

# Lecture-01 Introduction to numerical computation

Lecturer: Baijian Zhou

Email: [bjzhou@fudan.edu.cn](mailto:bjzhou@fudan.edu.cn)

Books:

- Numerical analysis (3rd), Timothy Sauer
- Matrix Computation (4th), Gene ....

Some notations:

$C(\mathbb{R})$ : set of all functions that are continuous on  $\mathbb{R}$ .

$C^1(\mathbb{R})$ : set of all  $f'$  continuous on  $\mathbb{R}$ .

$C^n[a, b]$ :  $f^{(n)}$  exists and continuous.

$\|\cdot\|$ : norm

$(a)_{10}$ : decimal numbers

$(a)_2$ : binary numbers

float: float point of a saved in comp.

problem 1: How to calculate  $\sqrt{2}$ .

Solution: Babylonian method

Let the numerical value of  $\sqrt{2}$  be  $x$ .

$$\Rightarrow \sqrt{2} \Rightarrow x = \sqrt{2}$$

$$1. \quad x = \sqrt{2} \Leftrightarrow x^2 = 2 \quad \Rightarrow \quad \frac{x}{2} = \frac{1}{x}$$

$$\Rightarrow \frac{x}{2} + \frac{x}{2} = \frac{x}{2} + \frac{1}{x} \Rightarrow x = \frac{x}{2} + \frac{1}{x}$$

We can guess a value of  $x$ , i.e.,  $x_0$ .

Hoping that  $\frac{x_0}{2} + \frac{1}{x_0}$  is getting closer to  $\sqrt{2}$ .

Let  $x_1 = \frac{x_0}{2} + \frac{1}{x_0}$ , repeat this ...

For example,

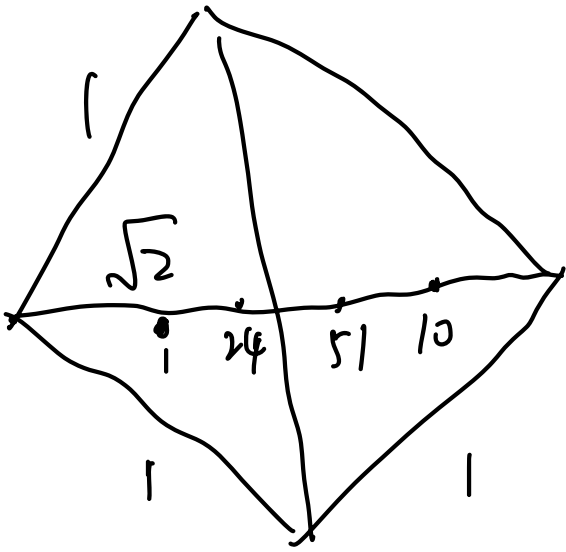
•  $x_0 = 1, \quad x_1 = \frac{1}{2} + 1 = 1.5 \quad x_1$  is better than  $x_0$

•  $x_2 = \frac{x_1}{2} + \frac{1}{x_1} = \frac{1.5}{2} + \frac{2}{3} \approx 1.4166\dots \quad x_2 \dots x_0$

⋮

★ Q1: why does  $x_{t+1} = \frac{x_t}{2} + \frac{1}{x_t}$  work?

# Illustration of Babylonian method



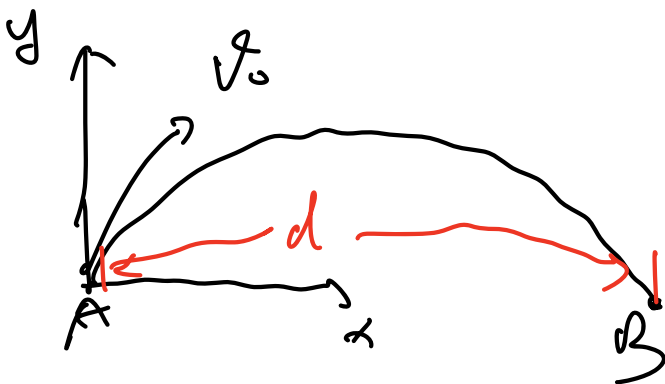
$$1 + \frac{24}{60} + \frac{51}{60^2} + \frac{10}{60^3}$$
$$= \frac{305470}{216000} = 1.41421 \overline{297}$$

six decimal digits accuracy!

