

Chapter 2 Computer Arithmetic

§2.1 Floating-Point Numbers and Roundoff Errors

decimal system

$$427.325 = 4 \times 10^2 + 2 \times 10^1 + 7 \times 10^0 + 3 \times 10^{-1} + 2 \times 10^{-2} + 5 \times 10^{-3}$$

binary system

$$(1001.11101)_2 = 1 \times 2^3 + 0 \times 2^2 + \dots + 1 \times 2^{-3} + 0 \times 2^{-4} + 1 \times 2^{-5} \\ = (9.90625)_{10}$$

rounding

$$\begin{matrix} x & \tilde{x} \\ 0.1735499 & \longrightarrow 0.1735 \\ 0.4321609 & \longrightarrow 0.4322 \end{matrix}$$

$$|x - \tilde{x}| \leq \frac{1}{2} \times 10^{-4}$$

chopped/truncated

$$\begin{matrix} \hat{x} \\ \longrightarrow 0.1735 \\ \longrightarrow 0.4321 \end{matrix}$$

$$|x - \hat{x}| \leq 10^{-4}$$

Normalized Scientific Notation

decimal system

$$732.5051 = 0.7325051 \times 10^3 \\ -0.005612 = -0.5612 \times 10^{-2}$$

$$x = \pm r \times 10^n$$

mantissa \swarrow exponent \nwarrow

$$\frac{1}{10} \leq r < 1$$

binary system

$$x = \pm \beta \times 2^m \quad \text{where} \quad \frac{1}{2} \leq \beta < 1$$

Hypothetical Computer Marc-32

a word length: 32 bits (binary digits)

a nonzero number $x = \pm f \times 2^m$

$= (-1)^s f \times 2^m$ normalized floating-point form

$f = (1.f)_2$ and $m = e - 127$

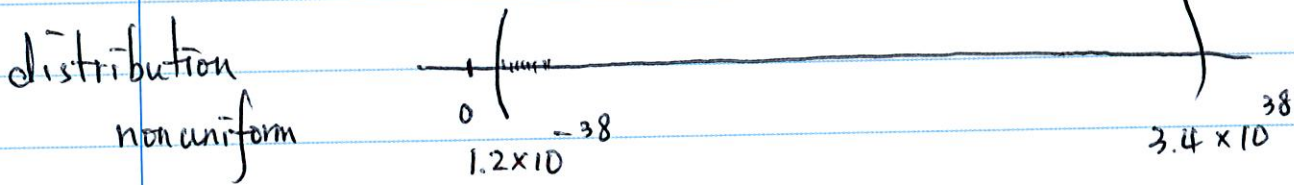
1-bit (pos. 0, neg. 1) *integer* *8-bit* 8-bit biased exponent

23-bit fractional part of real number x

$$0 < e < (11111111)_2 = 2^8 - 1 = 255$$

floating-point range of numbers $\Rightarrow -126 \leq m \leq 127$

$$\left(2^{-126}, (2 \cdot 2^{-23}) 2^{127} \right) \approx \left(1.2 \times 10^{-38}, 3.4 \times 10^{38} \right)$$



Nearby Machine Number (Floating-point number)

x — real number

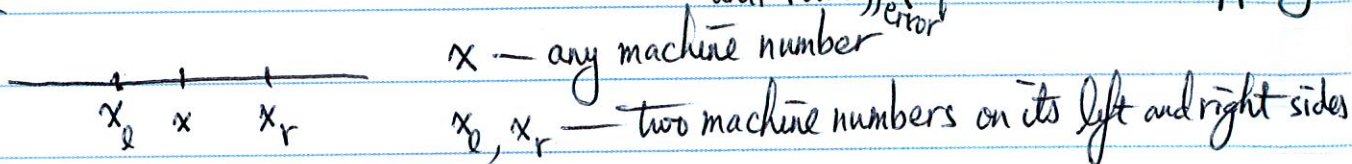
$f(x) = x^*$ — machine number closest to x satisfies $|\delta| = \left| \frac{x - x^*}{x} \right| \leq 2^{-24}$

unit roundoff error

floating-point number base β (decimal or binary)

$$fl(x) = x(1 + \delta) \text{ with } |\delta| \leq \epsilon = \begin{cases} \frac{1}{2}\beta^{1-n} & \text{rounding} \\ \beta^{1-n} & \text{chopping} \end{cases}$$

unit roundoff error



$$2^{-24} < \frac{x_r - x}{x} = \frac{x - x_l}{x} \leq 2^{-23} \quad \text{on Mac-32}$$

