



یادداشت های درسی آنالیز عددی پیشرفته ©

جواد فرضی^۱
دانشگاه صنعتی سهند، بخش ریاضی

بخش اول : خطا و انتشار آن

فهرست مندرجات

۳	خطا	۱
۳	نمایش اعداد	۱.۱
۵	حساب ممیز سیار و خطای گرد کردن	۲.۱
۹	انتشار خطا	۲
۹	مقدمه	۱.۲
۱۰	عدد حالت	۲.۲
۱۲	پایداری الگوریتم و خطای ذاتی	۳.۲
۱۷	جمع سری و واگرایی عددی	۴.۲
۱۸	مسائل	۳

^۱farzi@sut.ac.ir

۱ خطا

تعیین دقت نتایج عددی یکی از اهداف بسیار مهم در آنالیز عددی است. با مطالعه در این زمینه می توان چند نوع خطا که در محدود کردن دقت نتایج موثرند را برشمرد:

(۱) خطاهای داده های ورودی (errors in the input data)

(۲) خطاهای گرد کردن (roundoff errors)

(۳) خطاهای تقریب (approximation errors)

خطاهای ورودی از کنترل محاسبات خارج است. خطاهای گرد کردن مربوط به محاسبه با اعدادی است که نمایش آنها در ماشین به تعداد ارقام متناهی محدود شده است که گریزی از این واقعیت نیست. برای نوع سوم مساله P را در نظر می گیریم. بسیاری از روشها حتی با این فرض که محاسبات بدون خطای گرد کردن باشد به جواب دقیق این مساله منجر نمی شوند و به جای مساله P ساده تر \tilde{P} را حل می کنند. برای مثال مساله P، محاسبه سری نامتناهی

$$e = 1 + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \dots \quad (1)$$

را می توان با مساله \tilde{P} ، محاسبه سری تا یک تعداد متناهی از جملات سری جایگزین کرد. چنین خطای تقریبی خطای قطع (truncation error) نام دارد. البته خطای قطع در زمینه های دیگری مانند نمایش متناهی عدد نیز مطرح می شود. بسیاری از مسائل تقریب P با گسسته سازی (discretizing) مساله اصلی P بدست می آیند. انتگرالها با مجموع های متناهی و معادلات دیفرانسیل با معادلات تفاضلی تقریب می شوند. در چنین حالتی خطای تقریب با عنوان خطای گسسته سازی مطرح می شود. در اینجا تاثیر خطاهای ورودی و خطای گرد کردن بر روی نتایج عددی را بررسی می کنیم. به طور طبیعی خطای تقریب را در کنار روش تقریب ارائه شده در فصول بعدی بررسی خواهیم کرد.

۱.۱ نمایش اعداد

نمایش عدد x در مبنای ۲

$$x = \pm(\alpha_n 2^n + \alpha_{n-1} 2^{n-1} + \dots + \alpha_0 2^0 + \alpha_{-1} 2^{-1} + \alpha_{-2} 2^{-2} + \dots)$$
$$\alpha_i = 0 \quad \text{or} \quad 1.$$

برای مثال نمایش ۱۸.۵ در مبنای دو به صورت زیر است:

$$18.5 = 1 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0 + 1 \times 2^{-1}$$

مثال $4 = 3.999\dots$ نشان می دهد که نمایش ده دهی یک عدد ممکن است منحصر به فرد نباشد. همین در مورد نمایش دودویی هم درست است. برای دوری از چنین ابهامهایی نمایش متناهی را در نظر خواهیم گرفت. همچنین مثالهایی وجود دارد که نمایش یک عدد در یک دستگاه متناهی است ولی نمایش همان عدد در دستگاه دیگر متناهی نیست.

اعداد ماشینی ماشین‌های دیجیتالی با توجه به محدودیت حافظه در ذخیره کردن و محاسبات داخلی خود مکانهای ثابت و متناهی را برای اعداد اختصاص می‌دهند که به آن کلمه گفته می‌شود. این عدد n در برخی ماشین‌ها دو یا سه برابر قرار داده می‌شود. از این طول n به شیوه‌های مختلف می‌توان برای نمایش اعداد استفاده کرد.

• نمایش ممیز ثابت

در این نمایش تعداد n_1 رقم پیش و تعداد n_2 رقم پس از ممیز اعشار (دودویی) وجود دارد که $n = n_1 + n_2$. برای مثال با $n = 10$ ، $n_1 = 4$ و $n_2 = 6$ داریم

$$\begin{array}{l} 30.421 \rightarrow \overbrace{0030}^{n_1} \overbrace{421000}^{n_2} \\ 0.0437 \rightarrow \overbrace{0000}^{n_1} \overbrace{043700}^{n_2} \end{array}$$

• نمایش ممیز سیار

$$x = a \times 10^b (x = a \times 2^b), \quad |a| < 1 \quad (2)$$

در این نمایش عدد صحیح b جایگاه ممیز را نسبت به مانتیس a نشان می‌دهد. برای مثال نمایش عدد 30.421 به صورت ممیز سیار عبارتست از 30.421×10^2 و عدد 0.0437 را در مبنای دو با ممیز سیار به صورت 0.100101×2^{10} می‌توان نمایش داد.

در هر ماشینی تعداد مکانهای ثابت و متناهی t و e ، $n = t + e$ ، در حافظه به ترتیب برای نمایش مانتیس و توان وجود دارد. برای مثال با $t = 4$ و $e = 2$ نمایش ممیز سیار 5420 به صورت 0.5420×10^4 است. در صورتی که رقم اول مانتیس ناصفر باشد نمایش ممیز سیار به صورت نرمال است. در این صورت در رابطه (2) داریم $|a| \geq 10^{-1}$ ($|a| \geq 2^{-1}$). در این صورت ارقام با معنی یک عدد عبارتست از ارقام مانتیس بدون شمارش صفرهای پیشرو.

زیر مجموعه $A \subseteq R$ از اعداد حقیقی که قابل نمایش در یک ماشین مفروض باشند اعداد ماشینی نامیده می‌شوند. در یک ماشین ۳۲ و ۶۴ بیتی مکانهای حافظه برای ذخیره کردن یک عدد به صورت ممیز سیار (2)، به ترتیب، به این صورت اختصاص می‌یابد:

۳۲ بیتی		۶۴ بیتی	
۱	۱	۱	۱
۱	۱	۱	۱
۱۰	۷	۱۰	۷
۵۲	۲۳	۵۲	۲۳

۲.۱ حساب ممیزسیار و خطای گرد کردن

مجموعه اعداد A که قابل نمایش در ماشین باشند متناهی اند. این سؤال پیش می آید که چگونه می توان یک عدد $x \notin A$ که عدد ماشینی نیست را با یک عدد ماشینی مانند $g \in A$ نمایش داد. این سؤال نه تنها در خواندن اعداد از سوی رایانه، بلکه به هنگام محاسبات و نمایش نتایج میانی در رایانه نیز پیش می آید. به عبارتی با مثالهایی می توان نشان داد که مجموعه اعداد ماشینی نسبت به اعمال چهارگانه ریاضی بسته نیست. مثال زیر وقوع خطا به هنگام خواندن را نشان می دهد:

$$\frac{1}{10} = (0.0001100110011001\dots)_2 \quad (3)$$

در انتهای این بخش مثالهایی از حالت های دیگر ارائه می شود. طبیعی است بخواهیم که تقریب عدد $x \notin A$ با یک عدد ماشینی $rd(x) \in A$ در رابطه زیر صدق کند

$$|x - rd(x)| \leq |x - g| \quad \forall g \in A \quad (4)$$

در بیشتر مواقع چنین تقریبی را می توان با گرد کردن بدست آورد.
مثال (۴) $t = 4$

$$\begin{aligned} rd(0.14285 \times 10^0) &= 0.1429 \times 10^0, \\ rd(3.14159 \times 10^0) &= 0.3142 \times 10^1, \\ rd(0.142842 \times 10^2) &= 0.1428 \times 10^2. \end{aligned}$$

