

# Computer Arithmetics <sup>A</sup>

$\beta, N, L, U$

$$x = \pm (d_1 + d_2/\beta + \dots + d_N/\beta^{N-1}) \cdot \beta^E$$

$E, d_i$  is integer

$$0 \leq d_i \leq \beta - 1$$

$$L \leq E \leq U$$

Finite # of digits

$$\left| \frac{x - FE(x)}{x} \right| \leq \frac{\epsilon_{mach}}{2} \quad \text{normalized \#}$$

$$\epsilon_{mach} = \beta^{1-N}$$

$$f(x) = \frac{1 - (1-x)}{x}$$

~~$f(x)$~~ 

$$f\left(\frac{\epsilon_{mach}}{4} \cdot 0.99\right)$$

$$+ f\left(\frac{\epsilon_{mach}}{4} \cdot 1.01\right)$$

- loose digits

\*

/ -|| a lot of 1-1-1-1

integers

---

vpa

sym

# are represented

C

uniquely

$$13 = 1.3 \cdot 10 = 0.13 \cdot 10^2$$

Agreement  $d_i \neq \emptyset$

Binaries  $d_i = 1$

Subnormals

$d_i$  may be  $\emptyset$  if  $E = L$

$$N=3, \beta=?, L=-1, u=1$$

1.00

1.01

1.10

1.11

\* ~~10~~  $^{-1, \emptyset, 1}$

$$u - L = 2$$

D

Subnormals

$d_i$  may be  $\emptyset$

$$\text{IF } E = L$$

0.01

0.10  $\times 10^{-1}$

0.11

Subnormals

$1 \cdot 10^{-1}$  - smallest  
normalizer

$0.01 \cdot 10^{-1}$

What is smallest #

←

normalized #

realmin

$$X = 1.00 \cdot 10^4$$

$$X = \beta^k$$

Smallest subnormal

$$y = \beta^k \cdot \beta^{1-N}$$

Fit

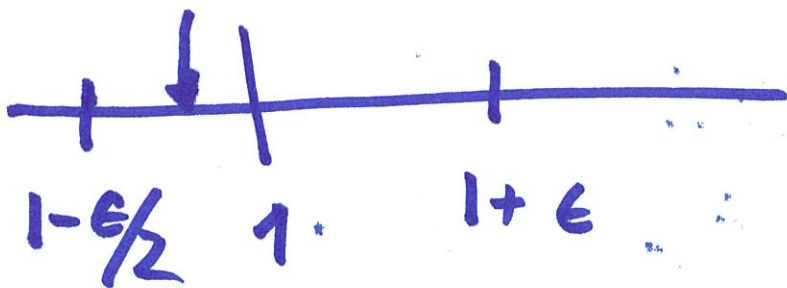
Smallest possible  
that can not be  
represented

$$\underbrace{11111}_{54} \approx \underbrace{1000}_{53}$$

$$f(x) = \frac{1 - (1-x)}{x}$$

F

- $x_1 = \epsilon^{1/4} \cdot 0.99$



$$1 - x \approx 1; f(x_1) = \emptyset$$

- $x_2 = \epsilon^{1/4} \cdot 1.01$

$$1 - x_2 \approx 1 - \epsilon^{1/2}$$

$$f(x_2) = \frac{1 - (1 - \epsilon^{1/2})}{\epsilon^{1/4} \cdot 1.01} \approx 2$$

G



double precision

H

8 bytes = 64 bits

$N = 53$  ← normalized

$E = 11$  bits  $d_1 = 1, 52$

1 bit for sign

$52 + 11 + 1$

∅.. 2047,  $\textcircled{?}$

-1023 ;



