

Numbers must be stored in computers and arithmetic operations must be performed on these numbers?

Two types of numbers:

1. fixed-point format: integer, long integer in C/Java/Python
2. floating-point format: float, double ...

We first consider the decimal floating-point:

For any $x \neq 0$, x can be uniquely represented as:

$$x = \sigma \cdot \bar{x} \cdot 10^e, \quad (\text{scientific notation})$$

where 1. $\sigma \in \{+1, -1\}$: sign

2. $\bar{x} \in [1, 10)$: significand (mantissa)

3. $e \in \mathbb{Z}$: exponent.

Similarly, we can define the binary floating-point:

For any $x \neq 0$, $x = \sigma \cdot \bar{x} \cdot 2^e$, (uniquely represented)

where 1. $\sigma \in \{+1, -1\}$: sign

2. $(1)_2 \leq \bar{x} < (10)_2$: significand (mantissa)

3. $e \in \mathbb{Z}$: integer, exponent

Examples of decimal:

$$x = 2018.401 \Rightarrow \sigma = +1, \bar{x} = 2.018401, e = 3$$

$$x = -0.02018 \Rightarrow \sigma = -1, \bar{x} = 2.018, e = -2$$

Examples of binary:

$$x = (1011.1)_2 \Rightarrow \sigma = +1, \bar{x} = (1.0111)_2, e = 3$$

$$x = -(0.000101)_2 \Rightarrow \sigma = -1, \bar{x} = (1.01)_2, e = -4$$

Question: How to store a binary floating-point number

into a computer? How to store σ, \bar{x}, e ?

IEEE 754 single-precision floating-point representation:

b_1	$b_2 \dots b_9$	$b_{10} b_{11} \dots b_{32}$
-------	-----------------	------------------------------

 (For 32-bit computer, it's a word).

$$1. \sigma = \begin{cases} +1 & \text{if } b_1 = 0 \\ -1 & \text{if } b_1 = 1 \end{cases} \quad (1 \text{ bit})$$

$$2. E = e + 127, \text{ where } E = (b_2 \dots b_9)_2 \quad (8 \text{ bits})$$

$$3. \bar{x} = 1.b_{10}b_{11} \dots b_{32} \quad (23 \text{ bits})$$

$$\begin{aligned} \text{So, } x &= \sigma \cdot \bar{x} \cdot 2^e = (-1)^{b_1} \cdot 2^{b_2 b_3 \dots b_9} \cdot 2^{-127} \cdot (1.b_{10}b_{11} \dots b_{32}) \\ &= (-1)^{b_1} \cdot 1.b_{10}b_{11} \dots b_{32} \cdot 2^{b_2 b_3 \dots b_9 - 127} \end{aligned}$$

$$\text{When } E = 0, \sigma = 0, b_{10}b_{11} \dots b_{32} = 00 \dots 0 \Rightarrow x = 0$$

$$E = 1, e = 1 - 127 = -126 : \text{minimum exponent}$$

$$E = 254, e = 254 - 127 = 127 : \text{maximum exponent}$$

$$E = 255, \text{ if } b_{10}b_{11} \dots b_{32} = 00 \dots 0 \Rightarrow x = \pm \infty, \text{ otherwise NaN. } \quad (2)$$

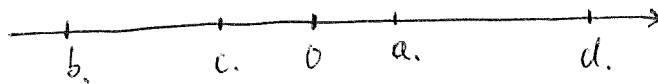
For single-precision format:

a. smallest positive number: $x = 1 \cdot (1.00 \dots 0) \cdot 2^{-126} \approx 1.2 \times 10^{-38}$

b. smallest number: $x = -1 \cdot (1.11 \dots 1) \cdot 2^{-127} \approx -3.4 \times 10^{-38}$

c. largest negative number: $x = -1 \cdot (1.00 \dots 0) \cdot 2^{-126} \approx -1.2 \times 10^{-38}$

d. largest number: $x = 1 \cdot (1.11 \dots 1) \cdot 2^{-127} \approx 3.4 \times 10^{-38}$



Question: How accurately a number x can be stored in single / double precision? / How to measure the accuracy of a floating-point format? machine epsilon / largest integer M

machine epsilon:

For any particular floating-point format, the difference between 1. and the next larger number that can be stored in that format.