Relative Roundoff Error in Adding

x_0, x_1, \dots, x_n — positive machine numbers in a computer with machine ϵ

the ~~total~~ relative roundoff error in computing $\sum_{i=0}^n x_i$

$$\leq (1 + \epsilon)^n - 1 \approx n\epsilon$$

§2.2 Absolute and Relative Errors: Loss of Significance

x — a given real number

x^* — an approx to x

absolute error $|x - x^*|$

relative error $\left| \frac{x - x^*}{x} \right|$

example $\left| \frac{x - f(x)}{x} \right| \leq \epsilon$

Loss of Significance

a decimal computer have a 5-digit mantissa

example $x = 0.3721478693$

$y = 0.3720230572$

$x - y = 0.0001248121$

$f(x) = 0.37215$

$f(y) = 0.37202$

$f(x) - f(y) = 0.00013 = 0.13000 \times 10^{-3}$

$$0.1248121 \times 10^{-3} \left| \frac{(x-y) - (f(x) - f(y))}{x-y} \right| = \left| \frac{0.0001248121 - 0.00013}{0.0001248121} \right| \approx 4\%$$

Subtraction of Nearly Equal Quantities

$$y = \sqrt{x^2 + 1} - 1 = \frac{x^2}{\sqrt{x^2 + 1} + 1}$$

(5)

Thm on Loss of Precision

$x > y$ — positive normalized floating-pt binary machine numbers
satisfy

$$2^{-p} \leq \frac{x-y}{x} \leq 2^{-p}$$

\Rightarrow at most q and at least p significant binary bits are lost in the subtraction $x-y$.

Proof on the lower bound

Let $x = r \times 2^n$

$\frac{1}{2} \leq r, s < 1$ i.e., $r = (0.1\dots)_2$

$y = s \times 2^m$

$\dots 2^{-1} 2^0 2^1 \dots$
 $\frac{1}{2} \quad \quad 1$

$x > y \Rightarrow y = (s \times 2^{m-n}) \times 2^n$

$x-y = (r - s \times 2^{m-n}) \times 2^n$

$\frac{1}{2} \cdot 2^{-q} \leq r - s \times 2^{m-n} = r \left(1 - \frac{s \times 2^m}{r \times 2^n} \right) = r \frac{x-y}{x} < 2^{-p}$
" $2^{-(q+1)}$

#

Example $y = x - \sin x$

(a) $\sin x \approx x$

$y = x - \sin x$
 $= \frac{x^3}{3!} - \frac{x^5}{5!} + \frac{x^7}{7!} - \dots$
 $\approx \frac{x^3}{6} \left\{ 1 - \frac{x^2}{20} \left[1 - \frac{x^2}{42} \left(1 - \frac{x^2}{72} \right) \right] \right\}$

6

(b) By the thrm, $\left| 1 - \frac{\sin x}{x} \right| \geq \frac{1}{2} = 2^{-1} \Rightarrow |x| \geq 1.9$

$y = x - \sin x$ loss at most of one bit

(c) $|x| < 1.9$

at the worst case $x = 1.9$

$$y = x - \sin x \approx \frac{x^3}{3!} - \frac{x^5}{5!} + \frac{x^7}{7!} - \frac{x^9}{9!} + \frac{x^{11}}{11!} - \frac{x^{13}}{13!}$$

has an error of at most 10^{-9} . (Prob 2.2.1)

Evaluation of Functions

for large x , for certain functions f
evaluation $f(x)$ could loss significant digits

range reduction

$$\cos(x + 2n\pi) = \cos x$$

$$\begin{aligned} \Rightarrow \cos 33278.21 &= \cos(33278.21 - 5296 \times 2\pi) \\ &= \cos(2.46) = -0.7765702835 \end{aligned}$$

