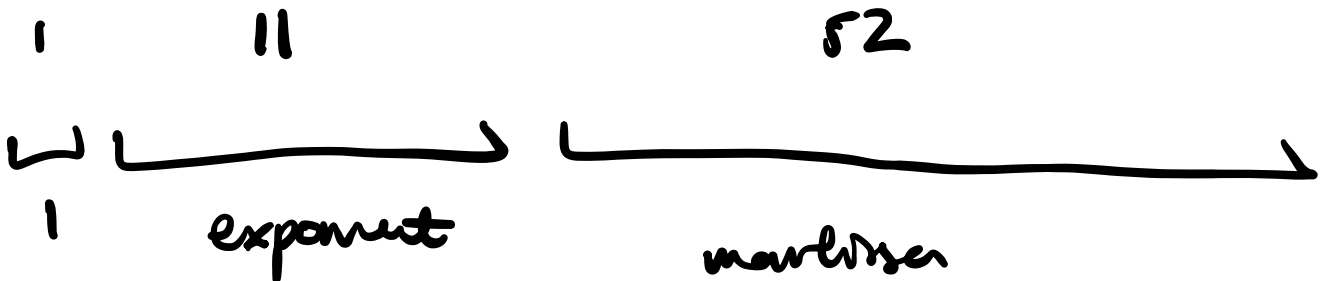


Floating Point Numbers

(64 bit machines)



$$\Rightarrow (-1)^{\text{sign}} \times 1.\underbrace{\dots\dots\dots}_{\text{mantissa}} \times 2^{\text{exponent} - 1023}$$

example 1.75 as Float 64

sign = 0

mantissa = 0.75 = $2^{-1} + 2^{-2}$

\Rightarrow 0.11 (in binary)

exponent = 1023 = 011111111111 (in binary)

$$\Rightarrow 1.75 = 0 \ 011111111 \ 1100 \dots 0$$

$\underbrace{\hspace{10em}}$
zeros

as a Float 64

Def. Machine Precision ϵ

→ The smallest representable floating point number s.t.

$$1 + \epsilon \neq 1$$

$$\begin{aligned} \text{In FP: } 1 &= 0 \ 011111111 \ 00 \dots 00 \\ 1 + \epsilon &= 0 \ 011111111 \ 00 \dots 01 \end{aligned}$$

$\downarrow +1$

$$\Rightarrow \epsilon = 2^{-52} \approx 10^{-16}$$

FP Operations

$$f(x \text{ op } y) = (x \text{ op } y) (1 + \delta)$$

\uparrow \uparrow

$\{+, \times, -, \div\}$ $|\delta| < \epsilon$

eg. Nested operations

\otimes is fp multiplication

$$(a_1 \otimes a_2) \otimes a_3$$

$$= (a_1 \times a_2) (1 + \delta_1) \otimes a_3, \quad |\delta_1| < \epsilon$$

$$= \left[(a_1 \times a_2) (1 + \delta_1) + a_3 \right] (1 + \delta_2), \quad |\delta_2| < \epsilon$$

$$= (a_1 \times a_2) (1 + \delta_1) (1 + \delta_2) + a_3 (1 + \delta_2)$$

$$= (a_1 \times a_2) (1 + \delta_1 + \delta_2 + \cancel{\delta_1 \delta_2}) + a_3 (1 + \delta_2)$$

$\sim \epsilon^2$, ignore

$= \delta_3, |\delta_3| < 2\epsilon$

$$= (a_1 \times a_2) (1 + \delta_3) + a_3 (1 + \delta_2)$$

- Two operators, get error of

$|\delta_3| < 2\epsilon$, makes sense

Ex 3 Notation and Order of Convergence

Def. If $f(h) \rightarrow q$ as $h \rightarrow 0$, and
 $|f(h) - q| \leq Kh^p$ for some $K \neq 0$,
then $f(h) \rightarrow q$ with order of convergence p .
(Equivalently, converges as $O(h^p)$)

eg. $\cos(2h) + 2h \sin(h)$

$$\cos(x) = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \dots$$

$$\sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots$$

$$\Rightarrow \cos(2h) = 1 - \frac{(2h)^2}{2!} + \frac{(2h)^4}{4!} - \dots$$

$$\cos(2h) = 1 - \frac{(2h)^2}{2!} + \frac{(2h)^4}{4!} - \dots$$

$$\Rightarrow \cos(2h) + 2h \sin(h)$$

$$= 1 - \frac{(2h)^2}{2!} + 2h^2 + \frac{(2h)^4}{4!} - \frac{2h^4}{3!} + \dots$$

$$= 1 + \left(\frac{2^4}{4!} - \frac{2}{3!} \right) h^4 + \dots$$

$$\Rightarrow \cos(2h) + 2h \sin(h) \rightarrow 1$$

with rate $\mathcal{O}(h^4)$