Example: For single-precision.

$$\begin{aligned} \text{eps} &= x_2 - x_1, \text{ where } \begin{cases} x_1 = 1.\underbrace{00 \dots 00}_{23 \text{ bits}} \\ x_2 = 1.\underbrace{00 \dots 01}_{22 \text{ bits}} \end{cases} \\ &= 1 + 2^{-23} - 1.0 \\ &= 2^{-23} \approx 1.19 \times 10^{-7} \end{aligned}$$

$$x = 1.0 + 2^{-23} \quad \checkmark \quad x = 1.0 + 2^{-24} \quad \times \text{ (cannot be exactly stored)}$$

For double-precision.

$$\text{eps} = 2^{-52} \approx 2.22 \times 10^{-16}$$

$$x = 1.0 + 2^{-52} \quad \checkmark \quad x = 1.0 + 2^{-53} \quad \times$$

largest integer M :

M has the following property:

For any x satisfying $0 \leq x \leq M$ can be stored exactly in the floating-point format.

$x = 0.\bar{x} \cdot 2^e$. If \bar{x} has n bits, then $\bar{x}_{\max} = \underbrace{1.11 \dots 1}_{n-1}$, $e = n-1$

, so $x = 2^n - 1$ can be stored exactly. Also, $x = 2^n$ can also be stored. But $x = 2^n + 1$ does NOT.

For single-precision, $M = 2^{24} = 16777216$

For double-precision, $M = 2^{53} \approx 9.0 \times 10^{15}$

Question: How to deal with numbers that cannot be stored exactly in a computer?

$$\left\{ \begin{array}{l} x = 1. + 2^{-24} \\ x = 1. + 2^{-53} \end{array} \right.$$

Rounding and chopping:

Suppose $\bar{x} = 1.b_1 \dots b_{n-1} b_n b_{n+1} \dots$, we have two ways:

- 1. truncate / chop: ignore $b_n b_{n+1} \dots$
- 2. round: $\left\{ \begin{array}{l} \text{if } b_n = 0, \text{ chop } \bar{x} \text{ to } n \text{ digits.} \\ \text{if } b_n = 1, \text{ add 1 to the last digit, } b_{n-1}. \end{array} \right.$

(Find out, other 5 rounding methods defined in IEEE 754)

Notation: $x = x_T$: true number,

$x_A, fl(x)$: machine floating-point number /
approximated number.

Let's assume that:

$$fl(x) = x \cdot (1 + \epsilon), \quad x = \sigma \cdot \bar{x} \cdot 2^e$$

$$\bar{x} = 1.b_1 b_2 \dots b_n b_{n+1} \dots$$

For chopping:

1. $\epsilon \leq 0$ For any x .

$$fl(x) - x = x \cdot \epsilon$$

if ~~$x < 0$~~ , then $fl(x) - x = x \cdot \epsilon$.

$$= \sigma \cdot \bar{x} \cdot 2^e \cdot \sigma \cdot (1.b_1 b_2 \dots b_n - 1.b_1 b_2 \dots b_n b_{n+1} \dots) \cdot 2^e$$

$$= \sigma^2 \cdot \bar{x} \cdot 2^{2e} \cdot (-0.\underbrace{00 \dots 0}_{n-1} b_n b_{n+1} \dots)$$

$$\leq 0.$$

2. Since $\epsilon \leq 0$, the errors could accumulate to be a large number.

3. $-2^{-n+1} \leq \epsilon \leq 0$.

