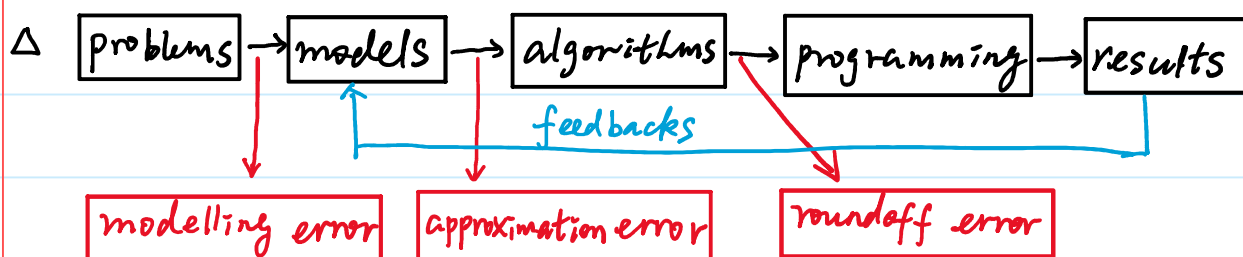


Week 1, Tuesday



⇒ core focus in this course.

△ roundoff error

e.g. $2+2=4$ $2.1 \times 3 = 6.3$ $7 - 6.15 = 0.85$ $4 \div 16 = 0.25$ $(\sqrt{3})^2 = 3$

exact arithmetic! but machine only has finite memory!

e.g. $\sqrt{3}$ is stored as 1.7321 in a calculator if it can only store 5 digits, in this calculator: $(\sqrt{3})^2 = 1.7321^2 = 3.0002$ → roundoff error

roundoff error: produced when perform finite-digit arithmetic

Q: how a number is stored in a machine?

△ binary machine numbers (IEEE 754-2008) 二进制

floating-point number: $(-1)^S 2^{C-1023} (1+f)$ 浮点数

$S = \begin{cases} 0 & \text{positive} \\ 1 & \text{negative} \end{cases}$, sign part

C : 11 digits, exponential part

f : 52 digits, mantissa part

64 digits (double)
双精度

e.g. $\begin{array}{ccc} 0 & 10000000011 & 1011001000100\dots0 \\ S & C & f \end{array}$

$S=0$, $C = 2^{10} + 2^1 + 2^0 = 1027$, $f = (\frac{1}{2})^1 + (\frac{1}{2})^3 + (\frac{1}{2})^4 + (\frac{1}{2})^5 + (\frac{1}{2})^8 + (\frac{1}{2})^{12}$

this number is: $(-1)^0 \cdot 2^{1027-1023} \cdot (1+f) = 27.56640625$

Rk: 1. range: $-1.79 \times 10^{308} \sim 1.79 \times 10^{308}$

Q: what is the smallest positive number and largest one?

2. machine epsilon: $2^{-53} \times (1/53) \sim 2^{-1023} \dots -16$

x: what is the smallest positive number and largest one?

$$2. \text{ machine error: } 2^{1023} \times \left(\frac{1}{2}\right)^{53} \approx 2^{-1023} \times 10^{-16}$$

only 16 significant digits (有效数字 later) in decimal number

△ Decimal machine number + 进制

normalized decimal floating-point form:

$$\pm 0.d_1 d_2 \dots d_k \times 10^n, \quad 1 \leq d_1 \leq 9, \quad 0 \leq d_i \leq 9, \quad i=2,3,\dots,k.$$

For any positive number $y = 0.d_1 d_2 \dots d_k d_{k+1} d_{k+2} \dots \times 10^n$

floating-point form:

① chopping: $f(y) = 0.d_1 d_2 \dots d_k \times 10^n$, chop off digits $d_{k+1} d_{k+2} \dots$

② rounding: $f(y) = 0.d_1 d_2 \dots d_k \times 10^n$, add $5 \times 10^{n-(k+1)}$ then chop.

$$(i) d_{k+1} \geq 5 \quad \tilde{d}_k = d_k + 1 \quad (ii) d_{k+1} < 5 \quad \tilde{d}_k = d_k \quad (\text{四舍五入})$$

$$\text{e.g. } \pi = 0.314159265 \dots \times 10^1$$

$$a) \text{ five-digit chopping } 0.31415 \times 10^1$$

$$b) \text{ five-digit rounding } 0.31416 \times 10^1$$

△ Def: suppose p^* is an approximation to p

actual error: $p - p^*$;

absolute error: $|p - p^*|$; relative error: $\frac{|p - p^*|}{|p|}$ if $|p| \neq 0$.