Problem 2: Predict the angle.

- Given
- $d$  distance  $A \rightarrow B$
  - $v_0$  speed of shell
  - gravitational acceleration  $g = 9.8$

Give the right angle so that shell can shell from A to B.



$t$ : total time

$$v_{0x} \cdot t = d.$$

$$v_y = v_{0y} - g t_{up}, \quad v_y = 0: \text{ reach the maximum.}$$

$$t_{up} = \frac{v_{0y}}{g}, \quad \text{total time } t = 2 t_{up}$$

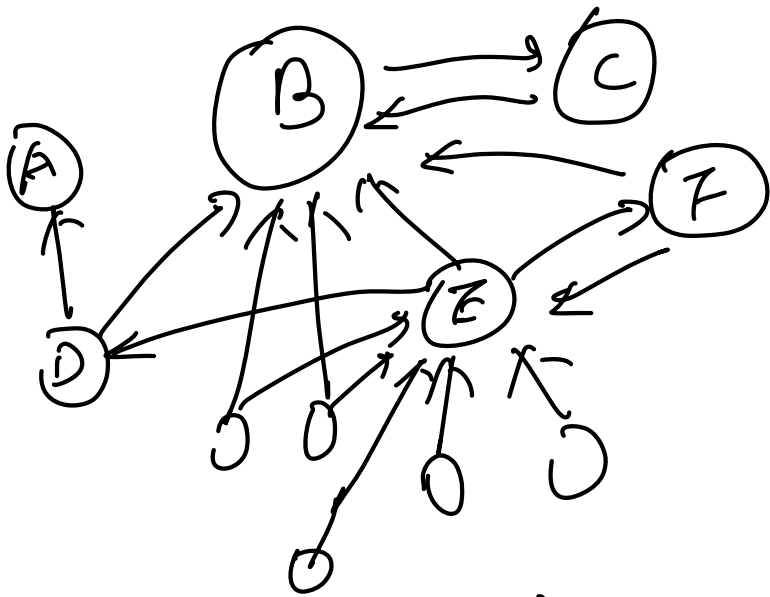
$$d = 2 \cdot \frac{v_{0y}}{g} \cdot v_{0x} \quad v_{0x} = v_0 \cos \theta$$

$$v_{0y} = v_0 \sin \theta$$

$$\Rightarrow 2 v_0^2 \sin \theta \cdot \cos \theta / g - d = 0. \quad \text{How to get } \theta?$$

# Problem 3: PageRank Problem.

Ranking Web pages.



directed graph  $\rightarrow$  adjacency matrix  $A$   
degree matrix  $D$ .

stochastic matrix  $A^T D^{-1}$ .

To solve:  $\pi = (2A^T D^{-1} + \frac{1-\alpha}{n} E) \pi$ .

