، به سادگی تعیین کرد و محاسبه‌ی صریح x_* لازم نیست.



۳- در روند کردن به بالا (یا به سمت $+\infty$), x با اولین عددی که بر سر راه حرکت آن به سمت $+\infty$ قرار دارد تقریب زده می‌شود (یعنی در این حالت، در شکل اخیر جهت تمام پیکان‌ها به سمت راست است). ضابطه‌ی این نگاشت به صورت $\text{fl}_u(x) = x_+$ است.

۴- در روند کردن به پایین (یا به سمت $-\infty$), x با اولین عددی که بر سر راه حرکت آن به سمت $-\infty$ قرار

دارد تقریب زده می‌شود (یعنی در این حالت، در شکل اخیر جهت تمام پیکان‌ها به سمت چپ است). ضابطه‌ی این نگاشت به صورت $\text{fl}_d(x) = x_-$ است.

مثال ۴۲-۱: سیستم نمایش ممیز شناور نرمال $\mathbb{F}(2,3, -1, 2)$ را در نظر بگیرید. در این صورت:

الف) اگر $w = 0.110101 \times 2^1$ ، آن‌گاه:

$$\begin{aligned} \text{fl}_c(w) &= 0.110 \times 2^1, & \text{fl}_r(w) &= 0.111 \times 2^1, \\ \text{fl}_u(w) &= 0.111 \times 2^1, & \text{fl}_d(w) &= 0.110 \times 2^1. \end{aligned}$$

ب) اگر $x = 0.10101 \times 2^{-1}$ ، آن‌گاه:

$$\begin{aligned} \text{fl}_c(x) &= 0.101 \times 2^{-1}, & \text{fl}_r(x) &= 0.101 \times 2^{-1}, \\ \text{fl}_u(x) &= 0.110 \times 2^{-1}, & \text{fl}_d(x) &= 0.101 \times 2^{-1}. \end{aligned}$$

ج) اگر $y = -0.1101 \times 2^1$ ، آن‌گاه:

$$\begin{aligned} \text{fl}_c(y) &= -0.110 \times 2^1, & \text{fl}_r(y) &= -0.110 \times 2^1, \\ \text{fl}_u(y) &= -0.110 \times 2^1, & \text{fl}_d(y) &= -0.111 \times 2^1. \end{aligned}$$

توجه کنید که در این‌جا، $y_+ = -0.110 \times 2^1$ ، $y_- = -0.111 \times 2^1$ بنابرین:

$$\begin{aligned} y_* &= - \left[0.110 \times 2^1 + \frac{1}{4} (0.001 \times 2^1) \right] \\ &= - \left[0.110 \times 2^1 + 0.001 \times 2^0 \right] \\ &= - \left[0.110 \times 2^1 + 0.0001 \times 2^1 \right] \\ &= -0.1101 \times 2^1. \end{aligned}$$

لذا $y_* = y$ و در نتیجه برای گرد کردن از قرارداد «گرد کردن به زوج» استفاده می‌کنیم.

د) اگر $z = 0.1111011 \times 2^1$ ، آن‌گاه:

$$\begin{aligned} \text{fl}_c(z) &= 0.111 \times 2^1, & \text{fl}_r(z) &= 0.100 \times 2^2, \\ \text{fl}_u(z) &= 0.100 \times 2^2, & \text{fl}_d(z) &= 0.111 \times 2^1. \end{aligned}$$

با توجه به مثال فوق و شکل‌های قبل به‌سادگی می‌توان دید که روش گرد کردن برای تقریب اعداد غیرماشینی، نسبت به سایر روش‌ها از دقت بالاتری برخوردار است (از این رو در ماشین‌های محاسباتی، به‌صورت پیش‌فرض از روش گرد کردن استفاده می‌شود). در ادامه، ضمن بررسی خطاهای مطلق و نسبی هر یک از روش‌ها، این مطلب را ثابت می‌کنیم. سیستم ممیز شناور نرمال زیر را برای نمایش اعداد در نظر بگیرید:

$$\mathbb{F} = \left\{ \pm (0.b_1 b_2 \dots b_m)_\beta \times \beta^e \mid b_i \in \{0, 1, \dots, \beta - 1\}, b_1 \neq 0, L \leq e \leq U \right\} \cup \{0\}$$

