# Floating-point numbers

Spring 2021    A. Donev

N — digits

Scientific / exponential format

$$\pm\ 3.15\ 78\ldots \cdot 10^{s}$$

$$\boxed{\pm}\ \underset{\text{zero}}{0}.\underset{\text{non zero}}{3}15\,78 \cdot 10^{6}$$

$$X = (-1)^{s} \cdot 0.\overset{m}{\underbrace{a_1 a_2 a_3 \ldots a_t}} \cdot \beta^{\overset{\text{exponent}}{e}}$$

$s = 0$ or $1$ — one bit

$t$ digits

$$\rightarrow \quad = (-1)^{s} \cdot \underset{\text{mantissa}}{m} \cdot \beta^{\underset{\text{fixed}}{e-t}}$$

On computer `floating point number

$$[ \; s \; | \; m \; | \; e \; ]$$

↑     ↑     ↑

sign    mantissa    exponent

$$\beta = 2$$

$$134_{10} = 1 \cdot 10^2 + 3 \cdot 10^1 + 4 \cdot 10^0$$

$$\rightarrow 100101_2 = 1 \cdot 2^5 + 1 \cdot 2^2 + 1 \cdot 2^0$$

# of bits = 32, 64, 128 bit

       ↑      ↑      ↑

4 bytes    8 bytes   16 byts

Standard    IEEE

— format

— rounding

— exceptions    $\frac{1}{0}, \sqrt{-1}$

# IEEE formats

## Single precision $= 4$ bytes
$= 32$ bits

$$\underbrace{1}_{\substack{\text{Sign} \\ S}} + \underbrace{8}_{\substack{\text{exponent} \\ e}} + \underbrace{23}_{\substack{\text{mantissa} \\ m}} = 32$$

## Double precision $= 64$ bits
$= 8$ bytes

$$1 + 11 + 52 = 64$$

## Rounding errors

$$\hat{x} = fl(x)$$

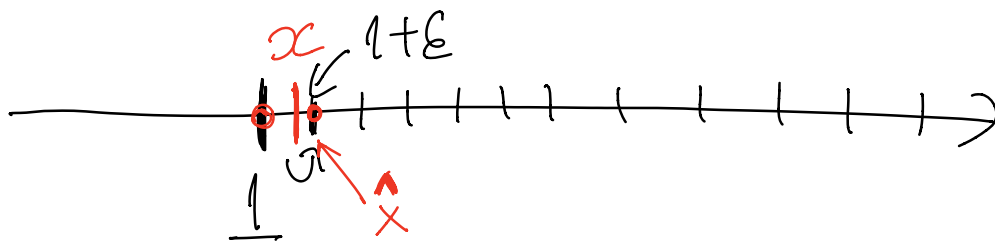$$x \in \mathbb{R}$$

$$\hat{x} \approx x$$

$$1 \quad 1+\varepsilon$$

spacing

eps in MATLAB

$$\varepsilon = 2^{-\#\text{ of bits in Mantissa}}$$

$$\varepsilon = \begin{cases} 2^{-23} & \text{single precision} \\ 2^{-53} & \text{double precision} \end{cases}$$

$$x \quad 1+\varepsilon$$

$$1 \qquad \hat{x}$$

If we round to nearest

$$\frac{|x - \hat{x}|}{|x|} \leq \frac{\varepsilon}{2} = u$$

roundoff unit
machine precision

relative error $= \begin{cases} 2^{-24} \approx 6 \cdot 10^{-8} & \text{SP} \\ 2^{-53} = 1 \cdot 10^{-16} & \text{DP} \end{cases}$

Double precision gives you
at most 16 digits of accuracy

There is a largest &
smallest number $\sim 10^{300}$

$$fl(x) \odot fl(y) = fl(x \circ y)$$

hardware

$+, -, *, /$

* and + not associative

$(x+y) + z \neq x + (y+z)$

$$1/0 = \underline{\underline{inf}}$$

$$\sqrt{-1} \quad \text{or} \quad 0/0 = NaN$$

not a number

**[Matlab uses double precision by default**

To know:

1) Do not compare floating-point numbers

$$x = y \quad ?$$

2) $x + y - z$    use parenthesis

$(x + y) - z$    to to)

$x + (y - z)$    (try control order

3) Built in "standard" functions sin, cos, ln, exp give you 16 digits.

# Propagation of roundoff error

Multiplication: $\underline{\underline{x}} \; \underline{\underline{y}}$

$$\frac{\left| (x + \underline{\underline{\delta x}})(y + \delta y) - xy \right|}{|xy|}$$

$$= \left| \underset{11}{\frac{\delta x}{x}} + \frac{\delta y}{y} + \overset{\text{small}}{\cancel{\left(\frac{\delta x}{x}\right)\left(\frac{\delta y}{y}\right)}} \right|$$

$$\sim 10^{-16}$$

$$\delta x / x \ll 1, \quad \delta y / y \ll 1$$

$$\leq \left| \frac{\delta x}{x} \right| + \left| \frac{\delta y}{y} \right)$$

$$\mathcal{E}_{xy} \leq \mathcal{E}_x + \mathcal{E}_y \quad \text{relative errors}$$

<u>Good</u>

Multiplication / division are "safe" ~ they don't loose digits

## <u>Addition</u>

$$\underline{1.0010} \qquad 5 \text{ digits}$$

$$\underline{0.00013\,678} \qquad 5 \text{ digits}$$

$$\underline{1.0013,\,678} \leftarrow \text{lost}$$

If round to nearest

1 . 0 0 1 4

1 3 $\boxed{6\ 7\ 8}$ ⟵ Lost 3 digits
1 4 0 0 0

Catastrophic cancellation

Add numbers of widely
different magnitude

Subtraction

$\underline{1 . 0 0 1 2}$ , 3 2 1
$-\underline{1 . 0 0 1 1}$ , 0 2 0

0 . 0 0 0 1 , 2 0 1
lost these digits

Subtraction of numbers
that are very close to each other

# Example

Harmonic sum

$$H(N) = \sum_{i=1}^{N} \frac{1}{i}$$

forward
or
backward
summation

$$\lim_{N \to \infty} \left( H(N) - \ln N \right) = \gamma =$$

Euler Constant
(not e)

## Backward is better

(Kahn summation)
see Wiki

## Example   Solve

$$x^2 - 2x + c = 0$$

$$x = 1 - \sqrt{1-c}$$

Assume $|c| \ll 1$

$\Rightarrow$ loss of digits

If $c = 10^{-9}$

$$1 - c = \underset{\text{9 zeros}}{(1).000000\,0\cdots\cdots\,1}\overbrace{\qquad\qquad}^{16}$$

I loose $16 - 9 = 7$

$$1 - \sqrt{1-c} = \frac{c}{1 + \sqrt{1-c}}$$

Good even for small $c$

$$\text{If } |c| \leq 1 \simeq \frac{c}{2}$$

$$1 - \sqrt{1-c} \underset{\substack{\text{Taylor} \\ \text{series} \\ \text{around } c=0}}{\approx} \frac{c}{2} + O(c^2)$$