

Numerical Methods

Derek Li

Contents

1	Scientific Computing	2
1.1	Approximations in Scientific Computation	2
1.1.1	Absolute Error and Relative Error	2
1.1.2	Data Error and Computational Error	2
1.1.3	Truncation Error	2
1.1.4	Forward Error and Backward Error	3
1.1.5	Conditioning	3
1.1.6	Stability of Algorithms	4
1.2	Computer Arithmetic	4
1.2.1	Floating-Point Numbers	4
1.2.2	IEEE Floating-Point Standard	5
1.2.3	Normalization	5
1.2.4	Properties of Floating-Point Systems	5
1.2.5	Rounding	5
1.2.6	Machine Epsilon	6
1.2.7	Subnormals and Gradual Underflow	6
1.2.8	Exceptional Values	7
1.2.9	Floating-Point Arithmetic	7
1.2.10	Cancellation	8

1 Scientific Computing

1.1 Approximations in Scientific Computation

1.1.1 Absolute Error and Relative Error

Let A be approximate value, T be true value. Absolute error and relative error are defined as follows:

$$\begin{aligned}\text{Absolute Error} &= A - T, \\ \text{Relative Error} &= \frac{A - T}{T} \text{ assuming } T \neq 0.\end{aligned}$$

If numbers written in scientific notation agree to p significant digits, then the magnitude of the relative error is about 10^{-p} (within a factor of 10).

Example 1.1. $A = 5.46729 \times 10^{-12}$, $T = 5.46417 \times 10^{-12}$. Thus, $A - T = 0.00312 \times 10^{-12}$ and $\frac{A-T}{T} = \frac{3.12 \times 10^{-3}}{5.46417}$, i.e., the relative error around 10^{-3} .

Example 1.2. $A = 1.00596 \times 10^{-10}$, $T = 0.99452 \times 10^{-10}$. Thus, $A - T = 0.01144 \times 10^{-10}$ and $\frac{A-T}{T} = \frac{1.144 \times 10^{-2}}{0.99452}$. A and T agree to 2 significant digits.

1.1.2 Data Error and Computational Error

The difference between exact function values due to error in the input and thus can be viewed as data error.

The difference between the exact and approximate functions for the same input and thus can be considered computational error.

1.1.3 Truncation Error

Truncation error is the difference between the true result (for the actual input) and the result that would be produced by a given algorithm using exact arithmetic. It is due to approximations such as truncating an infinite series, replacing derivatives by finite differences, or terminating an iterative sequence before convergence.

Example 1.3. $f(x) = \sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \cdots$, $\hat{f}(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!}$. The part $-\frac{x^7}{7!} + \cdots$ is the truncation error.

1.1.4 Forward Error and Backward Error

Suppose we want to compute the value of a function, $y = f(x)$, but we obtain instead an approximate value \hat{y} . The discrepancy between the computed and true values, $\Delta y = \hat{y} - y$, is called the forward error.

The quantity $\Delta x = \hat{x} - x$, where $f(\hat{x}) = \hat{y}$, is called the backward error.

Example 1.4. As an approximation to $y = \sqrt{2}$, the value $\hat{y} = 1.4$ has a forward error

$$\Delta y = \hat{y} - y = 1.4 - \sqrt{2} \approx -0.0142$$

or a relative forward error -1.004×10^{-2} . We find that $\sqrt{1.96} = 1.4$, so the backward error is

$$\Delta x = \hat{x} - x = 1.96 - 2 = -0.04$$

or the relative backward error -2×10^{-2} .

1.1.5 Conditioning

A problem is well-conditioned if a small change to the input produces a small change to the output; ill conditioned if there are some examples for which a small change in the input produces a large change in output.

Consider relative forward error $\frac{\hat{y}-y}{y} = \frac{f(\hat{x})-f(x)}{f(x)}$. Since

$$f(\hat{x}) - f(x) = f'(\tilde{x})(\hat{x} - x)$$

for some \tilde{x} between x and \hat{x} provided $f'(x)$ exists and is continuous between x and \hat{x} , then

$$\frac{\hat{y} - y}{y} = \frac{xf'(\tilde{x})}{f(x)} \frac{\hat{x} - x}{x} \approx \frac{xf'(x)}{f(x)} \frac{\hat{x} - x}{x}.$$

$\frac{xf'(x)}{f(x)}$ is called condition number.