Total # of #

I

3-i)

$$2 \times 3^{n-1} \times (4-L+1) + 1$$

zero

↑  
sign

↑  
mantissa

↑  
Exponent

Finding root of

3

$$f(x) = \emptyset$$

- Bisection  $r=1$
- Newton  $r=2$
- Secant  $1 < r < 2$

"golden ratio"

$$\epsilon_{k+1} = C \epsilon_k^r$$

↑ error at step  $k+1$

↘ error at step  $k$

$$m = \frac{a+b}{2};$$

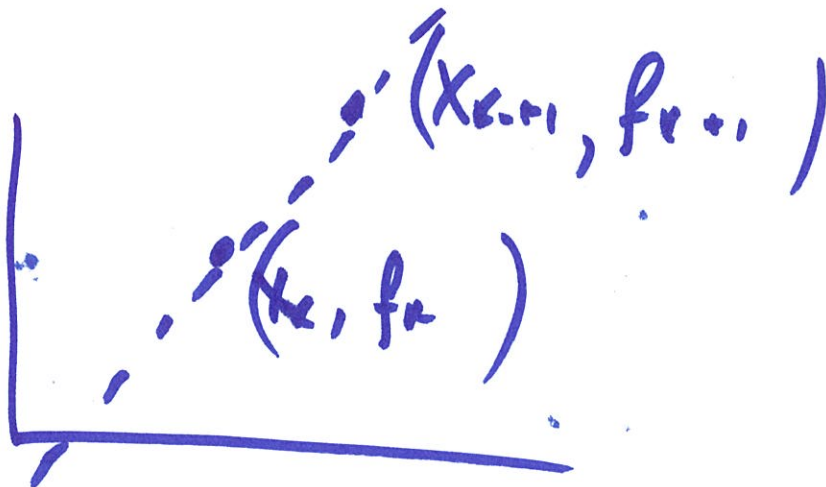
IF  $f(a)f(m) < \emptyset,$

$b = m$  else  $a = m$

END IF

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

$$x_{k+1} = x_k - \frac{f(x_k) \cdot (x_k - x_{k-1})}{f(x_k) - f(x_{k-1})}$$



1.12  $\epsilon_1 = 0.12$  Newton  
 1.01  $\epsilon_2 = 0.01$   $f'(x^*) \neq 0$   
 1.00  $\epsilon_3 = 0$   
 1.00 ←

1.16

1.08

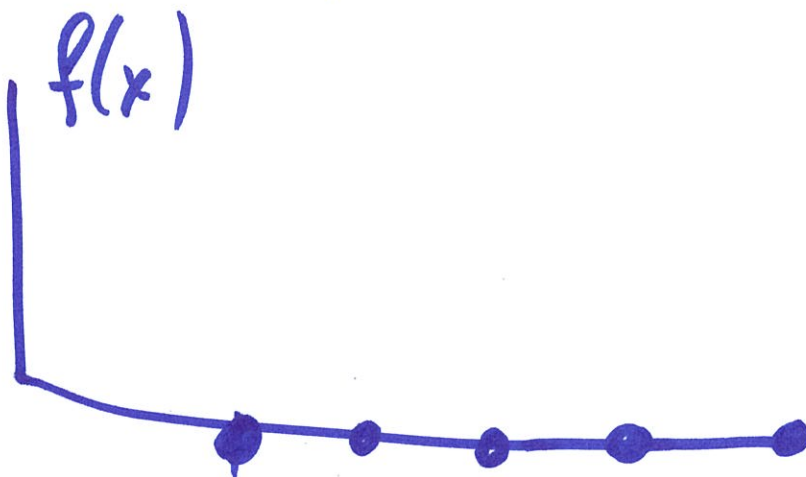
1.04

1.02

1.01

0.005

$$\epsilon_r = \epsilon_{r-1} / 2$$



$$f(x) = x^2$$

$$f'(x) = 2x$$

$$\epsilon_{r+1} = \epsilon_r / 2$$

$$x_{r+1} = x_r - \frac{x_r^2}{2x_r} = \frac{x_r}{2}$$

$$f(x)$$

$$f'(x) = \emptyset$$

∴ ∴ ↗ upper  
LL decomposition

M

↳ lower

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix}$$

$M_1$  eliminate below  $a_{11}$

$$M_1 A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \cancel{\phi} & a'_{22} & \dots & a'_{2n} \\ \cancel{\phi} & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ \cancel{\phi} & a'_{n2} & \dots & a'_{nn} \end{pmatrix}$$