# Machine representation of real numbers

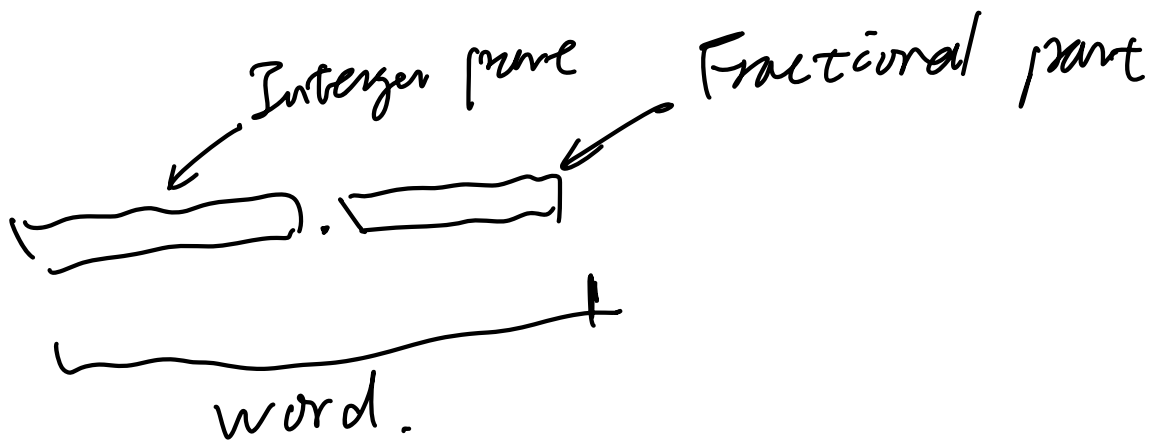
Examples:  $\pi$ ,  $e$ ,  $m_e \approx 9.1 \times 10^{-31}$ ,  $c = 2.9 \times 10^8$

atoms size  $10^{28} \sim 10^{82}$

Computers deal with numbers

→ store into a word.

Naïve idea: Fixed-point arithmetic



$10^{28}$  needs  $\approx 260$  bits of integer part

$\pi$  needs  $\approx 1$  bit but many digits in fractional part.

So, fixed-point is a bad idea!

Better method: float-point rep.

$$(x)_f = \pm m \times 10^n$$

↑ Sign      ↑ mantissa  
浮点      power/exponent

Double precision:

$$x = 1, \quad \pm 1.\overbrace{0\dots 0}^{s_2} \times 2^0$$

$$x = (1 \pm 2^{-s_2}) \pm 1.\overbrace{0\dots 0}^{s_1} \times 2^0, \quad \text{next float point}$$

$$\Rightarrow \text{We call } 2^{-s_2} = \epsilon_{\text{mach}}.$$

Rounding — chopping: biased.

$$x = \pm 1.(b_1 \dots b_{s_2} + 0 \dots 0 b_{s_3} \dots) \times 2^0$$

If  $x < 0$ , then  $b_{s_3}$  removal will make  $x \rightarrow 0$  always.

If  $x > 0$ , then  $x \rightarrow 0$  always

$\left| \begin{array}{c} \longrightarrow 0 \longleftarrow \\ x < 0 \qquad x > 0 \end{array} \right| : \text{chopping.}$

rounding :

$\left\langle \begin{array}{c} \longleftarrow \rightarrow 0 \longleftarrow \rightarrow \\ x < 0 \qquad x > 0 \end{array} \right\rangle$

If  $x = \pm 1. \boxed{b_1 \dots b_{s_1}} b_{s_2} b_{s_3} \dots \times 2^p$   
 $b_{s_2} = 1, b_{s_3} = 0, b_{s_4} = 0, \dots$

$x = \pm 1. \boxed{b_1 \dots b_{s_2}} \downarrow 00 \dots 0 \times 2^p$

- If  $b_{s_2} = 1$ ,  $\rightarrow$  round up  $\rightarrow b_{s_2} = 0$
- If  $b_{s_2} = 0$ ,  $\rightarrow$  round down  $\rightarrow b_{s_2} = 0$

$\left\langle \begin{array}{c} \longleftarrow \rightarrow 0 \longleftarrow \rightarrow \end{array} \right\rangle$

How to measure the error:

rounding error of  $x$ :

$x = 9.4, f(9.4) - 9.4 : \text{rounding error.}$



How to measure the error:

absolute error:  $|x_c - x|$

relative error:  $\frac{|x_c - x|}{|x|}$  or  $\frac{|x_c - x|}{|x_c|}$

If  $x \in \mathbb{R}^n$  or  $x \in \mathbb{R}^{n \times n}$ , use:

$\|x_c - x\|$  or  $\|x_c - x\|_{op}$ .

