

Lecture 1: Error Handling

Abhinav Jha

Indian Institute of Technology, Gandhinagar

9th March 2026



1 Why Numerical Methods?

2 Number System

2.1 Scientific Notation

2.2 IEEE Notation

3 Errors

3.1 Round Off Error

3.2 Measure of Error

3.3 Arithmetic Errors



Why Numerical Methods?

Why Numerical Methods? Number System Errors

- Existence
- Uniqueness



Why Numerical Methods?

Why Numerical Methods? Number System Errors

- Existence
- Uniqueness
- What is the solution?



Why Numerical Methods?

Why Numerical Methods? Number System Errors

- $\sum_{n=1}^{\infty} \frac{1}{n}$
- $\sum_{n=1}^{\infty} \frac{1}{n^2}$



Why Numerical Methods?

Why Numerical Methods? Number System Errors

- $\sum_{n=1}^{\infty} \frac{1}{n}$
- $\sum_{n=1}^{\infty} \frac{1}{n^2}$
-

$$\sum_{n=1}^{\infty} \frac{1}{n^2} = \frac{\pi^2}{6}$$



Why Numerical Methods?

n	Approximate Sum	Error	Time (s)
10^1	1.54976773	9.52×10^{-2}	0.0000
10^2	1.63498390	9.95×10^{-3}	0.0000
10^3	1.64393457	1.00×10^{-3}	0.0002
10^4	1.64483407	1.00×10^{-4}	0.0024
10^5	1.64492407	1.00×10^{-5}	0.0271
10^6	1.64493307	1.00×10^{-6}	0.2756
10^7	1.64493397	1.00×10^{-7}	2.8460

Table 1: Convergence of the series $\sum_{k=1}^n \frac{1}{k^2}$ towards $\pi^2/6$.



Why Numerical Methods?

Why Numerical Methods? Number System Errors

Numerical Methods

MA203

Numerical Analysis

MA 637



Number System

- Binary Number System (Base -2)

0 and 1

- Decimal Number System (Base -10)

$$619.916 = 6 \times 10^2 + 1 \times 10^1 + 9 \times 10^0 + 9 \times 10^{-1} + 1 \times 10^{-2} + 6 \times 10^{-3}$$



Number System

Why Numerical Methods? Number System Errors

$$x = 25_{10}$$

<u>Divisor</u>	<u>Quotient</u>	<u>Remainder</u>
25/2	12	1
12/2	6	0
6/2	3	0
3/2	1	1
1/2	0	1

$$25_{10} = 11001_2$$

$$\begin{aligned} 11001_2 &= 1 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 \\ &= 25_{10} \end{aligned}$$



Number System

Why Numerical Methods? Number System Errors

$0.125_{10} \rightarrow \text{Binary}$

<u>Multipl. with</u>	<u>Result</u>	<u>Integer Part</u>
0.125×2	0.25	0
0.25×2	0.5	0
0.5×2	1.0	1

$$0.125_{10} = 0.001_2$$

$$25.125_{10} = 11001.001_2$$



Scientific Notation

Why Numerical Methods? Number System Errors

$$x_1 = \pm x \cdot 10^n$$

$$0.1 \leq x < 1; \quad n \in \mathbb{Z} \text{ (Integers)}$$

$$\text{eg } 1301.122 = \underbrace{0.1301122}_{x} \times 10^4$$

$x, q_2 =$ Mantissa

$m, m =$ Exponent

$$x_2 = \pm q_1 \cdot 2^m$$

$$0.125 \leq q_1 < 1.02 \quad ; \quad \text{and } m \in \mathbb{Z}$$

$$x_2 \neq 0; \quad q_1 \neq 0$$



IEEE Notation

Why Numerical Methods? Number System Errors

64 bits

0.1 - - -

Sign

Exponent

Mantissa

s	c	f
1 bit	11 bits	52 bits

$$2^{-52} = 10^{-d} \quad \text{for some } d \geq 0$$
$$\Rightarrow \log_{10}(2^{-52}) = -d \log_{10}(10)$$
$$\Rightarrow d \approx 16$$



IEEE Notation

Why Numerical Methods? Number System Errors

Any number not represented by IEEE and is large leads to overflow
smaller underflow

