



UTTAR PRADESH  
RAJARSHI TANDON OPEN UNIVERSITY

# UGMM - 10

## Numerical Analysis

Block

# 1

### **SOLUTIONS OF NON-LINEAR EQUATIONS IN ONE VARIABLE**

---

#### **UNIT 1**

**Review of Calculus** **7**

---

#### **UNIT 2**

**Iteration Methods for Locating a Root** **30**

---

#### **UNIT 3**

**Chord Methods for Finding Roots** **46**

---

#### **UNIT 4**

**Approximate Roots of Polynomial Equations** **71**

---

---

## Course Design Committee

---

Prof. S.R.K. Iyengar (*Editor*)  
Dept. of Mathematics  
I.I.T., Delhi

Prof. J.P. Agarwal  
Dept. of Mathematics  
I.I.T., Kharagpur

Dr. U. Anantha Krishnaiah  
Dept. of Mathematics  
KRĒC, Surathkal

Prof. R.K. Jain  
Dept. of Mathematics  
I.I.T., Delhi

Prof. C. Prabhakara Rao  
Dept. of Mathematics and Humanities  
REC, Warangal

**Faculty Members**  
**School of Sciences, IGNOU**

Prof. R.K. Bose  
Dr. V.D. Madan  
Dr. Poornima Mital  
Dr. Manik Patwardhan  
Dr. Parvin Sinclair  
Dr. Sujatha Varma

---

## Block Preparation Team

---

Prof. S.R.K. Iyengar (*Editor*)  
Dept. of Mathematics  
I.I.T., Delhi

Prof. J.P. Agarwal  
Dept. of Mathematics  
I.I.T., Kharagpur

Dr. Sujatha Varma  
School of Sciences  
IGNOU

Prof. G.S. Rao (*Language Editing*)  
School of Humanities  
IGNOU

**Course Coordinator : Dr. Poornima Mital**

---

---

## Production

---

Mr. Balakrishna Selvaraj  
Registrar (PPD)  
IGNOU

Mr. M.P. Sharma  
Joint Registrar (PPD)  
IGNOU

---

May, 1993

© Indira Gandhi National Open University, 1993

ISBN-81-7263-338-6

*All rights reserved. No part of this work may be reproduced in any form, by mimeograph or any other means, without permission in writing from the Indira Gandhi National Open University.*

*Further information on the Indira Gandhi National Open University courses may be obtained from the University's office at Maidan Garhi, New Delhi-110 068.*

Reproduced and reprinted with the permission of Indira Gandhi National Open University  
by Dr.A.K.Singh, Registrar, U.P.R.T.Open University, Allahabad (May, 2013)  
Reprinted by : Nitin Printers, 1 Old Katra, Manmohan Park, Allahabad.

# NUMERICAL ANALYSIS

Mathematical modelling of physical/biological problems generally give rise to ordinary or partial differential equations or an integral equation or in terms of a set of such equation. A number of these problems can be solved exactly by mathematical analysis but most of them cannot be solved exactly. Thus, a need arises to devise numerical methods to solve these problems. These methods for solution of mathematical methods may give rise to a system of algebraic equations or a non-linear equation or system of non-linear equations. The numerical solution of these system of equations are quantitative in nature but when interpreted give qualitative results and are very useful. Numerical analysis deals with the development and analysis of the numerical methods. We are offering this course of numerical analysis to students entering the Bachelor's Degree Programme as an elective subject.

It was in the year 1624 that the English mathematician, Henry Briggs used a numerical procedure to construct his celebrated table of logarithms. The interpolation problem was first taken up by Briggs but was solved by the 17th century mathematicians and physicists, Sir Isaac Newton and James Gregory. Later on, other problems were considered and solved by more and more efficient methods. In recent years the invention and development of electronic calculators/computers have strongly influenced the development of numerical analysis.

This course assumes the knowledge of the course MTE-01 on calculus. It is a prerequisite for this course. Number of results from linear algebra are also used in this course. These results have been stated wherever required. For details of these results our linear algebra course (MTE-02) may be referred. This course is divided into 4 blocks. The first block, deals with the problem of finding approximate roots of a non-linear equation in one unknown. We have started the block with a recall of four important theorems from calculus which are referred to throughout the course. After introducing the concept of 'error' that arise due to approximations, we have discussed two basic approximation methods, namely, bisection and fixed point iteration methods and two commonly used methods, namely, secant and Newton-Raphson methods. In Block 2, we have considered the problem of finding the solution of system of linear equations. We have discussed both direct and iterative methods of solving system of linear equations.

Block 3 deals with the theory of interpolation. Here, we are concerned only with polynomial interpolation. The existence and uniqueness of interpolating polynomials are discussed. Several form of interpolating polynomials like Lagrange's and Newton's divided difference forms with error terms are discussed. This block concludes with a discussion on Newton's forward and backward difference form.

In Block 4, using interpolating polynomials we have obtained numerical differentiation and integration formulae together with their error terms. After a brief introduction to difference equations the numerical solution of the first order ordinary differential equation is dealt with. More precisely, Taylor series, Euler's and second order Runge Kutta methods are derived with error terms for the solution of differential equations.

Each block consists of 4 units. All the concepts given in the units are followed by a number of examples as well as exercises. These will help you get a better grasp of the techniques discussed in this course. We have used a scientific calculator for doing computations throughout the course. While attempting the exercises given in the units, you would also need a calculator which is available at your study centre. The solutions/answers to the exercises in a unit are given at the end of the unit. We suggest that you look at them only after attempting the exercises. A list of symbols and notations are also given in for your reference.

You may like to look up some more books on the subject and try to solve some exercises given in them. This will help you get a better grasp of the techniques discussed in this course. We are giving you a list of titles which will be available in your study centre for reference purposes.

## Some Useful Books

- 1) *Numerical Methods for Scientific and Engineering Computation* by M.K. Jain, S.R.K. Iyengar, R.K. Jain.
- 2) *Elementary Numerical Analysis* by Samuel D. Conte and Carl de Boor.

## NOTATIONS AND SYMBOLS

$\in$	belongs to
$\supset$	contains
$< (\leq)$	less than (less than or equal to)
$> (\geq)$	greater than (greater than or equal to)
$\mathbb{R}$	set of real numbers
$\mathbb{C}$	set of complex numbers
$n!$	$n(n-1) \dots 3 \cdot 2 \cdot 1$ (n factorial)
	closed interval
	open interval
$ x $	absolute value of a number x
i.e.	that is
$\sum_{i=1}^n a_i$	$a_1 + a_2 + \dots + a_n$
$x \rightarrow a$	x tends to a
$\lim_{x \rightarrow a} f(x)$	limit of f(x) as x tends to a
$P_n(x)$	nth degree polynomial
$f'(x)$	derivative of f(x) with respect to x
$f^{(n)}(x)$	nth derivative of f(x) with respect to x
$\approx$	approximately equal to
$\alpha$	alpha
$\beta$	beta
$\gamma$	gamma
$\epsilon$	epsilon
$\pi$	pi
$\Sigma$	capital sigma
$\zeta$	zeta

### Acknowledgements

Prof. R.K. Bose, Dr. Poornima Mital, Dr. Manik Patwardhan, Dr. Parvin Sinclair for comments on the manuscript.  
 Mrs. Manju Sharma for typing the manuscript.  
 Mrs. Wasima Shah for the artwork.



---

## BLOCK INTRODUCTION

---

This is the first of the four blocks which you will be studying in the Numerical Analysis course. In this block we shall be dealing with the problem of finding approximate roots of a non-linear equation in one unknown. In the Elementary Algebra course you have studied some methods for solving polynomial equations of degrees up to and including four. In this block we shall introduce you to some numerical methods for finding solutions of equation. These methods are applicable to polynomial and transcendental equations.

This block consists of four units. In **Unit 1**, we begin with a recall of four important theorems from calculus which are referred to throughout the course. We then introduce you to the concept of 'error' that arise due to approximation. We shall acquaint you with two types of errors that are common in numerical approximation. In **Unit 2**, we shall discuss two basic approximation methods, namely, bisection method and fixed point iteration method. Each of these methods involve a process that is repeated until an answer of required accuracy is achieved. These methods are known as iteration methods. We shall also discuss two accurate methods, namely, secant and Newton-Raphson methods in **Unit 3**. **Unit 4**, which is the last unit of this block, deals with the solutions of the most well-known class of equations, the polynomial equations. For finding the roots of polynomial equations we shall discuss Birge-Vieta and Graeffe's root squaring methods.

As already mentioned in the course introduction, we shall be using a scientific calculator for doing computations throughout the block. While attempting the exercises given in this block, you would also need a calculator which is available at your study centre. We therefore suggest you to go through the instructions manual, supplied with the calculator, before using it.

Lastly we remind you to go through the solved examples carefully, and to attempt all exercises in each unit. This will help you to gain some practice over various methods discussed in this block.



---

# UNIT 1 REVIEW OF CALCULUS

---

## Structure

- 1.1 Introduction
  - Objectives
- 1.2 Three Fundamental Theorems
  - Intermediate Value Theorem
  - Rolle's Theorem
  - Lagrange's Mean Value Theorem
- 1.3 Taylor's Theorem
- 1.4 Errors
  - Round-off Error
  - Truncation Error
- 1.5 Summary
- 1.6 Solutions/Answers

---

## 1.1 INTRODUCTION

---

The study of numerical analysis involves concepts from various branches of mathematics including calculus. In this unit, we shall briefly review certain important theorems in calculus which are essential for the development and understanding of numerical methods. You are already familiar with some fundamental theorems about continuous functions from your calculus course. Here we shall review three theorems given in that course, namely, Intermediate value theorem, Rolle's theorem and Lagrange's mean value theorem. Then we state another important theorem in calculus due to B. Taylor and illustrate the theorem through various examples.

Most of the numerical methods give answers that are approximations to the desired solutions. In this situation, it is important to measure the accuracy of the approximate solution compared to the actual solution. To find the accuracy we must have an idea of the possible errors that can arise in computational procedures. In this unit we shall introduce you to different forms of errors which are common in numerical computations.

The basic ideas and results that we have illustrated in this unit will be used often throughout this course. So we suggest you go through this unit very carefully.

### Objectives

After studying this unit you should be able to :

- apply
  - i) Intermediate value theorem
  - ii) Rolle's theorem
  - iii) Lagrange's mean value theorem
  - iv) Taylor's theorem;
- define the term 'error' in approximation;
- distinguish between rounded-off error and truncation error and calculate these errors as the situation demands.

---

## 1.2 THREE FUNDAMENTAL THEOREMS

---

In this section we shall discuss three fundamental theorems, namely, intermediate value theorem, Rolle's theorem and Lagrange's mean value theorem. All these theorems give properties of continuous functions defined on a closed interval  $[a, b]$ . We shall not prove them here, but we shall illustrate their utility with various examples. Let us take up these theorems one by one.

### 1.2.1 Intermediate Value Theorem

The intermediate value theorem says that a function that is continuous on a closed interval  $[a, b]$  takes on every intermediate value, i.e., every value lying between  $f(a)$  and  $f(b)$  if  $f(a) \neq f(b)$ .

Formally we can state the theorem as follows :

**Theorem 1 :** Let  $f$  be a function defined on a closed interval  $[a, b]$ . Let  $c$  be a number lying between  $f(a)$  and  $f(b)$  (i.e.  $f(a) < c < f(b)$  if  $f(a) < f(b)$  or  $f(b) < c < f(a)$  if  $f(b) < f(a)$ ). Then there exists at least one point  $x_0 \in [a, b]$  such that  $f(x_0) = c$ .

The following figure (Fig. 1) may help you to visualise the theorem more easily. It gives the graph of a function  $f$ .

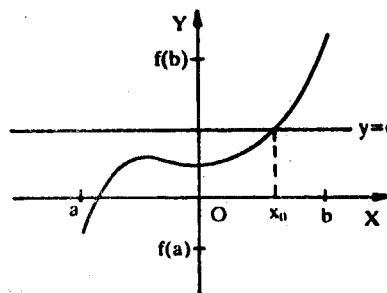


Fig. 1

In this figure  $f(a) < f(b)$ . The condition  $f(a) < c < f(b)$  implies that the points  $(a, f(a))$  and  $(b, f(b))$  lie on opposite sides of the line  $y = c$ . This, together with the fact that  $f$  is continuous, implies that the graph crosses the line  $y = c$  at some point. In Fig. 1 you see that the graph crosses the line  $y = c$  at  $(x_0, c)$ .

The importance of this theorem is as follows : If we have a continuous function  $f$  defined on a closed interval  $[a, b]$ , then the theorem guarantees the existence of a solution of the equation  $f(x) = c$ , where  $c$  is as in Theorem 1. However, it does not say what the solution is. We shall illustrate this point with an example.

**Example 1 :** Find the value of  $x$  in  $0 \leq x \leq \frac{\pi}{2}$  for which  $\sin(x) = \frac{1}{2}$ .

**Solution :** You know that the function  $f(x) = \sin x$  is continuous on  $\left[0, \frac{\pi}{2}\right]$ . Since  $f(0) = 0$  and  $f\left(\frac{\pi}{2}\right) = 1$ , we have  $f(0) < \frac{1}{2} < f\left(\frac{\pi}{2}\right)$ . Thus  $f$  satisfies all the conditions of Theorem 1.

Therefore, there exists at least one value of  $x$ , say  $x_0$ , such that  $f(x_0) = \frac{1}{2}$ , that is, the theorem guarantees that there exists a point  $x_0$  such that  $\sin(x_0) = \frac{1}{2}$ . Let us try to find this point from the graph of  $\sin x$  in  $\left[0, \frac{\pi}{2}\right]$  (see Fig. 2).

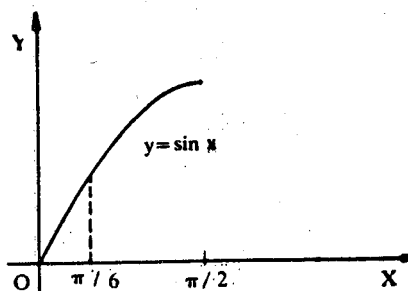


Fig. 2

From the figure, you can see that the line  $x = \frac{1}{2}$  cuts the graph at the point  $\left(\frac{\pi}{6}, \frac{1}{2}\right)$ . Hence there exists a point  $x_0 = \frac{\pi}{6}$  in  $\left[0, \frac{\pi}{2}\right]$  such that  $\sin(x_0) = \frac{1}{2}$ .

Let us consider another example.

