

## ME 261: Numerical Analysis

---

• 3.00 credit hours

Prereq.: ME 163/ME 171

• **Course content**

- Approximations and error types
- Roots of polynomials and transcendental equations
- Determinants and matrices
- Solution of linear and non-linear algebraic equations
- Eigenvalues and eigenvectors
- Interpolation methods
- Curve fitting
- Numerical differentiation and integration
- Solution of first-order differential equations
- Solving equations by finite differences

1

## ME 261: Numerical Analysis

---

• **Course teachers**

1. **Dr. S. Reaz Ahmed** (Sunday)

Professor

2. **Dr. M Zakir Hossain** (Saturday, Tuesday)

Assistant Professor

Dept. of Mechanical Engineering

BUET, Dhaka-1000, Bangladesh.

Room: M602 (EME Building)

email: zakir(at)me.buet.ac.bd, zakir92(at)yahoo.com

web: [teacher.buet.ac.bd/zakir](http://teacher.buet.ac.bd/zakir)

2

## ME 261: Numerical Analysis

- **Text book**

-Numerical Methods for Engineers  
by Steven C. Chapra and Raymond P. Canale  
Publishing company: Tata McGraw-Hill

- **Reference book**

1. Applied Numerical Analysis  
by Curtis F. Gerald and Patrik O. Wheatley  
Publishing company: Pearson
2. Numerical Analysis  
by Timothy Sauer  
Publishing company: Pearson

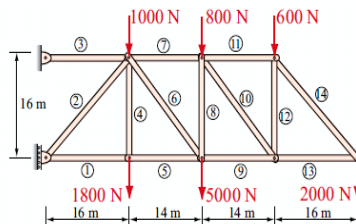
- **ACADEMIC OFFENCES**

Students must write their assignments in their own words.  
If students take an idea, or a passage of text from book, journal, web etc, they must acknowledge this by proper referencing such as footnotes or citations.  
**Plagiarism is a major academic offence.**

3

## Why Numerical Analysis ?

- Example: computing the forces in a TRUSS



- 14 equations with 14 unknowns

$$\begin{aligned}
 -F_1 + F_5 &= 0 \\
 -0.71071F_2 - F_3 + 0.6585F_6 + F_7 &= 0 \\
 -0.7071F_2 - F_4 - 0.7525F_6 &= 1000 \\
 F_4 &= 1800 \\
 -F_7 + 0.6585F_{10} + F_{11} &= 0 \\
 -F_8 - 0.7525F_{10} &= 800 \\
 -F_5 - 0.6585F_6 + F_9 &= 0 \\
 0.7525F_6 + F_8 &= 5000 \\
 -F_{11} + 0.7071F_{14} &= 0 \\
 -F_{12} - 0.7071F_{14} &= 600 \\
 -F_9 - 0.6585F_{10} + F_{13} &= 0 \\
 0.7525F_{10} + F_{12} &= 0 \\
 -F_{13} + 0.7071F_{14} &= 0 \\
 0.7071F_{14} &= 2000
 \end{aligned}$$

4

## Floating Point Arithmetic

---

- fractional quantities are typically represented in computers using floating point format
- this approach is very much similar to scientific notation
- for example, *fixed point number* 17.542 is the same as the floating point number  $.17542 \times 10^2$  which is often displayed as  $.17542e2$
- another example,  $-.004428$  is same as  $-.4428 \times 10^{-2}$

### General form of floating-point number in Computers:

$$\pm .d_1 d_2 d_3 \dots d_p \times B^e$$

- where,  $d_j$ 's are digits or bits with values from 0 to B-1  
B=the number base that is used, e.g. 2, 10,16,8  
P=number of significant bits (digits), that is, the **precision**.  
e=an integer exponent, ranging from  $E_{\min}$  to  $E_{\max}$
- $d_1 d_2 d_3 \dots d_p$  constitute the fractional part of the number

5

## Floating Point Arithmetic

---

### Normalization:

- the fractional digits are shifted and the exponents are adjusted so that  $d_1$  is **NONZERO**  
e.g.  
 $-.004428 \cong -.4428 \times 10^{-2}$  is the normalized form
- In base-2(binary) system normalization means *the first bit is always 1*