What if  $x=0$ , no worry for saving numbers. Since  $x \neq 0$ , there is no rounding error:  $\frac{0}{0}$ ; relative error is 0.

---

Theorem:

$$\frac{|f(x_c) - x|}{|x|} \leq \frac{1}{2} \epsilon_{mach}$$

Verify this:

$$|f(9.4) - 9.4| = 0.2 \times 2^{-49}$$

unit round off

$$\frac{|f(9.4) - 9.4|}{9.4} = \frac{0.2 \times 2^{-49}}{9.4} = \frac{8}{47} \times 2^{-52} \leq \frac{1}{2} \epsilon_{mach}$$

Proof:

W.L.O.G., Given  $x > 0$ , we want to measure

$$\left| \frac{x - f(x)}{x} \right|. \quad \text{We assume } x = q \times 2^m$$

$$q = (1.b_1 \dots b_{s_2} b_{s_3} b_{s_4} \dots)_2 \times 2^m \quad b_i \in \{0, 1\}$$

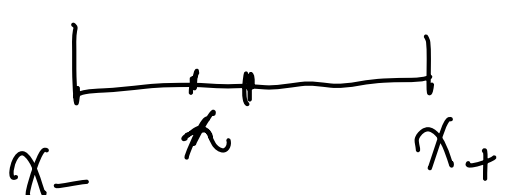
Two cases:

$$\text{rounding down: } x_- = (1.b_1 \dots b_{s_2})_2 \times 2^m.$$

$$\text{rounding up: } x_+ = \left[ (1.b_1 \dots b_{s_2})_2 + 2^{-s_2} \right] \times 2^m.$$

①. Case 1.  $f(x) = x_-$

$$(*) \quad \left| \frac{x - f(x)}{x} \right| = \left| \frac{x - x_-}{x} \right| \leq \frac{1}{2} \frac{|x_+ - x_-|}{|x|}$$

  $|x - x_-|$  must be no more than half of  $|x_+ - x_-|$

$$\frac{|x_+ - x_-|}{|x|} = \frac{2^{m-s_2}}{(1.b_1 \dots b_{s_2})_2 \times 2^m} = \frac{2^{-s_2}}{(1.b_1 \dots b_{s_2})_2} \leq 2^{-s_2}$$

$$\Rightarrow |x| \leq \frac{1}{2} \cdot 2^{-52}$$

③. Case 2;  $f(x) = x_+$

$$(*) \quad \frac{|x - f(x)|}{|x|} = \frac{|x - x_+|}{|x|} \leq \frac{1}{2} \frac{|x_+ - x_-|}{|x|}$$

$$\leq \frac{1}{2} \cdot 2^{-52} \quad \square$$

Remark: Let  $\delta = \frac{f(x) - x}{x}$ , then

①.  $f(x) = x(1 + \delta)$ , and  $|\delta| \leq 2^{-53}$

②. Generation:

$$\odot \in \{+, -, \times, \div\}$$

$x = f(x)$ ,  $y = f(y)$ : machine numbers

$x \odot y$ : computed and store

$\rightarrow f(x \odot y)$ .

$$f(x \odot y) = [x \odot y] \cdot (1 + \delta), \quad |\delta| \leq \epsilon.$$

log-sum-exp trick:

$$\log \pi_i = \log \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)}$$

$$\log \pi_i = x_i - \log \sum_{j=1}^n \exp(x_j)$$

$$\begin{aligned} \log \text{sumexp}(x) &= \log \sum_{j=1}^n \exp(x_j) + b - b \\ &= b + \log \sum_{j=1}^n \exp(x_j - b) \end{aligned}$$

$$b = \max \{ x_i, i = 1, 2, \dots, n \}.$$

---

this trick can avoid overflow.