**Example 2 :** Show that the equation  $2x^3 + x^2 - x + 1 = 5$  has a solution in the interval  $[1, 2]$ .

**Solution :** Let  $f(x) = 2x^3 + x^2 - x + 1$ . Since  $f$  is a polynomial in  $x$ ,  $f$  is continuous in  $[1, 2]$ . Also  $f(1) = 3$ ,  $f(2) = 19$  and 5 lies between  $f(1)$  and  $f(2)$ . Thus  $f$  satisfies all conditions of Theorem 1. Therefore, there exists a number  $x_0$  between 1 and 2 such that  $f(x_0) = 5$ . That is, the equation  $2x^3 + x^2 - x + 1 = 5$  has solution in the interval  $[1, 2]$ .

Thus we saw that the theorem enables us in establishing the existence of the solutions of certain equations of the type  $f(x) = 0$  without actually solving them. In other words, if you want to find an interval in which a solution (or root) of  $f(x) = 0$  exists, then find two numbers  $a, b$  such that  $f(a) \cdot f(b) < 0$ . Theorem 1, then states that the solution lies in  $]a, b[$ . We shall need some other numerical methods for finding the actual solution. We shall study the problem of finding solutions of the equation  $f(x) = 0$  more elaborately in Unit 2.

Why don't you try an exercise now.

---

E1) Show that the following equations have a solution in the interval given alongside.

a)  $x^3 - x - 5 = 0, [0, 2]$ .

b)  $\sin x + x = 1, \left[0, \frac{\pi}{6}\right]$ .

---

Let us now discuss another important theorem in calculus.

### 1.2.2 Rolle's Theorem

In this section we shall review the Rolle's theorem. The theorem is named after the seventeenth century French mathematician Michel Rolle (1652-1719).

**Theorem 2 (Rolle's Theorem) :** Let  $f$  be a continuous function defined on  $[a, b]$  and differentiable on  $]a, b[$ . If  $f(a) = f(b)$ , then there exists a number  $x_0$  in  $]a, b[$  such that  $f'(x_0) = 0$ .

Geometrically, we can interpret the theorem easily. You know that since  $f$  is continuous, the graph of  $f$  is a smooth curve (see Fig. 3).

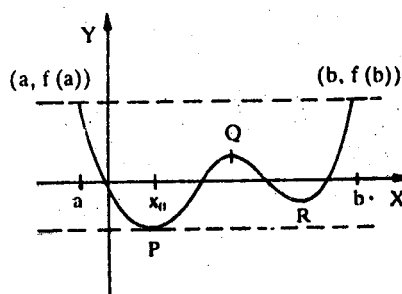


Fig. 3

You have already seen in your calculus course that the derivative  $f'(x_0)$  at some point  $x_0$  gives the slope of the tangent at  $(x_0, f(x_0))$  to the curve  $y = f(x)$ . Therefore the theorem states that if the end values  $f(a)$  and  $f(b)$  are equal, then there exists a point  $x_0$  in  $]a, b[$  such that the slope of the tangent at the point  $P(x_0, f(x_0))$  is zero, that is, the tangent is parallel to  $x$ -axis at that point (see Fig. 3). In fact we can have more than one point at which  $f'(x) = 0$  as shown in Fig. 3. This shows that the number  $x_0$  in Theorem 2 may not be unique.

The following example gives an application of Rolle's theorem.

**Example 3 :** Use Rolle's theorem to show that there is a solution of the equation  $\cot x = x$  in  $\left]0, \frac{\pi}{2}\right[$ .

**Solution :** Here we have to solve the equation  $\cot x - x = 0$ . We rewrite  $\cot x - x$  as  $\frac{\cos x - x \sin x}{\sin x}$ . Solving the equation  $\frac{\cos x - x \sin x}{\sin x} = 0$  in  $\left]0, \frac{\pi}{2}\right[$  is same as solving the equation  $\cos x - x \sin x = 0$ . Now we shall see whether we can find a function  $f$  which satisfies the conditions of Rolle's theorem and for which  $f'(x) = \cos x - x \sin x$ . Our experience in differentiation suggests that we try  $f(x) = x \cos x$ . This function  $f$  is continuous in  $\left]0, \frac{\pi}{2}\right[$ , differentiable in  $\left]0, \frac{\pi}{2}\right[$  and the derivative  $f'(x) = \cos x - x \sin x$ . Also  $f(0) = 0 = f\left(\frac{\pi}{2}\right)$ . Thus  $f$  satisfies all the requirements of Rolle's theorem. Hence, there exists a point  $x_0$  in  $]a, b[$  such that  $f'(x_0) = \cos x_0 - x_0 \sin x_0 = 0$ . This shows that a solution to the equation  $\cot x - x = 0$  exists in  $\left]0, \frac{\pi}{2}\right[$ .

You can try the following exercise on the same lines as Example 3.

---

E2) Using Rolle's theorem show that there is a solution to the equation  $\tan x - 1 + x = 0$  in  $]0, 1[$ .

---

Now, let us look at Fig. 3 carefully. We see that the line joining  $(a, f(a))$  and  $(b, f(b))$  is parallel to the tangent at  $(x_0, f(x_0))$ . Does this property hold when  $f(a) \neq f(b)$  also? In other words, does there exist a point  $x_0$  in  $]a, b[$  such that the tangent at  $(x_0, f(x_0))$  is parallel to the line joining  $(a, f(a))$  and  $(b, f(b))$ ? The answer to this question is the content of the well-known theorem, "Lagrange's mean value theorem", which we discuss next.

### 1.2.3 Lagrange's Mean Value Theorem

This theorem was first proved by the French mathematician Count Joseph Louis Lagrange (1736-1813).

**Theorem 3 :** Let  $f$  be a continuous function defined on  $[a, b]$  and differentiable in  $]a, b[$ . Then there exists a number  $x_0$  in  $]a, b[$  such that

$$f'(x_0) = \frac{f(b) - f(a)}{b - a} \quad \dots (1)$$

Geometrically we can interpret this theorem as given in Fig. 4.

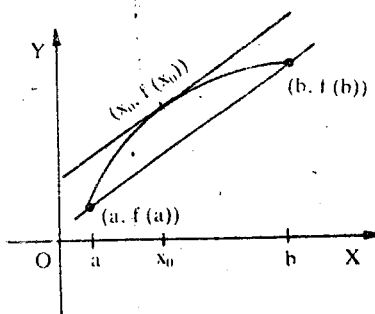


Fig. 4

In this figure you can see that the straight line connecting the end points  $(a, f(a))$  and  $(b, f(b))$  of the graph is parallel to some tangent to the curve at an intermediate point.

You may be wondering why this theorem is called 'mean value theorem'. This is because of the following physical interpretation.

Suppose  $f(t)$  denotes the position of an object at time  $t$ . Then the average (mean) velocity during the interval  $[a, b]$  is given by

$$\frac{f(b) - f(a)}{b - a}$$

Now Theorem 3 states that this mean velocity during an interval  $[a, b]$  is equal to the velocity  $f'(x_0)$  at some instant  $x_0$  in  $]a, b[$ .

We shall illustrate the theorem with an example.

**Example 4 :** Apply the mean value theorem to the function  $f(x) = \sqrt{x}$  in  $[0, 2]$  (see Fig. 5).

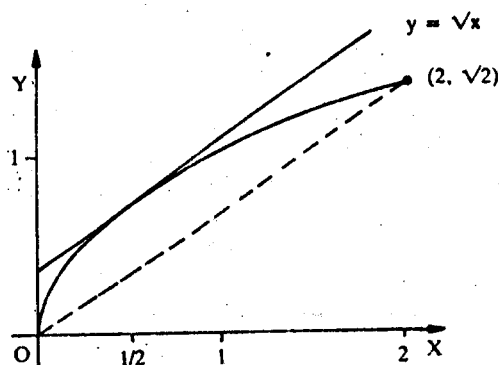


Fig. 5 : Graph of  $f(x) = \sqrt{x}$ .

**Solution :** We first note that the function  $f(x) = \sqrt{x}$  is continuous on  $[0, 2]$  and differentiable in  $]0, 2[$  and  $f'(x) = \frac{1}{2\sqrt{x}}$ .

Therefore by Theorem 3, there exists a point  $x_0$  in  $]0, 2[$  such that

$$f(2) - f(0) = f'(x_0)(2 - 0)$$

Now  $f(2) = \sqrt{2}$  and  $f(0) = 0$  and  $f'(x_0) = \frac{1}{2\sqrt{x_0}}$ .

Therefore we have

$$\sqrt{2} = \frac{1}{\sqrt{x_0}}$$

$$\text{i.e. } \sqrt{x_0} = \frac{1}{\sqrt{2}} \text{ and } x_0 = \frac{1}{2}.$$

Thus we get that the line joining the end points  $(0, 0)$  and  $(2, \sqrt{2})$  of the graph of  $f$  is parallel to the tangent to the curve at the point  $\left(\frac{1}{2}, \frac{1}{\sqrt{2}}\right)$ .

We shall consider one more example.

**Example 5 :** Consider the function  $f(x) = (x-1)(x-2)(x-3)$  in  $[0, 4]$ . Find a point  $x_0$  in  $]0, 4[$  such that

$$f'(x_0) = \frac{f(4) - f(0)}{4 - 0}.$$

**Solution :** We rewrite the function  $f(x)$  as

$$f(x) = (x-1)(x-2)(x-3) = x^3 - 6x^2 + 11x - 6$$

We know that  $f(x)$  is continuous on  $[0, 4]$ , since  $f$  is a polynomial in  $x$ . Also the derivative

$$f'(x) = 3x^2 - 12x + 11$$

exists in  $]0, 4[$ . Thus  $f$  satisfies all conditions of the mean value theorem. Therefore, there exists a point  $x_0$  in  $]0, 4[$  such that

$$f'(x_0) = \frac{f(4) - f(0)}{4 - 0}$$

$$\text{i.e., } 3x_0^2 - 12x_0 + 11 = \frac{6+6}{4-0} = 3$$

$$\text{i.e., } 3x_0^2 - 12x_0 + 8 = 0$$

This is a quadratic equation in  $x_0$ . The roots of this equation are

$$\frac{6 + 2\sqrt{3}}{8} \text{ and } \frac{6 - 2\sqrt{3}}{8}$$

Taking  $\sqrt{3} \approx 1.732$ , we see that there are two values for  $x_0$  lying in the interval  $]0, 4[$ .

The above example shows that the number  $x_0$  in Theorem 3 may not be unique. Again, as we mentioned in the case of Theorems 1 and 2, the mean value theorem guarantees the existence of a point only.

Why don't you try some exercises on your own?

E3) Let  $f(x) = \frac{1}{3}x^3 + 2x$ . Find a number  $x_0$  in  $]0, 3[$  such that

$$f'(x_0) = \frac{f(3) - f(0)}{3 - 0}$$

E4) Find all numbers  $x_0$  in the interval  $]-2, 1[$  for which the tangent to the graph of  $f(x) = x^3 + 4$  is parallel to the line joining the end points  $(-2, f(-2))$  and  $(1, f(1))$ .

E5) Show that Rolle's theorem is a special case of mean value theorem.

So far we have used the mean value theorem to show the existence of a point satisfying Eqn.1. Next we shall consider an example which shows another application of mean value theorem.

**Example 6 :** Find an approximate value of  $\sqrt[3]{26}$  using the mean value theorem.

**Solution :** Consider the function  $f(x) = x^{1/3}$ . Then  $f(26) = \sqrt[3]{26}$ . The number nearest to 26 for which the cube root is known is 27, i.e.  $f(27) = \sqrt[3]{27} = 3$ . Now we shall apply the mean value theorem to the function  $f(x) = x^{1/3}$  in the interval  $]26, 27[$ . The function  $f$  is continuous in  $[26, 27]$  and the derivative is

$$f'(x) = \frac{1}{3x^{2/3}}$$

Therefore, there exists a point  $x_0$  between 26 and 27 such that

$$\sqrt[3]{27} - \sqrt[3]{26} = \frac{1}{3x_0^{2/3}} (27 - 26)$$

$$\text{i.e. } \sqrt[3]{26} = 3 - \frac{1}{3x_0^{2/3}} \quad \dots (2)$$

Since  $x_0$  is close to 27, we approximate  $\frac{1}{3x_0^{2/3}}$  by  $\frac{1}{3(27)^{2/3}}$ , i.e.;

$$\frac{1}{3x_0^{2/3}} \approx \frac{1}{27}$$

The symbol  $\approx$  means approximately equal to.



Substituting this value in Eqn. (2) we get

$$\sqrt[3]{26} = 3 - \frac{1}{27} = 2.963.$$

Note that in writing the value of  $\sqrt[3]{26}$  we have rounded off the number after three decimal places. Using the calculator we find that the exact value of  $\sqrt[3]{26}$  is 2.9624961.

We have given this example just to illustrate the usefulness of the theorem. The mean value theorem has got many other applications which you will come across in later units.

Now we shall discuss another theorem in calculus.

### 1.3 TAYLOR'S THEOREM

You are already familiar with the name of the English mathematician Brook Taylor (1685-1731) from your calculus course. In this section we shall introduce you to a well-known theorem due to B. Taylor. Here we shall state the theorem without proof and discuss some of its applications.

You are familiar with polynomial equations of the form  $f(x) = a_0 + a_1 x + \dots + a_n x^n$  where  $a_0, a_1, \dots, a_n$  are real numbers. We can easily compute the value of a polynomial at any point  $x = a$  by using the four basic operations of addition, multiplication, subtraction and division. On the other hand there are functions like  $e^x$ ,  $\cos x$ ,  $\ln x$  etc. which occur frequently in all branches of mathematics which cannot be evaluated in the same manner. For example, evaluating the function  $f(x) = \cos x$  at 0.524 is not so simple. Now, to evaluate such functions we try to approximate them by polynomials which are easier to evaluate. Taylor's theorem gives us a simple method for approximating functions  $f(x)$  by polynomials.

Let  $f(x)$  be a real-valued function defined on  $R$  which is  $n$ -times differentiable (see MTE-01 Calculus Unit 6, Block 2). Consider the function

$$P_1(x) = f(x_0) + (x - x_0) f'(x_0)$$

where  $x_0$  is any given real number.

Now  $P_1(x)$  is a polynomial in  $x$  of degree 1 and  $P_1(x_0) = f(x_0)$  and  $P_1'(x_0) = f'(x_0)$ . The polynomial  $P_1(x)$  is called the **first Taylor polynomial of  $f(x)$  at  $x_0$** . Now consider another function

$$P_2(x) = f(x_0) + (x - x_0) f'(x_0) + \frac{(x - x_0)^2}{2!} f''(x_0).$$

Then  $P_2(x)$  is a polynomial in  $x$  of degree 2 and  $P_2(x_0) = f(x_0)$ ,  $P_2'(x_0) = f'(x_0)$  and  $P_2''(x_0) = f''(x_0)$ .  $P_2(x)$  is called the **second Taylor polynomial of  $f(x)$  at  $x_0$** .

