# Numerical Computing

# What is numerical computing?

- Numerical Analysis / computing is the branch of mathematics that provides tools and methods for solving mathematical problems in numerical form.
  - In numerical analysis we are mainly interested in implementation and analysis of numerical algorithms for finding an approximate solution to a mathematical problem
- Numerical computing is an approach for solving complex mathematical problems using only simple arithmetic operations. The approach involves formulation of mathematical models of physical situations that can be solved with arithmetic operations.

## What is error?

Error is a term used to denote the amount by which an approximation fails to equal the exact solution.

error = exact value – approximate value

#### **Sources of errors**

Numerical solutions are an approximation and since the computer program that executes the numerical method might have errors, a numerical solution needs to be examined closely. There are three major sources of error in computation: human errors, truncation errors, and round-off errors.

#### **Human errors**

Typical human errors are arithmetic errors, and/or programming errors: These errors can be very hard to detect unless they give obviously incorrect solution.

#### **■** Truncation errors

Error arise when approximations are used to estimate some quantity. These errors corresponding to the facts that

a finite (infinite)series of computational steps necessary to produce an exact result is "truncated" prematurely after a certain number of steps.

**Example:** the infinite Taylor series

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots$$

truncation error

# Three number systems.

- Beside the various bases for representing numbers binary, decimal and octal there are also three distinct number system that are used in computer machines.
- **1**) integers

$$0, \pm 1, \pm 2 \pm 3 \dots$$

**2**) Fixed point numbers

367.1432

-593.245678953

-0.00123675456

## ■ 3) Floating point numbers:

This number system differs in significant ways from the fixed point number system.

it has three parts,

- I. The sign
- II. The fractional part which often called mantissa.
- III. The exponent part often called the characteristic.

if fixed point number is 367.1432

floating point of this number is  $0.3671432 \times 10^3$ .

## Round off errors

- This is the most basic error in computer. All computing devices represent numbers, except for integers, with some imprecision.
- Digital computers will nearly always use floating point numbers of fixed word length; the true values are not expressed exactly by such representations.

# Round off errors

- Errors arise from the process of rounding off during computations.
- Chopping.

When computing devise drop without rounding is called chopping.

## Example

$$\pi = 0.314159265 \dots \times 10^{1}$$

The five digit floating point form of  $\pi$  using chopping is  $0.31415 \times 10^1 = 3.1415$ 

Is called chopped floating point representation of  $\pi$ .

Since the sixth digit of decimal expansion of  $\pi$  is 9 the floating point form of  $\pi$  using five digit.

 $(0.31415 + 0.00001) \times 10^1 = 3.1416$ 

Is called rounding floating point representation of  $\pi$ .

The error that results from replacing a number with its floating point is called round off error (regardless of whether the rounding or chopping method is used).

#### **■** Absolute error.

The **absolute error** is the magnitude of the difference between the exact value and the approximation.

absolute error = | exact value – approximate value |

#### Relative error.

The **relative error** is the absolute error divided by the magnitude of the exact value.

relative error = | exact value – approximate value | | exact value | | exact value |

■ Approximate error.

approximate error = present approx. - previous approx.

**Example: the infinite Taylor series** 

$$e^{x} = 1 + x + \frac{x^{2}}{2!} + \frac{x^{3}}{3!} + \frac{x^{4}}{4!} + \dots$$

Let x=1.2

$$= e^{1.2} = 1 + 1.2 + \frac{1.2^2}{2!} + \frac{1.2^3}{3!} + \frac{1.2^4}{4!} + \dots$$

