# Ch 1: Floating Point Numbers

Wednesday, August 27, 2025       6:35 PM

Learning objectives:
- Distinguish underflow vs. below machine epsilon
- There's a largest float
- Spacing between floats scales relatively

**How can we represent numbers on a computer?**

Integers are easy.  $6 = 1\cdot4 + 1\cdot2 + 0\cdot1$, ie., $110$ (binary)

What about fractions?

5 bits → $\dfrac{18}{19}$ + 3 bits → $\dfrac{22}{23}$ = $\dfrac{18\cdot23 + 22\cdot19}{19\cdot23}$ = $\dfrac{414 + 418}{437}$ = $\dfrac{832}{437}$ ) 10 bits / ) 9 bits

10 bits   10 bits

19 bits

**Problem:** memory increases

Could try decimals (fixed pt.)

$131.467$     or in binary   $2^2\ 2^1\ 2^0\ 2^{-1}\ 2^{-2}$
$\quad\ \underbrace{\ \ }_{3}\ \underbrace{\ \ }_{3}$        ... $100110.100110$ ...

Done in embedded system

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

Usually use IEEE 754 "double precision" → 64 bits = 8 bytes
(single precision = 32 bits)

Store numbers

$\mathbb{F}$ = Floating Pt.   significand
$\mathbb{R}$ = real numbers   $(-1)^s (1+f)\cdot 2^e$   $e = c - 1023$
                          radix or base

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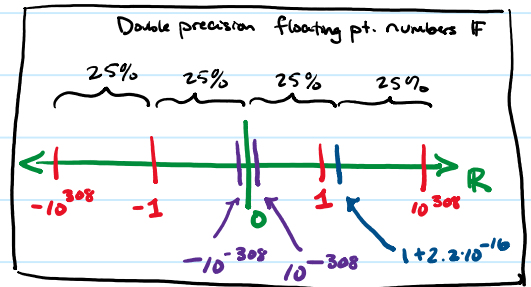$f$ = mantissa, 52 bits

For exponent, 11 bits, think of 1 sign bit so $\pm 2^{10}$ i.e. $\pm 1024$

(scientific notation)

$\mathbb{F}$ also includes $0$, NaN ($\frac{0}{0}$, $0\cdot\infty$, $\infty/\infty$), $\pm\infty$

Rule of thumb: precision $\simeq 2^{52} = 4.5\cdot10^{15}$
15 digits of precision in double
8 digits ...        in single



Double precision floating pt. numbers $\mathbb{F}$
25%  25%  25%  25%
$-10^{308}$  $-1$  $0$  $1$  $10^{308}$  $\mathbb{R}$
$-10^{-308}$  $10^{-308}$  $1 + 2.2\cdot10^{-16}$

**Implications**

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

Due to limitations in exponent

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

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

**Overflow** if not in range

**Note:** we're being a little loose w/ a few technicalities because we focus on the bigger message

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 $|x| < 2^{-1022} \approx 10^{-308}$, might be treated as $0$

# Ch 1: Floating Point Numbers, p. 2

← relative spacing

③ limit to spacing    "Machine epsilon",  SILENT ERROR

$1 \stackrel{?}{==} 1 + \varepsilon$ , ie., 1 and $1 + \varepsilon$ are indistinguishible

↗ due to limits in mantissa

true if $\varepsilon < \varepsilon_{machine} = 2^{-52} \approx 2.2 \cdot 10^{-16}$

NOTE: youtube video is wrong, it is relative

$2$ vs $2 + \varepsilon$ ,   $\varepsilon < 2 \cdot \varepsilon_{machine}$

**Def** of $\varepsilon_m$ :   the difference between 1 and the next largest floating pt. number  (aka "mainstream def" on wikipedia)

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

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

depends on precise definition

$\dfrac{|fl(x) - x|}{|x|}$ } relative error  = accuracy

$|fl(x) - x| = $ absolute error

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

Precision : #digits, need not be correct!

ˇ Accuracy : relative error, limited by precision

More implications:  lose associative rule

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

$$\left( 1 + \varepsilon_{machine}/2 \right) - 1 \stackrel{?}{\neq} 1 + \left( \varepsilon_m/2 - 1 \right)$$

⎵ 1 ⎵ → 0        ⎵ $\varepsilon_m/2 \approx 1.11 \cdot 10^{-16}$

**Ex** : softmax

$\log_{10} (10^{400} + 10^{400})$

$= \log_{10}(2) + 400$

but computer will overflow w/ naive implementation

**Catastrophic Cancellation**
Pretend we have 3 digits of precision.

$x = 1.23\underbrace{587325...}_{garbage}$
$y = 1.22\underbrace{3581926...}$

$x - y = 0.01\underbrace{~~~~~~~~}_{garbage}$

$= 1.\boxed{~~~~} \underline{~~~~} \cdot 10^{-2}$

you might think these are accurate since within our precision.

"Silent error"!

See demos

# Ch 1: Floating Point Numbers, extra diagram

$\mathcal{E}_m$ = machine epsilon

$\mathcal{E}$        $\mathcal{E}/2$        $\mathcal{E}/4$  $\mathcal{E}/8$...        $\mathcal{E}/4$        $\frac{\mathcal{E}}{2}$  $\leftarrow\mathcal{E}\rightarrow$

$-2$        $-1$        $-\frac{1}{2}$  $-\frac{1}{4}$        $2^{-2}$  $2^{-1}=\frac{1}{2}$        $2^0=1$        $1+\frac{\mathcal{E}}{2}$        $2^1=1$

$-1+\frac{\mathcal{E}}{2}$
(B)

(A)

(B')

$1-\frac{\mathcal{E}}{2}$

$0$

**(A)**

$$\left(1+\frac{\mathcal{E}}{2}\right)-1 \quad = \quad 1-1=0$$

$\underbrace{\quad\quad} = 1$ in floating pt.

**(B)**                    **(B')**

$$1+\left(\frac{\mathcal{E}}{2}-1\right) = 1-\underbrace{\left(1-\frac{\mathcal{E}}{2}\right)}_{} = \frac{\mathcal{E}}{2}$$

this is representable in floating pt.

Two definitions of $\mathcal{E}_{mach}$

① $1+\mathcal{E}_{mach}$ is next floating pt. number after **1**
   (what I drew in picture)            $\mathcal{E} = 2^{52} \approx 2.22\text{e-}16$
   This is "variant" in wikipedia

② Smallest number such that $1+\mathcal{E}_m \neq 1$
   i.e., "unit roundoff"                $\mathcal{E} = 2^{53} \approx 1.11\text{e-}16$
   if we "round to nearest" then
   this is $\frac{1}{2}$ the value from definition ①

In practice, we care that $\mathcal{E}_{mach} \approx 10^{-16}$,
specific numbers not always important

②ⁿ largest number $\mathcal{E}$ such that $1+\mathcal{E}$ rounds to 1

Take-home points:
   Floating pt. numbers have almost a constant
   <u>relative</u> accuracy, hence <u>absolute</u> accuracy
   depends on magnitude

   i.e., spacing of floating pt. #'s is
   <u>not</u> uniform.