Similarly we can define the  **$r$ th Taylor polynomial of  $f(x)$  at  $x_0$**  where  $1 \leq r \leq n$ . The  $r$ th Taylor polynomial at  $x_0$  is given by

$$P_r(x) = f(x_0) + (x - x_0) f'(x_0) + \dots + \frac{f^{(r)}(x_0)}{r!} (x - x_0)^r. \quad \dots (3)$$

You can check that  $P_r(x_0) = f(x_0)$ ,  $P_r'(x_0) = f'(x_0)$ ,  $\dots$

$$P_r^{(r)}(x_0) = f^{(r)}(x_0) \quad (\text{see E6})$$

Let us consider an example.

**Example 7 :** Find the fourth Taylor polynomial of  $f(x) = \ln x$  about  $x_0 = 1$ .

**Solution :** The fourth Taylor polynomial of  $f(x)$  is given by

$$P_4(x) = f(1) + (x - 1) f'(1) + \frac{(x - 1)^2}{2!} f''(1) + \frac{(x - 1)^3}{3!} f^{(3)}(1) + \frac{(x - 1)^4}{4!} f^{(4)}(1).$$

Now,  $f(1) = \ln 1 = 0$

$$f'(x) = \frac{1}{x}; f'(1) = 1$$

$$f''(x) = \left(-\frac{1}{x^2}\right); f''(1) = -1$$

$$f^{(3)}(x) = \frac{2}{x^3}; f^{(3)}(1) = 2$$

$$f^{(4)}(x) = \frac{-6}{x^4}; f^{(4)}(1) = -6$$

$$\text{Therefore, } P_4(x) = (x-1) - \frac{(x-1)^2}{2} + \frac{(x-1)^3}{3} - \frac{(x-1)^4}{4}$$

Now, you can try some exercises.

E6) If  $P_r$  denotes the  $r$ th Taylor polynomial as given by Eqn (3), then show that  
 $P_r(x_0) = f(x_0), P_r'(x_0) = f'(x_0), \dots, P_r^{(r)}(x_0) = f^{(r)}(x_0).$

E7) Obtain the third Taylor polynomial of  $f(x) = e^x$  about  $x = 0$ .

We are now ready to state the Taylor's theorem.

**Theorem 4 (Taylor's Theorem) :** Let  $f$  be a real valued function having  $(n + 1)$  continuous derivatives on  $]a, b[$  for some  $n \geq 0$ . Let  $x_0$  be any point in the interval  $]a, b[$ . Then for any  $x \in ]a, b[$ , we have

$$\begin{aligned} f(x) = & f(x_0) + \frac{(x-x_0)}{1!} f'(x_0) + \frac{(x-x_0)^2}{2!} f^{(2)}(x_0) + \dots \\ & + \dots + \frac{(x-x_0)^n}{n!} f^{(n)}(x_0) + \frac{(x-x_0)^{n+1}}{n+1!} f^{(n+1)}(c) \end{aligned} \quad \dots (4)$$

where  $c$  is a point between  $x_0$  and  $x$ .

The series given in Eqn. (4) is called the  **$n$ th Taylor's expansion of  $f(x)$  at  $x_0$** .

We rewrite Eqn. (4) in the form

$$f(x) = P_n(x) + R_{n+1}(x)$$

where  $P_n(x)$  is the  **$n$ th Taylor polynomial** of  $f(x)$  about  $x_0$  and

$$R_{n+1}(x) = \frac{(x-x_0)^{n+1}}{n+1!} f^{(n+1)}(c).$$

$R_{n+1}(x)$  depends on  $x, x_0$  and  $n$ .  $R_{n+1}(x)$  is called the **remainder (or error)** of the  $n$ th Taylor's expansion after  $n + 1$  terms.

Suppose we put  $x_0 = a$  and  $x = a + h$  where  $h > 0$ , in Eqn (4). Then any point between  $a$  and  $a + h$  will be of the form  $a + \theta h, 0 < \theta < 1$ .

Therefore, Eqn (4) can be written as

$$f(a+h) = f(a) + h f'(a) + \frac{h^2}{2!} f''(a) + \dots + \frac{h^n}{n!} f^{(n)}(a) + \frac{h^{n+1}}{n+1!} f^{(n+1)}(a + \theta h) \quad \dots (5)$$

Let us now make some remarks on the Taylor's theorem.

**Remark 1 :** Suppose that the function  $f(x)$  in Theorem 4 is a polynomial of degree  $m$ . Then  $f^{(r)}(x) = 0$  for all  $r > m$ . Therefore  $R_{n+1}(x) = 0$  for all  $n \geq m$ . Thus, in this case, the  $m$ th Taylor expansion of  $f(x)$  about  $x_0$  will be

$$f(x) = f(x_0) + \frac{(x-x_0)}{1!} f'(x_0) + \dots + \frac{(x-x_0)^m}{m!} f^{(m)}(x_0).$$

**Note that** the right hand side of the above equation is simply a polynomial in  $(x - x_0)$ .

Therefore, finding Taylor's expansion of a polynomial function  $f(x)$  about  $x_0$  is the same as expressing  $f(x)$  as a polynomial in  $(x - x_0)$  with coefficients from  $\mathbf{R}$ .

**Remark 2 :** Suppose we put  $x_0 = a$ ,  $x = b$  and  $n = 0$  in Eqn. (4). Then Eqn. (4) becomes

$$f(b) = f(a) + f'(c)(b - a)$$

or equivalently

$$f(b) - f(a) = f'(c)(b - a)$$

which is the Lagrange's mean value theorem. Therefore we can consider the mean value theorem as a special case of Taylor's theorem.

Let us consider some examples.

**Example 8 :** Expand  $f(x) = x^4 - 5x^3 + 5x^2 + x + 2$  in powers of  $(x - 2)$ .

**Solution :** The function  $f(x)$  is a polynomial in  $x$  of degree 4. Hence, derivatives of all orders exist and are continuous. Therefore by Taylor's theorem, the 4th Taylor expansion of  $f(x)$  about 2 is given by

$$f(x) = f(2) + \frac{x-2}{1!} f'(2) + \frac{(x-2)^2}{2!} f''(2) + \frac{(x-2)^3}{3!} f^{(3)}(2) + \frac{(x-2)^4}{4!} f^{(4)}(2).$$

Here  $f(2) = 0$

$$\begin{aligned} f'(x) &= 4x^3 - 15x^2 + 10x + 1, & f'(2) &= -7 \\ f''(x) &= 12x^2 - 30x + 10, & f''(2) &= -2 \\ f^{(3)}(x) &= 24x - 30, & f^{(3)}(2) &= 18 \\ f^{(4)}(x) &= 24, & f^{(4)}(2) &= 24 \end{aligned}$$

Hence the expansion is

$$\begin{aligned} f(x) &= -7(x-2) - \frac{2(x-2)^2}{2!} + \frac{18(x-2)^3}{3!} + \frac{24(x-2)^4}{4!} \\ &= -7(x-2) - (x-2)^2 + 3(x-2)^3 + (x-2)^4 \end{aligned}$$

**Example 9 :** Find the  $n$ th Taylor expansion of  $\ln(1+x)$  about  $x = 0$  for  $x \in ]-1, 1[$ .

**Solution :** We first note that the point  $x = 0$  lies in the given interval. Further, the function  $f(x) = \ln(1+x)$  has continuous derivatives of all orders. The derivatives are given by

$$\begin{aligned} f'(x) &= \frac{1}{1+x}, & f'(0) &= 1 \\ f''(x) &= \frac{-1}{(1+x)^2}, & f''(0) &= -1 \\ f^{(3)}(x) &= \frac{(-1)^2 2!}{(1+x)^3}, & f^{(3)}(0) &= 2 \\ &\dots & & \\ &\dots & & \\ &\dots & & \end{aligned}$$

$$f^{(n)}(x) = \frac{(-1)^{n-1} (n-1)!}{(1+x)^n}, \quad f^{(n)}(0) = (-1)^{n-1} (n-1)!.$$

Therefore by applying Taylor's theorem we get that for any  $x \in ]-1, 1[$

$$\ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots + \frac{(-1)^{n-1} x^n}{n} + \frac{(-1)^n n! x^{n+1}}{(n+1)! (1+c)^{n+1}}$$

where  $c$  is a point lying between  $0$  and  $x$ .

Now, let us consider the behaviour of the remainder in a small interval, say,  $[0, 0.5]$ . Then for  $x$  in  $[0, 0.5]$ , we have

$$\left| R_{n+1}(x) \right| = \left| \frac{(-1)^n n! x^{n+1}}{(n+1)! (1+c)^{n+1}} \right|$$

where  $0 < c < x$ .

Since  $|x| < 1$ ,  $|x|^{n+1} < 1$  for any positive integer  $n$ .

Also since  $c > 0$ ,  $\frac{1}{(1+c)^{n+1}} < 1$ . Therefore we have

$$\left| R_{n+1}(x) \right| < \frac{1}{n+1}$$

Now  $\frac{1}{n+1}$  can be made as small as we like by choosing  $n$  sufficiently large i.e.

$$\lim_{n \rightarrow \infty} \frac{1}{n+1} = 0. \text{ This shows that } \lim_{n \rightarrow \infty} \left| R_{n+1}(x) \right| = 0.$$

The above example shows that if  $n$  is sufficiently large, the value of the  $n$ th Taylor polynomial  $P_n(x)$  at any  $x_0$  will be approximately equal to the value of the given function  $f(x_0)$ . In fact, the remainder  $R_{n+1}(x)$  tell(s) us how close the value  $P_n(x_0)$  is to  $f(x_0)$ .

Now we shall make some general observations about the remainder  $R_{n+1}(x)$  in the Taylor's expansion of a function  $f(x)$ .

**Remark 3 :** Consider the  $n$ th Taylor expansion of  $f$  about  $x_0$  given by

$$f(x) = P_n(x) + R_{n+1}(x).$$

Then  $R_{n+1}(x) = f(x) - P_n(x)$ . If  $\lim_{n \rightarrow \infty} R_{n+1}(x) = 0$  for some  $x$ , then for that  $x$  we say that we can approximate  $f(x)$  by  $P_n(x)$  and we write  $f(x)$  as the infinite series.

$$f(x) = f_0(x) + f'(x)(x-x_0) + \frac{f^{(2)}(x_0)}{2!} (x-x_0)^2 + \dots + \frac{f^{(n)}(x_0)}{n!} (x-x_0)^n + \dots$$

$$= \sum_{n=0}^{\infty} \frac{f^{(n)}(x_0)}{n!} (x-x_0)^n \quad \dots (6)$$

You are already familiar with series of this type from your calculus course. This series is called Taylor's series of  $f(x)$ . If we put  $x_0 = 0$  in Eqn. (6) then the series

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(x_0)}{n!} x^n$$

is called Maclaurin's series.

**Remark 4 :** If the remainder  $R_{n+1}(x)$  satisfies the condition that  $\left| R_{n+1}(x) \right| < M$  for some  $n$  at some fixed point  $x = a$ , then  $M$  is called the **bound of the error at  $x = a$** .

In this case we have

$$|R_{n+1}(x)| = |f(x) - P_n(x)| < M$$

That is,  $f(x)$  lies in the interval  $]P_n(x) - M, P_n(x) + M[$ .

Now if  $M$  is considerably small for some  $n$ , then this interval becomes very small. In this case we say that  $f(x)$  is approximately equal to the value of the  $n$ th Taylor polynomial with error  $M$ . Thus the remainder is used to determine a bound for the accuracy of the approximation.

We shall explain these concepts with an example.

**Example 10 :** Find the 2nd Taylor's expansion of  $f(x) = \sqrt{1+x}$  in  $] -1, 1[$  about  $x = 0$ . Find the bound of the error at  $x = 0.2$ .

**Solution :** Since  $f(x) = \sqrt{1+x}$ , we have

$$f(0) = 1$$

$$f'(x) = \frac{1}{2\sqrt{1+x}}, \quad f'(0) = \frac{1}{2}$$

$$f''(x) = -\frac{1}{4}(1+x)^{-3/2}, \quad f''(0) = -\frac{1}{4}$$

$$f^{(3)}(x) = \frac{3}{8}(1+x)^{-5/2},$$

Applying Taylor's theorem to  $f(x)$ , we get

$$\sqrt{1+x} = 1 + \frac{1}{2}x - \frac{1}{8}x^2 + \frac{1}{16}x^3(1+c)^{-5/2}$$

where  $c$  is a point lying between 0 and  $x$ .

The error is given by  $R_3(x) = \frac{x^3}{16}(1+c)^{-5/2}$ .

When  $x = 0.2$ , we have

$$R_3(0.2) = \frac{(0.2)^3}{16(1+c)^{5/2}}$$

where  $0 < c < 0.2$ . Since  $c > 0$  we have

$$\left| \frac{1}{(1+c)^{5/2}} \right| < 1.$$

Hence,

$$|R_3(0.2)| \leq \frac{(0.2)^3}{16} = (0.5) 10^{-3}$$

Hence the bound of the error for  $n = 2$  at  $x = 0.2$  is  $(0.5) 10^{-3}$ .

Why don't you try some exercises now?

E8) Obtain the  $n$ th Taylor expansion of the function  $f(x) = \frac{1}{1+x}$  in  $]-\frac{1}{2}, 1[$  about  $x_0 = 0$ .

E9) Does  $f(x) = \sqrt{x}$  have a Taylor series expansion about  $x = 0$ ? Justify your answer.

E10) Obtain the 8th Taylor expansion of the function  $f(x) = \cos x$  in  $[-\frac{\pi}{4}, \frac{\pi}{4}]$  about  $x_0 = 0$ .

Obtain a bound for the error  $R_9(x)$ .

There are some functions whose Taylor's expansion is used very often. We shall list their expansions here.

$$e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \dots + \frac{x^n}{n!} + \frac{x^{n+1}}{(n+1)!} e^c \dots \quad \dots (7)$$

$$\begin{aligned} \sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots + \frac{(-1)^{n-1} x^{2n-1}}{(2n-1)!} \\ + \frac{(-1)^n x^{2n+1}}{(2n+1)!} \cos(c) \end{aligned} \quad \dots (8)$$

$$\begin{aligned} \cos x = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \dots + \frac{(-1)^n (x)^{2n}}{(2n)!} \\ + \frac{(-1)^{n+1} x^{2n+2}}{(2n+2)!} \cos(c) \end{aligned} \quad \dots (9)$$

$$\frac{1}{1-x} = 1 + x + x^2 + \dots + x^n + \frac{x^{n+1}}{(1-c)^{n+2}} \quad \dots (10)$$

where  $c$ , in each expansion, is as given in Taylor's theorem.