Example 1.5. The conditioning of $f(x) = \sqrt{x}$, $x \geq 0$.

Solution. We have $f'(x) = \frac{1}{2\sqrt{x}}$ and thus the condition number is

$$\frac{x}{2\sqrt{x}\sqrt{x}} = \frac{1}{2}.$$

Hence, $f(x)$ is well-conditioned.

Example 1.6. The conditioning of $f(x) = e^x$.

Solution. The condition number is x . Thus, e^x overflows or underflows if $|x|$ is large.

Example 1.7. The conditioning of $f(x) = \sin x$.

Solution. The condition number is $\frac{x \cos x}{\sin x}$.

(1) If $x \approx \pm\pi, \pm2\pi, \pm3\pi, \dots$, $\sin x$ overflows or underflows. Nevertheless, x could be 0 because

$$\left. \frac{\sin(x)}{x} \right|_{x=0} = \left. \frac{x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots}{x} \right|_{x=0} = \left. 1 - \frac{x^2}{3!} + \frac{x^4}{5!} - \dots \right|_{x=0} = 1.$$

(2) If $|x|$ is big and $\cos x \not\approx 0$, $\sin x$ overflows or underflows.

1.1.6 Stability of Algorithms

An algorithm is stable if small changes to the input result in small changes to the output. Note that stability of an algorithm does not by itself guarantee that the computed result is accurate.

1.2 Computer Arithmetic

1.2.1 Floating-Point Numbers

In a digital computer, the real number system \mathbb{R} of mathematics is represented by a floating-point number system.

Formally, a floating-point number system \mathbb{F} is characterized by four integers: β (base), p (precision), $[L, U]$ (exponent range). Any floating-point number $x \in \mathbb{F}$ has the form

$$x = \pm d_1.d_2d_3 \cdots d_p \times \beta^n,$$

where $0 \leq d_i < \beta, L \leq n \leq U$.

1.2.2 IEEE Floating-Point Standard

System	β	p	L	U
Single-Precision	2	24	-126	127
Double-Precision	2	53	-1022	1023

1.2.3 Normalization

A floating-point system is said to be normalized if the leading digit $d_1 \neq 0$ unless the number represented is zero.

1.2.4 Properties of Floating-Point Systems

There is a smallest positive normalized floating-point number, underflow limit

$$\text{UFL} = \underbrace{1.00 \cdots 0}_{p \text{ digits}} \times \beta^L.$$

There is a largest floating-point number, overflow limit

$$\text{OFL} = \underbrace{d.dd \cdots d}_{p \text{ digits}} \times \beta^U,$$

where $d = \beta - 1$ and thus

$$\text{OFL} = (\beta - \beta^{1-p})\beta^U = (1 - \beta^{-p})\beta^{U+1}.$$

1.2.5 Rounding

One of the commonly used rounding rules is round-to-nearest which is the default rounding rule in IEEE standard system.

Note 1. Most of the time, there is a unique closest p -digit number.

Note 2. In case of times, we round to the number with an even least significant digit (round to even).

Note 3. In binary, in case of times, we round to the number which has a 0 in its least significant digits.

Example 1.8. Consider a system with $\beta = 10, p = 3, L = -10, U = 10$.

$$1.54 \times 10^1 + 2.56 \times 10^{-1} = 1.5656 \times 10^1 \approx 1.57 \times 10^1.$$

1.2.6 Machine Epsilon

The accuracy of a floating-point system can be characterized by machine epsilon. We define machine epsilon

$$\varepsilon_{\text{mach}} = \beta^{1-p}.$$

There is same number of floating-point numbers in each interval $[\beta^n, \beta^{n+1})$, and numbers are evenly spaced with distance $\beta^n \cdot \varepsilon_{\text{mach}}$.

