

# Floating Point Numbers

Sunday, August 23, 2020

1:35 PM

Updated 8/22/25

How can we represent numbers on a computer?

Integers are easy.  $6 = 1 \cdot 4 + 1 \cdot 2 + 0 \cdot 1$ , ie., 110 (binary)

What about fractions?

$$\begin{array}{l} \text{5 bits} \rightarrow \frac{18}{19} + \frac{22}{23} = \frac{18 \cdot 23 + 22 \cdot 19}{19 \cdot 23} = \frac{414 + 418}{437} = \frac{832}{437} \end{array}$$

$\underbrace{\hspace{1.5cm}}_{10 \text{ bits}} \quad \underbrace{\hspace{1.5cm}}_{10 \text{ bits}} \quad \underbrace{\hspace{1.5cm}}_{19 \text{ bits}}$

Problem: memory increases

Could try decimals (fixed pt.)

$$\begin{array}{c} 131.467 \\ \hline 3 \quad 3 \end{array} \quad \text{or in binary} \quad \begin{array}{ccccccc} & 2^2 & 2^1 & 2^0 & 2^{-1} & 2^{-2} & \dots \\ \dots & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \dots \\ 1 & 0 & 0 & 1 & 1 & 0 & . & 1 & 0 & 0 & 1 & 1 & 0 \end{array}$$

Done in embedded system

Instead, the standard for numerical computing, is floating pt.

Usually use IEEE 754 "double precision"  $\rightarrow 64 \text{ bits} = 8 \text{ bytes}$   
(single precision  $\approx 32 \text{ bits}$ )

Store numbers

$$\mathbb{F} = \text{Floating Pt.} \quad (-1)^s (1 + f) \cdot 2^e \quad \begin{array}{l} e = e - 1023 \\ 2'' = 2048 \end{array}$$

$\mathbb{R}$  = real numbers  
 $s$  = sign bit (1 bit)  
 $e$  = exponent or characteristic, 11 bits  
 $f$  = mantissa, 52 bits

(scientific notation)

Also includes 0, NaN (%, 0.0,  $\infty/\infty$ ),  $\pm \infty$

Rule of thumb: precision  $\approx 2^{52} = 4.5 \cdot 10^{15}$

15 digits of precision in double  
8 digits ... in single

Implications

① We can't represent very large (or very negative) numbers

$$x \in \mathbb{F}, \text{ then } |x| \leq 2^{1024} = 10^{308}$$

ie.,  $x = -10^{400}$  is not in  $\mathbb{F}$  (it is  $-\infty$ )

Overflow if not in range

② we can't get too small in magnitude (close to 0)

$$x \in \mathbb{F}, |x| \geq 2^{-1022} \approx 10^{-308}$$

underflow if  $|x| < 2^{-1022}$

$$x = 2^{-1023}$$

$$x == 0 \rightarrow \text{True}$$

③ limit to spacing  $\leftarrow$  relative "Machine epsilon", SILENT ERROR

$1 \approx 1 + \epsilon$ , ie., 1 and  $1 + \epsilon$  are indistinguishable

$$\text{true if } \epsilon < \epsilon_{\text{machine}} = 2^{-52} \approx 2.2 \cdot 10^{-16}$$

$$2 \text{ vs } 2 + \epsilon, \quad \epsilon < 2 \cdot \epsilon_{\text{machine}}$$

NOTE: youtube video is wrong, it's relative

Notation:  $x \in \mathbb{R}$ ,  $\text{fl}(x)$  is nearest number to  $x$  that's in  $\mathbb{F}$

Notation:  $x \in \mathbb{R}$ ,  $f(x)$  is nearest number to  $x$  that's in  $\mathbb{F}$

$$\frac{|f(x) - x|}{|x|} \leq \frac{1}{2} \epsilon_{\text{machine}}, \quad f(x) = x \cdot (1 - u), \text{ some } |u| \leq \frac{1}{2} \epsilon_{\text{machine}}$$

depends on precise definition

$$\frac{|f(x) - x|}{|x|} \left\{ \begin{array}{l} \text{relative error} \\ \text{accuracy} \end{array} \right. \quad |f(x) - x| = \text{absolute error}$$

$$\text{digits of accuracy} \approx -\log_{10}(\text{rel. error})$$

**Precision**: # digits, need not be correct!

**Accuracy**: relative error, limited by precision

More implications: lose associative rule

$$(a+b)+c = a+(b+c)$$

$$\left( \underbrace{1 + \epsilon_{\text{machine}}/2}_1 \right) - 1 \neq 1 + \underbrace{(\epsilon_{\text{machine}}/2 - 1)}_{\epsilon_{\text{machine}}/2 \approx 1.11 \cdot 10^{-16}}$$

0

### Catastrophic Cancellation

Pretend we have 3 digits of precision.

$$x = 1.23 \overbrace{587325}^{\text{garbage}} \dots$$

$$y = 1.22 \overbrace{3581926}^{\text{garbage}} \dots$$

$$x - y = 0.01 \overbrace{\dots\dots\dots}^{\text{garbage}}$$

"Silent error"!

$$= 1 - \boxed{\dots\dots\dots} \cdot 10^{-2}$$

you might think these are accurate since within our precision.