به طور کلی، برای بدست آوردن $rd(x)$ در یک رایانه t -رقمی به صورت زیر عمل می کنیم: در آغاز $x \notin A$ را به صورت نرمال شده $x = a \times 10^b$ که $|a| \geq 10^{-1}$ نمایش می دهیم. فرض کنید نمایش اعشاری $|a|$ به صورت زیر باشد

$$|a| = 0.\alpha_1\alpha_2\dots\alpha_i\alpha_{i+1}\dots, \quad 0 \leq \alpha_i \leq 9, \quad \alpha_1 \neq 0.$$

در این صورت با معرفی

$$a' := \begin{cases} 0.\alpha_1\alpha_2\dots\alpha_t, & 0 \leq \alpha_{t+1} \leq 4, \\ 0.\alpha_1\alpha_2\dots\alpha_t + 10^{-t}, & \alpha_{t+1} \geq 5. \end{cases}$$

تقریب مورد نظر به صورت زیر ارائه می شود

$$\tilde{rd}(x) := \text{sign}(x).a' \times 10^b.$$

با توجه به $|a| \geq 10^{-1}$ خطای نسبی $\tilde{rd}(x)$ در نامساوی زیر صدق می کند:

$$\left| \frac{\tilde{rd}(x) - x}{x} \right| \leq \frac{5 \times 10^{-(t+1)}}{|a|} \leq 5 \times 10^{-t}.$$

با فرض $\text{eps} := 5 \times 10^{-t}$ این عبارت را می توان به صورت زیر نوشت

$$\tilde{rd}(x) = x(1 + \varepsilon), \quad |\varepsilon| \leq \text{eps}. \quad (5)$$

کمیت $eps := 5 \times 10^{-t}$ دقت ماشین نامیده می‌شود. در دستگاه دودویی به طور مشابه می‌توان $\tilde{rd}(x)$ را تعریف کرد: با شروع از $x = a \times 2^b$ که $1 < |a| \leq 2^{-1}$ و نمایش دودویی $|a|$ داریم

$$|a| = 0.\alpha_1\alpha_2\dots\alpha_t\alpha_{t+1}\dots, \quad \alpha_i = 0 \text{ or } 1, \quad \alpha_1 = 1.$$

با تشکیل a'

$$a' := \begin{cases} 0.\alpha_1\alpha_2\dots\alpha_t, & \alpha_{t+1} = 0, \\ 0.\alpha_1\alpha_2\dots\alpha_t + 2^{-t}, & \alpha_{t+1} = 1. \end{cases}$$

$$\tilde{rd}(x) := \text{sign}(x).a' \times 2^b.$$

دوباره رابطه (5) با $eps := 2^{-t}$ برقرار است.

اگر $\tilde{rd}(x) \in A$ یک عدد ماشینی باشد، آنگاه \tilde{rd} دارای خاصیت (4) است که روند گرد کردن درست را نشان می‌دهد، در این حالت داریم

$$\tilde{rd}(x) := rd(x) \quad \forall x \in A.$$

با توجه به اینکه مکانهای موجود برای نمایش توان تعداد محدود e است، همواره اعدادی وجود دارند که $x \notin A$ و $\tilde{rd}(x) \notin A$ (مثال $(e=2, t=4)$)

$$(a) \quad \tilde{rd}(0.31794 \times 10^{11}) = 0.3179 \times 10^{11} \notin A.$$

$$(b) \quad \tilde{rd}(0.99997 \times 10^{99}) = 0.1000 \times 10^{100} \notin A.$$

$$(c) \quad \tilde{rd}(0.012345 \times 10^{-99}) = 0.1235 \times 10^{-100} \notin A.$$

$$(d) \quad \tilde{rd}(0.54321 \times 10^{-110}) = 0.5432 \times 10^{-110} \notin A.$$

در حالت (a) و (b) نما یک مقدار بسیار بزرگ است که در فضای اختصاص داده شده نمی‌گنجد: این دو مورد مثالهایی از سرریز نما (exponent overflow) هستند. در حالت (b) سرریز نما بعد از گرد کردن اتفاق می‌افتد. حالتی (c) و (d) مثالهای از زیرریز نما (exponent overflow) هستند. در حالتی (c) و (d) با تعریف زیر می‌توان از زیرریز نما جلوگیری کرد

$$(c) \quad rd(0.012345 \times 10^{-99}) = 0.0123 \times 10^{-99} \in A,$$

$$(d) \quad rd(0.54321 \times 10^{-110}) = 0 \in A.$$

در این صورت rd در (5) صدق نخواهد کرد. چون خطای نسبی $rd(x)$ ممکن است از eps تجاوز کند. با صرف نظر از سرریز و زیرریز، گرد کردن را به صورت زیر تعریف می‌کنیم

$$rd(x) := \tilde{rd}(x).$$

با در نظر گرفتن مقیاس های مناسب و کنترل داده های ورودی و بررسی های مناسب به هنگام محاسبات با این مشکلات مواجه نمی شویم. بنابراین با توجه به اینکه این موارد در عمل چندان اتفاق نمی افتد در ادامه بحث فرض می کنیم $e = \infty$. با چنین شرایطی، $rd := \tilde{rd}$ دستوری برای گرد کردن فراهم می کند که در شرایط زیر صدق می کند

$$rd : R \rightarrow A, \quad (6)$$

$$\tilde{rd}(x) = x(1 + \varepsilon), \quad |\varepsilon| \leq eps \quad \forall x \in R. \quad (7)$$

مشاهده کردیم که در صورتی که نتیجه اعمال حسابی $x \pm y$, $x \times y$ و x/y لزوماً اعداد ماشینی نیستند، حتی اگر اعداد x و y اعداد ماشینی باشند. بنابراین انتظار نداریم که اعمال حسابی را در یک ماشین به دقت انجام شود. در اینجا می توان اعمال جایگزین $+^*$, $-^*$, \times^* و $/^*$ که اعمال ممیز سیار نام دارند را در نظر گرفت که اعمال حسابی را تا حد ممکن تقریب می کنند. بر اساس نگاشت گرد کردن rd می توان این اعمال را به صورت زیر تعریف کرد

$$\begin{aligned} x +^* y &:= rd(x + y), \\ x -^* y &:= rd(x - y), \quad \forall x, y \in A \\ x \times^* y &:= rd(x * y), \\ x /^* y &:= rd(x / y), \end{aligned} \quad (8)$$

لذا از (6) نتیجه می شود

$$\begin{aligned} x +^* y &:= (x + y)(1 + \varepsilon_1), \\ x -^* y &:= (x - y)(1 + \varepsilon_2), \quad |\varepsilon_i| \leq eps. \\ x \times^* y &:= (x * y)(1 + \varepsilon_3), \\ x /^* y &:= (x / y)(1 + \varepsilon_4), \end{aligned} \quad (9)$$

