

## Floating point numbers

Background: Scientific notation, e.g.

$$\text{Avogadro constant: } N_A = 6.022 \cdot 10^{23} \frac{1}{\text{mol}}$$

$$\text{Speed of light: } c = 3.00 \cdot 10^8 \frac{\text{m}}{\text{s}}$$

Advantages:

- Makes it easy to handle quantities of very different orders of magnitude
- "Number of significant digits" easy to trace through computation

Scientific notation can be thought of as "arbitrary precision floating point".

Floating point in hardware (IEEE 754):

$$\pm s \cdot 2^e = \pm \underbrace{d_0.d_1d_2 \dots d_{p-1}}_{\text{significant (or mantissa)}} \cdot 2^{\overset{\text{exponent}}{\downarrow} e}$$

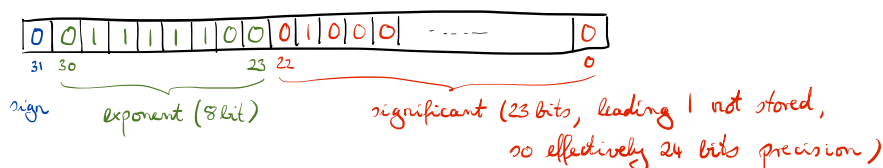
- significant has fixed length of  $p$  bits
- exponent has fixed number of bits s.t.  $e_{\min} \leq e \leq e_{\max}$ ; store  $e + \text{bias}$
- one sign bit
- if possible, "normalize" s.t.  $d_0 = 1$ , no need to store

E.g. layout for 32-bit IEEE floating point for

$$(0.15625)_{10} = 1.01 \cdot 2^{-3}$$

is computed as follows: Exponent bias per IEEE standard is 127,

$$\text{so } e + \text{bias} = -3 + 127 = 124 = (01111100)_2$$



Most common: "double precision" (64 bit) floating point with 11 bit exponent, effective 53 bit precision

## Features of the standard:

- Floating point addition  $\oplus$ , multiplication  $\odot$ , subtraction  $\ominus$ , division  $\oslash$  are exactly rounded, i.e.

$$\text{fl}(x) \oplus \text{fl}(y) = \text{fl}(\text{fl}(x) \cdot \text{fl}(y))$$

where "fl" denotes "round to floating point"

- Two zeros,  $+0$  and  $-0$ , which compare as equals
- NaN ("not a number"): propagate as NaN, except  $\text{NaN}^0 = 1$
- $\pm \text{inf}$  ("infinities"):  $1 \odot +0 = +\text{inf}$ ,  $1 \odot -0 = -\text{inf}$ ,  $0 \odot \text{inf} = \text{NaN}$
- Subnormal numbers  $0.d_1 \dots d_{p-1} \cdot 2^{e_{\min}}$

Needed to ensure that  $x = y \Leftrightarrow x - y = 0$  in floating point.

Expensive, even though handled in hardware on contemporary CPUs.

- Correct rounding (round to nearest, round to even in case of a tie) removes bias
- Lexical sorting is relatively easy

## Error measures:

absolute error: If  $\text{fl}(x) = 1.d_1d_2 \dots d_{p-1} \cdot 2^e$ , then

$$\begin{aligned} |x - \text{fl}(x)| &\leq 0.5 \cdot 2^{1-p} \cdot 2^e \\ &\equiv 0.5 \text{ ulp} \quad \text{"units in the last place"} \end{aligned}$$

relative error:  $\frac{|x - \text{fl}(x)|}{|x|} \leq \frac{0.5 \cdot 2^{1-p} \cdot 2^e}{2^e}$  (as  $|x| \geq 2^e$ )

$$= 0.5 \epsilon \quad \text{where } \epsilon = 2^{1-p} \text{ is the "machine epsilon"}$$

Remark:  $\epsilon$  is the smallest positive floating point number such that  $1 \oplus \epsilon > 1$ .

