

Floating Point Numbers

Sunday, August 23, 2020

1:35 PM

How can we represent numbers on a computer?

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

What about fractions?

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

$\underbrace{18}_{10 \text{ bits}} \quad \underbrace{22}_{10 \text{ bits}} \quad \underbrace{832}_{19 \text{ bits}}$

Problem: memory increases

Could try decimals (fixed pt.)

$$\begin{array}{c} 131.467 \\ \hline 3 \quad 3 \end{array}$$

or in binary 100110.100110

$\begin{array}{ccccccc} 2^2 & 2^1 & 2^0 & 2^{-1} & 2^{-2} & \dots \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \dots \end{array}$

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 = 32 bits)

Store numbers

\mathbb{F} = Floating Pt. $(-1)^s (1+f) \cdot 2^e$ $e = e - 1023$

\mathbb{R} = real numbers

s = sign bit (1 bit)
 e = exponent or characteristic, 11 bits $2^{11} = 2048$
 f = mantissa, 52 bits

(scientific notation)

Also includes 0, NaN (%, 0.00, 0/0), $\pm \infty$

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 ← relative complicated

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

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

$$2 \text{ vs } 2 + \epsilon, \quad \epsilon < \frac{1}{2} \epsilon_{\text{machine}}$$

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

$$\underbrace{|f(x) - x|}_{\text{absolute error}} \leq \frac{1}{2} \epsilon_{\text{machine}}, \quad f(x) = x \cdot (1 + \epsilon), \text{ some } |\epsilon| \leq \frac{1}{2} \epsilon_{\text{machine}}$$

$$\frac{|f(x) - x|}{|x|} \left. \vphantom{\frac{|f(x) - x|}{|x|}} \right\} \text{relative error} = \text{accuracy}$$

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

Precision: #digits, need not be correct!

✓ Accuracy: relative error, limited by precision

More implications: lose associative rule

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

$$\underbrace{\left(1 + \epsilon_{\text{machine}}/2\right)}_1 - 1 \neq 1 + \underbrace{\left(\epsilon_{\text{machine}}/2 - 1\right)}_{\epsilon_{\text{machine}}/2 \approx 1.1 \cdot 10^{-16}}$$