# Fundamentals of numerical analysis

Today: floating point arithmetic & rounding error analysis
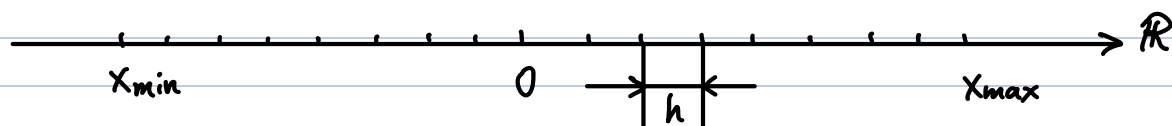
To do linear algebra on computers

**First step**: store numbers on computers & do arithmetic

**Challenge**: $\mathbb{R}$ is unbounded and forms a continuum
while computers are "discrete" & finite memory

## Idea 1 (Fixed point #s)

Discretize $\mathbb{R}$ into equally spaced points



Set $h = B^{-n}$

Denote the set of fixed point numbers
$$x = q \, B^{-n}, \qquad q \in \mathbb{Z}$$

On a (binary) computer, $\pm \underset{\uparrow}{1} \underbrace{0 \; 0 \; 1}.\underbrace{0 \; 1 \; 1 \; 0}$

sign    integer part    fraction part

m-digits    n-digits

$$x = \pm \sum_{i=-n}^{m-1} k_i \, B^i, \quad 0 \le k_i \le B-1$$

Nonzero fixed pt # range    $B^{-n} \le |x| \le B^m - B^{-n}$

Let $f_i(\cdot)$ map $\mathbb{R}$ to the nearest fixed point #

For $x$ in the range,
$$f_i(x) = \underbrace{x + \delta}_{\text{absolute error is small}}, \quad |\delta| \leq h$$

cons: · Less suitable for representing very large / small #s

· Values can overflow / underflow easily

## Idea 2 ( Floating point #s )

Mimics scientific notation $1.25 \times 10^{-1}$

Floating point #'s

$$x = \pm \frac{m}{B^t} B^e$$

· $t$ : precision

· $B$ : base ( usually $B = 2$ on a binary computer )

· $e$ : exponent $e_{min} \leq e \leq e_{max}$ (exponent range )

· $m$ : fraction $\quad B^{t-1} \leq m \leq B^t - 1$

   ↑
   "normalized"                    "0" is a special
   ensure unique representation    case ($m = 0$)

A more common way of expressing floating point # is

$$x = \pm B^e \times \left( \sum_{i=1}^{t} \frac{d_i}{B^i} \right) = \pm B^e \times . \underline{d_1} \underline{d_2} \cdots \underline{d_t}$$

each digit $0 \leq d_i \leq B-1$ . $d_1 \neq 0$ for normalized representation

- Decimal location "floats" depending on the size of #

  Less easy to overflow / underflow

  Range of nonzero floating point #s
  $$B^{e_{min}-1} \leq |x| \leq B^{e_{max}} (1 - B^{-t})$$

  example: IEEE 754 (1985, updated 2008)

| | B | t | $e_{min}$ | $e_{max}$ | $\epsilon_{mach}$ |
|---|---|---|---|---|---|
| single (FP32) prec. | 2 | 24 | -126 | 127 | $2^{-24} \approx 5.96 \times 10^{-8}$ |
| (FP64) double prec. | 2 | 53 | -1022 | 1023 | $2^{-53} \approx 1.11 \times 10^{-16}$ |

<span style="color:blue">↑ bits for fraction + 1 hidden bit (implicit digit $d_1 \equiv 1$)</span>

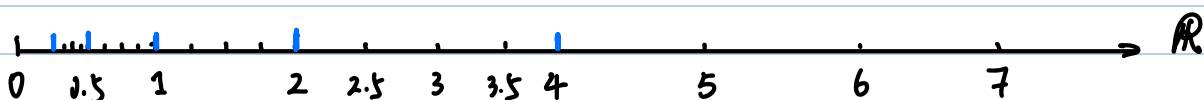Single precision: 32 bits = 1 + 8 + 23

double precision: 64 bits = 1 + 11 + 52

                                sign    exponent    fraction

Other precisions: FP8, FP16 (Half-prec.), single extended,
                     ↑                                  double extended, ...
             <span style="color:blue">multiple format ...</span>

