# Numerical Methods and Optimization

## Topic 1:

Introduction to Numerical Methods

### Lectures 1-4:

# Lecture 1
## Introduction to Numerical Methods

- What are **NUMERICAL METHODS**?
- Why do we need them?
- Topics covered in NMO.

**Reading Assignment:** Pages 3-10 of textbook

# Numerical Methods

**Numerical Methods:**

Algorithms that are used to obtain numerical solutions of a mathematical problem.

**Why do we need them?**

1. No analytical solution exists,
2. An analytical solution is difficult to obtain or not practical.

# What do we need?

**Basic Needs in the Numerical Methods:**

- Practical:

   Can be computed in a reasonable amount of time.

- Accurate:

   - Good approximate to the true value,
   - Information about the approximation error (Bounds, error order,… ).

# Outlines of the Course

- Number Representation
- Approximate solution of nonlinear Equations
- Solution of linear Equations (Direct methods)
- Solution of linear Equations (Iterative methods)
- Polynomial Interpolation
- Least Squares approximation
- Numerical Integration

# Solution of Nonlinear Equations

- Some simple equations can be solved analytically:

$$x^2 + 4x + 3 = 0$$

$$\text{Analytic solution } roots = \frac{-4 \pm \sqrt{4^2 - 4(1)(3)}}{2(1)}$$

$$x = -1 \ and \ \ x = -3$$

- Many other equations have no analytical solution:

$$\left. \begin{array}{c} x^9 - 2x^2 + 5 = 0 \\ x = e^{-x} \end{array} \right\} \text{ No analytic solution}$$

# Methods for Solving Nonlinear Equations

o **Bisection Method**

o **Newton-Raphson Method**

o **Secant Method**

# Solution of Systems of Linear Equations

$$x_1 + x_2 = 3$$

$$x_1 + 2x_2 = 5$$

We can solve it as :

$$x_1 = 3 - x_2, \qquad 3 - x_2 + 2x_2 = 5$$

$$\Rightarrow x_2 = 2, \; x_1 = 3 - 2 = 1$$

What to do if we have

1000 equations in 1000 unknowns.

# Cramer's Rule is Not Practical

Cramer's Rule can be used to solve the system :

$$x_1 = \frac{\begin{vmatrix} 3 & 1 \\ 5 & 2 \end{vmatrix}}{\begin{vmatrix} 1 & 1 \\ 1 & 2 \end{vmatrix}} = 1, \qquad x_2 = \frac{\begin{vmatrix} 1 & 3 \\ 1 & 5 \end{vmatrix}}{\begin{vmatrix} 1 & 1 \\ 1 & 2 \end{vmatrix}} = 2$$

But Cramer's Rule is not practical for large problems.

To solve N equations with N unknowns, we need $(N+1)(N-1)N!$ multiplica tions.

To solve a 30 by 30 system, $2.3 \times 10^{35}$ multiplica tions are needed.

A super computer needs more than $10^{20}$ years to compute this.
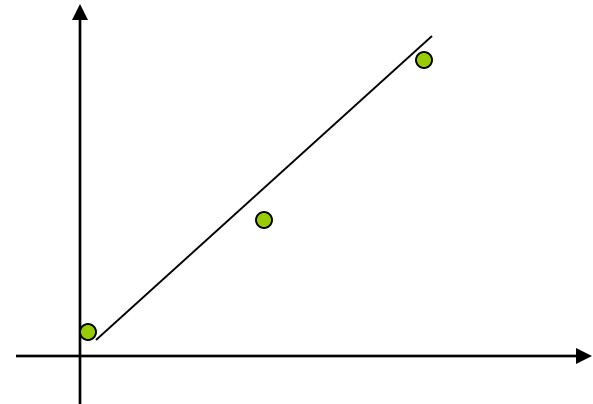
# Methods for Solving Systems of Linear Equations

o **Naive Gaussian Elimination**

o **Gaussian Elimination with Scaled Partial Pivoting**

o **Algorithm for Tri-diagonal Equations**

# Curve Fitting

- Given a set of data:

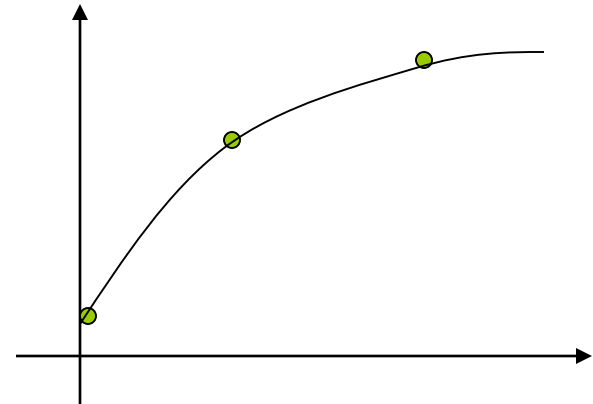| x | 0 | 1 | 2 |
|---|-----|------|------|
| y | 0.5 | 10.3 | 21.3 |

- Select a curve that best fits the data. One choice is to find the curve so that the sum of the square of the error is minimized.

# Polynomial Interpolation

□ Given a set of data:

| $x_i$ | 0 | 1 | 2 |
|---|---|---|---|
| $y_i$ | 0.5 | 10.3 | 15.3 |

□ Find a polynomial *P(x)* whose graph passes through all tabulated points.

$$y_i = P(x_i) \quad if \ x_i \ is \ in \ the \ table$$

# Methods for Curve Fitting

o **Least Squares**
- o **Linear Regression**
- o **Nonlinear Least Squares Problems**

o **Interpolation**
- o **Newton Polynomial Interpolation**
- o **Lagrange Interpolation**

# Integration

- Some functions can be integrated analytically:

$$\int_{1}^{3} x\,dx = \frac{1}{2}x^2\Big|_{1}^{3} = \frac{9}{2} - \frac{1}{2} = 4$$

But many functions have no analytical solutions :

$$\int_{0}^{a} e^{-x^2}\,dx = ?$$

# Methods for Numerical Integration

o **Upper and Lower Sums**

o **Trapezoid Method**

o **Romberg Method**

o **Gauss Quadrature**

# Summary

**Numerical Methods:** Algorithms that are used to obtain numerical solution of a mathematical problem.

**We need them when** No analytical solution exists or it is difficult to obtain it.

## Topics Covered in the Course

- Solution of Nonlinear Equations
- Solution of Linear Equations
- Curve Fitting
  - Least Squares
  - Interpolation
- Numerical Integration

# Lecture 2
# Number Representation and Accuracy

- Number Representation
- Normalized Floating Point Representation
- Significant Digits
- Accuracy and Precision
- Rounding and Chopping

**Reading Assignment:** Chapter 3

# Representing Real Numbers

- You are familiar with the decimal system:

$$312.45 = 3 \times 10^2 + 1 \times 10^1 + 2 \times 10^0 + 4 \times 10^{-1} + 5 \times 10^{-2}$$

- Decimal System:   Base = 10 , Digits (0,1,…,9)

- Standard Representations:

$$\pm \quad 3 \ 1 \ 2 \quad . \quad 4 \ 5$$

sign    integral        fraction

part            part

# Normalized Floating Point Representation

- Normalized Floating Point Representation:

$$\underset{\text{sign}}{\pm} \quad \underset{\text{mantissa}}{\underline{d.\ f_1\ f_2\ f_3\ f_4}} \times \underset{\text{exponent}}{10^{\pm n}}$$

$$d \neq 0, \qquad \pm n : \text{signed exponent}$$

- Scientific Notation: Exactly one non-zero digit appears before decimal point.

- Advantage: Efficient in representing very small or very large numbers.

# Binary System

- Binary System:     Base = 2, Digits {0,1}

$$\pm \quad \underline{1.\ f_1\ f_2\ f_3\ f_4} \quad \times\ 2^{\pm n}$$

sign     mantissa          signed exponent

$$(1.101)_2 = (1 + 1\times 2^{-1} + 0\times 2^{-2} + 1\times 2^{-3})_{10} = (1.625)_{10}$$

# Fact

- Numbers that have a finite expansion in one numbering system may have an infinite expansion in another numbering system:
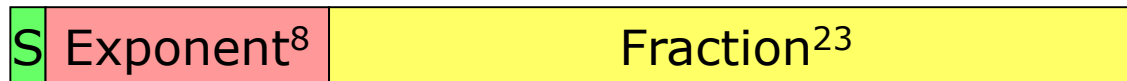
$$(1.1)_{10} = (1.00011001100110011 00...)_2$$