Now, let us consider some examples that illustrate the use of finding approximate values of some functions at certain points using truncated Taylor series.

**Example 11 :** Using Taylor's expansion for  $\sin x$  about  $x = 0$ , find the approximate value of  $\sin 10^\circ$  with error less than  $10^{-7}$ .

**Solution :** The  $n$ th Taylor expansion for  $\sin x$  given in Eqn. (9) is

$$\begin{aligned} \sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots + \frac{(-1)^{n-1} x^{2n-1}}{(2n-1)!} \\ + \frac{(-1)^n x^{2n+1}}{(2n+1)!} \cos c \end{aligned} \quad \dots (11)$$

where  $x$  is the angle measured in radians.

Now, in radian measure, we have

$$10^\circ = \frac{\pi}{18} \text{ radians.}$$

Therefore, by putting  $x = \frac{\pi}{18}$  in Eqn. (11) we get

$$\sin \frac{\pi}{18} = \frac{\pi}{18} - \frac{1}{3!} \left( \frac{\pi}{18} \right)^3 + \frac{1}{5!} \left( \frac{\pi}{18} \right)^5 + \dots + R_{n+1} \left( \frac{\pi}{18} \right)$$

where  $R_{n+1} \left( \frac{\pi}{18} \right)$  is the remainder after  $(n+1)$  terms.

Now

$$R_{n+1} \left( \frac{\pi}{18} \right) = \frac{(-1)^n}{(2n+1)!} \left( \frac{\pi}{18} \right)^{2n+1} \cos c.$$

If we approximate  $\sin \frac{\pi}{18}$  by  $P_n \left( \frac{\pi}{18} \right)$ , then the error introduced will be less than  $10^{-7}$  if

$$\begin{aligned} \left| \sin \left( \frac{\pi}{18} \right) - P_n \left( \frac{\pi}{18} \right) \right| &= \left| R_{n+1} \left( \frac{\pi}{18} \right) \right| \\ &= \left| \frac{(-1)^n}{(2n+1)!} \left( \frac{\pi}{18} \right)^{2n+1} \cos c \right| < 10^{-7}. \end{aligned}$$

Maximizing  $\cos c$ , we require that

$$\frac{1}{(2n+1)!} \left( \frac{\pi}{18} \right)^{2n+1} < 10^{-7}$$

Using the calculator, we find that the value of left hand side of Eqn. (12) for various n is

n	1	2	3
Left hand side	$.89 \times 10^{-3}$	$.13 \times 10^{-5}$	$.99 \times 10^{-9}$

From the table we find that the inequality in (12) is satisfied for  $n = 3$ . Hence the required approximation is

$$\sin\left(\frac{\pi}{18}\right) \approx \frac{\pi}{18} - \frac{1}{3!}\left(\frac{\pi}{18}\right)^3 + \frac{1}{5!}\left(\frac{\pi}{18}\right)^5 = 0.1745445$$

with error less than  $1.0 \times 10^{-7}$ .

Let us now find the approximate value of e using Taylor's theorem.

**Example 12 :** Using Maclaurin's series for  $e^x$ , show that  $e \approx 2.71806$  with error less than 0.001. (Assume that  $e < 3$ ).

**Solution :** The Maclaurin's series for  $e^x$  is

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

Putting  $x = 1$  in the above series, we get

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

Now we have to find n for which

$$|e - P_n(1)| = |R_{n+1}(1)| < 0.001.$$

$$\text{Now } |R_{n+1}(1)| \leq e^c \frac{1}{(n+1)!}$$

Since we have chosen  $x_0 = 0$  and  $x = 1$ , the value c lies between 0 and 1 i.e.  $0 < c < 1$ . Since  $e^c < e < 3$ , we get

$$|R_{n+1}(1)| \leq \frac{3}{(n+1)!}$$

The bound for  $R_{n+1}(1)$  for different n is given in the following table.

n	1	2	3	4	5	6
Bounds for $R_{n+1}$	1.5	.5	.1	.125	.004	.0006

From this table, we see that

$$R_{n+1} < .001 \text{ if } n = 6$$

Thus  $P_6(1)$  is the desired approximation to e. i.e.

$$e \approx 1 + 1 + \frac{1}{2} + \frac{1}{6} + \frac{1}{24} + \frac{1}{120} + \frac{1}{720} = \frac{1957}{720} \approx 2.71806$$

See if you can do the following exercises.

E11) Using Maclaurin's expansion for  $\cos x$ , find the approximate value of  $\cos \frac{\pi}{4}$  with the error bound  $10^{-5}$ .

E12) How large should n be chosen in Maclaurin's expansion for  $e^x$  to have

$$|e^x - P_n(x)| \leq 10^{-5}, \quad -1 \leq x \leq 1$$

In numerical analysis we are concerned with developing a sequence of calculations that will give a satisfactory answer to a problem. Since this process involves a lot of computations, there is a chance for the presence of some errors in these computations. In the next section we shall introduce you to the concept of 'errors' that arise in numerical computations.

## 1.4 ERRORS

In this section we shall discuss the concept of an 'error'. We consider two types of errors that are commonly encountered in numerical computations.

You are already familiar with the rounding off a number which has non-terminal decimal expansion from your school arithmetic. For example we use 3.1425 for  $22/7$ . These rounded off numbers are approximations of the actual values. In any computational procedure we make use of these approximate values instead of the true values. Let  $x_T$  denote the true value and  $x_A$  denote the approximate value. How do we measure the goodness of an approximation  $x_A$  to  $x_T$ ? The simplest measure which naturally comes to our mind is the difference between  $x_T$  and  $x_A$ . This measure is called the 'error'. Formally, we define error as a quantity which satisfies the identity.

$$\text{True value } x_T = \text{Approximate value } x_A + \text{error.}$$

Now if an 'error' in approximation is considerably small (according to some criterion), then we say that ' $x_A$  is a good approximation to  $x$ '.

Let us consider an example.

**Example 13 :** The true value of  $\pi$  is 3.14159265 ... In some mensuration problems the value  $22/7$  is commonly used as an approximation to  $\pi$ . What is the error in this approximation?

**Solution :** The true value of  $\pi$  is

$$\pi = 3.14159265 \dots \quad \dots (13)$$

Now, we convert  $22/7$  to decimal form, so that we can find the difference between the approximate value and true value. Then the approximate value of  $\pi$  is

$$\frac{22}{7} = 3.14285714 \dots \quad \dots (14)$$

Therefore,

$$\text{error} = \text{True value} - \text{approximate value} = -0.00126449 \dots \quad \dots (15)$$

Note that in this case the error is negative. Error can be positive or negative. We shall in general be interested in absolute value of the error which is defined as

$$|\text{error}| = |\text{True value} - \text{approximate value}|$$

For example, the absolute Error in Example 13 is

$$|\text{error}| = |-0.00126449 \dots| = 0.00126\dots$$

Sometimes, when the true value is very large or very small we prefer to study the error by comparing it with the true value. This is known as **Relative error** and we define this error as

$$|\text{Relative error}| = \left| \frac{\text{True value} - \text{approximate value}}{\text{True value}} \right|$$

In the case of Example 13,

$$|\text{Relative error}| = \frac{0.00126449 \dots}{3.14159265 \dots} = 0.00040249966 \dots$$

But note that in certain computations, the true value may not be available. In that case we replace the true value by the computed approximate value in the definition of relative error.

In numerical calculations, you will encounter mainly two types of errors: round-off error and truncation error. We shall discuss these errors in the next two subsections 1.4.1 and 1.4.2 respectively.



### 1.4.1 Round-off Error

Let us look at Example 13 again. You can see that the numbers appearing in Eqns (13), (14) and (15) consist of 8 digits after the decimal point followed by dots. The line of dots indicates that the digits continue and we are not able to write all of them. That is, these numbers cannot be represented exactly by a terminating decimal expansion. Whenever we use such numbers in calculations we have to decide how many digits we are going to take into account. For example, consider again the approximate value of  $\pi$ . If we approximate  $\pi$  using 2 digits after the decimal point (say), chopping off the other digits, then we have

$$\pi \approx 3.14$$

The error in this approximation is

$$\text{error} = 0.00159265 \dots \quad \dots (16)$$

If we use 3 digits after the decimal point, then using chopping we have

$$\pi \approx 3.141$$

In this case the error is given by

$$\text{error} = -0.00059265 \dots \quad \dots (17)$$

Now suppose we consider the approximate value rounded-off to three decimal places. You already know how to round off a number which has non-terminal decimal expansion. Then the value of  $\pi$  rounded-off to 3 digits is 3.142. The error in this case is

$$\text{error} = -0.00040734 \dots$$

which is smaller, in absolute value than 0.00059265 ... given in Eqn. (17). Therefore in general whenever we want to use only a certain number of digits after the decimal point, then it is always better to use the value rounded-off to that many digits because in this case the error is usually small. **The error involved in a process where we use rounding off method is called round-off error.**

We now discuss the concept of floating point arithmetic.

In scientific computations a real number  $x$  is usually represented in the form

$$x = \pm (.d_1 d_2 \dots d_n) 10^m$$

where  $d_1, d_2, \dots, d_n$  are natural numbers between 0 and 9 and  $m$  is an integer called exponent. Writing a number in this form is known as **floating point representation**. We denote this representation by  $fl(x)$ . Such a floating point number is said to be normalized if  $d_1 \neq 0$ . To translate a number into floating point representation we adopt any of the two methods — rounding and chopping. For example, suppose we want to represent the number 537 in the normalized floating point representation with  $n = 1$ , then we get

$$\begin{aligned} fl(537) &= .5 \times 10^3 \text{ chopped} \\ &= .5 \times 10^3 \text{ rounded} \end{aligned}$$

In this case we are getting the same representation in rounding and chopping. Now if we take  $n = 2$ , then we get

$$\begin{aligned} fl(537) &= .53 \times 10^3 \text{ chopped} \\ &= .54 \times 10^3 \text{ rounded} \end{aligned}$$

In this case, the representations are different.

Now if we take  $n = 3$ , then we get

$$\begin{aligned} fl(537) &= .537 \times 10^3 \text{ chopped} \\ &= .537 \times 10^3 \text{ rounded} \end{aligned}$$

The number  $n$  in the floating point representation is called **precision**.

The difference between the true value of a number  $x$  and rounded  $fl(x)$  is called round-off error. From the earlier discussion it is clear that the round-off error decreases when precision increases.

Mathematically we define these concepts as follows :

**Definition 2 :** Let  $x$  be a real number and  $x^*$  be a real number having non-terminal decimal expansion, then we say that  $x^*$  represents  $x$  rounded to  $k$  decimal places if

$$|x - x^*| \leq \frac{1}{2} 10^{-k}, \text{ where } k > 0 \text{ is a positive integer.}$$

Next definition gives us a measure by which we can conclude that the round-off error occurring in an approximation process is negligible or not.

**Definition 3 :** Let  $x$  be a real number and  $x^*$  be an approximation to  $x$ . Then we say that  $x^*$  is accurate to  $k$  decimal places if

$$\frac{1}{2} 10^{-(k+1)} \leq |x - x^*| \leq \frac{1}{2} 10^{-k} \quad \dots (18)$$

Let us consider an example.

**Example 14 :** Find out to how many decimal places the value of  $22/7$  obtained in Example 13 is accurate as an approximation to  $\pi = 3.14159265$  ?

**Solution :** We have already seen in Example 13 that

$$\left| \pi - \frac{22}{7} \right| = 0.00126449 \dots$$

Now  $.0005 < .00126 \dots < 0.005$

$$\text{or } \frac{1}{2} 10^{-3} < .00126 \dots < \frac{1}{2} 10^{-2}$$

Therefore the inequality (18) is satisfied for  $k = 2$ .

Hence, by Definition 3, we conclude that the approximation is accurate to 2 decimal places.

Here is an exercise for you.

E13) In some approximation problems where graphic methods are used, the value  $\frac{355}{113}$  is used as an approximation to  $\pi = 3.14159265 \dots$ . To how many decimal places the value  $\frac{355}{113}$  is accurate as an approximation to  $\pi$ ?

Now we make an important remark.

**Remark 5 :** Round-off errors can create serious difficulties in lengthy computations. Suppose we have a problem which involves a long calculation. In the course of these computations many rounding errors (some positive, and some negative) may occur in a number of ways. At the end of the calculations these errors will get accumulated and we don't know the magnitude of this error. Theoretically it can be large. But, in reality, some of these errors (between positive and negative errors) may get cancelled so that the accumulated error will be much smaller.

Let us now define another type of error called 'Truncation error'.

### 1.4.2 Truncation Error

We shall first illustrate this error with a simple example. In Sec 1.3, we have already discussed how to find approximate value of a certain function  $f(x)$ , for a given value of  $x$ , using Taylor's series expansion. Let

$$f(x) = \sum_{n=0}^{\infty} a_n (x - x_0)^n$$

denote the Taylor's series of  $f(x)$  about  $x_0$ . In practical situations, we cannot, in general, find

the sum of an infinite number of terms. So we must stop after a finite number of terms, say,  $N$ . This means that we are taking

$$f(x) = \sum_{n=0}^N a_n (x - x_0)^n,$$

and ignoring the rest of the terms, that is,  $\sum_{n=N+1}^{\infty} a_n (x - x_0)^n$ .

There is an error involved in this truncating process which arises from the terms which we exclude. This error is called the 'truncation error'. We denote this error by  $TE$ . Thus we have

$$TE = f(x) - \sum_{n=0}^N a_n (x - x_0)^n = \sum_{n=N+1}^{\infty} a_n (x - x_0)^n$$

You already know how to calculate this error from Sec. 1.3. There we saw that using Taylor's theorem we can estimate the error (or remainder) involved in a truncation process in some cases.

Let's see what happens if we apply Taylor's theorem to the function  $f(x)$  about the point  $x_0 = 0$ . We assume that  $f$  satisfies all conditions of Taylor's theorem. Then we have

$$f(x) = \sum_{n=0}^N a_n x^n + \frac{x^{N+1}}{N+1!} f^{N+1}(c) \quad \dots (19)$$

where  $a_n = \frac{f^{(n)}(0)}{n!}$  and  $0 < c < x$ .

Now, suppose that we want to approximate  $f(x)$  by  $\sum_{n=0}^N a_n x^n$ .