اگر x عددی حقیقی و غیر ماشینی با نمایش ممیز شناور نرمال $x = \pm (0.b_1 b_2 \dots b_m b_{m+1} \dots)_\beta \times \beta^e$ باشد،

آن‌گاه خطای مطلق بریدن آن برابر با قدرمطلق اختلاف x و $\text{fl}_c(x)$ و خطای مطلق گرد کردن آن برابر با قدرمطلق اختلاف x و $\text{fl}_r(x)$ است. بنابراین:

$$|x - \text{fl}_c(x)| = (\underbrace{\circ \dots \circ}_{\text{مرتبه } m} b_{m+1} \dots)_{\beta} \times \beta^e \leq (\underbrace{\circ \dots \circ}_{\text{مرتبه } m-1} 1)_{\beta} \times \beta^e = \beta^{e-m},$$

$$|x - \text{fl}_r(x)| \leq \frac{1}{\beta} (x_+ - x_-) = \frac{1}{\beta} \beta^{e-m}.$$

اکنون با توجه به این‌که $|x| \geq (\underbrace{\circ \dots \circ}_{\text{مرتبه } m-1} 1)_{\beta} \times \beta^e = \beta^{e-1}$ کران‌های بالای زیر را برای خطای نسبی $\text{fl}_c(x)$ و $\text{fl}_r(x)$ خواهیم داشت:

$$\frac{|x - \text{fl}_c(x)|}{|x|} \leq \frac{\beta^{e-m}}{\beta^{e-1}} = \beta^{1-m},$$

$$\frac{|x - \text{fl}_r(x)|}{|x|} \leq \frac{1}{\beta} \frac{\beta^{e-m}}{\beta^{e-1}} = \frac{1}{\beta} \beta^{1-m}.$$

همان‌طور که ملاحظه می‌کنید، کران بالای خطای نسبی به اندازه‌ی عدد (نما) بستگی ندارد و تنها به ارقام مانتیس وابسته است. بنابراین در سیستم ممیز شناور، اگرچه فاصله‌ی دو عدد متوالی بزرگ، زیاد است اما تأثیری در خطای نسبی ندارد. این یک مزیت برای اعداد ممیز شناور به شمار می‌رود.

در مورد خطاهای مطلق و نسبی روند کردن به بالا و پایین توجه کنید که اگر x عددی منفی (مثبت) باشد، آن‌گاه روند شده‌ی x به بالا (پایین) و بریده شده‌ی آن هر دو یکسان هستند. بنابراین می‌توان از تقارن اعداد ممیز شناور نسبت به مبدأ مختصات نتیجه گرفت که کران‌هایی که برای خطای مطلق و نسبی در حالت بریدن به دست آمدند، برای گرد کردن به سمت بالا و پایین نیز برقرارند. یعنی:

$$|x - \text{fl}_u(x)| \leq \beta^{e-m}, \quad \frac{|x - \text{fl}_u(x)|}{|x|} \leq \beta^{1-m},$$

$$|x - \text{fl}_d(x)| \leq \beta^{e-m}, \quad \frac{|x - \text{fl}_d(x)|}{|x|} \leq \beta^{1-m}.$$

اپسیلون ماشین و روند واحد

تعریف ۱-۱۴: به فاصله‌ی بین عدد یک و عدد ممیز شناور بلافاصله پس از آن در سیستم نمایش $\mathbb{F}(\beta, m, L, U)$ اپسیلون ماشین می‌گوییم و آن را با نماد ϵ_M یا eps نشان می‌دهیم. با توجه به این‌که عدد یک در این سیستم به صورت $1 = (\underbrace{\circ \dots \circ}_{\text{مرتبه } m-1} 1)_{\beta} \times \beta^1$ نمایش داده می‌شود، این فاصله عبارت است از:

$$\epsilon_M = 1_+ - 1 = (\underbrace{\circ \dots \circ}_{\text{مرتبه } m-2} 1)_{\beta} \times \beta^1 - (\underbrace{\circ \dots \circ}_{\text{مرتبه } m-1} 1)_{\beta} \times \beta^1 = \beta^{1-m}.$$

همان‌طور که ملاحظه می‌کنید اپسیلون ماشین تنها به مبنا و تعداد ارقام مانتیس بستگی دارد.