On standard 64-bit floating point,  $\epsilon \approx 2.22 \cdot 10^{-16}$

$$\Rightarrow \frac{1}{2} \epsilon \approx 1.11 \cdot 10^{-16}$$

## Error propagation:

$$(a) \quad \text{fl}(x) \odot \text{fl}(y) = \text{fl}(x) \cdot \text{fl}(y) (1 + \delta_1) \quad \text{with } |\delta_1| \leq \frac{\epsilon}{2}$$

$$\begin{aligned}
&= x(1+\delta_2) \cdot y(1+\delta_3) (1+\delta_1) \quad \text{with } |\delta_2|, |\delta_3| \leq \frac{\epsilon}{2} \\
&= x \cdot y (1 + (\delta_1 + \delta_2 + \delta_3) + \text{terms bounded by } \epsilon^2) \\
&= x \cdot y (1 + \delta) \quad \text{with } |\delta| \leq \frac{3}{2} \epsilon
\end{aligned}$$

$\Rightarrow$  Multiplication of floats is always well-behaved ("numerically stable")

$$\begin{aligned}
(b) \quad \text{fl}(x) \oplus \text{fl}(y) &= (\text{fl}(x) + \text{fl}(y)) (1 + \delta_1) \quad \text{with } |\delta_1| \leq \frac{\epsilon}{2} \\
&= (x(1+\delta_2) + y(1+\delta_3)) (1 + \delta_1) \quad \text{with } |\delta_2|, |\delta_3| \leq \frac{\epsilon}{2} \\
&= (x + y + x\delta_2 + y\delta_3) (1 + \delta_1) \\
&= x + y + (x+y)\delta_1 + x\delta_2 + y\delta_3 + \text{quadratic terms} \\
&= (x+y) \left( 1 + \delta_1 + \underbrace{\frac{x}{x+y}\delta_2 + \frac{y}{x+y}\delta_3}_{=: \delta} \right) + \dots \\
&\Rightarrow \text{Addition of same-signed numbers is well-behaved: } |\delta| \leq \frac{3}{2} \epsilon
\end{aligned}$$

Addition of different-signed numbers can have huge amplification of relative error. Worst case:  $x \approx -y \nabla$

(c) Suppose  $x, y, z$  are already floating point values. Then

$$\begin{aligned}
(x \oplus y) \oplus z &= ((x+y)(1+\delta_1) + z) (1 + \delta_2) \quad |\delta_1|, |\delta_2| \leq \frac{\epsilon}{2} \\
&= (x+y+z + (x+y)\delta_1) (1 + \delta_2) \\
&= x+y+z + (x+y+z)\delta_2 + (x+y)\delta_1 + \text{quadratic terms} \\
&= (x+y+z) \left( 1 + \delta_2 + \frac{x+y}{x+y+z} \delta_1 \right) + \dots
\end{aligned}$$

- If  $z \approx -(x+y)$ , catastrophic loss of accuracy may occur, as above.
- For same-signed numbers, repeated addition adds at most  $\frac{\epsilon}{2}$  error per term. But it helps to add smallest numbers first!

Problematic example:  $f'(x) \approx \frac{f(x+h) - f(x)}{h}$

(The  $\frac{1}{h}$  converts big relative errors to big absolute errors!)

## Mitigation:

① Rule of thumb: never test floating-point numbers for equality

E.g.  $a = 0.15 + 0.15$

$$b = 0.1 + 0.2$$

(i)  $a == b$  gives False in Python. Not a good idea.

(ii)  $\text{abs}(a-b) \leq \text{tol}$  Better, but absolute tolerance is wrong concept

(iii)  $\text{abs}(a-b)/\text{abs}(a) \leq \text{tol}$  Usually works, but there are edge cases

(iv) def almost-equal(a,b):

if  $a == b$ :

return True

(handles infinities!)

else if  $a == 0$  or  $b == 0$  or  $\text{abs}(a) + \text{abs}(b) < \text{min-normal}$ :

return  $\text{abs}(a-b) < \text{tol} * \text{min-normal}$

else return  $\text{abs}(a-b) / \min(\text{abs}(a) + \text{abs}(b), \text{max-float}) < \text{tol}$

② Rewrite formulas (often, but not always possible):

(i)  $a_0 + a_1x + a_2x^2 + \dots + a_nx^n = a_0 + x(a_1 + x(a_2 + \dots (a_nx)))$  (Horner's scheme)

(ii)  $x_{\pm} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$  (quadratic formula)

Problematic case:  $4ac \ll b^2$ . Rewrite, for  $b > 0$ ,

$$x_{+} = \frac{(-b + \sqrt{b^2 - 4ac})(-b - \sqrt{b^2 - 4ac})}{(-b - \sqrt{b^2 - 4ac})2a}$$

$$= \frac{b^2 - (b^2 - 4ac)}{(-b - \sqrt{b^2 - 4ac})2a} = \frac{2c}{b + \sqrt{b^2 - 4ac}}$$

No loss of significant digits!