- floating point numbers are not equally spaced



  0  0.5  1    2  2.5  3  3.5 4    5    6    7

If B=2, t=3, $e_{min}$ = -1, $e_{max}$ = 3

Floating point numbers:

$$2^3 \times .\underline{1}\,\underline{1}\,\underline{1} = 2^3 \times \left(\frac{1}{2} + \frac{1}{2^2} + \frac{1}{2^3}\right) = 7$$

$$2^3 \times .\underline{1}\,\underline{1}\,\underline{0} = 2^3 \times \left(\frac{1}{2} + \frac{1}{2^2} + \frac{0}{2^3}\right) = 6$$

$$2^3 \times .\underline{1}\,\underline{0}\,\underline{1} = 2^3 \times \left(\frac{1}{2} + \frac{0}{2^2} + \frac{1}{2^3}\right) = 5$$

$$2^3 \times .\underline{1}\,\underline{0}\,\underline{0} = 2^3 \times \left(\frac{1}{2} + \frac{0}{2^2} + \frac{0}{2^3}\right) = 4$$

$$2^3 \times .\underline{0}\,\underline{1}\,\underline{1} = 2^2 \times .\underline{1}\,\underline{1}\,\underline{1} = 2^2 \times \left(\frac{1}{2} + \frac{1}{2^2} + \frac{1}{2^3}\right) = 3.5$$

$$2^2 \times .\underline{1}\,\underline{1}\,\underline{0} = 3.0$$

$$2^2 \times .\underline{1}\,\underline{0}\,\underline{1} = 2.5$$

$$2^2 \times .\underline{1}\,\underline{0}\,\underline{0} = 2.0$$

$$2^2 \times .\underline{0}\,\underline{1}\,\underline{1} = 2^1 \times .\underline{1}\,\underline{1}\,\underline{1} = 1.75$$

$$\vdots$$

Set $\frac{m}{B^t} = \frac{1}{2^2} + \frac{1}{2^3}$

$m = 3 < m^{t-1}$

not a normalized representation

- So how to quantify accuracy of floating point # ?
machine epsilon (unit roundoff):

$\varepsilon_{mach}$ = half distant from 1.0 to the next larger float

$$1.0 = \frac{B^{t-1}}{B^t} B \quad , \quad \text{next } \# = \frac{B^{t-1}+1}{B^t} B$$

$\varepsilon_{mach} = \frac{1}{2} B^{1-t}$

$\varepsilon_{mach}$ = relative error of rounding $x \in \mathbb{R}$ to its nearest fp#

__Thm__ For every $x \in \mathbb{R}$ (in exponent range)

$$fl(x) = x(1+\delta), \qquad |\delta| \leq \varepsilon_{mach}$$

<span style="color:blue">relative error</span>

__Pf:__   w.l.o.g. assume that $x > 0$

We write $\qquad x = \mu \times B^{e-t}$,

where $B^{t-1} \leq \mu < B^{t}$, $\quad x \in [y_1, y_2]$

where $\quad y_1 = \lfloor \mu \rfloor B^{e-t}$, $\quad y_2 = \lceil \mu \rceil B^{e-t} = \dfrac{\lceil \mu \rceil}{B} B^{e-t+1}$

thus

$$\left| \frac{fl(x)-x}{x} \right| \leq \frac{1}{2} \left| \frac{y_2 - y_1}{x} \right| = \frac{1}{2} \frac{B^{e-t}}{\mu \times B^{e-t}} \leq \frac{1}{2} B^{1-t} \qquad \blacksquare$$