Def: Significant digits (figures) $\frac{|p - p^*|}{|p|} \leq 5 \times 10^{-t}$, take largest t

we said p^* approximates p to t significant digits.

△ floating-point: relative error = $\left| \frac{y - f(y)}{y} \right|$

$$y = 0.d_1 d_2 \dots d_k d_{k+1} \dots \times 10^n, \text{ chopping } f(y) = 0.d_1 d_2 \dots d_k \times 10^n$$

$$\left| \frac{y - f(y)}{y} \right| = \frac{0.d_{k+1} d_{k+2} \dots \times 10^{n-k}}{0.d_1 d_2 \dots \times 10^n} = \frac{0.d_{k+1} d_{k+2} \dots}{0.d_1 d_2 \dots} \times 10^{-k}$$

$$\leq \frac{1}{0.1} \times 10^{-k} = 10^{-k+1} \leq 5 \times 10^{-(k-1)}.$$

at least $(k-1)$ significant digits

if $d_{k+1} < 5$, then it has k significant digits, it means when $d_{k+1} < 5$, it more reasonable to chop.

△ finite-digit arithmetic: \oplus \ominus \otimes \oslash : \odot

store first, calculate second, store last

$$x \odot y = f(|f(x) \odot f(y)|)$$

e.g. $x = \frac{5}{7} = 0.\overline{714285}$, $y = \frac{1}{3} = 0.\overline{3}$, five-digit chopping

$$f(x) = 0.71428 \times 10^0, \quad f(y) = 0.33333 \times 10^0$$

$$x \oplus y = f(0.71428 \times 10^0 + 0.33333 \times 10^0) = f(1.04761 \times 10^0)$$

$$= 0.10476 \times 10^0$$

$$\text{absolute error: } \left| \frac{22}{21} - 0.10476 \times 10^0 \right| = 0.190 \times 10^{-4}$$

$$\text{relative error: } \frac{0.190 \times 10^{-4}}{22/21} = 0.182 \times 10^{-4}$$

Try $x \ominus y$, $x \otimes y$, $x \oslash y$

e.g. $u = 0.714251$, $v = 98765.9$, $w = 0.11111 \times 10^{-4}$

$$f(u) = 0.71425 \times 10^0, \quad f(v) = 0.98766 \times 10^5, \quad f(w) = 0.11111 \times 10^{-4}$$

$$x \ominus u = f(0.71428 \times 10^0 - 0.71425 \times 10^0) = f(0.00003 \times 10^0)$$

$$= 0.30000 \times 10^{-4}$$

$$\text{ab. error} = 0.471 \times 10^{-5}, \quad \text{re. error} = 0.136$$

$$(x \ominus u) \oslash w = f\left(\frac{0.30000 \times 10^{-4}}{0.11111 \times 10^{-4}}\right) = f(2.700027) = 0.27000 \times 10^1$$

$$\text{ab. error} = 0.424 \quad \text{re. error} = 0.136$$

$$(x \ominus u) \otimes v = f(0.30000 \times 10^{-4} \times 0.98766 \times 10^5) = f(0.296298 \times 10^1)$$

$$= 0.29629 \times 10^1$$

$$\text{ab. error} = 0.405 \quad \text{re. error} = 0.136$$

△ 4 cases to reduce roundoff error:

△ 4 cases to reduce roundoff error:

① avoid numerator \gg denominator

if $f(z) = z + \delta$, $z = 10^{-n}$

$$z \oplus \varepsilon = f\left(\frac{f(z)}{f(\varepsilon)}\right) = f\left(\frac{z + \delta}{10^{-n}}\right) = (z + \delta) \times 10^n$$

e.g. $x = 0.01$ $\frac{e^x - 1}{x} = 1.0050167 \dots$ $e^{0.01} = 1.010050$
 $= 1.0000$ (five digit chopping)