منظور از $fl(E)$ محاسبه عبارت E با حساب ممیز سیار است. اگر E یک عبارت مقدماتی باشد، $fl(E) = rd(E)$. اعمال ممیز سیار در قوانین شناخته شده اعمال حسابی صدق نمی کنند. برای مثال،

$$x +^* y = x \quad \text{if} \quad |y| < \frac{eps}{B}|x|, \quad x, y \in A,$$

که در آن B مبنای دستگاه اعداد است. دقت ماشین را می توان به صورت کوچکترین عدد مثبت ماشینی g که $1 +^* g > 1$ تعریف کرد:

$$eps = \min\{g \in A \mid 1 +^* g > 1, g > 0\}.$$

اعمال ممیز سیار لزوماً جابجایی و توزیع پذیر نیستند.
مثال (8) با فرض

$$\begin{aligned} a &:= 0.23371258 \times 10^{-4}, \\ b &:= 0.33678429 \times 10^2, \\ c &:= -0.33677811 \times 10^2, \end{aligned}$$

داریم

$$\begin{aligned} a +^* (b +^* c) &= 0.23371258 \times 10^{-4} +^* 0.61800000 \times 10^{-3}, \\ &= 0.64137126 \times 10^{-3}, \\ (a +^* b) +^* c &= 0.33678452 \times 10^2 -^* 0.33677811 \times 10^2, \\ &= 0.64100000 \times 10^{-3}. \end{aligned}$$

مقدار دقیق عبارتست از

$$a + b + c = 0.641371258 \times 10^{-3}.$$

در تفریق دو عدد $x, y \in A$ با علامت یکسان باید مواظب خنثی سازی (cancellation) بود. این وقتی پیش می آید که x و y یک یا چند رقم پیشرو مشترک با یک نما داشته باشند، برای مثال

$$\begin{aligned} x &= 0.315876 \times 10^1, \\ y &= 0.314289 \times 10^1. \end{aligned}$$

تفریق موجب از بین رفتن ارقام پیشرو مشترک می شود. مقدار دقیق $x - y$ یک عدد ماشینی است، لذا خطای گرد کردن جدیدی پیش نمی آید: $x -^* y = x - y$. بنابراین، تفریق در حالت خنثی سازی یک عمل کاملاً بی خطر است. با این حال، خنثی سازی به دلیل انتشار خطاهای قبلی بسیار خطرناک است.

اگر نمایش عدد ماشینی نرمال $\pm 0.\alpha_1\alpha_2\dots\alpha_s \times \beta^e$ را در ماشین های مختلف در نظر بگیریم با فرض $m \leq e \leq M$ جدول زیر را خواهیم داشت

ماشین	β	s	m	M	eps
VAX	2	24	-128	127	6.0×10^{-08}
VAX	2	56	-128	127	1.4×10^{-17}
CRAY-1	2	48	-16384	16383	3.6×10^{-15}
IBM 3081	16	6	-64	63	9.5×10^{-07}
IBM 3081	16	14	-64	63	2.2×10^{-16}
IEEE					
single	2	24	-125	128	6.0×10^{-08}
double	2	53	-1021	1024	1.1×10^{-16}

۲ انتشار خطا

۱.۲ مقدمه

در عمل وقتی توابع مرکب از چندین تابع را با داده‌های اولیه خطادار محاسبه می‌کنیم در هر مرحله‌ای خطایی مرکب می‌شویم تا جایی که در بعضی از موارد نتیجه حاصل به طور کامل غلط است. بنابراین لازم است قبل از محاسبه یک تابع از رفتار آن آگاه باشیم. اینجاست که مفهوم ناپایداری عددی ظاهر می‌شود؛ با یک بیان ساده و نه عبارت دقیق، گوییم یک روند عددی ناپایدار است در صورتی که خطاهای کوچک حاصل از یک مرحله این روند محاسباتی دقت محاسبات در مراحل بعدی را به صورت محسوس کاهش دهد.

مثال. دنباله عددی زیر را در نظر می‌گیریم

$$\begin{cases} x_0 = 1, & x_1 = \frac{1}{3} \\ x_{n+1} = \frac{13}{3}x_n - \frac{4}{3}x_{n-1}, & n \geq 1. \end{cases} \quad (10)$$

با حل این رابطه بازگشتی داریم

$$x_n = \left(\frac{1}{3}\right)^n \quad (11)$$

اگر برای محاسبه این دنباله از (۱۰) استفاده کنیم در یک ماشین ۳۲ بیتی نتایج زیر بدست می‌آید که بعضی از نتایج بدست آمده به صورت کامل نادرست است.

$x_0 = 1.0000000$		$x_8 = 0.0003757$	رقم بامعنا
$x_1 = 0.3333333$	رقم بامعنا	$x_9 = 0.0009437$	
$x_2 = 0.1111112$	رقم بامعنا	$x_{10} = 0.0035887$	
$x_3 = 0.0370373$	رقم بامعنا	$x_{11} = 0.0142927$	
$x_4 = 0.0123466$	رقم بامعنا	$x_{12} = 0.0571502$	
$x_5 = 0.0041187$	رقم بامعنا	$x_{13} = 0.2285939$	
$x_6 = 0.0013857$	رقم بامعنا	$x_{14} = 0.9143735$	
$x_7 = 0.0005131$	رقم بامعنا	$x_{15} = 3.657493$	نادرست! با خطای نسبی 10^8

این نتایج را می‌توان با کد ساده زیر در MATLAB محاسبه کرد. براساس این برنامه $x_{29} = -14.273$:

```
x(1) = 1;
x(2) = 1/3;
for n=2:30
    x(n+1) = 13/3*x(n)-4/3*x(n-1);
end
x'
```

این الگوریتم ناپایدار است. خطای موجود در x_n در محاسبه x_{n+1} با عامل $\frac{1}{3}$ افزایش می‌یابد. یعنی خطا در x_1 با عامل $(\frac{1}{3})^{14}$ به x_{15} انتشار می‌یابد. با توجه به اینکه خطای مطلق x_1 تقریباً برابر 10^{-8} است و در x_{15} خطا به عامل $(\frac{1}{3})^{14}$ که تقریباً برابر است با 10^9 ضرب می‌شود لذا، خطا در x_{15} تقریباً برابر است با 10 . بنابراین، نتایج بدست آمده از نرم افزارهای پیشرفته در صورتی قابل اعتماد است که از نظر تئوری رفتار خطا بررسی شده باشد.

مثال. مثال دیگر برای ناپایداری عددی عبارتست از محاسبه انتگرال

$$y_n = \int_0^1 x^n e^x dx, \quad n \geq 0 \quad (12)$$

اگر از انتگرال گیری جزء به جزء برای y_{n+1} استفاده کنیم داریم:

$$y_{n+1} = e - (n+1)y_n \quad (13)$$

با استفاده از این رابطه و $y_0 = e - 1$ داریم:

$$y_1 = e - y_0 = e - (e - 1) = 1$$

با شروع از $y_1 = 1$ بر روی یک ماشین ۳۲ بیتی مقادیر y_2, y_3, \dots, y_{15} با رابطه (۱۳) محاسبه می‌کنیم. بعضی از این نتایج عبارتند از

$$y_2 = .7182817$$

$$y_{11} = 1.422453$$

$$y_{15} = 39711.43$$

این نتایج نمی‌تواند درست باشد. از (۱۲) واضح است که $y_1 > y_2 > \dots > 0$ و $\lim_{n \rightarrow \infty} y_n = 0$.