The machine epsilon is important because it bounds the relative error in representing any nonzero real number x within the normalized range of a floating-point system:

$$|\text{fl}(x) - x| \leq \frac{1}{2}\beta^n \varepsilon_{\text{mach}} \Rightarrow \left| \frac{\text{fl}(x) - x}{x} \right| \leq \frac{\frac{1}{2}\beta^n \varepsilon_{\text{mach}}}{\beta^n} = \frac{1}{2}\varepsilon_{\text{mach}}.$$

In IEEE, $\text{fl}(a \text{ op } b)$ is closest floating-point number to $a \text{ op } b$ (round to nearest and no overflows or underflows), where op means basic operations: $+$, $-$, \times , $/$, $\sqrt{\cdots}$. We also have

$$\left| \frac{\text{fl}(a \text{ op } b) - (a \text{ op } b)}{a \text{ op } b} \right| \leq \frac{1}{2}\varepsilon_{\text{mach}}.$$

1.2.7 Subnormals and Gradual Underflow

There is a noticeable gap around zero in floating-point system because of normalization. Subnormal numbers have $d_1 = 0$, which can fill in the gap, but the inequality

$$\left| \frac{\text{fl}(x) - x}{x} \right| \leq \frac{1}{2}\varepsilon_{\text{mach}}$$

might not hold.

Such an augmented floating-point system is sometimes said to exhibit gradual underflow, since it extends the lower range of magnitudes representable rather than underflowing to zero as soon as the minimum exponent value would otherwise be exceeded.

Example 1.9 (Underflow to a Subnormal). Consider a system with $\beta = 10, p = 3, L = -10, U = 10$.

$$1.01 \times 10^{-5} \times 2.02 \times 10^{-6} = 2.0402 \times 10^{-11} = 0.20402 \times 10^{-10} \approx 0.20 \times 10^{-10}.$$

Example 1.10 (Underflow to Zero).

$$1.01 \times 10^{-6} \times 2.02 \times 10^{-7} = 2.0402 \times 10^{-13} = 0.0020402 \times 10^{-10} \approx 0.00 \times 10^{-10}.$$

1.2.8 Exceptional Values

The IEEE floating-point standard provides two additional special values to indicate exceptional situations:

- Inf, which stands for infinity, results from dividing a finite number by zero or $\text{Inf} + \text{Inf}$.
- NaN, which stands for not a number, results from an undefined or indeterminate operation such as $\frac{0}{0}$, $0 \times \text{Inf}$, $\frac{\text{Inf}}{\text{Inf}}$, or $\text{Inf} - \text{Inf}$.

Example 1.11 (Overflow).

$$3.56 \times 10^5 \times 5.41 \times 10^5 = 19.2596 \times 10^{10} \rightarrow \text{Inf}.$$

1.2.9 Floating-Point Arithmetic

IEEE Standard guarantees that if a and b are IEEE floating-point numbers, then $\text{fl}(a \text{ op } b)$ is the correctly rounded value of $a \text{ op } b$ (could be $\pm\text{inf}$, NaN, underflow).

If a and b are normalized floating-point numbers and if no overflows or underflows occur in computing $a \text{ op } b$, then

$$\text{fl}(a \text{ op } b) = (a \text{ op } b)(1 + \delta),$$

where $|\delta| \leq \frac{1}{2}\varepsilon_{\text{mach}}$, $\delta = \frac{\text{fl}(a \text{ op } b) - (a \text{ op } b)}{a \text{ op } b}$ is the relative error.

If there is an underflow, we may get $|\delta| > \frac{1}{2}\varepsilon_{\text{mach}}$.

Example 1.12. Suppose $L = -10$, $p = 3$ approximate $\text{fl}(2.02 \times 10^{-6} \times 1.11 \times 10^{-6})$.

The exact value is $2.2422 \times 10^{-12} = 0.022422 \times 10^{-10} \approx 0.02 \times 10^{-10}$. We know $\delta \approx -0.108$, but $\frac{1}{2}\varepsilon_{\text{mach}} = 0.005$.

1.2.10 Cancellation

Here is an example of catastrophic cancellation.

Example 1.13. Suppose $p = 3$, $\text{fl}(1.00 \times 10^3 + 1.00 \times 10^7 - 1.00 \times 10^7)$.