but $\frac{e^x - 1}{x} \approx 1 + \frac{1}{2}x = 1.0050$

② avoid two nearly equal numbers $x \approx y$ doing $x \ominus y$

$$f(x) = 0.d_1d_2 \dots d_p \alpha_{p+1} \dots \alpha_k \times 10^n$$

$$f(y) = 0.d_1d_2 \dots d_p \beta_{p+1} \dots \beta_k \times 10^n$$

$$x \ominus y = 0.\delta_{p+1}\delta_{p+2} \dots \delta_k \times 10^{n-p}, \quad 0.\delta_{p+1}\delta_{p+2} \dots \delta_k = f(0.\alpha_{p+1} \dots \alpha_k - \beta_{p+1} \dots \beta_k)$$

only $k-p$ digits of significance

e.g. $10 \ominus \sqrt{99}$ (3-digit chopping) $10 - \sqrt{99} \approx 0.0501$

$$f(\sqrt{99}) = 0.994 \times 10^1, \quad f(10) = 0.100 \times 10^2$$

$$10 \ominus \sqrt{99} = f(0.100 \times 10^2 - 0.994 \times 10^1) = f(0.100 \times 10^2 - 0.099 \times 10^2) \\ = f(0.001 \times 10^2) = 0.100 \times 10^0 \quad \text{error: } 0.0499$$

but $1 \oplus (10 \oplus \sqrt{99}) = f\left(\frac{0.100 \times 10^1}{0.199 \times 10^2}\right) = f(0.5025 \times 10^{-1}) = 0.503 \times 10^{-1}$
 error: 0.0002

e.g. $A = 10^7(1 - \cos 2^\circ)$ 4-digit rounding: $A = 0.6092 \times 10^4$

$$f(\cos 2^\circ) = 0.9994 \times 10^0 \quad f(1) = 0.1000 \times 10^1 \quad 10^7 = 0.1000 \times 10^8$$

$$A = f(0.1000 \times 10^8 \times f(0.1000 \times 10^1 - 0.9994 \times 10^1)) = 0.1000 \times 10^4$$

but: $1 - \cos \theta = 2 \sin^2 \frac{\theta}{2}$ $f(\sin 1^\circ) = 0.1745 \times 10^{-1}$

$$f(0.1000 \times 10^8 \times 0.2 \times 10^1 \times (0.1745 \times 10^{-1})) = 0.6090 \times 10^4$$

③ avoid $x \oplus y$ for $x \gg y$

e.g. $A = 52492 + \delta_1 + \delta_2 + \dots + \delta_{1000}$, $\delta_i = 0.1 = 52592 = 0.52592 \times 10^5$

5-digit rounding:

$$52492 = 0.52492 \times 10^5, \quad 0.1 = 0.10000 \times 10^1$$

$$52492 \oplus 0.1 = f1(0.52492 \times 10^5 + 0.00000 \times 10^5) = 0.52492 \times 10^5$$

then $52492 \oplus 0.1 \oplus 0.1 \oplus \dots \oplus 0.1 = 0.52492 \times 10^5$ error = 100

but: $52492 \oplus (0.1 \oplus 0.1 \oplus \dots \oplus 0.1) = f1(0.52492 \times 10^5 + 0.1 \times 10^3)$
 $= f1(0.52492 \times 10^5 + 0.00100 \times 10^5) = 0.52592 \times 10^5$ no error

④ reduce steps:

e.g. evaluation of polynomials

$$f(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$$

direct calculation: multiplications: $n + n-1 + \dots + 1 = \frac{n(n+1)}{2}$

additions: n

(text book)

起源于《算数九章》(汉代)

but with nested arithmetic e.g. 秦九韶算法 (1202年-1261年)

$$f(x) = (\dots (a_n x + a_{n-1}) x + a_{n-2}) x + \dots + a_1) x + a_0$$

n multiplications + n additions

RK: also called Horner's algorithm ≈ 1800

Pseudo code: Input: $a_n, a_{n-1}, \dots, a_0, x$

Output: $f(x)$

Step 1: $S = a_n$

Step 2: For $k = n-1, n-2, \dots, 1, 0$

Set $S = xS + a_k$

Step 3: set $f(x) = S$, Output ($f(x)$)

Stop.

HW1, Sec. 1.2: 1, 3, 5 ac 14

HW1, Sec. 1.2: 1, 3, 5 ac 14
15, 19 20 25 9 28, 29
Only Question (a)