$M_2$

↳ "eliminates"  $a'_{22}$

$$M_2 M_1 A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \cancel{\phi} & a'_{22} & \dots & a'_{2n} \\ \cancel{\phi} & \cancel{\phi} & a''_{33} & \dots & a''_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \cancel{\phi} & \cancel{\phi} & \cancel{\phi} & \dots & a''_{nn} \end{pmatrix}$$

$$(AB)^{-1} = B^{-1}A^{-1} \quad \mathcal{N}$$

$$LI = (M_{n-1} \dots M_3 M_2 M_1) A$$

$$M = M_{n-1} M_{n-2} \dots M_3 M_2 M_1$$

$$M^{-1}(LI = MA)$$

$$A = M^{-1}LI$$

$$M^{-1} = (M_{n-1} M_{n-2} \dots M_3 M_2 M_1)^{-1}$$

$$= \underbrace{M_1^{-1} M_2^{-1} M_3^{-1}} = L$$

$$A = LLI$$



$$\underline{A} \underline{x} = \underline{b}$$

$$A = LU$$

$$\underline{L} \underline{U} \underline{x} = \underline{b}$$

z

$\underline{L} \underline{z} = \underline{b} \Leftarrow$  Forward substitution

$\underline{U} \underline{x} = \underline{z} \Leftarrow$  Backward substitution



# Example

p

$$\begin{pmatrix} 3 & 4 \\ 6 & 9 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 11 \\ 24 \end{pmatrix}$$

A

$$\begin{pmatrix} 3 & 4 \\ 6 & 9 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 11 \\ 24 \end{pmatrix}$$

$$M = \begin{pmatrix} 1 & 0 \\ -2 & 1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 0 \\ -2 & 1 \end{pmatrix} \begin{pmatrix} 3 & 4 \\ 6 & 9 \end{pmatrix} = \begin{pmatrix} 3 & 4 \\ 0 & 1 \end{pmatrix}$$

$$M \cdot A = U$$

$$\begin{pmatrix} 3 & 4 \\ 6 & 9 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} 3 & 4 \\ 0 & 1 \end{pmatrix}$$

$$\begin{pmatrix} 3 & 4 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 11 \\ 2 \end{pmatrix}$$

$$x_2 = 2$$

$$3x_1 + 4x_2 = 11$$

$$x_1 = \frac{1}{3}(11 - 4x_2) = \frac{1}{3}(11 - 8) = 1$$

$$x_1 = 1$$

$$x_2 = 2$$

$$\underline{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} 3 & 4 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 11 \\ 24 \end{pmatrix} \quad Q$$

Z

$$\begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} 11 \\ 24 \end{pmatrix}$$

$$z_1 = 11$$

$$z_2 = 24 - 2z_1 = 24 - 22 = 2$$

$$\underline{z} = \begin{pmatrix} 11 \\ 2 \end{pmatrix}$$

~~$$\begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} 3 & 4 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} 11 \\ 24 \end{pmatrix}$$~~

$$\begin{pmatrix} 3 & 4 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 11 \\ 2 \end{pmatrix}$$

~~$$x_2 = 2; \quad x_1 = \frac{1}{3}(11 - 4x_2) = \frac{1}{3}(11 - 8) = 1$$~~



# Example

$$\begin{pmatrix} 1 & 2 & 3 \\ 2 & 5 & 8 \\ 3 & 10 & 20 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 10 \\ 24 \\ 49 \end{pmatrix}$$

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 5 & 8 \\ 3 & 10 & 20 \end{pmatrix}$$

$$M_1 = \begin{pmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ -3 & 0 & 1 \end{pmatrix}$$

$$M_1 A = \begin{pmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ -3 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 & 3 \\ 2 & 5 & 8 \\ 3 & 10 & 20 \end{pmatrix}$$

$$= \begin{pmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 4 & 11 \end{pmatrix}; M_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -4 & 1 \end{pmatrix}$$

$$M_2 M_1 A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 4 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 4 & 11 \end{pmatrix} S$$

$$= \begin{pmatrix} 1 & 2 & 3 \\ 0 & -1 & 2 \\ 0 & 0 & 3 \end{pmatrix} = I$$

$$A = (M_2 M_1)^{-1} I = M_1^{-1} M_2^{-1} I$$

$$= \begin{pmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 3 & 4 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 3 \end{pmatrix}$$

$$= \begin{pmatrix} 1 & 2 & 3 \\ 2 & 5 & 8 \\ 3 & 10 & 20 \end{pmatrix}$$