

Machine precision is the distance between 1 and the next fp number. In this

case,

$$1 = (0 | E=0 | \overbrace{1.0 \dots 0}^{23 \text{ bits}})$$
$$1 + 2^{-23} = (0 | E=0 | 1.0 \dots 01)$$
$$\hookrightarrow \approx 1.2 \times 10^{-7}$$

Usually calculations are done in double precision:

sign: 1 bit
exponent: 11
mantissa: 52

machine precision: $\epsilon_{\text{ps}} = 2^{-52} \approx 2.2 \times 10^{-16}$

How to calculate machine precision? Demo

Floating point representations were standardized by

the IEEE in 1980s

\hookrightarrow Inst. of Electrical & Electronics Engin.

Consistent Rules

- Representation (# bits)
- Correct / consistent rounding
- sensible treatment of exceptions: $\frac{1}{6}$

Hidden bit rep.: 0, NaN, $\pm\infty$ (over/underflow)

single precision: 32 bits ($\epsilon \sim 10^{-7}$)

double precision: 64 bits ($\epsilon \sim 10^{-16}$)

extended precision: 80 bits ($\epsilon \sim 10^{-20}$) (no hidden bits)

Rounding

Example: Decimal notation: $\frac{1}{10} = 0.1$

$$\frac{1}{3} - \frac{1}{4} = \frac{1}{12} = \frac{1}{16}$$

Binary notation: $\frac{1}{10} = (1.100\overline{1100})_2 \times 2^{-4}$

$$\frac{1}{12} - \frac{1}{16} = \frac{4}{48} - \frac{3}{48}$$

$$= \frac{1}{48} = \frac{1}{34}$$

How is this rounded?

4 options: Up
Down
to 0
nearest.

$$\sqrt{10} = 1.100 \overline{1100} \times 2^{-4}$$

$$= 1.100 1100 1100 1100 \dots 1100 1$$

(down, toward 0)

$$= 1.100 1100 1100 1100 \dots 1101 0$$

(up, nearest)

↓ default

Absolute rounding error = $|\text{round}(x) - x|$

Relative rounding error = $\frac{|\text{round}(x) - x|}{x}$

Most important rule; IEEE says

$$\text{Round}(a+b) = a \oplus b$$

↑
exact
addition

↑
computer
addition

$$= (a+b) / (1+\delta)$$

↑

$\leq \epsilon$

machine
precision