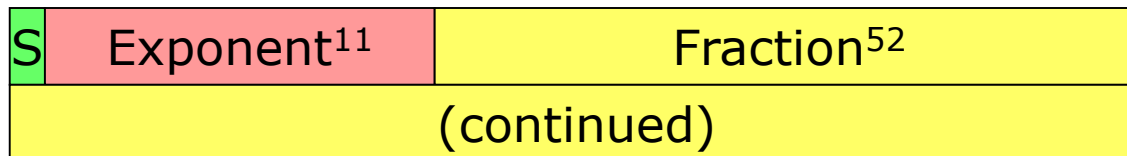- You can never represent 1.1 exactly in binary system.

# IEEE 754 Floating-Point Standard

□ Single Precision (32-bit representation)

  ■ 1-bit Sign + 8-bit Exponent + 23-bit Fraction

| S | Exponent$^8$ | Fraction$^{23}$ |
|---|---|---|

□ Double Precision (64-bit representation)

  ■ 1-bit Sign + 11-bit Exponent + 52-bit Fraction

| S | Exponent$^{11}$ | Fraction$^{52}$ |
|---|---|---|
| (continued) | | |

# Significant Digits

- Significant digits are those digits that can be used with confidence.

- Single-Precision: 7 Significant Digits

  $1.175494\ldots \times 10^{-38}$ to $3.402823\ldots \times 10^{38}$

- Double-Precision: 15 Significant Digits

  $2.2250738\ldots \times 10^{-308}$ to $1.7976931\ldots \times 10^{308}$

# Remarks

- Numbers that can be exactly represented are called machine numbers.