۲.۲ عدد حالت

این بحث را با یک سوال آغاز می‌کنیم. در محاسبه تابع $f(x)$ اگر x کمی پریشیده شود تاثیر آن بر $f(x)$ چه خواهد بود؟ اگر در این سوال هدف خطای مطلق باشد از قضیه مقدار میانگین داریم:

$$f(x+h) - f(x) = f'(\xi)h \approx hf'(x) \quad (14)$$

بنابراین در صورتی که $f'(x)$ زیاد بزرگ نباشد تاثیر این پریشیدگی روی $f(x)$ کم است. با این وجود، خطای نسبی وضعیت را بهتر نمایان می‌کند. وقتی x به میزان h پریشیده می‌شود $f(x)$ به $f(x+h)$ پریشیده می‌شود و خطای نسبی برابر است با

$$\frac{f(x+h) - f(x)}{f(x)} \approx \frac{hf'(x)}{f(x)} = \left[\frac{xf'(x)}{f(x)} \right] \left(\frac{h}{x} \right) \quad (15)$$

عامل $\frac{xf'(x)}{f(x)}$ عدد حالت این مساله است.

مثال. عدد حالت برای محاسبه تابع $f(x) = \sin^{-1} x$ چیست؟

$$\frac{xf'(x)}{f(x)} = \frac{x}{\sqrt{1-x^2} \sin^{-1} x} \quad (16)$$

اگر x نزدیک ۱ باشد، $\sin^{-1} x \approx \pi/2$ و عدد حالت به بی‌نهایت میل خواهد کرد. لذا، در نزدیکی $x = 1$ خطاهای نسبی کوچک به خطاهای نسبی بزرگ در $\sin^{-1} x$ منجر خواهد شد.

مثال. در این مثال مساله پیدا کردن ریشه یک تابع را در نظر می‌گیریم. این مساله در بخش‌های بعدی به صورت الگوریتمی بررسی خواهد شد. فرض کنید توابع f و g در همسایگی r متعلق به رده توابع C^2 باشند، که r ریشه f است. فرض می‌کنیم r یک ریشه ساده باشد یعنی؛ $f'(r) \neq 0$. اگر تابع f را به $F \equiv f + \varepsilon g$ پریشیده کنیم ریشه جدید در کجا خواهد بود؟ فرض کنیم ریشه جدید $r + h$ باشد؛ می‌خواهیم یک فرمول تقریبی برای h بدست آوریم. پریشیدگی h در معادله $F(r + h) = 0$ ، به طور معادل در رابطه زیر، صدق می‌کند:

$$f(r + h) + \varepsilon g(r + h) = 0$$

با توجه به اینکه f و g در C^2 هستند می‌توان بسط تیلور را برای $F(r + h) = 0$ نوشت:

$$[f(r) + hf'(r) + \frac{1}{2}h^2 f''(\xi)] + \varepsilon[g(r) + hg'(r) + \frac{1}{2}h^2 g''(\eta)] = 0$$

با صرف نظر از جملات شامل h^2 و توجه به این نکته که $f(r) = 0$ داریم

$$h \approx -\varepsilon \frac{g(r)}{f'(r) + \varepsilon g'(r)} \approx -\varepsilon \frac{g(r)}{f'(r)}$$

برای بررسی بهتر مثال زیر را در نظر می‌گیریم

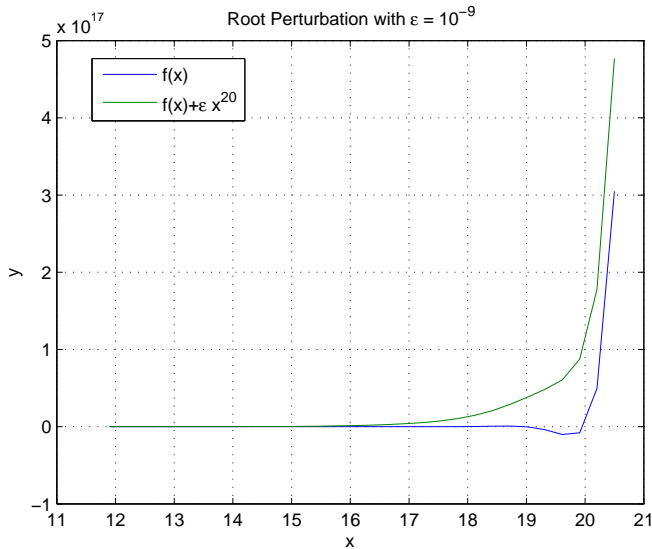
$$f(x) = \prod_{k=1}^{20} (x - k) = (x - 1)(x - 2) \dots (x - 20)$$

$$g(x) = x^{20}$$

به وضوح ریشه‌های f برابر است با $1, 2, \dots, 20$. اگر f را با $f + \varepsilon g$ جایگزین کنیم ریشه $r = 20$ چه تاثیری روی ریشه $r = 20$ خواهد داشت؟ جواب عبارتست از

$$h \approx -\varepsilon \frac{g(20)}{f'(20)} = -\varepsilon \frac{20^{20}}{19!} \approx -10^9 \varepsilon$$

بنابراین، یک تغییر ε در ضریب x^{20} در $f(x)$ یک پریشیدگی به اندازه $10^9 \varepsilon$ در ریشه 20 ایجاد می‌کند. در نتیجه، ریشه‌های این چندجمله‌ای به پریشیدگی در ضرایب بسیار حساس هستند.



شکل فوق انحراف ریشه $x = 20$ را با پریشیدگی $\varepsilon = 10^{-9}$ نشان می‌دهد. با این پریشیدگی تغییر قابل توجهی در مکان این ریشه بوجود می‌آید.

۳.۲ پایداری الگوریتم و خطای ذاتی

تابع مفروض $y = \phi(x)$, $\phi: D \subseteq R^n \rightarrow R^m$ را در نظر می‌گیریم و آن را به دنباله‌ای از توابع مقدماتی تجزیه می‌کنیم. برای راحتی نمادهایی را هم معرفی می‌کنیم

$$\begin{aligned} \phi^{(i)}: D_i &\rightarrow D_{i+1}, \quad i = 0, 1, \dots, r, \quad D_0 = D, D_i \subseteq R^{n_i}, D_{r+1} \subseteq R^{n_{r+1}} = R^m \\ \phi &= \phi^{(r)} \circ \phi^{(r-1)} \circ \dots \circ \phi^{(0)} \\ x^{(i)} &= \begin{pmatrix} x_1^{(i)} \\ \vdots \\ x_{n_i}^{(i)} \end{pmatrix}, \quad \phi^{(i)}(x^{(i)}) = x^{(i+1)} \\ x &= x^{(0)} \rightarrow \phi^{(0)}(x^{(0)}) = x^{(1)} \rightarrow \dots \rightarrow \phi^{(r)}(x^{(r)}) = x^{(r+1)} = y \end{aligned}$$

چنین تجزیه‌ای منحصر بفرد نیست. مثال زیر را در نظر می‌گیریم
 مثال. برای محاسبه $\phi(a, b, c) = a + b + c$ از دو الگوریتم زیر می‌توان استفاده کرد

$$\begin{aligned} \phi^{(0)}(a, b, c) &= \begin{pmatrix} a + b \\ c \end{pmatrix} \in R^2, \quad \phi^{(1)}(u, v) = u + v \in R \\ \phi^{(0)}(a, b, c) &= \begin{pmatrix} a \\ b + c \end{pmatrix} \in R^2, \quad \phi^{(1)}(u, v) = u + v \in R. \end{aligned}$$

همانطور که قبلاً هم بیان شد این دو الگوریتم از نظر عددی یکسان نیستند. در این بخش با بررسی خطاهای گرد کردن در روند محاسبه یک تابع بر اساس توابع مقدماتی به یک معیار کلی برای ارزیابی الگوریتم‌های مختلف برای حل یک مسأله می‌رسیم.

