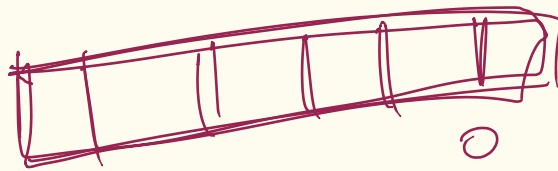


IEEE: 64 bits

1 sign bit;
52 bits for mantissa
11 bits for exp

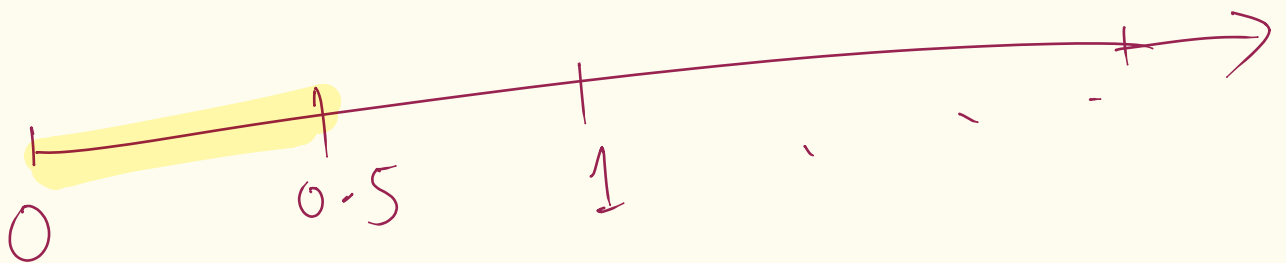
$$\pm (0.1 \underbrace{d_1 d_2 \dots d_{52}}_{\text{mantissa}})_2 \times 2^e$$



0 . . . 1
1 1 . . . 1

0
2047
0

smallest positive: $(0.100\dots0) \times 2^0$
 $= 0.5$



exp range $[0, 2047]$
bias $-1022 \rightarrow$ $[-1022, +1025]$
0 ∞

example:

$x = \text{distance}$
 $f(x) = x \quad \text{distance}$

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

unitless

$$|8| < 10^{-10}$$

error is less than 10^{-10}

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

$$< \frac{10^{-10}}{10^{-10}}$$

$$\frac{0.1 + 0.2}{0.3}$$

$$0.3$$

$$(0.1 + 0.2 - 0.3) \times 10^{-17}$$

$$= 5.55111 \dots \times 10^{-17}$$

$$= 5.55 \times 10^{-17}$$

Bounding the relative rounding error:

$$x = (0.\underbrace{d_1 d_2 d_3 \dots d_m}_{\text{exact}} d_{m+1}) \times \beta^e$$

$$\text{round} \downarrow$$

$$fl(x) = (0.\underbrace{d_1 d_2 \dots d_m}_{\text{exact}}) \times \beta^e$$

$$x = (0.\underbrace{314159265}_{\text{round to 9}}) \times 10^1$$

$$(0.\underbrace{3142}) \times 10^1$$

$$fl(x) - x = \left((0.d_1 d_2 \dots d_m) - (0.d_1 d_2 \dots) \right) \times \beta^e$$

$$= \left(0.\underbrace{000 \dots 0}_{m \text{ digits}} \boxed{} \right) \times \beta^e$$

$$\leq \left(0.\underbrace{000 \dots 0}_{m \text{ digits}} \frac{\beta}{2} \right) \times \beta^e$$

$$= \frac{\beta}{2} \times \beta^{-(m+1)} \times \beta^e$$

$$= \frac{1}{2} \times \beta^e \times \beta^{-m}$$

$x = 0.31415$
 round to 4 digits
 0.3142

$\Rightarrow \text{error} = 0.00005$

always

base 10

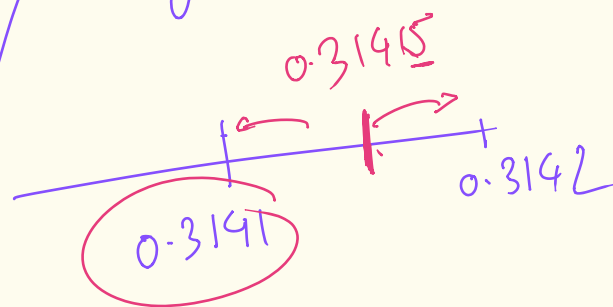
β

≤ 5

next digit ≤ 5

$\dots \leq \frac{\beta}{2}$

can this be larger than 5?