Then Eqn (19) tells us that the truncation error in approximating  $f(x)$  by  $\sum_{n=0}^N a_n x^n$  is given by

$$TE = R_{N+1}(x) = \frac{x^{N+1}}{N+1!} f^{N+1}(c) \quad \dots (20)$$

Theoretically we can use this formula for truncation error for any sufficiently differentiable function. But practically it is not easy to calculate the  $n$ th derivative of many functions. Because of the complexity in differentiation of such functions, it is better to obtain indirectly their Taylor polynomials by using one of the standard expansions we have listed in Sec. 1.3.

For example consider the function  $f(x) = e^{x^2}$ . It is difficult to calculate the  $n$ th derivative of this function. Therefore, for convenience, we obtain Taylor's expansion of  $e^{x^2}$  using Taylor's expansion of  $e^y$  by putting  $y = x^2$ . We shall illustrate this in the following example.

**Example 15 :** Calculate a bound for the truncation error in approximating  $e^{x^2}$  by

$$e^{x^2} \approx 1 + x^2 + \frac{x^4}{2!} + \frac{x^6}{3!} + \frac{x^8}{4!} \text{ for } x \in ]-1, 1[.$$

**Solution :** Put  $u = x^2$ . Then  $e^{x^2} = e^u$ . Now we apply the Taylor's theorem to function  $f(u) = e^u$  about  $u = 0$ . Then, we have

$$e^u = 1 + u + \frac{u^2}{2!} + \frac{u^3}{3!} + \frac{u^4}{4!} + R_5(u) \text{ where}$$

$$R_5(u) = \frac{e^c u^5}{5!}$$

and  $0 < c < u$ . Since  $|x| < 1$ ,  $u = x^2 < 1$  i.e.  $c < 1$ . Therefore,  $e^c < e < 3$ . Thus

$$|R_5(u)| \leq \left| \frac{3x^{10}}{5!} \right| < \frac{3}{5!} = \frac{1}{40} = .025$$

Hence the truncation error in approximating  $e^{x^2}$  by the above expression is less than  $.25 \times 10^{-1}$ .

If the absolute value of the TE is less, then we say that the approximation is good.

Now, in practical situations we should be able to find out the value of  $n$  for which the summation  $\sum a_n x^n$  gives a good approximation to  $f(x)$ . For this we always specify the accuracy (or error bound) required in advance. Then we find  $n$  using formula (20) such that the absolute error  $|R_{n+1}(x)|$  is less than the specified accuracy. This gives the approximation within the prescribed accuracy.

Let us consider an example.

**Example 16 :** Find an approximate value of the integral

$$\int_0^1 e^{x^2} dx$$

with an error less than 0.025.

**Solution :** In Example 15 we observed that

$$e^{x^2} \approx 1 + \frac{x^2}{1!} + \frac{x^4}{2!} + \frac{x^6}{3!} + \frac{x^8}{4!}$$

$$\text{with TE} = \frac{e^{x^2} x^{10}}{5!}$$

Now we use this approximation to calculate the integral. We have

$$\int_0^1 e^{x^2} dx \approx \int_0^1 \left( 1 + x^2 + \frac{x^4}{2!} + \frac{x^6}{3!} + \frac{x^8}{4!} \right) dx \quad \dots (21)$$

with the truncation error

$$\text{TE} = \int_0^1 \frac{e^{x^2} x^{10}}{5!} dx.$$

We have

$$|\text{TE}| \int_0^1 \frac{e^{x^2} |x|^{10}}{5!} \leq \frac{3}{5!} = .25 \times 10^{-1}$$

Integrating the right hand side of (21), we get

$$\begin{aligned} \int_0^1 e^{x^2} dx &\approx \int_0^1 \left( 1 + x^2 + \frac{x^4}{2!} + \frac{x^6}{3!} + \frac{x^8}{4!} \right) dx = \left[ x + \frac{x^3}{3} + \frac{x^5}{5 \times 2!} + \frac{x^7}{7 \times 3!} + \frac{x^9}{9 \times 4!} \right]_0^1 \\ &= \left[ x + \frac{x^3}{3} + \frac{x^5}{10} + \frac{x^7}{42} + \frac{x^9}{216} \right]_0^1 \\ &= 1 + \frac{1}{3} + \frac{1}{10} + \frac{1}{40} + \frac{1}{216} \\ &= 0.0048 \end{aligned}$$

Here is an important remark.

**Remark :** The magnitude of the truncation error could be reduced within any prescribed accuracy by retaining sufficiently large number of terms. Likewise the magnitude of the round-off error could be reduced by retaining additional digits.

You can now try the following exercises.

E14) a) Calculate a bound for the truncation error in approximating  $f(x) = \sin x$  by

$$\sin x \approx 1 - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} \text{ where } -1 \leq x \leq 1.$$

b) Using the approximation in (a), calculate an approximate value of the integral

$$\int_0^1 \frac{\sin x}{x} dx$$

with an error  $10^{-4}$ .

E15) a) Calculate the truncation error in approximating

$$e^{-x^2} \text{ by } 1 - x^2 + \frac{x^4}{2}, \quad -1 \leq x \leq 1.$$

b) Using the approximation in (a) calculate an approximate value of  $\int_0^{0.1} e^{-x^2} dx$

within an error bound of  $10^{-7}$ .

We end this unit by summarising what we have learnt in this unit.

## 1.5 SUMMARY

In this unit we have :

- recalled three important theorems in calculus, namely
  - i) Intermediate value theorem
  - ii) Rolle's theorem
  - iii) Lagrange's mean value theorem
- state Taylor's theorem and demonstrated it with the help of examples.  
The  $n$ th Taylor's expansion :

$$f(x) = f(x_0) + \frac{(x-x_0)}{1!} f'(x_0) + \frac{(x-x_0)^2}{2!} f^{(2)}(x_0) + \dots$$
$$\dots + \frac{(x-x_0)^n}{n!} f^{(n)}(x_0) + \frac{(x-x_0)^{n+1}}{(n+1)!} f^{(n+1)}(c)$$

- defined the term 'error' occurring in numerical computations.
- discussed two types of errors namely
  - i) **Round-off error** : Error occurring in computations where we use rounding off method to represent a number is called round-off error.
  - ii) **Truncation error** : Error occurring in computations where we use truncation process to represent the sum of an infinite number of terms.
- explained how Taylor's theorem is used to calculate the truncation error.

## 1.6 SOLUTIONS/ANSWERS

- E1) a) The given equation is of the form  $f(x) = c$  where  $f(x) = x^3 - x - 5$  and  $c = 0$ .  
 $f$  is a continuous function in the interval  $[0, 2]$  and  $f(0) = -5$  and  $f(2) = 1$ . Then 0 lies between  $f(0)$  and  $f(2)$ . Therefore by IV theorem, the equation  $f(x) = 0$  has a solution in the interval  $[0, 2]$ .

b) Here the equation is of the form  $f(x) = c$  where  $f(x) = \sin x + x$  and  $c = 1$ .

$f$  is a continuous function defined on  $\left[0, \frac{\pi}{6}\right]$  and  $f(0) = 0$  and  $f\left(\frac{\pi}{6}\right) = \frac{1}{2} + \frac{\pi}{6}$

$= 0.5 + 0.523 = 1.023$ . Thus  $f(0) < 1 < f\left(\frac{\pi}{6}\right)$ . Therefore by IV theorem, the result follows.

E2) Let  $f(x) = (x - 1) \sin x$

$$= x \sin x - \sin x$$

Then  $f'(x) = x \cos x + \sin x - \cos x$

$$= (x - 1) \cos x + \sin x$$

Now  $f'(x) = 0$  implies that  $(x - 1) \cos x + \sin x = 0$ . That is  $(x - 1) + \tan x = 0$ .

This shows that there exists a function  $f(x) = (x - 1) \sin x$  such that  $f$  is continuous in  $[0, 1]$  and differentiable on  $]0, 1[$  and  $f'(x) = \tan x - 1 + x$ .

Therefore by Rolle's theorem there exists a point  $x_0$  in  $]0, 1[$  such that  $f'(x_0) = \tan x_0 - 1 + x_0 = 0$ .

E3) Note that  $f$  is a continuous function in  $[0, 3]$  and differentiable in  $]0, 3[$  and

$$f'(x) = x^2 + 2$$

Therefore by Lagrange's mean value theorem there exists a point  $x_0$  in  $]0, 3[$  such that

$$f'(x_0) = \frac{f(3) - f(0)}{3}$$

But  $f'(x_0) = x_0^2 + 2$  and  $f(3) = 15$  and  $f(0) = 0$ . Thus

$$x_0^2 + 2 = \frac{15}{3} = 5$$

$$\text{i.e. } x_0^2 = 3$$

$x_0 = \sqrt{3}$ , since  $x_0$  is a point in  $]0, 3[$ , we consider only the positive value.

E4)  $f(x) = x^3 + 4$  satisfies the requirements of Lagrange's mean value theorem in the interval  $]-2, 1[$ . Therefore there exists a point  $x_0$  in  $]-2, 1[$  such that the slope  $f'(x_0)$  of the tangent line at  $x_0$  is the same as the slope of the line joining  $(-2, f(-2))$  and  $(1, f(1))$ .

$$\text{i.e. } f'(x_0) = \frac{f(-2) - f(1)}{-2 - 1}$$

$$\text{But } f'(x_0) = 3x_0^2$$

Therefore we get

$$3x_0^2 = \frac{5 + 4}{3} = 3$$

$$\text{i.e. } x_0^2 = 1$$

Since  $x_0$  lies in  $]-2, 1[$ , we don't consider the positive value. Therefore there exists only one point  $x_0 = -1$  satisfying the theorem.

E5) Suppose  $f$  is a function defined on  $[a, b]$  which satisfies all the requirements of Lagrange's mean value theorem. Then there exists a point  $x_0$  in  $]a, b[$  such that

$$f'(x_0) = \frac{f(b) - f(a)}{b - a}$$

Suppose in particular  $f$  satisfies the condition that  $f(a) = f(b)$ , i.e.  $f(b) - f(a) = 0$ , then we get  $f'(x_0) = 0$ . This is what the Rolle's theorem states. Hence we deduce that, in the statement of Lagrange's mean value theorem, if we put the extra condition that  $f(a) = f(b)$ , then we get the Rolle's theorem.

E6) Put  $x = x_0$  in Eqn. (3), then we get

$$P_r(x_0) = f(x_0)$$

To calculate,  $P'_r(x_0)$ , we differentiate both sides of Eqn. (3). Then we have

$$P'_r(x) = f'(x_0) + \frac{2f''(x_0)(x-x_0)}{2!} + \frac{3f'''(x_0)(x-x_0)^2}{3!} + \dots (*)$$

Putting  $x = x_0$  on both sides of the above expression, we get

$$P'_r(x_0) = f'(x_0) + 0 + 0 \dots = f'(x_0).$$

Note that apart from the first term, all other terms in the R.H.S. contain the factor  $(x - x_0)$  and therefore when we put  $x = x_0$ , these terms vanish.

By differentiating \* further, we get  $P_r^{(i)}(x_0) = f^{(i)}(x_0)$ ,  $i = 2, \dots, r$ .

E7) The 3rd Taylor polynomial of  $f(x) = e^x$  about  $x = 0$  is  $P_3(x)$

$$= f(0) + x f'(0) + \frac{x^2}{2!} f''(0) + \frac{x^3}{3!} f'''(0).$$

Here  $f(0) = e^0 = 1$

$f'(x) = e^x$ ,  $f'(0) = 1$

Similarly  $f''(0) = 1 = f'''(0)$

Therefore  $P_3(x) = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!}$ .

E8)  $f(x) = \frac{1}{1+x}$ ,  $f(0) = 1$

$$f'(x) = \frac{-1}{(1+x)^2}, f'(0) = -1$$

$$f''(x) = \frac{(-1)(-2)}{(1+x)^3}, f''(0) = 2$$

$$f'''(x) = \frac{(-1)(-2)(-3)}{(1+x)^4}, f'''(0) = (-1)^3 3!$$

$$f^{(n)}(x) = \frac{(-1)(-2)\dots(-n)}{(1+x)^{n+1}}, f^{(n)}(0) = (-1)^n n!$$

The function  $f(x)$  and its derivatives of different orders are continuous in  $\left] -\frac{1}{2}, 1 \right[$ .

Therefore by Taylor's theorem

$$\begin{aligned} f(x) &= 1 - x + \frac{2}{2!} x^2 + \dots + \frac{(-1)^n n!}{n!} x^n + \frac{(-1)^{n+1} (n+1)!}{(n+1)! (1+c)^{n+1}} \\ &= 1 - x + x^2 - x^3 + \dots + (-1)^n x^n + \frac{(-1)^{n+1}}{(1+c)^{n+1}} \end{aligned}$$

where  $c$  is a point lying between 0 and  $x$ .

E9) No. Because the derivatives of  $f(x)$  are not defined at  $x = 0$ .

E10) 8th Taylor expansion of  $f(x) = \cos x$  about  $x_0 = 0$  is

$$\begin{aligned} \cos x &= \cos 0 - \frac{x}{1!} \sin 0 - \frac{x^2}{2!} \cos 0 + \frac{x^3}{3!} \sin 0 + \frac{x^4}{4!} \cos 0 \\ &\quad - \frac{x^5}{5!} \sin 0 - \frac{x^6}{6!} \cos 0 + \frac{x^7}{7!} \sin 0 + \frac{x^8}{8!} \cos 0 - \frac{x^9}{9!} \sin c \\ &= 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \frac{x^8}{8!} - \frac{x^9}{9!} \sin c \end{aligned}$$

The remainder is given by

$$R_9(x) = -\frac{x^9}{9} \sin c.$$

Now since  $x$  lies in  $\left[-\frac{\pi}{4}, \frac{\pi}{4}\right]$  we have  $|x| \leq \frac{\pi}{4} < 1$

Therefore, we get

$$|R_9(x)| \leq \frac{1}{9!} = 0.00000275573.$$

E11) We have seen in E 10 that the remainder in the 8th Taylor expansion of  $\cos x$  is such that

$$|R_9(x)| \leq 10^{-5}$$

$$\text{i.e. } |\cos x - P_9(x)| \leq 10^{-5}$$

$$\text{where } P_9(x) = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \frac{x^8}{8!}$$

This shows that we can approximate  $f(x)$  by 8th Taylor polynomial with error bound  $10^{-5}$  i.e.,

$$\cos x \approx 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \frac{x^8}{8!}.$$

Putting  $x = \frac{\pi}{4}$ , we get

$$\cos \frac{\pi}{4} \approx 1 - \frac{\left(\frac{\pi}{4}\right)^2}{2!} + \frac{\left(\frac{\pi}{4}\right)^4}{4!} - \frac{\left(\frac{\pi}{4}\right)^6}{6!} + \frac{\left(\frac{\pi}{4}\right)^8}{8!}$$

$$E12) |e^x - P_n(x)| = \left| \frac{x^{n+1}}{(n+1)!} e^c \right|$$

where  $c$  lies between 0 and  $x$ . Since  $|x| \leq 1$ , we get that  $|e^c| \leq e$ . Therefore

$$|R_{n+1}| = \left| \frac{x^{n+1} e^c}{(n+1)!} \right| \leq \frac{e}{(n+1)!}$$

Now, we have to find an integer  $n$  such that  $\frac{e}{(n+1)!} \leq 10^{-5}$ .