اکنون فرض کنید x عددی مخالف صفر و متعلق به Ω باشد و ننگاشت نمایش ساز fl را در نظر بگیرید. قرار می‌دهیم:

$$\delta = \frac{\text{fl}(x) - x}{x},$$

در این صورت $\text{fl}(x) = x(1 + \delta)$ و $|\delta| \leq u$ که در آن:

$$u = \begin{cases} \beta^{1-m} = \epsilon_M; & \text{اگر } \text{fl}(x) = \text{fl}_c(x) \\ \frac{1}{4}\beta^{1-m} = \frac{\epsilon_M}{4}; & \text{اگر } \text{fl}(x) = \text{fl}_r(x) \end{cases}$$

عدد u در این جا، روند واحد نامیده می‌شود و علاوه بر مبنا و تعداد ارقام مانتیس به نوع ننگاشت نمایش ساز بستگی دارد. بنابراین، اگر u روند واحد وابسته به ننگاشت نمایش ساز fl باشد، به ازای هر $x \in \Omega, x \neq 0$ می‌توان نوشت:

$$\boxed{\text{fl}(x) = x(1 + \delta); \quad |\delta| \leq u}$$

تعریف $\text{fl}(x)$ بدین صورت و استفاده از آن، نخستین بار توسط ویلکینسون* انجام شده است. این رابطه هم‌چنین در تحلیل خطای تولید شده توسط اعمال حسابی در کامپیوتر بسیار اهمیت دارد. نکته‌ی مهمی که از این رابطه استنتاج می‌شود این است که هرچه ϵ_M کوچک‌تر باشد، u و در نتیجه خطای نسبی روند کردن کمتر است. یعنی کوچک بودن اندازه‌ی اپسیلون ماشین به‌گونه‌ای بیانگر دقت آن ماشین است و چون اپسیلون ماشین تنها به دو عامل مبنا و تعداد ارقام مانتیس بستگی دارد، با افزایش تعداد ارقام مانتیس می‌توان دقت ماشین را بالا برد. به همین دلیل است که وقتی دقت بالاتری در محاسبات مد نظر است، از دقت مضاعف استفاده می‌شود. در دقت مضاعف، برای نمایش یک عدد ممیز شناور از دو کلمه‌ی کامپیوتر استفاده می‌شود و در این حالت مانتیس حداقل دو برابر بیت نسبت به دقت معمولی در اختیار دارد. برخی از کامپیوترها قادرند اعداد را حتی تا ۴ برابر دقت معمولی ذخیره کنند.

نکته‌ی دیگر این‌که با توجه به تعریف $\text{fl}(x)$ و با قرار دادن $x = 1$ داریم:

$$1 = \text{fl}(1) = 1 + \delta; \quad |\delta| \leq u,$$

بنابراین:

$$\text{fl}(1 + \delta) = \text{fl}(\text{fl}(1)) = \text{fl}(1) = 1; \quad |\delta| \leq u,$$

یعنی در یک سیستم نمایش ممیز شناور نرمال به‌همراه ننگاشت نمایش ساز fl ، روند واحد بزرگ‌ترین مقداری است که اگر با عدد ۱ جمع شود، هم‌چنان عدد ۱ به عنوان حاصل جمع نمایش داده می‌شود. به عبارت دیگر:

$$\boxed{u = \sup\{\delta | \text{fl}(1 + \delta) = 1\}}$$

نتیجه: اگر $|\delta| \leq u$ ، آنگاه نمایش $1 + \delta$ برابر با ۱ است. هم‌چنین با توجه به این‌که

$$x + y = x\left(1 + \frac{y}{x}\right)$$

اگر داشته باشیم $u \leq \left|\frac{y}{x}\right|$ آنگاه نمایش $x + y$ با نمایش x یکسان خواهد بود.