$$\begin{aligned} \epsilon &= \frac{fl(x) - x}{x} = \frac{\sigma \cdot (-0.00 \dots 0 b_n b_{n+1} \dots) \cdot 2^e}{\sigma \cdot 1.b_1 b_2 \dots b_n b_{n+1} \dots \cdot 2^e} \\ &= - \frac{0.00 \dots 0 b_n b_{n+1} \dots}{1.b_1 b_2 \dots b_n b_{n+1} \dots} \quad (-1 \leq \epsilon) \end{aligned}$$

$$\text{So, } -0.00 \dots 0 b_n b_{n+1} \dots \leq \epsilon.$$

$$\text{where } -0.00 \dots 0 b_n b_{n+1} \dots = -\sum_{i=n}^{\infty} 2^{-i} = -2^{-n+1}$$

$$\text{therefore, } -2^{-n+1} \leq \epsilon \leq 0.$$

For rounding:

$$f(x) = \begin{cases} \sigma \cdot 1.b_1 b_2 \dots b_{n-1} \cdot 2^e & \text{if } b_n = 0 \\ \sigma \cdot (1.b_1 b_2 \dots b_{n-1} + 2^{-n+1}) \cdot 2^e & \text{if } b_n = 1 \end{cases}$$

Case 1: $b_n = 0$

$$\begin{aligned} |\epsilon| &= \left| \frac{f(x) - x}{x} \right| = \left| \frac{\sigma \cdot (0.00\dots 0 \underset{n-1}{0} b_{n+1} b_{n+2} \dots) \cdot 2^e}{\sigma \cdot 1.b_1 b_2 \dots b_{n-1} b_n \dots \cdot 2^e} \right| \\ &= \frac{0.00\dots 0 b_{n+1} b_{n+2} \dots}{1.b_1 b_2 \dots b_{n-1} b_n \dots} \leq 0.00\dots 0 b_{n+1} b_{n+2} \dots \\ &= \sum_{i=n+1}^{\infty} 2^{-i} < 2^{-n}, \quad \text{so } |\epsilon| \leq 2^{-n}. \end{aligned}$$

Case 2: $b_n = 1$

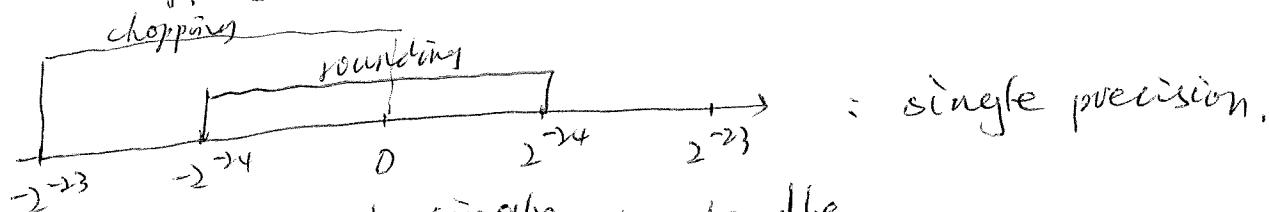
$$\begin{aligned} |\epsilon| &= \left| \frac{f(x) - x}{x} \right| = \left| \frac{\sigma \cdot (2^{-n+1} - 0.00\dots 0 1 b_{n+1} b_{n+2} \dots) \cdot 2^e}{\sigma \cdot 1.b_1 b_2 \dots b_{n-1} b_n \dots \cdot 2^e} \right| \\ &= \frac{2^{-n+1} - 0.00\dots 0 1 b_{n+1} b_{n+2} \dots}{1.b_1 b_2 \dots b_{n-1} 1 b_{n+1} \dots} = \frac{2^{-n} - 0.00\dots 0 b_{n+1} b_{n+2} \dots}{1.b_1 b_2 \dots b_{n-1} 1 b_{n+1} \dots} \\ &\leq \frac{2^{-n} - 0.00\dots 0 b_{n+1} b_{n+2} \dots}{1. + 2^{-n}} \leq \frac{2^{-n}}{1. + 2^{-n}} \leq 2^{-n}. \end{aligned}$$

For both cases, we have, $|\epsilon| \leq 2^{-n}$.

$$-2^{-n} \leq \epsilon \leq 2^{-n}.$$

Summary for rounding and chopping:

1. chopping could be as twice ^(large) as rounding



	single	double
rounding	$[-2^{-24}, 2^{-24}]$	$[-2^{-53}, 2^{-53}]$
chopping	$[-2^{-23}, 0]$	$[-2^{-52}, 0]$

$E \epsilon_r = 0$, $E \epsilon_c = -2^{-24}$, rounding is better.