- Floating point arithmetic

To carry out rounding analysis, we need to make some assumptions about the accuracy of the basic arithmetic operation. The most common assumptions are embodied in the following model

Let $\mathbb{F}$ be the set of all floating point numbers

Let $*$ be one of the operations $+, -, \times$ or $\div$

Let $\circledast$ be its floating point analogue

Standard model
(Fundamental Axiom of Floating Point Arithmetic)

$$\forall x, y \in \mathbb{F}, \quad x \circledast y = fl(x * y)$$

that is, $\exists \delta$ with $|\delta| \leq \epsilon_{mach}$ s.t.

$$x \circledast y = (x * y)(1 + \delta)$$

This model is valid for most computers, including IEEE standard arithmetic

example: $f(x) = \big(\big((x - 0.5) + x\big) - 0.5\big) + x$

in exact arithmetic, $f\left(\frac{1}{3}\right) = 0$

in double precision, $f(x) \neq 0$, $\forall x \in \mathbb{F}$

$\left(* \text{ Hint: Show that } f(x) = 3x - 1 \text{ for } x \equiv fl(x) \text{ near } \frac{1}{3}\right)$

- Catastrophic cancellation:

Substracting two nearly equal numbers cancel the most significant digits but the result can have large relative error

ex1. evaluate $\dfrac{1}{1-x} - 1$ for $|x| << 1$, $x \in F$

Method 1: Direct evaluation

$$\text{Output}_1 = \left[ \frac{1}{(1-x)(1+\delta_1)}(1+\delta_2) - 1 \right](1+\delta_3) \qquad |\delta_i| \leq \varepsilon_{mach}$$

$$= \frac{\left[1+\delta_2 - (1-x)(1+\delta_1)\right](1+\delta_3)}{(1-x)(1+\delta_1)}$$

$$= \frac{\delta_2 - \delta_1 + x(1+\delta_1)}{1+x} \cdot \frac{1+\delta_3}{1+\delta_1}$$

when $x \sim O(\delta_2 - \delta_1)$, relative error $\sim O\left(\dfrac{\delta_2-\delta_1}{x}\right) = O(1)$

Method 2: Rearrange calculation

from $\dfrac{1}{1-x} - 1 = \dfrac{x}{1-x}$

$$\text{Output}_2 = \frac{x(1+\delta_1)}{(1-x)(1+\delta_2)}(1+\delta_2)$$

relative error $\sim O(\delta)$ even when $x \sim O(\delta)$

ex 2. $\dfrac{e^x - 1}{x}$, $|x| << 1$

assume the exp and log function are both computed with a relative error not excceeding $\varepsilon_{mach}$

from Taylor expansion

$$\frac{e^x - 1}{x} = \frac{1 + x + \frac{1}{2}x^2 + \cdots - 1}{x} \approx 1 + \frac{1}{2}x + O(x^2)$$

Method 1:  Direct evaluation

$$Output_1 = \frac{\left[e^x(1+\delta_1) - 1\right](1+\delta_2)}{x(1+\delta_3)}(1+\delta_4)$$

$$= \frac{\left(1 + x + \frac{1}{2}x^2 + \cdots\right)(1+\delta_1) - 1}{x} \frac{1+\delta_2}{1+\delta_3}(1+\delta_4)$$

$$\approx \left(\frac{\delta_1}{x} + 1 + \frac{1}{2}x\right)\frac{1+\delta_2}{1+\delta_3}(1+\delta_4)$$

$\uparrow$

relative  error  $\sim O\left(\frac{\delta_1}{x}\right)$

Method 2:  Rearrange  calculation

First  compute   $\hat{y} = e^x(1+\delta_1)$

then   $Output_2 = \frac{\hat{y} - 1}{\log \hat{y}}(1+\delta_2)$

exercise: while the relative errors of numerator and denominator are $O(1)$ for $x \sim O(\varepsilon_{mach})$, $Output_2$ has $O(\varepsilon_{mach})$ relative error and is accurate.