فرض کنید

$$\phi(x) = \begin{pmatrix} \phi_1(x_1, \dots, x_n) \\ \vdots \\ \phi_m(x_1, \dots, x_n) \end{pmatrix}, \quad \phi^{(i)}(u) = \begin{pmatrix} \phi_1^{(i)}(u) \\ \vdots \\ \phi_{n_i+1}^{(i)}(u) \end{pmatrix}$$

اگر \tilde{x} تقریب x باشد برای محاسبه $y = \phi(x)$ در عمل به دلیل وجود خطاهای گردکردن $\tilde{y} = \phi(\tilde{x})$ محاسبه می شود. بنابراین می توان تاثیر خطاهای اولیه $\Delta x_i = \tilde{x}_i - x_i$ بر y را بدست آورد. به عبارت دیگر می توان تقریبی برای خطای $\Delta y_i = \tilde{y}_i - y_i$ تعیین کرد. با بسط تیلور و صرف نظر کردن از جملات مرتبه بالا داریم

$$\begin{aligned} \Delta y_i &= \tilde{y}_i - y_i = \phi_i(\tilde{x}) - \phi_i(x) = \sum_{j=1}^n (\tilde{x}_j - x_j) \frac{\partial \phi_i(x)}{\partial x_j} \\ &= \sum_{j=1}^n \frac{\partial \phi_i(x)}{\partial x_j} \Delta x_j, \quad i = 1, \dots, m, \end{aligned} \quad (17)$$

اگر به صورت ماتریسی بنویسیم داریم

$$\Delta y = \begin{bmatrix} \Delta y_1 \\ \vdots \\ \Delta y_m \end{bmatrix} = \begin{bmatrix} \frac{\partial \phi_1}{\partial x_1} & \dots & \frac{\partial \phi_1}{\partial x_n} \\ \vdots & \dots & \vdots \\ \frac{\partial \phi_m}{\partial x_1} & \dots & \frac{\partial \phi_m}{\partial x_n} \end{bmatrix} \begin{bmatrix} \Delta x_1 \\ \vdots \\ \Delta x_m \end{bmatrix} = D\phi(x)\Delta x. \quad (18)$$

که در آن $D\phi(x)$ ماتریس ژاکوبی است. با تقسیم طرفین (17) به $y_i = \phi_i(x)$ داریم

$$\varepsilon_{y_i} = \sum_{j=1}^n \frac{x_j}{\phi_i(x)} \frac{\partial \phi_i(x)}{\partial x_j} \varepsilon_{x_j}, \quad i = 1, \dots, m. \quad (19)$$

که ε_{y_i} و ε_{x_i} به ترتیب خطای نسبی y_i و x_i هستند. برای ادامه بحث به مفهوم نگاشت مانده نیاز خواهیم داشت که به صورت زیر تعریف می شود

$$\psi^{(i)} = \phi^{(r)} \circ \phi^{(r-1)} \circ \dots \circ \phi^{(i)} : D_i \rightarrow R^m, \quad i = 0, 1, 2, \dots, r. \quad (20)$$

در این صورت $\psi^{(0)} = \phi$ و $D\psi^{(i)}$ و $D\phi^{(i)}$ ماتریسهای ژاکوبی نگاشت های $\psi^{(i)}$ و $\phi^{(i)}$ هستند. با توجه به

$$D(f \circ g)(x) = Df(g(x)).Dg(x),$$

داریم

$$D\phi(x) = D\phi^{(r)}(x^{(r)}) . D\phi^{(r-1)}(x^{(r-1)}) . \dots . \phi^{(0)}(x), \quad (21)$$

$$D\psi^{(i)}(x^{(i)}) = D\phi^{(r)}(x^{(r)}) . D\phi^{(r-1)}(x^{(r-1)}) . \dots . \phi^{(i)}(x^{(i)}), \quad i = 0, 1, 2, \dots, r. \quad (22)$$

خطاهای ورودی و خطاهای گردکردن مقادیر $x^{(i)}$ را پریشیده می کنند و مقادیر تقریبی $\tilde{x}^{(i)}$ با $\tilde{x}^{(i+1)} = fl(\phi^{(i)}(\tilde{x}^{(i)}))$ محاسبه می شوند. با داشتن خطاهای $\Delta x^{(i)} = \tilde{x}^{(i)} - x^{(i)}$ داریم

$$\Delta x^{(i+1)} = [fl(\phi^{(i)}(\tilde{x}^{(i)}) - \phi^{(i)}(\tilde{x}^{(i)}))] + [\phi^{(i)}(\tilde{x}^{(i)}) - \phi^{(i)}(x^{(i)})] \quad (23)$$

با صرف نظر کردن از جملات مرتبه بالا، طبق (۱۸) داریم

$$\phi^{(i)}(\tilde{x}^{(i)}) - \phi^{(i)}(x^{(i)}) = D\phi^{(i)}(x^{(i)})\Delta x^{(i)}. \quad (24)$$

با توجه به اینکه نگاشت $\phi^{(i)}$ یک نگاشت مقدماتی است داریم

$$fl(\phi^{(i)}(u)) = rd(\phi^{(i)}(u)). \quad (25)$$

اگر (۲۵) را به صورت مولفه‌ای بنویسیم داریم

$$fl(\phi_j^{(i)}(u)) = rd(\phi_j^{(i)}(u)) = (1 + \varepsilon_j)\phi_j^{(i)}(u), \quad |\varepsilon_j| < eps, \quad j = 1, 2, \dots, n_{i+1} \quad (26)$$

لذا (۲۵) را می‌توان به صورت زیر نوشت

$$fl(\phi^{(i)}(u)) = (I + E_{i+1})\phi^{(i)}(u) \quad (27)$$

که در آن I ماتریس همانی و E_{i+1} ماتریس خطای قطری است

$$E_{i+1} = \begin{bmatrix} \varepsilon_1 & & & \circ \\ & \varepsilon_2 & & \\ & & \ddots & \\ \circ & & & \varepsilon_{n_{i+1}} \end{bmatrix}, \quad |\varepsilon_j| \leq eps. \quad (28)$$

بنابراین برای عبارت اول در طرف راست رابطه (۲۳) داریم

$$fl(\phi^{(i)}(\tilde{x}^{(i)})) - \phi^{(i)}(\tilde{x}^{(i)}) = E_{i+1} \cdot \phi^{(i)}(\tilde{x}^{(i)}).$$

علاوه بر آن، اگر عبارت $\phi^{(i)}(\tilde{x}^{(i)})$ را حول $x^{(i)}$ بسط داده و به E_{i+1} ضرب کنیم جملات مرتبه بالای خطا پدیدار می‌شوند که با چشم‌پوشی از آنها نتیجه می‌شود $E_{i+1} \cdot \phi^{(i)}(\tilde{x}^{(i)}) \approx E_{i+1} \cdot \phi^{(i)}(x^{(i)})$. لذا، داریم

$$fl(\phi^{(i)}(\tilde{x}^{(i)})) - \phi^{(i)}(\tilde{x}^{(i)}) = E_{i+1} \cdot \phi^{(i)}(x^{(i)}) = E_{i+1} \cdot x^{(i+1)} =: \alpha_{i+1}. \quad (29)$$