Sign	Exponent	Mantissa
s	c	f
1 bit	11 bits	52 bits

$$2^{11} = 2048, \quad -1023 \text{ to } 1024$$

$$(-1)^c 2^{c-1023} (1+f); \quad 0 \leq f < 1$$



Round Off Errors

Why Numerical Methods? Number System Errors

$$y = 0.d_1 d_2 d_3 \dots d_k d_{k+1} \dots \times 10^n$$

- Chopping: Chop a number to k -digits

$$f_c(y) = 0.d_1 d_2 \dots d_k \times 10^n$$

- Rounding:

$$f_r(y) = \begin{cases} 0.d_1 d_2 \dots (d_k + 1) \times 10^n & ; d_{k+1} \geq 5 \\ 0.d_1 d_2 \dots d_k \times 10^n & ; \text{else} \end{cases}$$

eg: $\pi = 0.3141592 \dots \times 10^1$

Chopping five-digit $f_c(\pi) = 0.31415 \times 10^1$

Rounding $f_r(\pi) = 0.31416 \times 10^1$



Measure of Error

Why Numerical Methods? Number System Errors

p = Known value

p^* = Numerical value

Absolute error : $|p - p^*|$

Relative error : $\frac{|p - p^*|}{|p|}$, $p \neq 0$



Arithmetic Errors

Why Numerical Methods? Number System Errors

$$\oplus \quad \ominus \quad \otimes \quad \odot$$

$$2 + 3 = 5$$

$$x \oplus y = fe(fe(x) + fe(y))$$

$$x \otimes y = fe(fe(x) \times fe(y))$$

$$\text{eg: } x = \frac{5}{7}; \quad y = \frac{1}{3}$$

$$\text{five-digit chopping: } fe(x) = 0.71428 \times 10^0$$

$$fe(y) = 0.33333 \times 10^0$$

$$x \oplus y = fe(fe(x) + fe(y))$$

$$= fe(0.71428 \times 10^0 + 0.33333 \times 10^0)$$

$$= fe(1.04761 \times 10^0) = fe(0.104761 \times 10^1) \\ = 0.10476 \times 10^1$$

$$x + y = \frac{5}{7} + \frac{1}{3} = \frac{22}{21} = 1.04761904$$

$$|(x+y) - (x \oplus y)| \\ = 0.19 \times 10^{-4}$$

Relative error

$$= 0.182 \times 10^{-4}$$



Arithmetic Errors

Why Numerical Methods? Number System Errors

Example: $x_0^R = 1$

$$x_1^R = \frac{1}{3}$$

$$x_{n+1}^R = \frac{13}{8} x_n^R - \frac{4}{8} x_{n-1}^R \quad ; \quad n \geq 1$$

True solution: $x_n^T = \left(\frac{1}{3}\right)^n$

7 digit chopping in simulations

10



Arithmetic Errors

Why Numerical Methods? Number System Errors

n	x_n^R	x_n^T	Error
0	1.0000000	1.0000000	0.0000000
1	0.3333333	0.3333333	0.0000000
2	0.1111110	0.1111111	0.0000001
3	0.0370366	0.0370370	0.0000004
4	0.0123439	0.0123457	0.0000018
5	0.0041081	0.0041152	0.0000071
6	0.0013432	0.0013717	0.0000285
7	0.0003431	0.0004572	0.0001141
8	-0.0003042	0.0001524	0.0004566
9	-0.0017757	0.0000508	0.0018265
10	-0.0072891	0.0000169	0.0073060
11	-0.0292185	0.0000056	0.0292241
12	-0.1168947	0.0000019	0.1168966
13	-0.4675857	0.0000006	0.4675863
14	-1.8703451	0.0000002	1.8703453
15	-7.4813812	0.0000001	7.4813813

Table 2: Computed values, exact values, and absolute error

Conditioning!



Taylor Series

$$f(x) \in C^\infty[0, b]$$

$$f(x+h) = f(x) + \frac{h}{1!} f'(x) + \frac{h^2}{2!} f''(x) + \dots + \frac{h^n}{n!} f^{(n)}(x) + \dots$$

Truncate it \rightarrow Some n

Mathematical errors

Root finding Methods

$$x^2 - 1 = 0$$

$$x \sin(x) - e^x \cos(x) + \tan(x) \geq 0$$

