

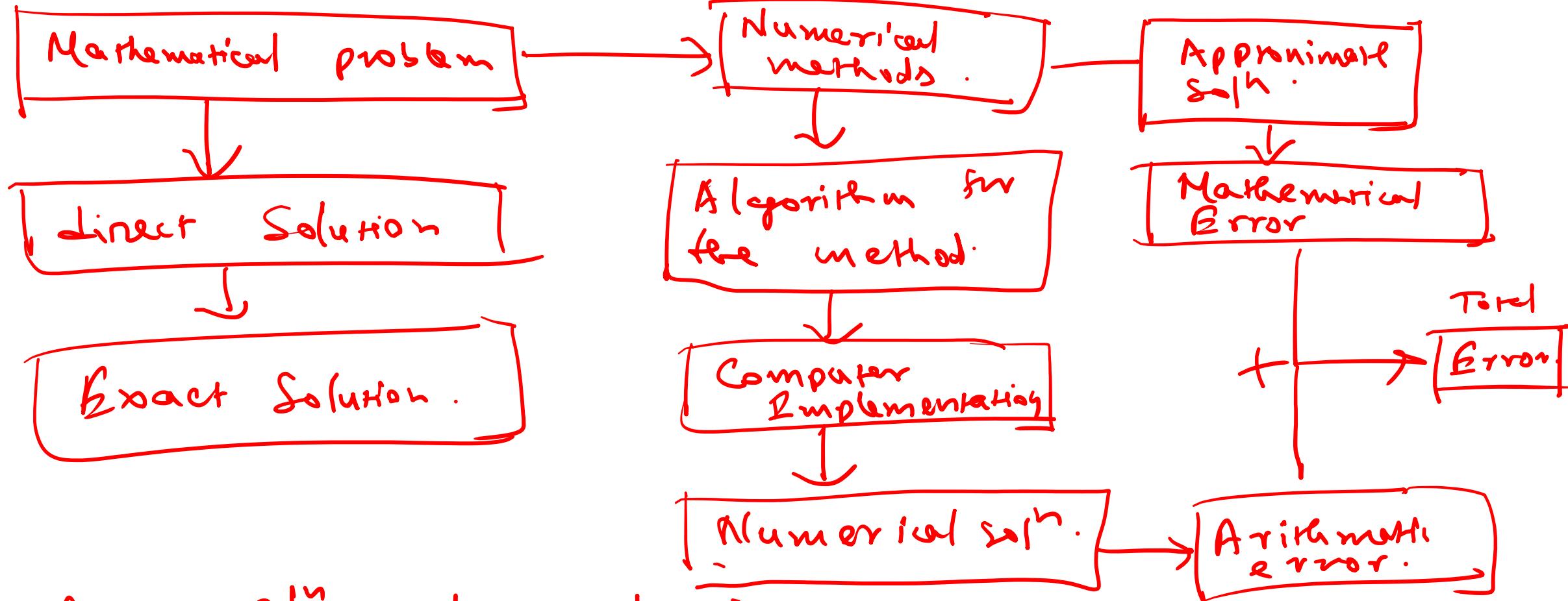
Computational Numerical Methods

CS 374

Prosenjit Kundu

Computer Arithmetic & Errors.

1. floating point representation of number



$$\text{exact soln} = \text{numerical soln} + \text{error.}$$

$$\frac{1}{3} =$$

$$0.\overline{33333}$$

$$0.\overline{333333}$$

Floating point representation of numbers

Let $\beta \neq 2$ let $n \in \mathbb{N}$ & $\beta \geq 2$

Any real number can be represented exactly in in
in base β as

$$(-1)^s \times \underbrace{(d_1 d_2 \dots d_n d_{n+1} \dots)}_{\beta} \times \beta^e.$$

with $d_i \in \{0, 1, \dots, \beta-1\}$ $\underline{d_1 \neq 0}.$ or $d_1 = d_2 = \dots = d_n = 0$

$s = 0$ or 1 is called the sign.

an appropriate integer e is called exponent.