This is satisfied if  $n = 8$  because  $\frac{e}{9!} \approx 0.749 \times 10^{-5}$

Therefore  $n = 8$  is the required number.

E13)  $\frac{355}{113} = 3.14159292 \dots$  (using a scientific calculator) and  $\pi = 3.14159265 \dots$

$$|\text{error}| = \left| \pi - \frac{355}{113} \right| = 0.00000027$$

$$\text{Then } \frac{1}{2} 10^{-7} < 0.00000027 < \frac{1}{2} \times 10^{-6}$$

Therefore the approximation is accurate to 6 decimal places.

E14) a) We apply Taylor's theorem to the function  $f(x) = \sin x$  in  $]-1, 1[$  about  $x = 0$ .

Then for  $n = 7$ , we have

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + R_8(x)$$

$$\text{where } |R_8(x)| = \left| \frac{x^8}{8!} \sin(c) \right| \leq \frac{1}{8!} = 0.000024802.$$

Therefore, the truncation error T. E. =  $R_8(x)$  is  $.24802 \times 10^{-4}$



$$b) \sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + R_8(x)$$

Hence

$$\frac{\sin x}{x} = 1 - \frac{x^2}{3!} + \frac{x^4}{5!} - \frac{x^6}{7!} + \frac{R_8(x)}{x}$$

$$\text{where } \frac{R_8(x)}{x} = \frac{x^8 \sin c}{8! x} = \frac{x^7}{8!} \sin(c)$$

Thus

$$\int_0^1 \frac{\sin x}{x} dx = \int_0^1 \left( 1 - \frac{x^2}{3!} + \frac{x^4}{5!} - \frac{x^6}{7!} \right) dx + \int_0^1 \frac{R_8(x)}{x} dx$$

$$\begin{aligned} \text{Now } \left| \int_0^1 \frac{R_8(x)}{x} dx \right| &= \left| \int_0^1 \frac{x^7}{8!} \sin(c) dx \right| \\ &\leq \frac{1}{8!} \int_0^1 x^7 dx = \frac{1}{8 \times 8!} \leq 0.24802 \times 10^{-4} \end{aligned}$$

$$\text{Therefore we have } \int_0^1 \frac{\sin x}{x} dx \approx x - \frac{x^3}{3! \cdot 3} + \frac{x^5}{5! \cdot 5} - \frac{x^7}{7! \cdot 7} \Bigg|_0^1$$

$$= 1 - \frac{1}{3! \cdot 3} + \frac{1}{5! \cdot 5} - \frac{1}{7! \cdot 7}$$

$$= 0.946$$

with an error less than  $0.25 \times 10^{-4}$

E15) a) Put  $u = -x^2$ . Then  $e^{-x^2} = e^u$ . We consider the 2nd Taylor expansion of  $e^u$  given by

$$e^u = 1 + u + \frac{u^2}{2} + R_3(u)$$

$$\text{where } R_3(u) = \frac{e^c u^3}{3!}$$

$$|R_3(u)| = \left| \frac{e^c u^3}{3!} \right| \leq \frac{e^c |u|^3}{3!}$$

$$\text{Since } u \leq 0, e^c \leq 1. \text{ Hence } |R_3(u)| \leq \frac{1}{3!} = \frac{1}{6}$$

$$b) \text{ From (a) we have } e^{-x^2} = 1 - x^2 + \frac{x^4}{2} + R_3(-x^2)$$

Hence

$$\int_0^{0.1} e^{-x^2} dx = \int_0^{0.1} \left( 1 - x^2 + \frac{x^4}{2} \right) dx + \int_0^{0.1} R_3(-x^2) dx$$

$$\text{Now, } \left| \int_0^{0.1} R_3(-x^2) dx \right| \leq \int_0^{0.1} \frac{e^{-x^2} x^6}{3!} dx$$

$$\leq \int_0^{0.1} \frac{x^6}{3!} dx = \frac{x^7}{3! \cdot 7} \Bigg|_0^{0.1} = \frac{x^7}{3! \cdot 7} \Bigg|_0^{0.1}$$

$$= \frac{(0.1)^7}{3! \cdot 7} 10^{-7}$$

$$\text{Therefore } \int_0^{0.1} e^{-x^2} dx \approx \int_0^{0.1} \left( 1 - x^2 + \frac{x^4}{2} \right) dx = x - \frac{x^3}{3} + \frac{x^5}{10} \Bigg|_0^{0.1} = .099667666.$$

## UNIT 2 ITERATION METHODS FOR LOCATING A ROOT

### Structure

- 2.1 Introduction
  - Objectives
- 2.2 Initial Approximation to a Root
  - Tabulation Method
  - Graphical Method
- 2.3 Bisection Method
- 2.4 Fixed Point Iteration Method
- 2.5 Summary
- 2.6 Solutions/Answers

### 2.1 INTRODUCTION

We often come across equations of the forms  $x^4 + 3x^2 + 2x + 1 = 0$  or  $e^x = x - 2$  or  $\tanh x = x$  etc. Finding one or more values of  $x$  which satisfy these equations is one of the important problems in Mathematics. From your elementary algebra course (MTE-04), you are already familiar with some methods of solving equations of degrees 1, 2, 3 and 4. Equations of degrees 1, 2, 3 and 4 are called linear, quadratic, cubic and biquadratic respectively. There you might have realised that it is very difficult to use the methods available for solving cubic and biquadratic equations. In fact no formula exists for solving equations of degree  $n \geq 5$ . In these cases we take recourse to approximate methods for the determination of the solution of equations of the form

$$f(x) = 0 \quad \dots (1)$$

The problem of finding approximate values of roots of polynomial equations of higher degree was initiated by Chinese mathematicians. The methods of solution in various forms appeared in the 13th century work *che' in kiu-shoo*. The first noteworthy work in this direction was done in Europe by the English mathematician Fibonacci. Later in the year 1600 Vieta and Isaac Newton made significant contributions to the theory.

In this unit as well as in the next two units we shall discuss some numerical methods which gives an approximate solution of an equation  $f(x) = 0$ . We can classify the methods of solution into two types, namely (i) Direct methods and (ii) Iteration methods. Direct methods produce solutions in a finite number of steps whereas iteration methods give an approximate solution by repeated application of a numerical process. As we said earlier, direct methods you have done in MTE-04. You will find later that for using iteration methods we have to start with an approximate solution. Iteration methods improve this approximate solution. We shall begin this unit by first discussing methods which enable us to determine an initial approximate solution and then discuss iteration methods to refine this approximate solution.

#### Objectives

After studying this unit you should be able to :

- find an initial approximation of the root using (1) tabulation method (2) graphical method.
- use bisection method for finding approximate roots.
- use fixed point iteration method for finding approximate roots.

## 2.2 INITIAL APPROXIMATION TO A ROOT

You know that in many problems of engineering and physical sciences you come across equations in one variable of the form  $f(x) = 0$ .

For example, in Physics, the pressure—volume—temperature relationship of real gases can be described by the equation

$$PV = RT + \frac{\beta}{V} + \frac{r}{V^2} + \frac{s}{V^3} \quad \dots (2)$$

where  $P, V, T$  are pressure, volume and temperature respectively.  $R, \beta, r, s$  are constants. We can rewrite Eqn. (2) as

$$PV^4 - RTV^3 - \beta V^3 - rV - s = 0 \quad \dots (3)$$

Therefore the problem of finding the specific volume of a gas at a given temperature and pressure reduces to solving the biquadratic equation Eqn. (3) for the unknown variable  $V$ .

Consider another example in life sciences, the study of genetic problem of recombination of chromosomes can be described in the form

$$p(1 - p) = p^2 - p + k = 0,$$

where  $p$  stands for the recombination fraction with the limitation  $0 \leq p \leq \frac{1}{2}$  and  $(1 - p)$  stands for the non-recombination fraction. The problem of finding the recombination fraction of a gene reduces to the problem of finding roots of the quadratic equation  $p^2 - p + k = 0$ .

In these problems we are concerned with finding value (or values) of the unknown variable  $x$  that satisfies the equation  $f(x) = 0$ . The function  $f(x)$  may be a polynomial of the form

$$f(x) = a_0 + a_1 x + \dots + a_n x_n$$

or it may be a combination of polynomials, trigonometric, exponential or logarithmic functions. By a root of this equation we mean a number  $x_0$  such that  $f(x_0) = 0$ . The root is also called a zero of  $f(x)$ .

If  $f(x)$  is linear, then Eqn. (1) is of the form  $ax + b = 0$ ,  $a \neq 0$  and it has only one root given by  $x = -\frac{b}{a}$ . Any equation which is not linear is called a **non-linear equation**. In this unit we shall discuss some methods for finding roots of the equation  $f(x) = 0$  where  $f(x)$  is a non linear function. You are already familiar with various methods for calculating roots of quadratic, cubic and biquadratic equations (see MTE-04, Unit 3). But there is no such formula for solving polynomial equations of degree more than 4 or even for a simple equation like

$$x - \cos x = 0$$

Here we shall discuss some of the numerical approximation methods. These methods involve two steps :

**Step 1 :** To find an initial approximation of a root.

**Step 2 :** To improve this approximation to get a more accurate value.

We first consider step 1. Finding an initial approximation to a root means locating (or estimating) a root of an equation approximately. There are two ways for achieving this—tabulation method and graphical method.

Let us start with Tabulation method.

### 2.2.1 Tabulation Method

This method is based on the intermediate value theorem (IV Theorem), (see Theorem 1, Unit 1). Let us try to understand the various steps involved in the method through an example.

**Solutions of Non-linear Equations  
in one Variable**

Suppose we want to find a root of the equation

$$2x - \log_{10} x = 7.$$

We first compute values of  $f(x) = 2x - \log_{10} x - 7$  for different values of  $x$ , say  $x = 1, 2, 3$ , and 4.

When  $x = 1$ , we have  $f(1) = 2 - \log_{10} 1 - 7 = -5$

Similarly, we have

$$f(2) = 4 - \log_{10} 2 - 7 = -3.301$$

(Note that  $\log_{10} 2$  is computed using a scientific calculator.)

$$f(3) = 6 - \log_{10} 3 - 7 = -1.477$$

$$f(4) = 8 - \log_{10} 4 - 7 = -0.3977$$

These values are given in the following table :

**Table 1**

x	1	2	3	4
f(x)	-5	-3.301	-1.477	0.397

We find that  $f(3)$  is negative and  $f(4)$  is positive. Now we apply IV Theorem to the function  $f(x) = 2x - \log_{10} x - 7$  in the interval  $I_1 = [3, 4]$ . Since  $f(3)$  and  $f(4)$  are of opposite signs, by IV theorem there exists a number  $x_0$  lying between 3 and 4 such that  $f(x_0) = 0$ . That is, a root of the function lies in the interval  $]3, 4[$ . Note that this root is positive.

Let us now repeat the above computations for some values of  $x$  lying in  $]3, 4[$  say  $x = 3.5, 3.7$  and 3.8. In the following table we report the values of  $f(x)$ .

**Table 2**

x	3.5	3.7	3.8
f(x)	-0.544	-0.168	0.0202

We find that  $f(3.7)$  and  $f(3.8)$  are of opposite signs. By applying IV theorem again to  $f(x)$  in the interval  $I_2 = [3.7, 3.8]$ , we find that the root of  $f(x)$  lies in the interval  $]3.7, 3.8[$ . Note that this interval is smaller than the previous interval. We call this interval a refinement of the previous interval. Let us repeat the above procedure once again for the interval  $I_2$ . In Table 3 we give the values of  $f(x)$  for some  $x$  between 3.7 and 3.8.

**Table 3**

x	3.75	3.78	3.79
f(x)	-0.074	-0.017	0.00137

Table 3 shows that the root lies within the interval  $]3.78, 3.79[$  and this interval is much smaller compared to the original interval  $]3, 4[$ . The procedure is terminated by taking any value of  $x$  between 3.78 and 3.79 as an approximate value of the root of the equation  $f(x) = 2x - \log_{10} x - 7 = 0$ .

The method illustrated above is known as **Tabulation method**. Let us write the steps involved in the method.

**Step 1 :** Select some numbers  $x_1, x_2, \dots, x_n$  and calculate  $f(x_1), f(x_2), \dots, f(x_n)$ . If  $f(x_i) = 0$  for some  $i$ , then  $x_i$  is a root of the equation. If none of the  $x_i$ s are zero, then proceed to step 2.

**Step 2 :** Find values  $x_i$  and  $x_{i+1}$  such that  $f(x_i)$  and  $f(x_{i+1})$  are of opposite signs i.e.  $f(x_i) f(x_{i+1}) < 0$ . Rename  $x_i = a_j$  and  $x_{i+1} = b_j$ . Then by the IV Theorem a root lies in between  $a_j$  and  $b_j$ . Test for all values of  $f(x_j), j = 1, 2, \dots, n$  and determine other intervals, if any, in which some more roots may lie.

We will talk about the choice of  $x_1, x_2, \dots, x_n$  later.

**Step 3 :** Repeat Step 1 by taking some numbers between  $a_1$  and  $b_1$ . Again, if  $f(x_j) = 0$  for some  $x_j$  between  $a_1$  and  $b_1$ , then we have found the root  $x_j$ . Otherwise, continue step 2.

Continue the steps 1, 2, 3 till we get a sufficiently small interval  $[a, b]$  in which the root lies. Then any value between  $[a, b]$  can be chosen as an initial approximation to the root. You may have noticed that the test values  $x_j, j = 1, 2, \dots, n$  chosen are dependent on the nature of the function  $f(x)$ .

We can always gather some information regarding the root either from the physical problem in which the equation  $f(x) = 0$  occur, or it is specified in the problem. For example, we may ask for the smallest positive root or a root closest to a given number etc.

For a better understanding of the method let us consider one more example.