α_{i+1} را می‌توان به عنوان خطای مطلق گرد کردن که در محاسبه $\phi^{(i)}$ تازه تولید شده، تعبیر کرد و همچنین اعضای قطری E_{i+1} خطاهای نسبی متناظر هستند. در نتیجه، بنابر (۲۳)، (۲۴) و (۲۹) عبارت $\Delta x^{(i+1)}$ را می‌توان به صورت تقریب مرتبه اول زیر بیان کرد

$$\Delta x^{(i+1)} = \alpha_{i+1} + D\phi^{(i)}(x^{(i)}) \cdot \Delta x^{(i)} = E_{i+1} \cdot x^{(i+1)} + D\phi^{(i)}(x^{(i)}) \cdot \Delta x^{(i)},$$

$$i \geq \circ, \quad \Delta x^{(\circ)} := \Delta x.$$

ولذا داریم

$$\begin{aligned} \Delta y = \Delta x^{(i+1)} &= \alpha_{r+1} + D\phi^{(r)}(x^{(r)}) \cdot \Delta x^{(r)} \\ &= \alpha_{r+1} + D\phi^{(r)}(x^{(r)}) \cdot (\alpha_r + D\phi^{(r-1)}(x^{(r-1)}) \cdot \Delta x^{(r-1)}) \\ &\vdots \\ &= \alpha_{r+1} + D\phi^{(r)}(x^{(r)}) \cdot \alpha_r + D\phi^{(r)}(x^{(r)}) \cdot D\phi^{(r-1)}(x^{(r-1)}) \cdot \alpha_{r-1} + \dots \\ &\quad + D\phi^{(r)}(x^{(r)}) \cdot \dots \cdot D\phi^{(1)}(x^{(1)}) \cdot \alpha_1 + D\phi^{(r)}(x^{(r)}) \cdot \dots \cdot D\phi^{(\circ)}(x^{(\circ)}) \cdot \Delta x. \end{aligned}$$

بنابراین با استفاده از نگاشت مانده داریم

$$\begin{aligned}\Delta y &= \alpha_{r+1} + D\psi^{(r)}(x^{(r)}) \cdot \alpha_r + \dots + D\psi^{(1)}(x^{(1)}) \cdot \alpha_1 + D\phi(x) \cdot \Delta x \\ &= E_{r+1} \cdot y + D\psi^{(r)}(x^{(r)}) \cdot E_r x^{(r)} + \dots + D\psi^{(1)}(x^{(1)}) \cdot E_1 x^{(1)} + D\phi(x) \cdot \Delta x.\end{aligned}\quad (30)$$

اگر الگوریتم‌های مختلفی برای محاسبه ϕ وجود داشته باشد، در عبارت خطا $D\phi$ تغییر نخواهد کرد. با این وجود، ماتریسهای ژاکوبی $D\psi^{(i)}$ که انتشار خطای گردکردن را ارزیابی می‌کنند تغییر خواهند کرد، و لذا تاثیر کلی خطای گردکردن به صورت زیر خواهد بود

$$\alpha_{r+1} + D\psi^{(r)}(x^{(r)}) \cdot \alpha_r + \dots + D\psi^{(1)}(x^{(1)}) \cdot \alpha_1.\quad (31)$$

برای محاسبه $\phi(x)$ ، به ازای مجموعه‌ای از داده‌های x ، یک الگوریتم را در مقایسه با الگوریتم دیگر به طور عددی قابل اعتماد (trustworthy) گوئیم اگر تاثیر کلی خطاهای گردکردن (31) برای الگوریتم اول در مقایسه با الگوریتم دوم کمتر باشد.

در عبارت خطای (30) جمله اول، بدون وابستگی به الگوریتم استفاده شده، دارای کران زیر است

$$|E_{r+1}y| \leq |y|eps,$$

لذا اندازه خطای $|y|eps$ برای هر الگوریتمی قابل انتظار است. علاوه بر آن، با استفاده از اعداد با مفسر t -رقمی داده‌های ورودی دارای خطای زیر هستند (جز در حالتی که داده‌های اولیه اعداد ماشینی باشند)

$$\Delta^{(\circ)}x \leq |x|eps,$$

در نتیجه برای هر الگوریتمی خطای زیر همواره پیش‌بینی می‌شود

$$\Delta^{(\circ)}y := [|D\phi(x)| \cdot |x| + |y|]eps\quad (32)$$

$\Delta^{(\circ)}y$ خطای ذاتی y نامیده می‌شود. چون این خطا همواره وجود دارد، لذا نباید انتظار داشت که خطاهای گردکردن در محاسبات میانی از $\Delta^{(\circ)}y$ کمتر باشند. بنابراین خطاهای گردکردن α_i یا E_i در (30) را بی‌ضرر گوئیم اگر خطای کلی Δy حداکثر از مرتبه خطای ذاتی $\Delta^{(\circ)}y$ باشد:

$$|D\psi^{(i)}(x^{(i)}) \cdot \alpha_i| = |D\psi^{(i)}(x^{(i)}) \cdot E_i x^{(i)}| \approx \Delta^{(\circ)}y.$$

اگر خطاهای گردکردن یک الگوریتم بی‌ضرر باشند، در این صورت الگوریتم را خوش‌رفتار یا به طور عددی پایدار گوئیم. مثال. محاسبه $\phi(a, b) = a^2 - b^2 = (a+b)(a-b)$ از دو الگوریتم زیر را در نظر می‌گیریم

الگوریتم ۱: $\phi^{(\circ)}(a, b) = \begin{bmatrix} a^2 \\ b \end{bmatrix}$, $\phi^{(1)}(u, v) = \begin{bmatrix} u \\ v^2 \end{bmatrix}$, $\phi^{(2)}(u, v) = u - v$

الگوریتم ۲: $\phi^{(\circ)}(a, b) = \begin{bmatrix} a \\ b \\ a+b \end{bmatrix}$, $\phi^{(1)}(a, b, u) = \begin{bmatrix} u \\ a-b \end{bmatrix}$, $\phi^{(2)}(u, v) = u \cdot v$

برای محاسبه خطاهای کلی گردکردن و خطای ذاتی داریم

$$x = x^{(\circ)} = \begin{bmatrix} a \\ b \end{bmatrix}, \quad x^{(1)} = \begin{bmatrix} a^2 \\ b \end{bmatrix}, \quad x^{(2)} = \begin{bmatrix} a^2 \\ b^2 \end{bmatrix}, \quad x^{(3)} = y = a^2 - b^2,$$

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

$$D\phi(x) = (2a, -2b),$$

$$D\psi^{(1)}(x^{(1)}) = (1, -2b), \quad D\psi^{(2)}(x^{(2)}) = (1, -1)$$

$$\text{داریم } fl(\phi^{(o)}(x^{(o)})) - \phi^{(o)}(x^{(o)}) = \begin{bmatrix} a \times^* a \\ b \end{bmatrix} - \begin{bmatrix} a^2 \\ b \end{bmatrix} \text{ با توجه به}$$

$$\alpha_1 = \begin{bmatrix} \varepsilon_1 a^2 \\ \circ \end{bmatrix}, \quad E_1 = \begin{bmatrix} \varepsilon_1 & \circ \\ \circ & \circ \end{bmatrix},$$

و به طور مشابه،

$$\alpha_2 = \begin{bmatrix} \circ \\ \varepsilon_2 b^2 \end{bmatrix}, \quad E_2 = \begin{bmatrix} \circ & \circ \\ \circ & \varepsilon_2 \end{bmatrix},$$

$$\alpha_3 = \varepsilon_3(a^2 - b^2), \quad |\varepsilon_i| \leq eps, \quad i = 1, 2, 3.$$

$$\text{با } \Delta x = \begin{bmatrix} \Delta a \\ \Delta b \end{bmatrix} \text{ از (۳۰) داریم}$$

$$\Delta y = 2a\Delta a - 2b\Delta b + a^2\varepsilon_1 - b^2\varepsilon_2 + (a^2 - b^2)\varepsilon_3. \quad (33)$$