### IEEE Standard for the floating-point numbers:

- consists of a set of binary representations of real numbers
- normalization gives the advantage that the **first bit does not need to be stored**
- Hence, nonzero binary floating point numbers can be expressed as,

$$\pm (1 + f) \times 2^e$$

where  $f$  =the mantissa (or the fraction)  
e= exponent

6

## Floating Point Arithmetic

- e.g. the decimal number 9, which is 1001 in binary, would be stored as,

$$+1.001 * 2^3$$

although the original number has 4 significant digits, we only have to store the 3 fractional bits: .001

- There are three commonly used levels of precision for floating point numbers:
  - single precision (32 bits)
  - double precision (64 bits)
  - long double (80 bits)

The bits are divided as follows:

precision	sign	exponent	mantissa
single	1	8	23
double	1	11	52
long double	1	15	64

7

## Floating Point Arithmetic

- By default, MATLAB has adopted IEEE double precision format



Figure: The manner in which a floating point number is stored in IEEE double precision format

Note: because of normalization, 53 bits can be stored in mantissa.

- Sign bit is 0 for positive numbers and 1 for negative numbers

8







## Machine Epsilon $\epsilon_{\text{mach}}$

---

- Double precision numbers can be computed reliably

up to machine epsilon.

- MATLAB has a library function `eps` to display  $\epsilon_{\text{mach}}$

15

## Absolute and relative error

---

Let  $X_{\text{comp}}$  be a computed version of the exact quantity  $X_{\text{exact}}$

$$\text{Absolute error} = |X_{\text{comp}} - X_{\text{exact}}| \qquad \text{Relative error} = \frac{|X_{\text{comp}} - X_{\text{exact}}|}{|X_{\text{exact}}|}$$

- **absolute error** does not account the order of magnitude of the value under examination.
- **Relative error** can be expressed in %
- For example, we want to measure the length of bridge and a rivet.
- We measured the length of bridge is 9999 cm and the length of the rivet is 9 cm.
- The exact values are 10000 and 10 cm, respectively.
- The absolute error for both cases is 1 cm.
- The percentage relative errors are 0.01% and 10%.
- That means, measurement of the length of the bridge is quite accurate, but the measurement of the rivet has significant error.

16



## Arithmetic accuracy in Computers

---

- *Roundoff errors* arise because digital computers cannot represent some quantities exactly.
- Aside from the limitations of a computer's number system, the actual arithmetic manipulation involving this numbers can also result in *roundoff error*.
- To understand how this occurs, let's look at how the computer performs simple addition and subtraction.
- Because of their familiarity, normalized base-10 numbers will be employed to illustrate the effect of roundoff errors on simple arithmetic manipulations.
- Other number bases would behave in a similar fashion.
- To simplify the discussion, we assume a hypothetical computer with a 3-digit mantissa, and a 1-digit exponent, i.e.  $B=10$ ,  $p=3$ ,  $-9 \leq e \leq 9$ .

17

## Arithmetic accuracy in Computers

---

- When two floating-point numbers are added /or subtracted,
  - the numbers are first expressed so that they have *the same exponents*.

examples,

**Addition:**

- if we want to add  $1.371 + 0.02693$
- the computer would express the numbers as

$$\begin{array}{r} 0.1371 \quad * 10^1 \\ + \quad 0.002693 * 10^1 \\ \hline 0.139793 * 10^1 \end{array}$$

← Align the decimal points

- Because this hypothetical computer only carries 3-digit mantissa
  - if rounded,  $\longrightarrow 0.140 * 10^1$
  - if chopped,  $\longrightarrow 0.139 * 10^1$

18

## Arithmetic accuracy in Computers

---

### Subtraction:

- compute  $76.4 - 76.3$  on a 3-decimal-digit computer.
- The computer would express the numbers as

$$\begin{array}{r} 0.764 * 10^2 \\ - 0.763 * 10^2 \\ \hline 0.001 * 10^2 \end{array}$$

← 3 significant digits

- After normalization,  $0.100 * 10^0$

Two non-significant zeros are appended

### Subtractive cancellation/loss of significant:

- subtraction of two nearly equal numbers.
- a major source of error in floating-point operations.