Bisection method:

Given:  $[a, b]$ ,  $f$ ,  $\epsilon$ . Such that

Bisection goes as the following:  $f(a) \cdot f(b) < 0$

for  $t = 0, 1, 2, \dots$

$$c = \frac{a+b}{2}$$

if  $f(c) = 0$  return  $c$

if  $f(a) \cdot f(c) < 0$  then

$$b = c$$

else

$$a = c$$

return  $c$

Q. Error analysis: with  $\begin{cases} a_0 = a \\ b_0 = b \end{cases}$

Bisection generates:

$[a_0, b_0], [a_1, b_1], \dots, [a_n, b_n],$

where

$$a_0 \leq a_1 \leq a_2 \leq \dots \leq a_n \leq b_n \leq b_{n-1} \leq \dots \leq b_2 \leq b_1 \leq b_0$$

and

$$b_{n+1} - a_{n+1} = \frac{1}{2} (b_n - a_n) \quad (n \geq 0)$$

(Recall every time, it cuts  $[a_n, b_n]$  in half.)

Recursively,

$$b_n - a_n = 2^{-n} (b_0 - a_0).$$

$$\text{Thus, } \lim_{n \rightarrow \infty} b_n - \lim_{n \rightarrow \infty} a_n$$

$$= \lim_{n \rightarrow \infty} 2^{-n} (b_0 - a_0)$$

$$= 0.$$

$$\text{If we put } r = \lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} b_n,$$

by taking the limit in the inequality

$$f(a_n) \cdot f(b_n) \leq 0.$$

$$\Rightarrow \lim_{n \rightarrow \infty} f(a_n) \cdot f(b_n) \leq 0 \Rightarrow f(r) \cdot f(r) \leq 0$$

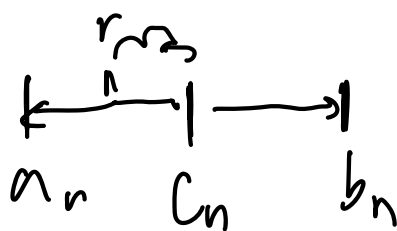
$\Rightarrow f(r) = 0.$

$|a_n - b_n| < \epsilon$  stopped:

$r$  must be in  $[a_n, b_n]$ ,

$c_n = \frac{a_n + b_n}{2}$  is the estimate of  $r$ .

So,  $|r - c_n| \leq \frac{1}{2}(b_n - a_n)$


$$= 2^{-n-1}(b_0 - a_0)$$

Finally,  $|r - c_n| \leq 2^{-(n+1)} \cdot (b_0 - a_0).$

Thm. Error analysis of Bisection:

Let  $[a_n, b_n]$  be intervals used in Bisection, then

$\lim_{n \rightarrow \infty} a_n$  and  $\lim_{n \rightarrow \infty} b_n$  exist, are equal, and

represent a zero of  $f$ . Def.  $c_n = \frac{a_n + b_n}{2}$

and  $r = \lim_{n \rightarrow \infty} c_n$ , then

$$|r - c_n| \leq 2^{-cn+1} (b_0 - a_0).$$

---

②. Time complexity:  $O(n)$ .

③. Example 3.4 of using bisection.

$$\text{We know } |r - c_n| \leq 2^{-cn+1} (b_0 - a_0)$$

$$\text{Let } 2^{-cn+1} (b_0 - a_0) \leq 0.5 \times 10^{-p}$$

$$2^{-n} \leq 10^{-6}$$

$$a_0 = 0$$

$$b_0 = 1$$

$$p = 6$$

$$-n \cdot \log_2 2 \leq -6 \log_2 10$$

$$n \geq \left\lceil \frac{6 \log_2 10}{\log_2 2} \right\rceil = 20.$$

④. Exercise: Suppose that the bisection method is started with the interval  $[50, 63]$ . How many steps should be taken to compute a root with relative accuracy of one part in  $10^{12}$ ?