First step: $\text{fl}(1.0000 \times 10^7 + 0.0001 \times 10^7) = \text{fl}(1.0001 \times 10^7) \rightarrow 1.00 \times 10^7$.
(Note: it is not a underflow.)

Second step: $\text{fl}(1.00 \times 10^7 - 1.00 \times 10^7) = 0$.

If you compute a sum and some of the intermediate values are much larger in magnitude than final result, then the relative error in the computed sum may be very inaccurate.

One exception in sum: suppose $a \geq 0, b \geq 0, c \geq 0$ (or $a \leq 0, b \leq 0, c \leq 0$) and assume no overflows or underflows.

$$\begin{aligned} \text{fl}((a+b)+c) &= [(a+b)(1+\delta_1)+c](1+\delta_2) \\ &:= (a+b+c)(1+\hat{\delta}_1)(1+\delta_2) := (a+b+c)(1+\tilde{\delta}), \end{aligned}$$

where $|\delta_1|, |\delta_2| \leq \frac{1}{2}\varepsilon_{\text{mach}}$. We can show that $|\hat{\delta}_1| \leq \frac{1}{2}\varepsilon_{\text{mach}}$ and $|\tilde{\delta}| \leq 1.01\varepsilon_{\text{mach}}$.

Proof. We have

$$(a+b)(1+\delta_1)+c \leq (a+b)(1+\frac{1}{2}\varepsilon_{\text{mach}})+c \leq (a+b+c)(1+\frac{1}{2}\varepsilon_{\text{mach}}).$$

Similarly,

$$(a+b+c)(1-\frac{1}{2}\varepsilon_{\text{mach}}) \leq (a+b)(1+\delta_1)+c := (a+b+c)(1+\hat{\delta}_1).$$

Therefore, $|\hat{\delta}_1| \leq \frac{1}{2}\varepsilon_{\text{mach}}$.

Hence,

$$\begin{aligned} |\tilde{\delta}| &\leq |\hat{\delta}_1| + |\delta_2| + |\hat{\delta}_1\delta_2| \leq \frac{1}{2}\varepsilon_{\text{mach}} + \frac{1}{2}\varepsilon_{\text{mach}} + \frac{1}{4}\varepsilon_{\text{mach}}^2 \\ &= (1 + \frac{1}{4}\varepsilon_{\text{mach}})\varepsilon_{\text{mach}} \leq 1.01\varepsilon_{\text{mach}}, \end{aligned}$$

when p is large. □

In multiplication, the situation could be better. Assume no overflows or underflows.

$$\begin{aligned}\text{fl}((a * b) * c) &= [(a * b)(1 + \delta_1)] * c(1 + \delta_2) = (a * b * c)(1 + \delta_1)(1 + \delta_2) \\ &= (a * b * c)(1 + \hat{\delta}),\end{aligned}$$

where $|\delta_1|, |\delta_2| \leq \frac{1}{2}\varepsilon_{\text{mach}}, |\hat{\delta}| \leq 1.01\varepsilon_{\text{mach}}$.

Example 1.14 (Catastrophic Cancellation). $f(x) = \sqrt{1 + x^2} - 1$.

(1) $f(x)$ is well conditioned because

$$\frac{xf'(x)}{f(x)} = \frac{x^2}{\sqrt{1 + x^2}(\sqrt{1 + x^2} - 1)} = 1 + \frac{1}{\sqrt{1 + x^2}} \in [1, 2].$$

(2) $\text{fl}(\sqrt{1 + x^2} - 1)$ does not always give an accurate result in the relative error sense. If x is small enough, $\text{fl}(\sqrt{1 + x^2} - 1) \rightarrow 0$. Recall that if $A = 0, T \neq 0, \frac{A-T}{T} = -1$, which is a bad relative error because there is no same digit.

(3) We could change the function to a mathematically equivalent value that has a much smaller floating-point error: $\frac{x^2}{\sqrt{1 + x^2} + 1}$.

