# Numerical Analysis

GOAL: solve maths/physics problems with a COMPUTER

Problem → | MATHEMATICAL MODEL | → (ALGORITHM) | NUMERICAL MODEL | → | COMPUTER IMPLEMENTATION AND RUN | → Solution

## SOURCES of ERRORS

- Physical approximation (MATH MODEL)     $\nabla \times \bar{H} = \bar{J} + \frac{\partial \bar{D}}{\partial t} \simeq \bar{J}$

- TRUNCATION ERRORS (NUMERICAL MODEL)     conversion of analytical operators

- ROUNDOFF ERRORS (COMPUTER IMPLEMENT.)     finite # of digits to represent numbers

## NUMBER REPRESENTATION

### POSITIONAL REPRESENTATION

EX:  $(3012)_{10} = 3 \cdot 10^3 + 0 \cdot 10^2 + 1 \cdot 10^1 + 2 \cdot 10^0$

Pos. 2      Pos. 0

digits          BASE

### INTEGERS

$$q = a_m \beta^m + a_{m-1} \beta^{m-1} + \cdots + a_1 \beta^1 + a_0 \beta^0$$

DIGITS, $a_k \in \mathbb{N}$

$0 \le a_k \le \beta - 1$

$a_m \ne 0$

BASE, $\beta \in \mathbb{N}$

$\beta \ge 2$

REALS    $X = \lfloor X \rfloor + frac(x)$

                ↑                    ↑
          integer part      fractional part

EX:  $(3012\overset{\text{Pos}=-1}{\underset{}{.}}\,4\overset{\text{Pos}=-3}{01})_{10}$

RADIX POINT

$3 \cdot 10^3 + 0 \cdot 10^2 + 1 \cdot 10^1 + 2 \cdot 10^0 + 4 \cdot 10^{-1} + 0 \cdot 10^{-2} + 1 \cdot 10^{-3}$

$$(X)_\beta = \underbrace{a_m \beta^m + a_{m-1}\beta^{m-1} + \cdots + a_0 \beta^0}_{\lfloor X \rfloor} + \underbrace{b_1 \beta^{-1} + b_2 \beta^{-2} + \cdots + b_m \beta^{-m}}_{frac(x)}$$

$b_K \in \mathbb{N}$

$0 \le b_K \le \beta - 1$

$b_m \neq 0$

## FIXED POINT REPRESENTATION

○ Positional representation ⟨ fixed number of digits
                              fixed position RADIX POINT

## FIXED POINT SET

$$X(\beta, t, q) = \left\{ x \in \mathbb{R} = sign(x) \left[ \sum_{K=0}^{t-(q+1)} a_K \beta^K + \sum_{K=1}^{q} b_K \beta^{-K} \right] \right\}$$

base

number of digits

number of digits for fractional part
$0 \le q \le t$, $q \in \mathbb{N}$

$$X(\beta, t, q) = \left\{ x \in \mathbb{R} = \text{sign}(x) \left[ \sum_{k=0}^{t-(q+1)} a_k \beta^k + \sum_{k=1}^{q} b_k \beta^{-k} \right] \right\}$$

EX: $X(\beta = 10, t = 4, q = 1)$

$\underbrace{\qquad}$ 3 digits for $\lfloor x \rfloor$

1 " frac($x$)

$\max(x) = 9 \cdot 10^2 + 9 \cdot 10^1 + 9 \cdot 10^0 + 9 \cdot 10^{-1} = 999.9$

$\min(x) = 0 \cdot 10^2 + 0 \cdot 10^1 + 0 \cdot 10^0 + 1 \cdot 10^{-1} = 0.1$



• # of POSITIVE ELEMENTS

$$\beta^t - 1 = 10000 - 1 = 9999$$

○ maximum element

$$(\beta^t - 1) \beta^{-q} = 9999 \cdot 10^{-1} = 999.9$$
$$\qquad\qquad\qquad |$$
$$\qquad\qquad\qquad 1$$

○ minimum element

$$\beta^{-q} = 10^{-1} = 0.1$$

⊡ SPACING : $\Delta = \beta^{-q}$

⟹ CONSTANT

ERRORS of Fixed Point representation    Hp: $X(\beta = 10, t = 4, q = 1)$

ABS. ERROR    $E(x) = x - fip(x)$     $fip(x)$

$x_1 = \frac{10^3}{3} = 333.\overline{3}$     $E(x_1) = 333.\overline{3} - 333.3 = 0.0\overline{3}$

$x_2 = \frac{10^0}{3} = 0.\overline{3}$     $E(x_2) = 0.\overline{3} - 0.3 = 0.0\overline{3}$

→ ABS ERROR is CONSTANT

RELATIVE ERROR: $e(x) = \left| \dfrac{E(x)}{x} \right|$

$e(x_1) = \left| \dfrac{0.0\overline{3}}{10^3/3} \right| = \left| \dfrac{10^{-1}/3}{10^3/3} \right| = 10^{-4} \rightarrow$ "WRONG" by $^1/_{10000}$

