

University of IdahoFinite Precision ArithmeticNumber systems

$$\begin{aligned}
 x &= \pm (d_n d_{n-1} \dots d_1 d_0 \cdot d_{-1} d_{-2} \dots)_{\beta} = \\
 &= \pm d_n \cdot \beta^n + d_{n-1} \cdot \beta^{n-1} + \dots + d_1 \cdot \beta^1 + d_0 \cdot \beta^0 + \\
 &\quad + d_{-1} \cdot \beta^{-1} + d_{-2} \cdot \beta^{-2} + \dots, \quad d_n \neq 0
 \end{aligned}$$

β : base, d_i : digits, $0 \leq d_i \leq \beta - 1$

Ex $\beta = 10$: decimal system

$$(2015)_{10} = 2 \cdot 10^3 + 0 \cdot 10^2 + 1 \cdot 10^1 + 5 \cdot 10^0$$

University of Idaho

$\beta = 2$: binary system

$$(10101.01)_2 = 1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 + \\ + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} = 16 + 4 + 1 + \frac{1}{4} = (21.25)_{10}$$

Q: how are real numbers represented in a computer?

Answer: floating point representation

$$x = \pm (0.d_1d_2\dots d_n)_\beta \cdot \beta^e \quad ; \quad x \text{ has } n \text{ significant digits}$$

$d_1 \neq 0$

$0.d_1d_2\dots d_n$: mantissa

e : exponent $-M \leq e \leq M$

University of Idaho

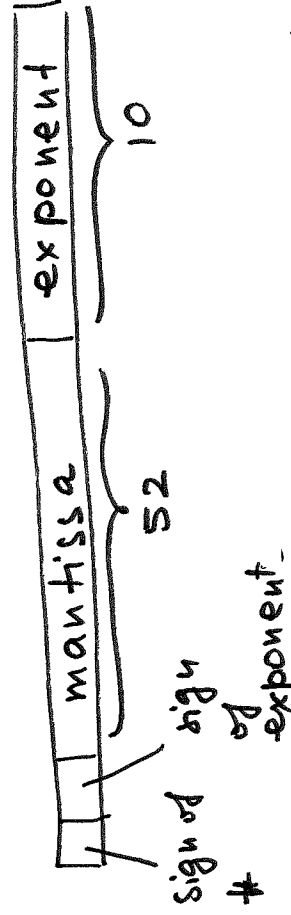
Ex Consider a computer with $\beta=2$, $n=4$, $M=3$

$$x_{\max} = (0.1111)_2 \cdot 2^3 = (111.1)_2 = 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 + 1 \cdot 2^{-1} = \\ = (4+2+1+0.5)_{10} = (7.5)_{10}$$

$$x_{\min} = (0.1000)_2 \cdot 2^{-3} = (0.0001)_2 = 1 \cdot 2^{-4} = (0.0625)_{10}$$

Note

1. In IEEE double precision arithmetic, each number is stored in memory as a string of 64 bits



The first two bits are used to store signs of a number and of exponent, 52 bits for mantissa,

and remaining 10 for exponent.

Hence, we have $\beta = 2$, $n = 52$

$$M = (1111111111)_2 = 2^{10} - 1 = 1023$$

Aside

$$\begin{aligned} 11_2 &= 1 \cdot 2^1 + 1 \cdot 2^0 \\ &= 3 = 2^2 - 1 \end{aligned}$$

2. If x is a real number and $f(x)$ is its floating point representation, then $x - f(x)$ is called the roundoff error.

Number x can be chopped or rounded.

$$f|_{\text{chop}}(x) = \pm (0.d_1d_2 \dots d_n)_\beta \cdot \beta^e$$

or number x can be rounded

$$f|_{\text{round}}(x) = \begin{cases} \pm (0.d_1d_2 \dots d_n)_\beta \cdot \beta^e & \text{if } d_{n+1} < \frac{\beta}{2} \\ \pm [(0.d_1d_2 \dots d_n)_\beta + \beta^{-n}] \cdot \beta^e & \text{if } d_{n+1} > \frac{\beta}{2} \end{cases}$$

x : exact

\tilde{x} : approximation

$x - \tilde{x}$: error

$|x - \tilde{x}|$: absolute error

$\frac{|x - \tilde{x}|}{|x|}$: relative error

Note: it can be shown that when a number is rounded, the bounds on both absolute and relative errors due to roundoff are one-half the bounds when a number is chopped.

Ex $\pi = (3.14159265358979\dots)_{10} = 1.\underline{2}^1 + 1.\underline{2}^0 + 0.2^{-1} + 0.2^{-2} + 1.2^{-3} + 1.2^{-6} + 1.2^{-12} + \dots = 2^{-1} = \frac{1}{2} = 0.5$
 $\frac{1}{8} \quad \frac{1}{64} \quad \frac{1}{4096}$

University of Idaho

$$= (11.0010010000010\dots)_2 =$$
$$= (0.\overbrace{110010010000010}^{\text{with rounding}})_2 \cdot 2^2$$

$$n=4 \Rightarrow fl(x) = (0.1101)_2 \cdot 2^2 = 3.25 : \text{closest 4-bit representation of } \pi$$

In reality, $n=52$, the roundoff error in $fl(x)$

$$\text{is } 2^{-52} \approx 10^{-15}$$