Other important points:
   (See iJulia notebook)
   · Input/output rounding
   · Non associativity
   · Catastrophic cancellation. etc.

- Rounding error analysis

Consider the inner product

$$x^T y \quad , \quad x, y \in \mathbb{F}^n \qquad \text{special case: summation}$$

Naive summation algorithm

> $$S_1 = fl(x_1 y_1)$$
> $$S_i = fl(S_{i-1} + fl(x_i y_i)) \qquad i = 2, \cdots, n$$

$$S_1 = x_1 y_1 (1+\delta_1) \quad , \qquad |\delta_1| \leq \varepsilon_{mach}$$

$$S_2 = \left(S_1 + x_2 y_2 (1+\delta_2)\right)(1 + \delta_2') \qquad\qquad |\delta_2|, |\delta_2'| \leq \varepsilon_{mach}$$

$$= x_1 y_1 (1+\delta_1)(1+\delta_2') + x_2 y_2 (1+\delta_2)(1+\delta_2')$$

$$S_n = x_1 y_1 (1+\delta_1) \prod_{i=2}^{n} (1+\delta_i') + \sum_{j=2}^{n} x_j y_j (1+\delta_j) \prod_{i=j}^{n} (1+\delta_i')$$

$$|\delta_i'|, |\delta_i| \leq \varepsilon_{mach}$$

Lemma: If $|\delta_i| \leq \varepsilon_{mach}$, and $n \varepsilon_{mach} < 1$, then

$$\prod_{i=1}^{n} (1+\delta_i) = 1 + \theta_n$$

with $\quad |\theta_n| \leq \dfrac{n \varepsilon_{mach}}{1 - n \varepsilon_{mach}} =: \gamma_n \qquad \leftarrow$ linear in $n$

Pf: By induction                          ▨

By this Lemma, we obtain

$$S_n = x_1 y_1 (1+\theta_n') + \sum_{j=2}^{n} x_j y_j (1 + \theta_j)$$

with $|\theta_n'| \leq \gamma_n$, $|\theta_j| \leq \gamma_{n-j+2}$

Let $\Delta x = (\theta_n' x_1, \theta_2 x_2, \cdots, \theta_n x_n)^T$

or $\Delta y = (\theta_n' y_1, \theta_2 y_2, \cdots, \theta_n y_n)^T$

then $S_n = (x + \Delta x)^T y = x^T (y + \Delta y)$

Note that $|\Delta x| \leq \gamma_n |x|$, $|\Delta y| \leq \gamma_n |y|$,

$$|x| := (|x_i|)_{i=1}^n$$

$$|x|^T |y| := \sum_{i=1}^n |x_i| |y_i|$$

$|x^T y - S_n| = |\Delta x^T y| \leq \gamma_n |x|^T |y|$

**Remark 1:** Better algorithm can do better

For example,

1) $\quad x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \begin{matrix} n/2 \\ n/2 \end{matrix} \qquad y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \begin{matrix} n/2 \\ n/2 \end{matrix}$

do $S_{n/2}^{(1)}$ for $x_1, y_1$, $S_{n/2}^{(2)}$ for $x_2, y_2$, then $S_{n/2}^{(1)} + S_{n/2}^{(2)}$

Error: $\left| fl(S_{n/2}^{(1)} + S_{n/2}^{(2)}) - x^T y \right| \leq \gamma_{n/2 + 1} |x|^T |y|$

(Hint: perform the induction as above)

2) The error can be further improved by partitioning into $k$ pieces

Error: $\quad Error \leq \gamma_{\frac{n}{k} + k - 1} |x|^T |y|$

take $k = \sqrt{n}$ minimize the dependence in $n$

3) Use pairwise summation recursively can improve to

Error: $\quad Error \leq \gamma_{\lceil \log_2 n \rceil} |x|^T |y| \quad \Big\uparrow$

— exercise

**Remark 2:** Worst case bound is often pessimistic

can be improved by rounding probabilistically

(Higham–Mary, 2019)