19

## Subtractive cancellation/loss of significant

---

- subtraction of two nearly equal numbers
- a major source of error in floating-point operations

### Example-1

### Example-2

### Lesson

- It is important to find ways to avoid subtracting nearly equal numbers in calculations

20

## Arithmetic accuracy in Computers

---

- Roundoff error can happen in several circumstances other than just storing numbers - for example:
  - *Large computations* - if a process performs a large number of computations, roundoff errors may build up to become significant
    - Example
  - *Adding a Large and a Small Number* - Since the small number's mantissa is shifted to the right to be the same scale as the large number, digits are lost
  - *Smearing* - Smearing occurs whenever the individual terms in a summation are larger than the summation itself.
    - $(x+10^{-20})-x = 10^{-20}$  mathematically, but  
 $x=1$ ;  $(x+1e^{-20})-x$  gives a 0 in MATLAB!

21

## Adding a Large and a Small Number

---

- Suppose we want to add a small number 0.001 to a large number 4000 using a hypothetical computer with the 4-digit mantissa and the 1-digit exponent

$$\begin{array}{r}
 4000 \quad 0.400 \quad * 10^4 \\
 0.001 \quad 0.000\ 000\ 1 \quad * 10^4 \\
 \hline
 0.400\ 000\ 1 \quad * 10^4
 \end{array}$$

- The final result will be  $0.4000 * 10^4$  (since we have 4-digit mantissa)
- It seems that we have not performed any addition !!!
- This type of error may occur in the computation of an infinite series

### Example

- One way to mitigate this type of error is to sum the series in reverse order
- In this way each new term will be of comparable magnitude to the accumulated sum

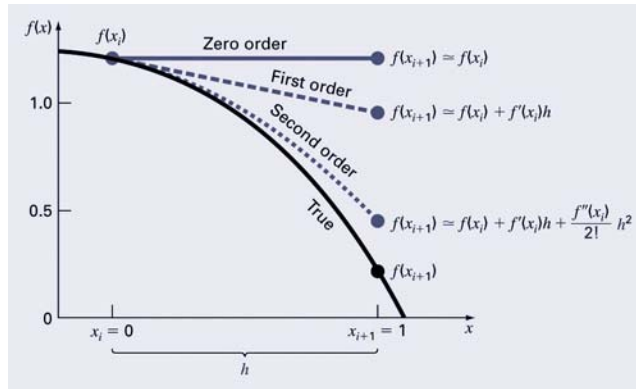
22

## Truncation Error

Truncation errors are those that result from using an approximation in place of an exact mathematical procedure.

Example: The Taylor Series

$$f(x_{i+1}) = f(x_i) + f'(x_i)h + \frac{f''(x_i)}{2!}h^2 + \frac{f^{(3)}(x_i)}{3!}h^3 + \dots + \frac{f^{(n)}(x_i)}{n!}h^n + R_n$$



23

## Truncation Error

$R_n = O(h^{n+1})$ , meaning:

- The more terms are used, the smaller the error, and
- The smaller the spacing, the smaller the error for a given number of terms.

Example 2: Numerical differentiation

- The first order Taylor series can be used to calculate approximations to derivatives:

$$f'(x_i) = \frac{f(x_{i+1}) - f(x_i)}{h} + O(h)$$

If  $h$  is decreased ( $\downarrow$ ) truncation error decreases ( $\downarrow$ )

24

## Error Propagation

---

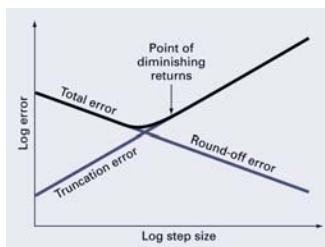
- Functions of a single variable
  - Derivation using Taylor series for a single variable
  - Example
- Functions of more than one variable
  - Derivation using multivariable version of Taylor series
  - Example

25

## Total Numerical Error

---

- The *total numerical error* is the summation of the truncation and roundoff errors.
- The truncation error generally *increases* as the step size increases,
- To minimize the roundoff error, increase the number of precision (bits in mantissa) in computers
- Round off error increases due to subtractive cancellation
  - Or due to an increase in the number of computations in an analysis.
- A decrease in stepsize can lead to subtractive cancellation, or to an increase in computations,  $\Rightarrow$  the round off errors are increased.



26

## Other Errors

---

- Blunders - errors caused by malfunctions of the computer or human imperfection.
- Model errors - errors resulting from incomplete mathematical models.
- Data uncertainty - errors resulting from the inaccuracy and/or imprecision of the measurement of the data.