**Example 1 :** Find the approximate value of the real root of the equation

$$2x - 3 \sin x - 5 = 0.$$

**Solution :** Let  $f(x) = 2x - 3 \sin x - 5$ .

Since  $f(-x) = -2x + 3 \sin x - 5 < 0$  for  $x > 0$ , the function  $f(x)$  is negative for all negative real numbers  $x$ . Therefore the function has no negative real root. Hence the roots of this equation must lie in  $[0, \infty[$ . Now following step 1, we compute values of  $f(x)$ , for  $x = 0, 1, 2, 3, 4, \dots$

We have

$$f(0) = -5.0,$$

$$f(1) = 2 - 3 \sin 1 - 5 = -5.5224$$

using the calculator. Note that  $x$  is in radians. The values  $f(0), f(1), f(2)$  and  $f(3)$  are given in Table 4.

**Table 4**

$x$	0	1	2	3
$f(x)$	-5.0	-5.51224	-3.7278	0.5766

Now we follow step 2. From the table we find that  $f(2)$  and  $f(3)$  are of opposite signs. Therefore a root lies between 2 and 3. Now, to get a more refined interval, we evaluate  $f(x)$  for some values between 2 and 3. The values are given in Table 5.

**Table 5**

$x$	2	2.5	2.8	2.9
$f(x)$	-3.7278	-1.7954	-0.4049	0.0822

This table of values shows that  $f(2.8)$  and  $f(2.9)$  are of opposite signs and hence the root lies between 2.8 and 2.9. We repeat the process once again for the interval  $[2.8, 2.9]$  by taking some values as given in Table 6.

**Table 6**

$x$	2.8	2.85	2.88	2.89
$f(x)$	-0.4049	-1.1624	-0.0159	0.0232

From Table 6 we find that the root lies between 2.88 and 2.89. This interval is small, therefore we take any value between 2.88 and 2.89 as an initial approximation of the root. Since  $f(2.88)$  is near to zero than  $f(2.89)$ , we can take any number near to 2.88 as an initial approximation to the root.

Why don't you try some exercises now.

**E1)** Find an initial approximation to a root of the equation  $3x - \sqrt{1 + \sin x} = 0$  using tabulation method.

**E2)** Find an initial approximation to a positive root of the equation  $2x - \tan x = 0$  using tabulation method.

You might have realized that the tabulation method is a lengthy process for finding an initial approximation of a root. However, since only a rough approximation to the root is required, we normally use only one application of the tabulation method. In the next sub-section we shall discuss the graphical method.

### 2.2.2 Graphical Method

In this method, we draw the approximate graph of  $y = f(x)$ . The points where the curve cuts the x-axis are taken as the required approximate values of the roots of the equation  $f(x) = 0$ . Let us consider an example.

**Example 2 :** Find an approximate value of a root of the biquadratic equation

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

using graphical method.

**Solution :** We first sketch the fourth degree polynomial  $f(x) = x^4 + 4x^3 + 4x^2 - 2$ . This graph is given in Fig. 1.

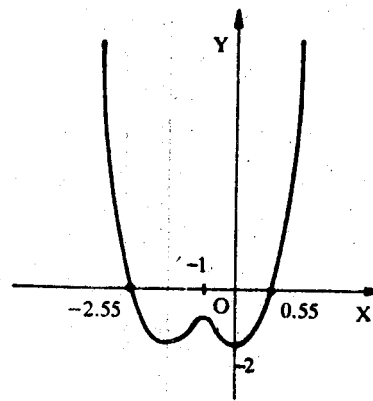


Fig. 1 : Graph of  $f(x) = x^4 + 4x^3 + 4x^2 - 2$ .

The figure shows that the graph cuts the x-axis at two points  $-2.55$ , and  $0.55$ , approximately. Hence  $-2.55$  and  $0.55$  are taken as the approximate roots of the equation  $x^4 + 4x^3 + 4x^2 - 2 = 0$ .

Now go back for a moment to Unit 1 and see Example 1 in Sec. 1.2. There we applied graphical method to find the roots of the equation  $\sin x = \frac{1}{2}$ .

Let us consider another example.

**Example 3 :** Find the approximate value of a root of

$$x^2 - e^x = 0$$

using graphical method.

**Solution :** First thing to do is to draw the graph of the function  $f(x) = x^2 - e^x$ . It is not easy to graph this function. Now if we split the function as

$$f(x) = f_1(x) - f_2(x)$$

where  $f_1(x) = x^2$  and  $f_2(x) = e^x$ , then we can easily draw the graphs of the functions  $f_1(x)$  and  $f_2(x)$ . The graphs are given in Fig. 2.

The figure shows that the two curves  $y = x^2$  and  $y = e^x$  intersect at some point P. From the figure, we find that the approximate point of intersection of the two curves is  $-0.7$ . Thus we

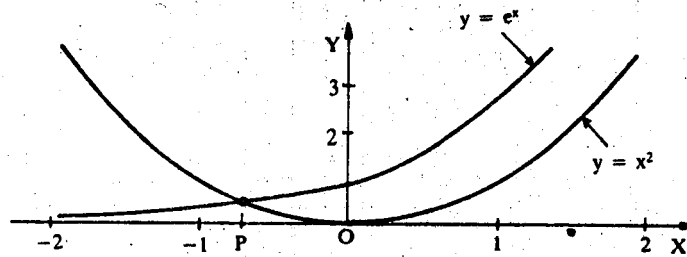


Fig. 2 : Graphs of  $f_1(x) = x^2$  and  $f_2(x) = e^x$ .

have  $f_1(-0.7) \approx f_2(-0.7)$ , and therefore  $f(-0.7) = f_1(-0.7) - f_2(-0.7) \approx 0$ . Hence  $-0.7$  is an approximate value of the root of the equation  $f(x) = 0$ .

From the above example we observe the following : Suppose we want to apply the graphic method for finding an approximate root of  $f(x) = 0$ . Then we may try to simplify the method by splitting the equation as

$$f(x) = f_1(x) - f_2(x) = 0 \quad \dots (4)$$

where the graphs of  $f_1(x)$  and  $f_2(x)$  are easy to draw. From Eqn.(4), we have  $f_1(x) = f_2(x)$ .

The x-coordinate of the point at which the two curves  $y_1 = f_1(x)$  and  $y_2 = f_2(x)$  intersect gives an approximate value of the root of the equation  $f(x) = 0$ . Note that we are interested only in the x-coordinate, we don't have to worry about the point of intersection of the curves.

Often we can split the function  $f(x)$  in the form (4) in a number of ways. But we should choose that form which involves minimum calculations and the graphs of  $f_1(x)$  and  $f_2(x)$  are easy to draw. We illustrate this point in the following example.

**Example 4 :** Find an approximate value of the positive real root of  $3x - \cos x - 1 = 0$  using graphic method.

**Solution :** Since it is easy to plot  $3x - 1$  and  $\cos x$ , we rewrite the equation as  $3x - 1 = \cos x$ . The graphs of  $y = f_1(x) = 3x - 1$  and  $y = f_2(x) = \cos x$  are given in Figure 3.

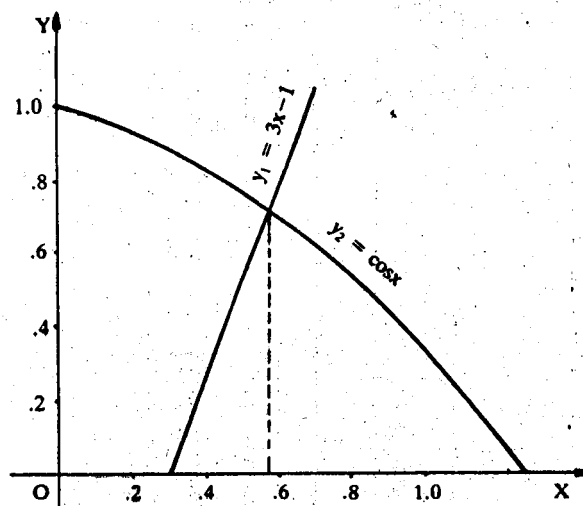


Fig. 3 : Graphs of  $f_1(x) = 3x - 1$  and  $f_2(x) = \cos x$ .

It is clear from the figure that the x-coordinate of the point of intersection is approximately 0.6. Hence  $x = 0.6$  is an approximate value of the root of the equation  $3x - \cos x - 1 = 0$ .

We now make a remark.

**Remark 1 :** You should take some care while choosing the scale for graphing. A magnification of the scale may improve the accuracy of the approximate value.

Here is an exercise for you.

E3) Find the approximate location of the roots of the following equations in the regions given using graphic method.

a)  $f(x) = e^{-x} - x = 0$ , in  $0 \leq x \leq 1$

b)  $f(x) = e^{0.4x} - 0.4x - 9 = 0$ , in  $0 \leq x \leq 7$

We have discussed two methods, namely, tabulation method and graphical method which help us in finding an initial approximation to a root. But these two methods give only a rough approximation to a root. Now to obtain more accurate results, we need to improve these crude approximations. In the tabulation method we found that one way of improving the process is refining the intervals within which a root lies. A modification of this method is known as bisection method. In the next section we discuss this method.

### 2.3 BISECTION METHOD

In the beginning of the previous section we have mentioned that there are two steps involved in finding an approximate solution. The first step has already been discussed. In this section we consider the second step which deals with refining an initial approximation to a root.

Once we know an interval in which a root lies, there are several procedures to refine it. The bisection method is one of the basic methods among them. We repeat the steps 1, 2, 3 of the tabulation method given in subsection 2.2.1 in a modified form. For convenience we write the method as an algorithm.

An algorithm is a complete and unambiguous set of instructions leading to the solution of a problem.

This method is also called as Bolzano method, Bracketing Method.

Suppose that we are given a continuous function  $f(x)$  defined on  $[a, b]$  and we want to find the roots of the equation  $f(x) = 0$  by bisection method. We describe the procedure in the following steps :

**Step 1 :** Find points  $x_1, x_2$  in the interval  $[a, b]$  such that  $f(x_1) \cdot f(x_2) < 0$ . That is, those points  $x_1$  and  $x_2$  for which  $f(x_1)$  and  $f(x_2)$  are of opposite signs—(see Step 1 of subsection 2.2.1). This process is called “finding an initial bisecting interval”. Then by IV theorem a root lies in the interval  $] x_1, x_2 [$ .

**Step 2 :** Find the middle point  $c$  of the interval  $] x_1, x_2 [$  i.e.  $c = \frac{x_1 + x_2}{2}$ . If  $f(c) = 0$ , then  $c$  is the required root of the equation and we can stop the procedure. Otherwise we go to Step 3.

**Step 3 :** Find out if

$$f(x_1) f(c) < 0$$

If it holds, then the root lies in  $] x_1, c [$ . Otherwise the root lies in  $] c, x_2 [$  (see Fig. 4). Thus in either case we have found an interval half as wide as the original interval that contains the root.

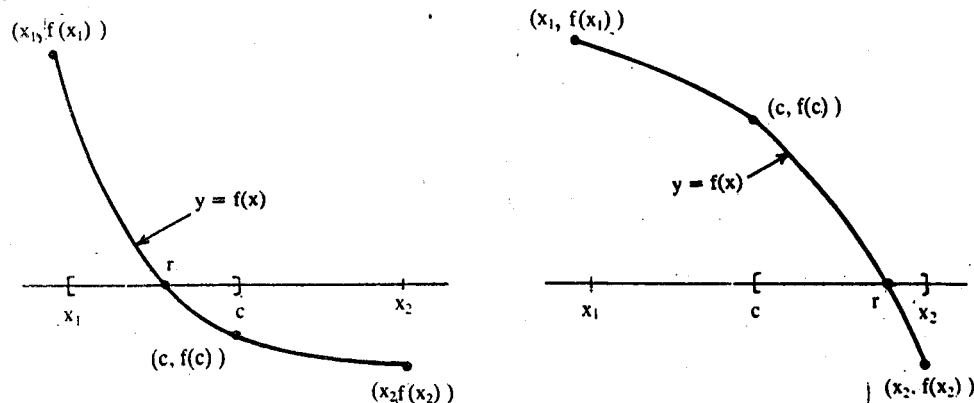


Fig. 4 : The decision process for the bisection method.



**Step 4 :** Repeat Steps 2 and 3 with the new interval. This process either gives you the root or an interval having width  $1/4$  of the original interval  $]x_1, x_2[$  which contains the required root.

**Step 5 :** Repeat this procedure until the interval width is as small as we desire. Each bisection halves the length of the preceding interval. After  $N$  steps, the original interval length will be reduced by a factor  $1/2^N$ .

Now we shall see how this method helps in refining the initial intervals in some of the problems we have done in subsection 2.2.1.

**Example 5 :** Consider the equation  $2x - \log_{10} x = 7$  lies in  $]3.78, 3.79[$ . Apply bisection method to find an approximate root of the equation correct to three decimal places.

**Solution :** Let  $f(x) = 2x - \log_{10} x - 7$ . From Table 2 in subsection 2.2.1, we find that  $f(3.78) = -0.01749$  and  $f(3.79) = 0.00136$ . Thus a root lies in the interval  $]3.78, 3.79[$ .

Then we find the middle point of the interval  $]3.78, 3.79[$ . The middle point is  $c = (3.78 + 3.79)/2 = 3.785$  and  $f(c) = f(3.785) = -0.00806 \neq 0$ . Now, we check the condition in Step 3. Since  $f(3.78) f(3.785) > 0$ , the root does not lie in the interval  $]3.78, 3.785[$ . Hence the root lies in the interval  $]3.785, 3.79[$ . We have to refine this interval further to get better approximation. Further bisections are shown in the following Table.

**Table 7**

Number of Bisections	Bisected value $x_i$	$f(x_i)$	Improved Interval
1	3.785	-0.00806	]3.785, 3.79[
2	3.7875	$-3.3525 \times 10^{-3}$	]3.7875, 3.79[
3	3.78875	$-9.9594 \times 10^{-4}$	]3.78875, 3.79[
4	3.789375	$1.824 \times 10^{-4}$	]3.78875, 3.789375[
5	3.7890625	$-4.068 \times 10^{-4}$	]3.78906, 3.789375[

The table shows that the improved interval after 5 bisections is  $]3.78906, 3.789375[$ . The width of this interval is  $3.789375 - 3.78906 = 0.000315$ . If we stop further bisections, the maximum absolute error would be 0.000315. The approximate root can therefore be taken as  $(3.78906 + 3.789375)/2 = 3.789218$ . Hence the desired approximate value of the root rounded off to three decimal places is 3.789.