$$\underbrace{d_1 d_2 \dots d_n d_{n+1} \dots}_{\text{called mantissa}} = \frac{d_1}{\beta} + \frac{d_2}{\beta^2} + \dots + \frac{d_n}{\beta^n} + \dots$$

$$0.3145 = (-1)^0 (\cdot 3145)_{10} 10^0$$

$$-0.003145 = (-1)^1 (\cdot 3145)_{10} 10^{-2}$$

$$0.3145 = \cancel{(-1)} \frac{3}{10} + \frac{1}{100} + \frac{4}{1000} + \frac{5}{10000}$$

$$= \frac{3}{10} + \frac{1}{10^2} + \frac{4}{10^3} + \frac{5}{10^4}.$$

If $B = 2$.

$$(-1)^{\text{sgn}} \times (\cdot d_1 d_2 \dots d_n d_{n+1} \dots)_2 \times 2^e.$$

is called the binary floating point representation.

For $B = 10$

is called decimal floating point representations.

$$\text{Ex } x = -0.0017$$

$$= (-1)^{\text{sgn}} 0.17 \times 10^{-2}.$$

Note

There are only finite number of digit if ~~it is~~ in the n digit f.p. representation.

But

$$\begin{aligned}\frac{1}{3} &= 0.3333\ldots \\ &= (-1)^0 (0.333\ldots) \cdot (10^0)\end{aligned}$$

Underflow & Overflow in memory

- # When the value of the exponent e in floating point representation exceeds the maximum limit of (M) the memory we encounter the overflow of memory.
 - # When the value undergoes below of the minimum (m) range we will encounter underflow.

$$m \leq e \leq M.$$

- " If $e > m$ then the number is considered as infinity
if $e < m$ " " "

Chopped & Rounded numbers

Let x be a real number given in the ~~FP~~ FP representation as

$$x = (-1)^s \times (d_1.d_2 \dots d_n d_{n+1} \dots)_{\beta} \xrightarrow{\text{to}} \beta^e.$$

The chopped approximation of x is

$$f(x) = (-1)^s (d_1 d_2 \dots d_n)_{\beta} \beta^e.$$

Let x be a real number given in the FP representation

as

$$x = (-1)^s (\cdot d_1 d_2 \dots d_n d_{n+1} \dots)_B \times B^e.$$

The rounded approximation of x is given by

$$f(x) = \begin{cases} (-1)^s (\cdot d_1 d_2 \dots d_n)_B \times B^e & 0 \leq d_{n+1} < \beta / 2 \\ (-1)^s (\cdot d_1 d_2 \dots (d_n + 1))_B \times B^e & \frac{\beta}{2} \leq d_{n+1} \leq \beta \end{cases}$$

Performing arithmetic operations

Let $*$ denote any of the basic arithmetic operations.

$$* = \{ +, -, \times, \div \}$$

Let x, y be two real numbers.

The process of computing $x * y$ is as follows.

Step 1

Create two n-digit FP ~~approximation~~ approximation

$f(x)$ & $f(y)$ of the numbers x, y respectively

Step 2

Perform the calculations $f(x) * f(y)$ using exact arithmetic.

Step 3 Get the n digit FP representation ~~f(x)~~

g $f_1(f(x) * f(y))$

Example Consider $f(n) = n(\sqrt{n+1} - \sqrt{n})$

Evaluate $f(100000)$ using a 8 digit rounding.

$$f_1(100000) =$$

$$f_1(100001)$$

$$f(100000) = \frac{100000}{100001} (\sqrt{100001} - \sqrt{100000})$$

$$\sqrt{100001} \approx 316229347 = (-1)^0 (0.316229347)$$

$$= (-1)^0 0.316229 \times 10^3$$

$$f(\sqrt{100000}) = (-1)^0 \cdot 316228 \cdot 10^3.$$

$$\therefore f(f(f(x+1)) - f(x)$$

$$f(\sqrt{100001} - f(100000)) = 0.1 \times 10^{-2}$$

Finally

$$f(f(100000)) = f(100000) \cancel{\times} (0.1 \times 10^{-2})$$
$$= 00$$