So far,

$$|fl(x) - x| \leq \frac{1}{2} \beta^e \cdot \beta^{-m}$$

divide by $|x|$.

$$x = (0.d_1 d_2 d_3 \dots) \times \beta^e$$

lowest

$$x = (0.1 0 0 0 0) \times \beta^e$$

$$= 1 \times \beta^{-1} \times \beta^e$$

$$x \geq \beta^{-1} \times \beta^e$$

①

②

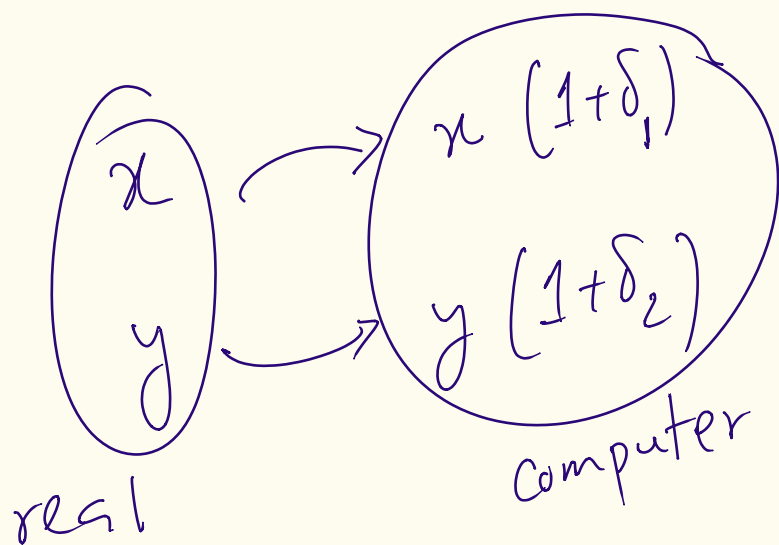
$$\left| \frac{fl(x) - x}{x} \right| \leq \frac{\frac{1}{2} \beta^e \beta^{-m}}{\beta^{-1} \beta^e} = \frac{1}{2} \beta^{1-m}$$

maximum
rounding
error

Machine epsilon $\epsilon_M = \frac{1}{2} \beta^{1-m}$

IEEE: $\beta = 2, m = 52$

$$\epsilon_M = \frac{1}{2} \cdot 2^{1-52} = 2^{-52} = 2.22 \times 10^{-16}$$



$$\begin{aligned} \frac{fl(x) - x}{x} &= \delta \\ fl(x) - x &= \delta x \\ fl(x) &= x + \delta x \\ &= x(1 + \delta) \end{aligned}$$

$$\begin{aligned} &x(1 + \delta_1) + y(1 + \delta_2) \\ &= x + y + x\delta_1 + y\delta_2 \end{aligned}$$

relative error:

$$\frac{(x+y+x\delta_1+y\delta_2) - (x+y)}{(x+y)}$$

$$= \frac{x\delta_1 + y\delta_2}{x+y}$$

subtraction:

$$\frac{x\delta_1 - y\delta_2}{x-y}$$

0.00001
0.00002

loss of significance:
when the relative error $\frac{x\delta_1 \pm y\delta_2}{x \pm y}$
becomes larger than ϵ_M .

$$x^2 - 56x + 1 = 0$$

$$x = 28 \pm \sqrt{783}$$

$$28 + \sqrt{783}$$

$$= 55.9821$$

$$= 55.98 = 0.5598 \times 10^2$$

$$28 - \sqrt{783}$$

$$= 0.0178628$$

$$= 0.01786$$

$$= 0.1786 \times 10^{-1}$$

$$28 + \sqrt{783}$$

$$\rightarrow 27.9821$$

$$\downarrow$$

$$27.98$$

$$28 + 27.98$$

$$= 55.98$$

~~$$28 - \sqrt{783}$$

$$= 28 - 27.98$$

$$= 0.02$$~~

$$x_1 \times x_2 = 1$$

$$x_2 = \frac{1}{x_1}$$

$$= \frac{1}{55.98}$$

$$= \underline{\underline{0.01786}}$$

2 digit

$$5.9 + 5.5 + 0.4$$

$$(5.9 + 5.5) + 0.4$$

$$= 11 + 0.4$$

$$= 11$$

$$5.9 + (5.5 + 0.4)$$

$$= 5.9 + 5.9$$

$$= 12$$

5.01, 5.02

(3 digit)

avg:
$$\frac{5.01 + 5.02}{2} = \frac{10.0}{2} = 5.0$$

Assignment 1:

Deadline: 19 February

Quiz 1: 19 February

in-person

→ na17.pdf chapter 1,
floating point arithmetic note [error]

$$\left. \begin{array}{l} x = 0.31415 \\ \downarrow \\ \bar{x} = 0.3142 \end{array} \right\} \begin{array}{l} \bar{x} - x \\ = (0.3142) - (0.3141\cancel{0}) \\ = 0.0000\textcircled{5} \end{array}$$

$$\text{error} = 0.00005$$

$$\begin{array}{l} x = 0.31414 \\ \downarrow \\ \text{round } 0.3141 \end{array}$$

$$\begin{array}{lcl} \text{base } 10 & \longrightarrow & 5 \\ \text{" } \beta & \longrightarrow & \beta/2 \end{array}$$