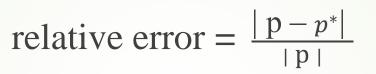
Actual value of  $e^{1.2}=3.320126923$ 

n	Approximate value	Approximate error
1	1	
2	2.2	1.2
3	2.92	0.72

#### **Example:**

Determine the absolute and relative errors when approximation p by  $p^*$  when

a) 
$$p = 0.3000 \times 10^{1}$$
 and  $p^{*} = 0.3100 \times 10^{1}$   
absolute error =  $|p - p^{*}|$   
=  $|0.3000 \times 10^{1} - 0.3100 \times 10^{1}|$   
=  $|-0.1|$   
= 0.1



$$= \frac{\left|0.3000 \times 10^{1} - 0.3100 \times 10^{1}\right|}{\left|0.3000 \times 10^{1}\right|}$$

$$= 0.3333333... \times 10^{-1}$$

b) 
$$p = 0.3000 \times 10^{-3}$$
 and  $p^* = 0.3100 \times 10^{-3}$   
absolute error =  $|p - p^*|$   
=  $|0.3000 \times 10^{-3} - 0.3100 \times 10^{-3}|$   
=  $0.1 \times 10^{-4}$   
relative error =  $\frac{|p - p^*|}{|p|}$   
=  $\frac{|0.3000 \times 10^{-3} - 0.3100 \times 10^{-3}|}{|0.3000 \times 10^{-3}|}$   
=  $0.3333333 \dots \times 10^{-1}$ 

⇒ c) p = 0.3000 × 10<sup>4</sup> and 
$$p^* = 0.3100 \times 10^4$$
  
absolute error =  $|p - p^*|$   
=  $|0.3000 \times 10^4 - 0.3100 \times 10^4|$   
=  $0.1 \times 10^3$   
relative error =  $\frac{|p - p^*|}{|p|}$   
=  $\frac{|0.3000 \times 10^4 - 0.3100 \times 10^4|}{|0.3000 \times 10^4|}$   
= 0.3333333 ... ... × 10<sup>-1</sup>

This example shows that the same relative error occur for widely varying absolute error, As a measure of accuracy the absolute error can be misleading and the relative error is more meaningful because the relative error takes into consideration the size of the value.

#### **NUMERICAL ALGORITHM**

- A complete set of procedures which gives an approximate solution to a mathematical problem.
- An algorithm is a procedure that describes, in an unambiguous manner, a finite sequence of steps to be performed in a specified order. The object of the algorithm is to implement a procedure to solve a problem or approximate a solution to the problem.

- One criterion we will impose on an algorithm whenever possible is that small changes in the initial data produce correspondingly small changes in the final results. An algorithm that satisfies this property is called **stable**; otherwise it is **unstable**.
- Some algorithms are stable only for certain choices of initial data, and are called conditionally stable.

■ The two cases that arise most often in practice are defined as follows.

Suppose that  $E_0 > 0$  denotes an error introduced at some stage in the calculations and  $E_n$  represents the magnitude of the error after n subsequent operations.

- If  $E_n \approx CnE_0$ , where C is a constant independent of n, then the growth of error is said to be **linear**.
- If  $E_n \approx C^n E_0$ , for some C > 1, then the growth of error is called **exponential**.

#### **NUMERICAL ITERATION METHOD**

A mathematical procedure that generates a sequence of improving approximate solution for a class of problems i.e. the process of finding successive approximations.

In the problem of finding the solution of an equation, an iteration method uses as initial guess to generate successive approximation to the solution.

# CONVERGENCE CRITERIA FOR A NUMERICAL COMPUTATION

- If the method leads to the value close to the exact solution, then we say that the method is convergent otherwise the method is divergent.
- ► Why we use numerical iterative methods for solving equations?

As analytic solutions are often either too tiresome or simply do not exist, we need to find an approximate method of solution. This is where numerical analysis comes into picture.

#### LOCAL CONVERGENCE

An iterative method is called locally convergent to a root, if the method converges to root for initial guesses sufficiently close to root.

#### RATE OF CONVERGENCE OF AN ITERATIVE METHOD

Suppose that the sequence  $(x_k)$  converges to "r" then the sequence  $(x_k)$  is said to converge to "r" with order of convergence "a" if there exist a positive constant "p" such that

$$\lim_{k\to\infty}\frac{|x_{k+1}-r|}{|x_k-r|^a}=\lim_{k\to\infty}\frac{\epsilon_{k+1}}{\epsilon_k^a}=p$$

Thus if a=1, the convergence is linear. If a=2, the convergence is quadratic and so on.Where the number "a" is called convergence factor.