و تاثیر کلی خطاهای گردکردن عبارتست از

$$|a^2\varepsilon_1 - b^2\varepsilon_2 + (a^2 - b^2)\varepsilon_3| \leq (a^2 + b^2 + |a^2 - b^2|)eps. \quad (34)$$

به طور مشابه برای الگوریتم ۲ داریم

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

$$\psi^{(1)}(u, v) = u.v,$$

$$D\phi(x) = (2a, -2b), \quad D\psi^{(1)}(x^{(1)}) = (a-b, a+b),$$

$$\alpha_1 = \begin{bmatrix} \varepsilon_1(a+b) \\ \varepsilon_2(a-b) \end{bmatrix}, \quad \alpha_2 = \varepsilon_3(a^2 - b^2), \quad E_1 = \begin{bmatrix} \varepsilon_1 & \circ \\ \circ & \varepsilon_2 \end{bmatrix}, \quad |\varepsilon_i| \leq eps,$$

دوباره از (۳۰) داریم

$$\Delta y = 2a\Delta a - 2b\Delta b + (a^2 - b^2)(\varepsilon_1 + \varepsilon_2 + \varepsilon_3). \quad (35)$$

$$|(a^2 - b^2)(\varepsilon_1 + \varepsilon_2 + \varepsilon_3)| \leq 3|a^2 - b^2|eps. \quad (36)$$

با مقایسه (۳۴) و (۳۶) به راحتی می توان نشان داد که اگر $3 < |\frac{a}{b}|^2 < \frac{1}{\frac{1}{3}}$ الگوریتم ۲ و در غیر این صورت الگوریتم ۱ قابل اعتمادتر است.

۴.۲ جمع سری و واگرایی عددی

بنا بر تعریف، سری $\sum_{n=1}^{\infty} x_n$ به عدد s همگراست هرگاه، به ازای $\varepsilon > 0$ دلخواه عدد صحیح $N > 0$ وجود داشته باشد به طوری که

$$\left| s - \sum_{n=1}^N x_n \right| < \varepsilon \quad (37)$$

بنابراین با جمع تعداد کافی از جملات این سری می‌توان به اندازه دلخواه به مقدار واقعی سری دست یافت. اما این موضوع با اضافه کردن جملات گرد شده درست نیست. الگوریتم جمع کردن این سری با ترتیب معمولی عبارتست از

$$s_{k+1} = s_k + x_{k+1}$$

که $s_k = \sum_{n=1}^k x_n$ مجموع جزئی سری مورد نظر است. دو منبع اصلی خطا در جمع سریها وجود دارد:

- در صورتی که همه جمله‌ها مثبت باشند، خطاها انباشته می‌شوند و در مجموع های جزئی خطاهای بزرگی را پدید می‌آورند.

- در سری‌های متناوب که دارای جملات با علامت متناوب هستند، براساس تحلیل ریاضی

$$s_{k+1} = s_k + x_{k+1} = s_{k-1} + x_{k+1} + x_k$$

(علامت x_k با x_{k+1} متفاوت است) نشان می‌دهد که در صورتی که $|x_{k+1}| \approx |x_k|$ خطای خنثی سازی ارقام بامعنی اتفاق می‌افتد.

در اینجا دو مثال از واگرایی عددی ناشی از گرد کردن ارائه می‌کنیم.

مثال ۱. سری زیر را در نظر می‌گیریم [۴]

$$\begin{aligned} & 0.5 - 0.4 \quad (38) \\ & + \underbrace{0.05 - 0.049 + 0.005 - 0.0049 + \dots + 0.005 - 0.0049}_{\text{جمله } 2^0} \\ & + \underbrace{0.005 - 0.00499 + 0.0005 - 0.000499 + \dots + 0.0005 - 0.000499}_{\text{جمله } 2^{00}} \\ & + \dots \end{aligned}$$

به وضوح جمع این سری عبارتست از

$$s = 0.1 + 10 \times 0.001 + 100 \times 0.00001 + \dots = \frac{1}{9} \quad (39)$$

اگر جملات سری تا n رقم اعشار گرد شوند، جملات سطر $(n+1)$ ام به بعد ناپدید می‌شوند. در سطرهای باقی مانده شرایط زیر حاکم است: جملات از سطر ۱ تا سطر $\lceil \frac{n+1}{3} \rceil$ تاثیری از گرد کردن نخواهند داشت. جملات از سطر $1 + \lceil \frac{n+1}{3} \rceil$ تا n دوتا دوتا همدیگر را حذف می‌کنند (خطای خنثی سازی). در نهایت در سطر $n+1$ مجموع برابر است

با $1 = 10^n \times 10^{-n}$. بنابراین مجموع جملات با گرد کردن تا n رقم اعشار حداقل ۱ خواهد شد. قضیه. فرض کنید x_0, x_1, \dots, x_n اعداد مثبت ماشینی در یک رایانه باشند که خطای گرد کردن آن ε باشد. آنگاه خطای نسبی گرد کردن در محاسبه $\sum_{k=0}^n x_k$ به صورت معمولی، حداکثر برابر است با $(1 + \varepsilon)^n - 1$. این کمیت تقریباً برابر $n\varepsilon$ است.

مثال ۲. محاسبه $\sin(4^\circ)$ را با سری زیر در نظر می‌گیریم

$$\sin x = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} x^{2n+1}$$

با در نظر گرفتن 8° جمله از این سری با دقت معمولی برابر است با $0.8 + 5.2344314E$ ، اما مقدار دقیق آن برابر است با $\sin(4^\circ) = 0.0697564737441252623456789$.

۳ مسائل

(۱) در محاسبه سری نامتناهی $\sum_{n=1}^{\infty} x_n$ ، جواب با خطای مطلق کمتر از ε مورد نیاز است. آیا توقف جمع کردن تا جملاتی که اندازه آنها کمتر از ε است روش مطمئنی است؟ پاسخ را با سری $\sum_{n=1}^{\infty} (0.99)^n$ نشان دهید.

(۲) مساله قبلی را در حالتی که جملات سری به طور متناوب مثبت و منفی باشند و $|x_n|$ یکنوا به صفر همگرا شود بررسی کنید. (از قضیه همگرایی سریهای متناوب از ریاضیات عمومی استفاده کنید).

(۳) چرا از دست رفتن ارقام بامعنی در تقریب زیر مهم نیست

$$x - \sin x \approx \left(\frac{x^3}{6}\right) \left(1 - \left(\frac{x^2}{40}\right) \left(1 - \frac{x^2}{44}\right)\right).$$

(۴) ضرب نقطه‌ای دو بردار زیر را به چند طریق محاسبه کنید:

$$x = [2.718281828, -3.141592654, 1.414213562, 0.5772156649, 0.3010299957]$$

$$y = [1486.2497, 878366.9879, -22.27492, 4773714.647, 0.000185049]$$

الف) $\sum_{i=1}^n x_i y_i$

ب) $\sum_{i=n}^1 x_i y_i$

ج) با ترتیب بزرگترین - به - کوچکترین (اعداد مثبت را به ترتیب بزرگترین به کوچکترین جمع کنید و سپس اعداد منفی را از کوچکترین به بزرگترین جمع کنید و دو مجموع جزئی را با هم جمع کنید).