- Difference between machine numbers is not uniform

- Sum of machine numbers is not necessarily a machine number

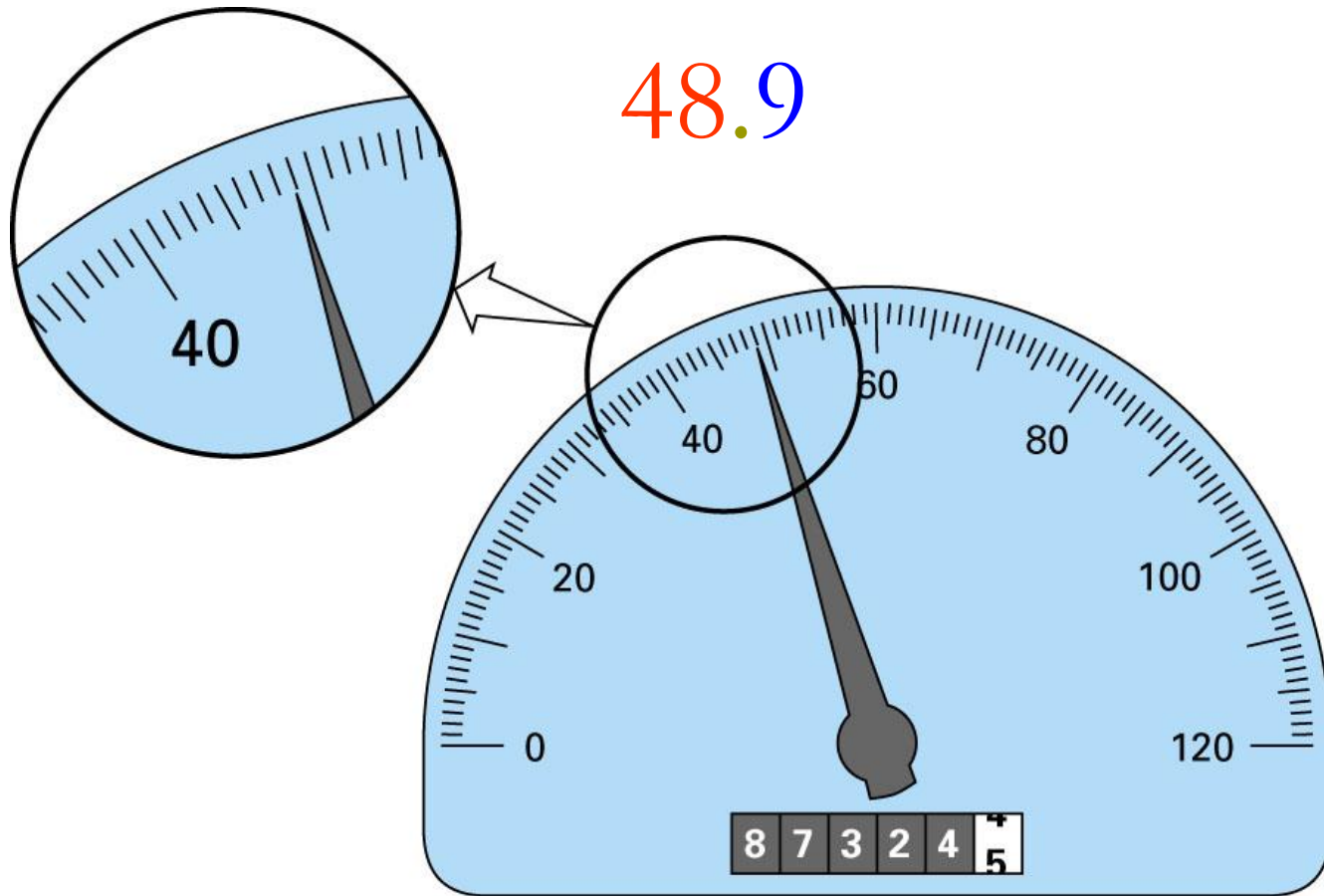# Calculator Example

□ Suppose you want to compute:

3.578 * 2.139

using a calculator with two-digit fractions
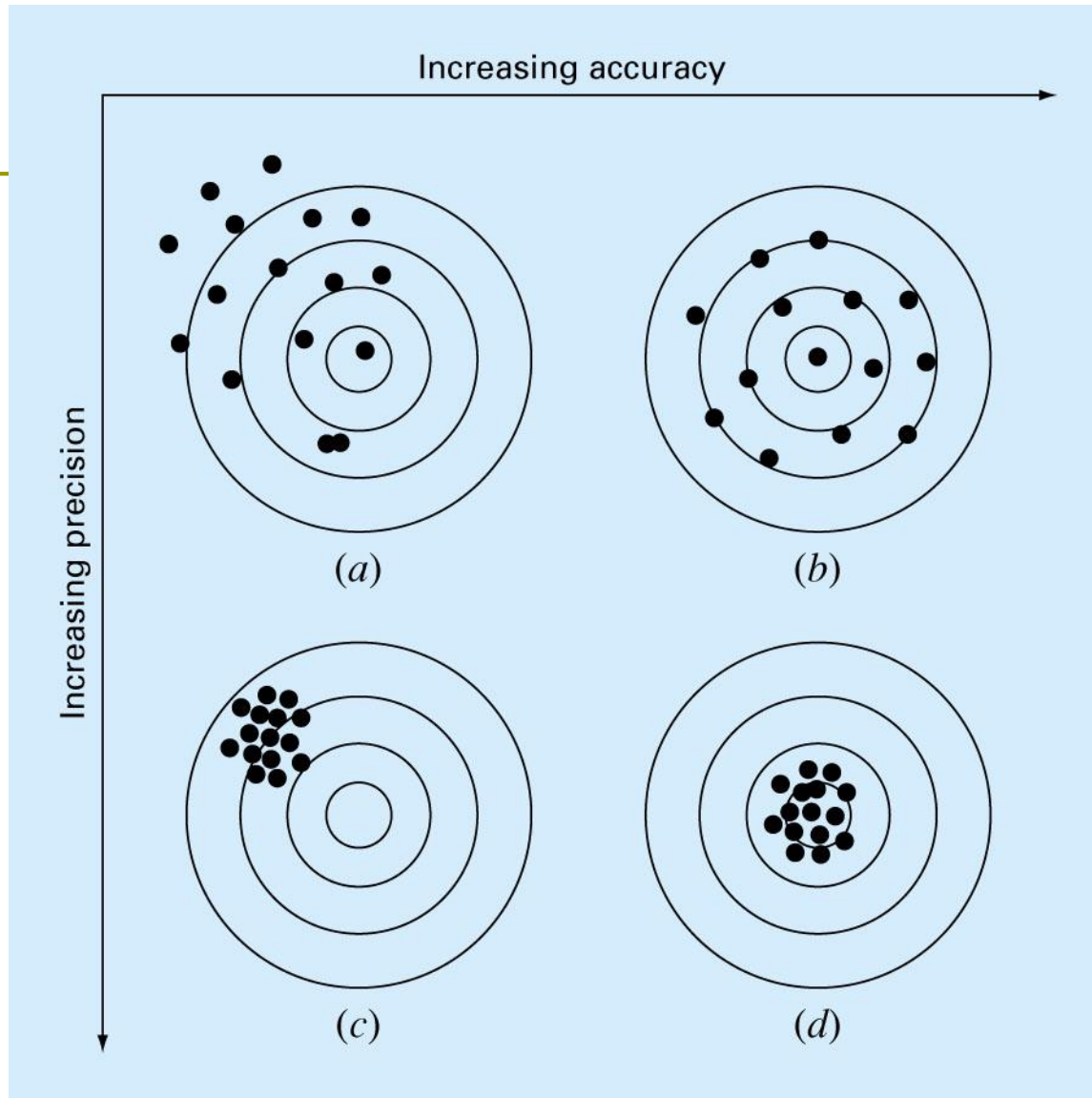
3.57  *  2.13  =  7.60