Proof. Want to show $\text{fl}\left(\frac{x^2}{\sqrt{1 + x^2} + 1}\right) = \frac{x^2}{\sqrt{1 + x^2} + 1}(1 + \tilde{\delta})$. We have

$$\begin{aligned}\text{fl}\left(\frac{x^2}{\sqrt{1 + x^2} + 1}\right) &= \frac{x^2(1 + \delta_1)(1 + \delta_5)}{\left[\sqrt{[1 + x^2(1 + \delta_1)](1 + \delta_2)(1 + \delta_3) + 1}\right](1 + \delta_4)} \\ &:= \frac{x^2(1 + \delta_1)(1 + \delta_5)}{\left[\sqrt{(1 + x^2)(1 + \hat{\delta}_1)(1 + \delta_2)(1 + \delta_3) + 1}\right](1 + \delta_4)} \\ &:= \frac{x^2(1 + \delta_1)(1 + \delta_5)}{\left(\sqrt{1 + x^2}(1 + \tilde{\delta}_1)(1 + \tilde{\delta}_2)(1 + \delta_3) + 1\right)(1 + \delta_4)} \\ &:= \frac{x^2(1 + \delta_1)(1 + \delta_5)}{(\sqrt{1 + x^2} + 1)(1 + \tilde{\delta}_1)(1 + \tilde{\delta}_2)(1 + \tilde{\delta}_3)(1 + \delta_4)} \\ &:= \frac{x^2}{\sqrt{1 + x^2} + 1}(1 + \delta_1)(1 + \delta_5)(1 + \hat{\delta}_1)(1 + \hat{\delta}_2)(1 + \hat{\delta}_3)(1 + \hat{\delta}_4) \\ &:= \frac{x^2}{\sqrt{1 + x^2} + 1}(1 + \tilde{\delta}),\end{aligned}$$

where $|\delta_i| \leq \frac{1}{2}\varepsilon_{\text{mach}}$.

Note that (1) $\left|\widehat{\delta}_1\right| \leq \frac{1}{2}\varepsilon_{\text{mach}}$ since $\widehat{\delta}_1 = \frac{x^2}{1+x^2}\delta_1 \leq \frac{x^2}{2(1+x^2)}\varepsilon_{\text{mach}} \leq \frac{1}{2}\varepsilon_{\text{mach}}$.

(2) Suppose $|\delta| \leq \frac{1}{2}\varepsilon_{\text{mach}}$, we have

$$1 - \frac{1}{2}\varepsilon_{\text{mach}} \leq \sqrt{1 - \frac{1}{2}\varepsilon_{\text{mach}}} \leq \sqrt{1 + \delta} \leq \sqrt{1 + \frac{1}{2}\varepsilon_{\text{mach}}} \leq 1 + \frac{1}{2}\varepsilon_{\text{mach}},$$

provided $\frac{1}{2}\varepsilon_{\text{mach}} \leq 2$. Therefore, $\exists \widehat{\delta}$ s.t. $\sqrt{1 + \delta} = 1 + \widehat{\delta}$ for $\left|\widehat{\delta}\right| \leq \frac{1}{2}\varepsilon_{\text{mach}}$.

(3) Suppose $|\delta| \leq \frac{1}{2}\varepsilon_{\text{mach}}$, we have

$$\frac{1}{1 + \delta} \leq \frac{1}{1 - \frac{1}{2}\varepsilon_{\text{mach}}} = \frac{2}{2 - \varepsilon_{\text{mach}}} = 1 + \frac{1}{2 - \varepsilon_{\text{mach}}}\varepsilon_{\text{mach}} \leq 1 + \frac{1.01}{2}\varepsilon_{\text{mach}},$$

for $\varepsilon_{\text{mach}} \leq 0.02$. Similarly,

$$1 - \frac{1.01}{2}\varepsilon_{\text{mach}} \leq \frac{1}{1 + \delta}$$

and thus $\left|\widehat{\delta}\right| \leq \frac{1.01}{2}\varepsilon_{\text{mach}}$.

Therefore, we have $|\delta_1|, |\delta_5| \leq \frac{1}{2}\varepsilon_{\text{mach}}$, $\left|\widehat{\widehat{\delta}}_i\right| \leq \frac{1.01}{2}\varepsilon_{\text{mach}}$, and $\left|\widetilde{\delta}\right|$ is kind of small. □