د) کوچکترین - به - بزرگترین (ترتیب عکس روش قبلی)

دقت معمولی و دقت مضاعف را برای جوابها بکار ببرید. نتایج را با جواب درست تا ۷ رقم اعشار $10^{-9} \times 1.006571$ مقایسه کنید.

(۵) در یک رایانه محاسبات طوری انجام می‌شود که خطای گرد کردن دارای کران زیر است

$$x_1 \text{ op } x_2 \rightarrow (x_1 \text{ op } x_2)(1 + \varepsilon), \quad |\varepsilon| \leq 2^{-36}$$

که op نمایانگر +، -، × یا / است. فرض کنید خطای نسبی محاسبه تابع $\exp(x)$ برابر 10^{-10} است.

(الف) برای خطای نسبی محاسبه $\sinh x$ با $\frac{1}{4}(e^x - e^{-x})$ ، کران بالایی، به صورت تابعی از x ، ارائه کنید.

(ب) برای مقادیر کوچک x ، محاسبه $\sinh x$ به صورت $x + \frac{x^3}{6}$ در مقایسه با روش (الف) دارای خطای نسبی کمتری است. برای این خطا کران بالایی به صورت تابعی از x ارائه کنید.

(ج) برای دو روش x را طوری محاسبه کنید که هر دو دارای خطای یکسان باشند. مقدار این کران را محاسبه کنید.

جوابها:

(الف) $|R/\sinh x| \leq v \coth x + u$ که $u = 2^{-36}$ و $v = 10^{-10}$

(ب) $|R/\sinh x| \leq (u(x + 2x^3/3) + x^5/120)/\sinh x$

(ج) $x_0 = 0.026$ که نتیجه می‌دهد $|R/\sinh x| \lesssim 0.38 \times 10^{-8}$.

(۶) محاسبه تابع $f(x) = e^x - 1 - 0.99x$ را برای مجموعه‌ای بزرگ از مقادیر $x > 0$ در نظر می‌گیریم. به دلیل خنثی سازی برای مقادیر کوچک x باید از سری تیلور استفاده شود. برای $x > 0$ و $n \geq 2$ رابطه $g_0(x) < f(x) < g_1(x)$ برقرار است، که در آن

$$g_0(x) = 0.01x + \frac{x^2}{2!} + \dots + \frac{x^{n-1}}{(n-1)!} + \frac{x^n}{n!} e^{\theta x}$$

برنامه‌ای بنویسید که به صورت زیر محدوده مقادیر $x \leq x_m$ را تعیین کند که بسط سری باید مورد استفاده قرار گیرد (مقادیر محاسبه شده با * مشخص شده‌اند):

(i) مقادیر $g_0^*(x_i)$ ، $f^*(x_i)$ و $g_1^*(x_i)$ را برای $x_i = 10^{-10+0.1i}$ ، $i = 1, 2, 3, \dots$ محاسبه کنید. تعداد جملات n_i در سری را طوری انتخاب کنید که

$$g_1^*(x_i) - g_0^*(x_i) \leq \frac{1}{4} 10^{-7} (g_1^*(x_i) + g_0^*(x_i)).$$

(ii) هرگاه رابطه $g_0^*(x_i) < f^*(x_i) < g_1^*(x_i)$ سه مرتبه متوالی، برای مثال برای $i = m-1, m, m+1$ صدق کرد، آنگاه مقادیر x_m و تعداد جملات n_m مورد استفاده را چاپ کرده و خاتمه پیدا کند.

(۷) سری زیر را تا ۵ رقم اعشار درست محاسبه کنید

$$\sum_{n=1}^{\infty} (-1)^{n-1} \frac{1}{n^2 - \frac{1}{3^6}}$$

جواب: ۰.۸۴۹۵۵

۸) مقدار $\sum_{n=1}^{10^6} n^{-\frac{1}{2}}$ را با حداکثر خطای 10^{-5} محاسبه کنید. تخمین خطاها را ارائه کنید.
 جواب: $\sum = 1998.540144 \pm 10^{-6}$

۹) مقدار $\prod_{n=1}^{\infty} \frac{(2n)^{\frac{1}{2n}}}{(2n+1)^{\frac{1}{(2n+1)}}}$ را تا چهار رقم اعشار محاسبه کنید.
 جواب: 1.173357 ± 10^{-6}

(جواب دقیق عبارتست از $(2 \log 2)(\gamma - \frac{1}{2} \log 2)$ که γ ثابت اویلر است).

۱۰) می دانیم

$$\lim_{n \rightarrow \infty} \frac{\sin \frac{1}{n}}{\frac{1}{n}} = 1,$$

یعنی برای هر $\varepsilon > 0$ عدد طبیعی n_0 وجود دارد به طوری که به ازای هر $n \geq n_0$

$$\left| \frac{\sin \frac{1}{n}}{\frac{1}{n}} - 1 \right| < \varepsilon$$

کوچکترین مقدار ممکن n_0 را برای $\varepsilon = 10^{-6}$ بدست آورید.
 جواب: $n_0 = 409$.

۱۱) برای تقریب e^x در $0 \leq x \leq 1$ دو جمله از بسط تیلور e^x حول $x = a$ را در نظر می گیریم. a را چگونه باید انتخاب کرد تا بیشترین خطا کمینه باشد؟ a را تا دو رقم اعشار درست محاسبه کنید.
 جواب: $a \approx 0.54$

۱۲) با رابطه بازگشتی زیر مقدار a_{100} را تا ۲ رقم بامعنی درست محاسبه کنید

$$\begin{cases} a_n - \frac{2}{3}a_{n-1} - a_{n-2} = 0 \\ a_0 = 0 \\ a_1 = \frac{5}{3} \end{cases} \quad (40)$$

جواب: $a_{100} \approx 1.3 \times 10^{30}$.

۱۳) تا ۳ رقم اعشار کوچکترین مقدار مثبت x را پیدا کنید به طوری که

$$\frac{\tan x}{x} \geq 3.$$

جواب: $x = 1.324$.

۴ بررسی مراجع

برای محاسبه

$$s = a_1 + a_2 + \dots + a_n$$

با جمعوندهای نامنفی به شیوه‌های مختلف می‌توان عمل کرد. در مرجع [۲] نشان داده شده است که جمع این مقادیر با ترتیب صعودی دارای خطای کمتری است و کران خطای نسبی عبارتست از

$$\frac{|s - fl(s)|}{s} \leq (\log_2 n + 1)eps,$$

که eps دقت ماشین است. با این حال می‌توان از روشهای دقیق‌تری هم استفاده کرد که دارای محاسبات بیشتر است و نیاز به حافظه بیشتری دارد. استفاده از فرمول اویلر-ماکلورن برای جمع سریها در مرجع [۲] بحث شده است.

مراجع

- [1] D . A. Macdonald, *A Note On The Summation Of Slowly Convergent Alternating Series*, BIT 36 :4 (1996), 766-774 .
- [2] O. Caprani, *Roundoff Errors In Floating-Point Summation*, BIT 15 (1975), 5-9.
- [3] Stoer. *Introduction to Numerical Analysis*
- [4] Riesel H. *A case of numerical divergence*, springer BIT, Volume 1 No. 2 June, 1961.
- [5] Kincaid D., Cheney W. *Numerical Analysis*, (Brooks/Cole, 1991).
- [6] R. L. Burden and J. D. Faires. *Numerical Analysis*, (Brooks/Cole, 1997).

[۷] هوسکینگ، نخستین گامها در آنالیز عددی، ترجمه دکتر بابلیان و دکتر میرنیا، نشر دانشگاهی.