Question: How to measure errors?

Let x_T : true number, x_A : machine number.

Absolute error: $\text{Err}(x_A) = x_T - x_A$

Relative error: $\text{Rel}(x_A) = \frac{x_T - x_A}{x_T}$

Example: $x_T = \pi$, $x_A = 22/7$.

$$\text{Err}(x_A) \approx -0.000126$$

$$\text{Rel}(x_A) \approx -0.000403$$

Question: Which one is better?

Ex 1: Albany $\xrightarrow{\text{NYC}}$

$$x_T = 152 \text{ mi} \quad \text{Err}(x_A) = 1$$

$$x_A = 151 \text{ mi} \quad \text{Rel}(x_A) = 1/151 \approx 0.658\%$$

Albany $\xrightarrow{\text{Crossgates}}$

$$x_T = 2 \text{ miles}$$

$$x_A = 1 \text{ mile}$$

⑦

There are 8 types of errors. Check: example-02.ipynb

1. modeling errors:

$$N(t) = N_0 e^{kt}, \quad N_0 = 3.929000 \times e^{-1790k}, \quad k = 0.02975$$

$$1790 \leq t \leq 1860.$$

2. blunders and mistakes.

3. physical measurement : speed of light in a vacuum

$$c = (2.997925 \pm 2) \cdot 10^{10} \text{ cm/sec}, \quad |e| \leq 0.000003.$$

4. machine representation error : rounding / chopping.

5. mathematical approximation errors.

6. loss-of-significant errors.

$$\cos(2\theta) = 2\cos^2(\theta) - 1 = 1 - 2\sin^2(\theta).$$

$$f(x) = \frac{1 - \cos x}{x^2} = \frac{2\sin^2(x/2)}{x^2} = \frac{1}{2} \left[\frac{\sin(x/2)}{x/2} \right]^2.$$

7. noise in function evaluation :

8. overflow / underflow.

Errors in arithmetic operations:

Propagated error: $E = (x_T \text{ w } y_T) - (x_A \text{ w } y_A)$

1. x_T, y_T are true numbers.
2. x_A, y_A are machine numbers.
3. $w \in \{+, -, \times, \div\}$.
4. Suppose we know the rounding error ϵ_x, ϵ_y .

$$x_T = x_A + \epsilon_x, \quad y_T = y_A + \epsilon_y$$

①. For multiplication: $w = \times$

We try to use $\text{Rel}(x_A \times y_A)$ to bound the error.

$$\begin{aligned}\text{Rel}(x_A \times y_A) &= \frac{x_T y_T - x_A y_A}{x_T \cdot y_T} = \frac{x_T y_T - (x_T - \epsilon_x) \cdot (y_T - \epsilon_y)}{x_T y_T} \\&= \frac{x_T \epsilon_y + y_T \epsilon_x - \epsilon_x \epsilon_y}{x_T y_T} = \frac{\epsilon_y}{y_T} + \frac{\epsilon_x}{x_T} - \frac{\epsilon_x}{x_T} \cdot \frac{\epsilon_y}{y_T} \\&= \text{Rel}(x_A) + \text{Rel}(y_A) - \underbrace{\text{Rel}(x_A) \cdot \text{Rel}(y_A)}_{\text{small}} \\&\approx \text{Rel}(x_A) + \text{Rel}(y_A).\end{aligned}$$

Numerical stable.