**True answer:**  7.653342

# Significant Digits - Example



48.9

# Accuracy and Precision

- Accuracy is related to the closeness to the true value.

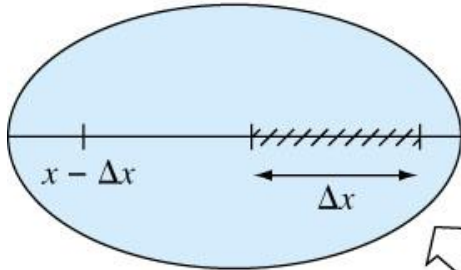- Precision is related to the closeness to other estimated values.
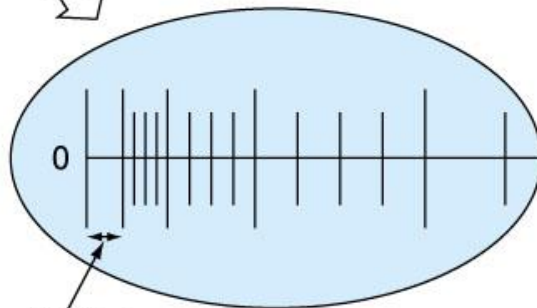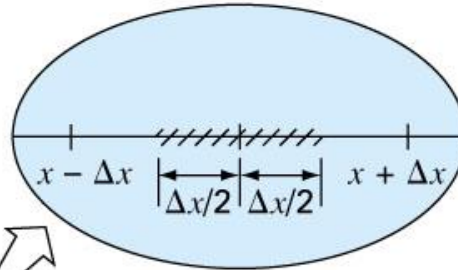
# Rounding and Chopping

- Rounding: Replace the number by the nearest machine number.

- Chopping: Throw all extra digits.

# Rounding and Chopping

# Error Definitions – True Error

Can be computed if the true value is known:

Absolute True Error

$$E_t = \left| \text{true value} - \text{approximation} \right|$$

Absolute Percent Relative Error

$$\varepsilon_t = \left| \frac{\text{true value} - \text{approximation}}{\text{true value}} \right| * 100$$

# Error Definitions – Estimated Error

When the true value is not known:

Estimated  Absolute  Error

$$E_a = \left| \text{current estimate} - \text{previous estimate} \right|$$

Estimated  Absolute  Percent  Relative  Error

$$\varepsilon_a = \left| \frac{\text{current estimate} - \text{previous estimate}}{\text{current estimate}} \right| *100$$

# Notation

We say that the estimate is correct to *n* decimal digits if:

$$|\text{Error}| \leq 10^{-n}$$

We say that the estimate is correct to *n* decimal digits **rounded** if:

$$|\text{Error}| \leq \frac{1}{2} \times 10^{-n}$$

# Summary

- **Number Representation**

  Numbers that have a finite expansion in one numbering system may have an infinite expansion in another numbering system.

- **Normalized Floating Point Representation**

  - Efficient in representing very small or very large numbers,

  - Difference between machine numbers is not uniform,

  - Representation error depends on the number of bits used in the mantissa.