$e(x_2) = \left| \dfrac{0.0\overline{3}}{10^0/3} \right| = \left| \dfrac{10^{-1}/3}{10^0/3} \right| = 10^{-1} \rightarrow$ "WRONG" by $^1/_{10}$

$\rightarrow$ VARIABLE relative error $\rightarrow$ relative ACCURACY is smaller for small numbers

PRO :

- SIMPLE $\rightarrow$ ALLOWS for \
  FAST ARITHMETICS

$\Rightarrow$ VIDEOGAMES

CONS:

- NON-CONSTANT REL. ERROR

## FLOATING POINT REPRESENTATION

any $x$ can be represented:

$$x = \text{sign}(x) \left[ \sum_{k=0}^{\infty} d_k \beta^{-k} \right] \beta^{+p} \quad \leftarrow \text{EXPONENTIAL PART, } p \in \mathbb{N}$$

MANTISSA $m$
$1 \le m < \beta$

digits

$0 \le d_k \le \beta - 1$
$d_0 \neq 0$

## FLOATING POINT SET

$$F(\beta, t, L, U) = \{0\} \cup \left\{ x \in \mathbb{R} = \text{sign}(x) \left[ \sum_{k=0}^{t-1} d_k \beta^{-k} \right] \beta^p \right\}$$

# of
digits of
$m$

$\in \mathbb{N}$

$p \in [L, U]$

$$X(\beta=10, t=4, q=1)$$

$$\max(X) = \mathbf{999.9}$$

$$\min(X) = 0,1$$

RANGE $X \rightarrow 0,1 \div 999.9$

RANGE $F \rightarrow 1 \div 999.9 \cdot 10^9$

DISCRETIZATION

$\Delta \rightarrow$ NOT CONSTANT

$$\left| \frac{E(X)}{X} \right|$$
$\downarrow$

ERRORS:

$$e(X) = \left| \frac{X - fl_P(X)}{X} \right| = \left| \frac{\text{sign}(x)\left[\sum_{K=0}^{\infty} d_K \beta^{-K}\right]\beta^P - \text{sign}(x)\left[\sum_{K=0}^{t-1} d_K \beta^{-K}\right]\beta^P}{\underbrace{\text{sign}(x)\left[\sum_{K=0}^{\infty} d_K \beta^{-K}\right]\beta^P}_{m(x)}} \right|$$

$$= \frac{\sum_{K=t}^{\infty} d_K \beta^{-K}}{m} \leq \frac{\sum_{K=t}^{\infty} d_K \beta^{-K}}{1} \cdot \Rightarrow e(x) < \beta^{1-t}$$

$$m \geq 1$$

$$\sum_{k=t}^{\infty} d_k \beta^{-k} = d_t \beta^{-t} + d_{t+1}\beta^{-(t+1)} + \cdots + d_{t+m}\beta^{-(t+m)} + \cdots$$

$$\downarrow < \beta^{-t} \cdot \beta$$

$$= \beta^{-t}\left[\underbrace{d_t \beta^0}_{9} + \underbrace{d_{t+1}\beta^{-1}}_{0,9} + \underbrace{\cdots}_{0,09} \cdots + d_{t+m}\beta^{-m} + \cdots\right] < \beta^{1-t}$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxx}}_{< \beta}$$

$$d_t \leq \beta - 1$$

$$F(\beta=10, t=3, L=0, U=9)$$
$$\underbrace{\phantom{xxx}}_{P \in [0,9]}$$

$$\max(F) = \left[\underbrace{d_0\beta^0 + d_{-1}\beta^{-1} + d_{-2}\beta^{-2}}_{m}\right]\cdot \beta^P$$

$$\max(m) = 9,99$$

$$\left\{ \max(\beta^P) = 10^9 \right.$$

$$\max(F) = 9,99 \cdot 10^9 \sim \mathbf{10^{10}}$$

$$\min(F) = \left[\underset{9}{\overset{1}{d_0}}\beta^0 + \underset{0}{d_{-1}}\beta^{-1} + \underset{0}{d_{-2}}\beta^{-2}\right]\cdot \beta^P$$
$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxx}}_{m}$$

$$\min(F) = 1.00 \cdot 10^0 = 1$$

$$e(x) < \beta^{1-t}$$

$$\rightarrow \text{ACCOUNT for ROUNDING}$$

$$e(x) < K \beta^{1-t}$$

$\downarrow$

$1/2$ for ==ROUNDING TO NEAREST== floating point

$$\boxed{K \beta^{1-t} \Rightarrow \varepsilon \quad \text{MACHINE PRECISION}}$$

$$flp(x) = x(1 \pm \varepsilon) \qquad \text{im MATLAB} \quad \varepsilon \sim 10^{-16}$$

0,2  0,3

$x = 0,26 \rightarrow \nearrow^{0,3}_{0,2}$

flp

ROUNDING TOWARD ZERO

$$x_1 = -2/3 = -0,\overline{6}$$
$$x_2 = +2/3 = 0,\overline{6}$$

$x_1$   $\beta^{1-t}$

$-0,7$   $-0,6$   $0$

ROUND T. 0 : $flp(x_1) = -0,6$      ROUND TO NEAREST : $flp(x_1) = -0,7$