②. For division:

$$\begin{aligned}\text{Rel}\left(\frac{x_A}{y_A}\right) &= \frac{x_T/y_T - x_A/y_A}{x_T/y_T} = \frac{\left(x_T/y_T - \frac{x_A}{y_A}\right) \cdot \frac{y_A}{x_T}}{\left(\frac{x_T}{y_T}\right) \cdot \frac{y_A}{x_T}} \\&= \frac{\frac{y_A}{y_T} - \frac{x_A}{x_T}}{\frac{y_A}{y_T}} = \frac{1 - \frac{x_A}{x_T} - \left(1 - \frac{y_A}{y_T}\right)}{1 - \left(1 - \frac{y_A}{y_T}\right)} = \frac{\text{Rel}(x_A) - \text{Rel}(y_A)}{1 - \text{Rel}(y_A)}\end{aligned}$$

$$\text{As } \lim_{n \rightarrow \infty} \sum_{n=0}^{\infty} [\text{Rel}(y_A)]^n = \frac{1}{1 - \text{Rel}(y_A)}$$

$$\text{So } \text{Rel}\left(\frac{x_A}{y_A}\right) = [\text{Rel}(x_A) - \text{Rel}(y_A)] \cdot \sum_{n=0}^{\infty} [\text{Rel}(y_A)]^n$$

$$\approx \text{Rel}(x_A) - \text{Rel}(y_A) + \underbrace{\text{Rel}(x_A) \cdot \text{Rel}(y_A) - \text{Rel}(y_A)^2 + \dots}_{\text{small}}$$

$$\approx \text{Rel}(x_A) - \text{Rel}(y_A).$$

So, division is also numerical stable.

Addition / Subtraction:

$$\begin{aligned} \text{Rel}(x_A + y_A) &= \frac{x_T + y_T - (x_A + y_A)}{x_T + y_T} = \frac{x_T \cdot \frac{x_T - x_A}{x_T} + y_T \cdot \frac{y_T - y_A}{y_T}}{x_T + y_T} \\ &= \frac{x_T \cdot \text{Rel}(x_A) + y_T \cdot \text{Rel}(y_A)}{x_T + y_T} \end{aligned}$$

• If x_T, y_T have same sign:

$$|\text{Rel}(x_A + y_A)| \leq |\text{Rel}(x_A)| + |\text{Rel}(y_A)|$$

• If x_T, y_T have different sign:

Consider $x_T = x+1, y_T = -x, x \geq 0$.

$$\text{Rel}(x_A + y_A) = \frac{(x+1) \cdot \text{Rel}(x_A) - x \cdot \text{Rel}(y_A)}{x+1 - x}$$

$$= \underbrace{x \cdot [\text{Rel}(x_A) - \text{Rel}(y_A)]}_{\text{this could be large}} + \text{Rel}(x_A)$$

So, addition is not numerical stable.

Propagated error in function evaluation:

$f(x)$ differentiable $[a, b]$.

$$f(x_T) - f(x_A) = f'(\eta) \cdot (x_T - x_A).$$

$$\eta \in [x_T, x_A] \text{ or } [x_A, x_T]$$

Since x_T and x_A are close to each other.

$$f(x_T) - f(x_A) \approx f'(x_T) \cdot (x_T - x_A).$$

$$\text{Rel}(f(x_A)) = \frac{f(x_T) - f(x_A)}{f(x_T)} = \frac{f'(\eta) \cdot (x_T - x_A)}{f(x_T)}$$

$$\approx \frac{f'(x_T) \cdot x_T \cdot (x_T - x_A)}{f(x_T) \cdot x_T}$$

$$= \frac{f'(x_T) \cdot x_T}{f(x_T)} \cdot \text{Rel}(x_A)$$

Example: $f(x) = b^x$, $f'(x) = (\log b) \cdot b^x$.

$$\text{So, } \text{Rel}(f(x_A)) \approx \frac{(\log b) \cdot b^x}{b^x}.$$

$$b^{x_T} - b^{x_A} \approx (\log b) \cdot b^{x_T} \cdot (x_T - x_A)$$

$$\text{Rel}(b^{x_A}) \approx (\log b) \cdot x_T \cdot \text{Rel}(x_T).$$

$K = (\log b) \cdot x_T$. condition number.

It could be large!