محاسبه‌ی تقریبی روند واحد با استفاده از برنامه‌ی زیر امکان‌پذیر است. این برنامه کوچک‌ترین عدد ممیز شناور x را محاسبه می‌کند که به ازای آن $\text{fl}(1+x) > 1$. همچنین از این برنامه برای تقریب eps ماشین می‌توان استفاده کرد. توجه کنید که eps ماشین و روند واحد دو مفهوم مجزا هستند. با این وجود، اختلاف این مقادیر در بسیاری از کامپیوترها به قدری کوچک است که معمولاً تمایزی بین آن‌ها گذارده نمی‌شود.

```
x := 1;
while 1 + x > 1
    x := x/2;
end
```

۱-۶ تولید و انتشار خطا

همان‌طور که پیش‌تر گفته شد در ماشین‌های محاسباتی عددهای غیر ماشینی با اعداد ماشینی تقریب زده می‌شوند و فرآیند محاسبات بر روی تقریب آن‌ها صورت می‌گیرد. در این بخش می‌خواهیم این مطلب را به‌طور دقیق‌تر بررسی کنیم و ببینیم که خطاهای موجود در مراحل یک فرآیند محاسباتی چگونه به‌وجود می‌آیند و آیا می‌توان برای کاهش آن‌ها یا حداقل جلوگیری از افزایش آن‌ها در مراحل بعدی محاسبات کاری انجام داد یا خیر؟ برای این منظور ابتدا به‌طور مختصر نحوه‌ی انجام محاسبات بر روی اعداد ممیز شناور را تشریح می‌کنیم و سپس به بررسی خطای موجود در حاصل یک محاسبه‌ی عددی می‌پردازیم.

حساب ممیز شناور

انجام محاسبات بر روی اعداد ممیز شناور را حساب ممیز شناور می‌نامند. ماشینی را در نظر بگیرید که از سیستم ممیز شناور نرمال $\mathbb{F}(\beta, m, L, U)$ استفاده می‌کند. همچنین فرض کنید $*$ بیان‌گر یکی از چهار عمل حسابی اصلی باشد، یعنی $*$ $\in \{+, -, \times, \div\}$. حال اگر X و Y دو عدد حقیقی باشند و بخواهیم $X * Y$ را به کمک این ماشین محاسبه کنیم، آن‌گاه مراحل زیر توسط ماشین انجام می‌شود:

۱- ممکن است X یا Y ماشینی نباشند، بنابراین در نخستین مرحله آن‌ها با $x = \text{fl}(X)$ و $y = \text{fl}(Y)$ جایگزین می‌شوند. لذا محاسبه‌ی $X * Y$ به $x * y$ تغییر می‌یابد.

۲- $x * y$ با استفاده از حساب ممیز شناور و با دقت مضاعف محاسبه می‌شود (در عمل، بسیاری از کامپیوترها اعمال حسابی را در ثبات‌های خاص که تعداد بیت‌های بیشتری نسبت به اعداد معمولی ماشین دارند انجام می‌دهند. این بیت‌ها که اصطلاحاً بیت پشتیبان نامیده می‌شوند به اعداد اجازه می‌دهند که موقتاً دقت بیشتری داشته باشند و تعداد آن‌ها در حالت کلی از ماشینی به ماشین دیگر متفاوت است). در این‌جا نیز ممکن است که نتیجه‌ی عمل عددی ماشینی نباشد. لذا پس از انجام عمل حسابی، نتیجه‌ی عمل با استفاده از نگاهت نمایش ساز ماشین به یک عدد ماشینی تبدیل می‌شود. در نتیجه $y * x$ با $y \otimes x = \text{fl}(x * y)$ جای‌گزین می‌شود.