§2.3 Stable and Unstable Computations: Conditioning

Numerical Instability

a numerical process is unstable

⇔ small errors made at one stage of the process are magnified in subsequent stages and seriously degrade the accuracy of the overall calculation.

example 1

$$\begin{cases} x_0 = 1, x_1 = \frac{1}{3} \\ x_{n+1} = \frac{13}{3}x_n - \frac{4}{3}x_{n-1} \quad \text{for } n \geq 1 \end{cases} \Rightarrow x_n = \left(\frac{1}{3}\right)^n$$

on Marc-32 $x_0 = 1, x_1 = \dots, x_{15} = 3.657493$ (see p65)

⇒ This algorithm is unstable.

general solution $x_n = A\left(\frac{1}{3}\right)^n + B4^n$

$$\left\{ \begin{array}{l} 1 = A + B \\ \frac{1}{3} = \frac{A}{3} + 4B \end{array} \right\} \Rightarrow A = 1 \text{ and } B = 0$$

$$\begin{cases} x_0 = 1, x_1 = 4 \\ x_{n+1} = \frac{13}{3}x_n - \frac{4}{3}x_{n-1} \end{cases} \Rightarrow x_n = 4^n \quad \underline{\text{this is stable}}$$

$$\begin{aligned} x_1 &= 4.000006 \\ x_{10} &= 1.048576 \times 10^6 \\ x_{20} &= 1.099512 \times 10^{12} \end{aligned}$$

example 2

$$y_n = \int_0^1 x^n e^x dx$$

$$u = x^n, v = e^x$$

$$u' = nx^{n-1}, v = e^x$$

$$= x^n e^x \Big|_0^1 - n \int_0^1 x^{n-1} e^x dx$$

$$\Rightarrow \begin{cases} y_n = e^{-n} y_{n-1} \\ y_0 = e^{-1} \end{cases}$$

$$y_2 = 0.7182817$$

$$y_{11} = 1.422453$$

$$y_{15} = 39711.43$$

but these cannot be correct, because

$$1 = y_1 > y_2 > y_3 > \dots > 0 \text{ and } \lim_{n \rightarrow \infty} y_n = 0$$

Conditioning

- evaluation of a function f at a pt x

$$f(x+h) - f(x) = f'(\xi)h \approx hf'(x) \text{ for small } h$$

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

relative perturbation in evaluation of f
condition number
relative perturbation in x

example

$$f(x) = \arcsin x$$

$$\frac{x f'(x)}{f(x)} = \frac{x}{\sqrt{1-x^2} \arcsin x}$$

is large for x near 1.

- compute a simple root of f

$$f(r) = 0 \text{ and } f'(r) \neq 0$$

perturbation of f $F \equiv f + \epsilon g$

root for F : $r+h$ ~~$r+h$~~ $0 = F(r+h) = f(r+h) + \epsilon g(r+h)$
 $= [f(r) + hf'(r) + \frac{1}{2}h^2 f''(\xi)] + \epsilon [g(r) + hg'(r) + \frac{1}{2}h^2 g''(\eta)]$

$$\approx hf'(r) + \epsilon g(r) + hg'(r)$$

$$\Rightarrow h \approx -\epsilon \frac{g(r)}{f'(r) + \epsilon g'(r)} \approx -\epsilon \left(\frac{g(r)}{f'(r)} \right) \sim \text{condition number}$$

$$h = (r+h) - r$$

$$\epsilon = (F - f) / g$$

example

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

$g(x) = x^{20}$ means an approximation error in the coefficient of x^{20}

\Rightarrow for root $r=20$

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

Condition Number for $A_{n \times n}$

$$\kappa(A) = \|A\| \cdot \|A^{-1}\|$$

Hilbert matrix $H_n = (h_{ij})_{n \times n}$ $h_{ij} = \frac{1}{i+j-1}$, $1 \leq i, j \leq n$.

$$\kappa_{\infty}(H_n) = \|H_n\|_{\infty} \|H_n^{-1}\|_{\infty} = c e^{3.5n}$$

$$\|A\|_{\infty} = \max_{i,j} |a_{ij}|$$