**Example 6 :** Apply bisection method to find an approximation to the positive root of the equation

$$2x - 3 \sin x - 5 = 0$$

rounded off to three decimal places.

**Solution :** Let  $f(x) = 2x - 3 \sin x - 5$ .

In Example 1, we had shown that a positive root lies in the interval  $]2.8, 2.9[$ . Now we apply bisection method to this interval. The results are given in the following table.

Table 8

Number of bisection	Bisected value $x_i$	$f(x_i)$	Improved interval
1	2.85	-0.1624	]2.85, 2.9[
2	2.875	-0.0403	]2.875, 2.9[
3	2.8875	0.02089	]2.875, 2.8875[
4	2.88125	$-9.735 \times 10^{-3}$	]2.88125, 2.8875[
5	2.884375	$5.57781 \times 10^{-3}$	]2.88125, 2.884375[
6	2.8828125	$-2.0795 \times 10^{-3}$	]2.8828125, 2.884375[
7	2.8835938	$1.7489 \times 10^{-3}$	]2.8828125, 2.8835938[
8	2.8832031	$-1.6539 \times 10^{-4}$	]2.8832031, 2.8835938[

After 8 bisections the width of the interval is  $2.8835938 - 2.8832031 = 0.0003907$ . Hence, the maximum possible absolute error to the root is 0.0003907. Therefore the required approximation to the root is 2.883.

Now let us make some remarks.

**Remark 2 :** While applying bisection method we must be careful to check that  $f(x)$  is continuous. For example, we may come across functions like  $f(x) = \frac{1}{x-1}$ . If we consider the interval  $]1.5, 1.5[$ , then  $f(1.5) f(1.5) < 0$ . In this case we may be tempted to use bisection method. But we cannot use the method here because  $f(x)$  is not defined at the middle point  $x = 1$ . We can overcome these difficulties by taking  $f(x)$  to be continuous throughout the initial bisecting interval. (Note that if  $f(x)$  is continuous, by IV theorem  $f(x)$  assumes all values between the interval.)

Therefore you should always examine the continuity of the function in the initial interval before attempting the bisection method.

**Remark 3 :** It may happen that a function has more than one root in an interval. The bisection method helps us in determining one root only. We can determine the other roots by properly choosing the initial intervals.

You can try some exercises now.

- E4) Starting with the interval  $[a_0, b_0]$ , apply bisection method to the following equations and find an interval of width 0.05 that contains a solution of the equations
- $e^x - 2 - x = 0, [a_0, b_0] = [1.0, 1.8]$
  - $\ln x - 5 + x = 0, [a_0, b_0] = [3.2, 4.0]$
- E5) Using bisection method find an approximate root of the equation  $x^3 - x - 4 = 0$  in the interval  $]1, 2[$  to two places of decimal.

While applying bisection method we repeatedly apply steps 2, 3, 4 and 5. You recall that in the introduction we classified such a method as an **Iteration method**. As we mentioned in the beginning of Sec. 2.2, a numerical process starts with an initial approximation and iteration improves this approximation until we get the desired accurate value of the root.

Let us consider another iteration method now.

## 2.4 FIXED POINT ITERATION METHOD

The bisection method we have described earlier depends on our ability to find an interval in which the root lies. The task of finding such intervals is difficult in certain situations. In such cases we try an alternate method called **Fixed Point Iteration Method**. We shall discuss the advantage of this method later.

Iteration means repeated application of a numerical process or a pattern of action.

The first step in this method is to rewrite the equation  $f(x) = 0$  as

$$x = g(x) \quad \dots (5)$$

For example consider the equation  $x^2 - 2x - 8 = 0$ . We can write it as

$$x = \sqrt{2x + 8} \quad \dots (6)$$

$$x = \frac{2x + 8}{x} \quad \dots (7)$$

$$x = \frac{x^2 - 8}{2} \quad \dots (8)$$

We can choose the form (5) in several ways. Since  $f(x) = 0$  is the same as  $x = g(x)$ , finding a root of  $f(x) = 0$  is the same as finding a root of  $x = g(x)$  i.e. a fixed point of  $g(x)$ . Each such  $g(x)$  given in (6), (7) or (8) is called an **iteration function** for solving  $f(x) = 0$ .

A fixed point of a function  $g$  is a point  $\alpha$  such that  $g(\alpha) = \alpha$ .

Once an iteration function is chosen, our next step is to take a point  $x_0$ , close to the root, as the initial approximation of the root.

Starting with  $x_0$ , we find the first approximation  $x_1$  as

$$x_1 = g(x_0)$$

Then we find the next approximation as

$$x_2 = g(x_1)$$

Similarly we find the successive approximations  $x_2, x_3, x_4, \dots$  as

$$x_3 = g(x_2)$$

$$x_4 = g(x_3)$$

$$\vdots$$

$$x_{n+1} = g(x_n)$$

Each computation of the type  $x_{n+1} = g(x_n)$  is called an iteration. Now, two questions arise (i) when do we stop these iterations? (ii) Does this procedure always give the required solution?

To ensure this we make the following assumptions on  $g(x)$ :

#### Assumption\*

The derivative  $g'(x)$  of  $g(x)$  exists,  $g'(x)$  is continuous and satisfies  $|g'(x)| < 1$  in an interval containing  $x_0$ . (That would mean that we require  $|g'(x_i)| < 1$  at all iterates  $x_i$ .)

The iteration is usually stopped whenever  $|x_{i+1} - x_i|$  is less than the accuracy required.

In Unit 3 you will prove that if  $g(x)$  satisfies the above conditions, then there exists a unique point  $\alpha$  such that  $g(\alpha) = \alpha$  and the sequence of iterates approach  $\alpha$ , provided that the initial approximation is close to the point  $\alpha$ .

Now we shall illustrate this method with the following example.

**Example 7 :** Find an approximate root of the equation

$$x^2 - 2x - 8 = 0$$

using fixed point iteration method, starting with  $x_0 = 5$ . Stop the iteration whenever  $|x_{i+1} - x_i| < 0.001$ .

**Solution :** Let  $f(x) = x^2 - 2x - 8$ . We saw that the equation  $f(x) = 0$  can be written in three forms (6), (7) and (8). We shall take up the three forms one by one.

**Case 1 :** Suppose we consider form (5). In this form the equation is written as

$$x = (2x + 8)^{1/2}$$

Here  $g(x) = (2x + 8)^{1/2}$ . Let's see whether Assumption (\*) is satisfied for this  $g(x)$ . We have

$$g'(x) = \frac{1}{(2x + 8)^{1/2}}$$

Then  $|g'(x)| < 1$  whenever  $(2x + 8)^{1/2} > 1$ . For any positive real number  $x$ , we see that the inequality  $(2x + 8)^{1/2} > 1$  is satisfied. Therefore, we consider any interval on the positive side of  $x$ -axis. Since the starting point is  $x_0 = 5$ , we may consider the interval  $I = [3, 6]$ . This contains the point 5. Now,  $g(x)$  satisfies the condition that  $g'(x)$  exists on  $I$ ,  $g'(x)$  is continuous on  $I$  and  $|g'(x)| < 1$  for every  $x$  in the interval  $[3, 6]$ . Now we apply fixed point iteration method to  $g(x)$ .

We get

$$x_1 = g(5) = \sqrt{18} = 4.243$$

$$x_2 = g(4.243) = 4.060$$

$$x_3 = 4.015$$

$$x_4 = 4.004$$

$$x_5 = 4.001$$

$$x_6 = 4.000.$$

Since  $|x_6 - x_5| = |-0.001| = 0.001$ , we conclude that an approximate value of a root of  $f(x) = 0$  is 4.

Case 2 : Let us consider the second form,

$$x = \frac{2x + 8}{x}$$

Here  $g(x) = \frac{2x + 8}{x}$  and  $g'(x) = \frac{-8}{x^2}$ . The  $|g'(x)| < 1$  for any real number  $x \geq 3$ . Hence  $g(x)$  satisfies Assumption (\*) in the interval  $[3, 6]$ . Now we leave it as an exercise for you to complete the computations (See E6).

Case 3 : Here we have  $x = \frac{x^2 - 8}{2}$ . Then  $g(x) = \frac{x^2 - 8}{2}$  and  $g'(x) = x$ . In this case  $|g'(x)| < 1$  only if  $|x| < 1$  i.e. if  $x$  lies in the interval  $]-1, 1[$ . But this interval does not contain 5. Therefore  $g(x)$  does not satisfy the Assumption (\*) in any interval containing the initial approximation. Hence, the iteration method cannot provide approximation to the desired root.

Note : This example may appear artificial to you. You are right because in this case we have got a formula for calculating the root. This example is taken to illustrate the method in a simple way.

Let us consider another example.

**Example 8 :** Use fixed point iteration procedure to find an approximate root of  $2x - 3 \sin x - 5 = 0$  starting with the point  $x_0 = 2.8$ . Stop the iteration-whenver  $|x_{i+1} - x_i| < 10^{-5}$ .

**Solution :** We can rewrite the equation in the form,

$$x = \frac{3}{2} \sin x + \frac{5}{2}.$$

Here  $g(x) = \frac{3}{2} \sin x + \frac{5}{2}$  and  $g'(x) = \frac{3}{2} \cos x$ .

Now at  $x_0 = 2.8$ , we have

$$|g'(2.8)| = 1.413$$

which is greater than 1. Thus  $g(x)$  does not satisfy Assumption (\*) and therefore in this form the iteration method fails.

Use a calculator to evaluate the square root.

Let us now rewrite the equation in another form. We write

$$x = x - \frac{2x - 3 \sin x - 5}{2 - 3 \cos x}$$

$$\text{Then } g(x) = x - \frac{2x - 3 \sin x - 5}{2 - 3 \cos x}$$

You may wonder how did we get this form. Note that here  $g(x)$  is of the form  $g(x) = x - \frac{f(x)}{f'(x)}$ . You will find later that the above equation is the iterated formula for another popular iteration method.

$$\begin{aligned} \text{Then } g'(x) &= 1 - \left[ \frac{(2 - 3 \cos x)(2 - 3 \cos x) - (2x - 3 \sin x + 5) 3 \sin x}{(2 - 3 \cos x)^2} \right] \\ &= \frac{2x - 3 \sin x + 5}{(2 - 3 \cos x)^2} 3 \sin x \end{aligned}$$

$$\text{At } x_0 = 2.8, |g'(x_0)| = 0.0669315 \text{ (or } 0.02174691) < 1$$

Therefore  $g(x)$  satisfies the Assumption (\*). Using the initial approximation as  $x_0 = 2.8$ , we get the successive approximation as

$$x_1 = 2.8839015$$

$$x_2 = 2.8832369$$

$$x_3 = 2.8832369$$

Since  $|x_2 - x_3| < 10^{-5}$  we stop the iteration here and conclude that 2.88323 is an approximate value of the root.

Next we shall use another form

$$x = \sin^{-1} \left( \frac{2x - 5}{3} \right)$$

$$\text{Here } g(x) = \sin^{-1} \left( \frac{2x - 5}{3} \right) \text{ and } g'(x) = \frac{2}{\sqrt{9 - (2x - 5)^2}}$$

At  $x_0 = 2.8$ ,  $g'(x_0) = 0.6804 < 1$ . In fact, we can check that in any small interval containing 2.8,  $|g'(x)| < 1$ . Thus  $g(x)$  satisfies the Assumption (\*). Applying the iteration method, we have

$$x_1 = \sin^{-1} \left( \frac{2(2.8) - 5}{3} \right) = 0.201358$$

We find that there are two values which satisfy the above equation. One value is 0.201358 and the other is  $\pi - 0.201358 = 2.940235$ . In such situations, we take a value close to the initial approximation. In this case the value close to the initial approximation is 2.940235. Therefore we take this value as the starting point of the next approximation.

$$x_1 = 2.940235$$

Next we calculate

$$\begin{aligned} x_2 &= \sin^{-1} \left( \frac{2(2.940235) - 5}{3} \right) \\ &= 0.297876 \text{ or } 2.843717 \end{aligned}$$

Continuing like this, it needed 17 iterations to obtain the value  $x_{17} = 2.88323$ , which we got from the previous form. This means that in this form the convergence is very slow.

From examples 7 and 8, we learn that if we choose the form  $x = g(x)$  properly, then we can get the approximate root provided that the initial approximation is sufficiently close to the root. The initial approximation is usually given in the problem or we can find using the IV theorem.

Now we shall make a remark here.

**Remark :** The Assumption (\*) we have given for an iteration function, is a stronger assumption. In actual practice there are a variety of assumptions which the iteration function  $g(x)$  must satisfy to ensure that the iterations approach the root. But, to use those assumptions you would require a lot of practice in the application of techniques in mathematical analysis. In this course, we will be restricting ourselves to functions that satisfies Assumption (\*). If you would like to know about the other assumptions, you may refer to 'Elementary Numerical Analysis' by Samuel D Conte and Carl de Boor.

To get some practice over this method, you can try the following exercises.

- E5) Apply fixed point iteration method to the form  $x = \frac{2x+8}{x}$  starting with  $x_0 = 5$ , to obtain a root of  $x^2 - 2x - 8 = 0$ .
- E7) a) Apply fixed point iteration method to the following equations with the initial approximation given alongside. In each case find an approximate root rounded off to 4 decimal places.
- i)  $x = -45 + \frac{2}{x}$ ,  $x_0 = -20$ .
  - ii)  $x = \frac{1}{2} + \sin x$ ,  $x_0 = 1$ .
- b) Compute the exact roots of the equation  $x^2 + 45x - 2 = 0$  using quadratic formula and compare with the approximate root obtained in (a) (i).

Let us now briefly recall what we have done in this unit.

## 2.5 SUMMARY

In this unit we have covered the following points :

- We have seen that the methods for finding an approximate solution of an equation involve two steps :
  - i) Find an initial approximation to a root.
  - ii) Improve the initial approximation to get a more accurate value of the root.
- We have described the following iteration methods for improving an initial approximation of a root.
  - i) Bisection method
  - ii) Fixed point iteration method.

## 2.6 SOLUTIONS/ANSWERS

E1) Let  $f(x) = 3x - \sqrt{1 + \sin x}$

Since  $f(-x) = -3x - \sqrt{1 + \sin(-x)} = -3x - \sqrt{1 + \sin x} < 0$  for  $x > 0$ ,  $f(x)$  has no negative real root.

Computing values of  $f(x)$  for  $x = 0, 1, 2, \dots$  radians, we get

$$f(0) = 3 \times 0 - \sqrt{1} = -1$$

and  $f(1) = 3 - \sqrt{1 + \sin 1} = 3 - \sqrt{1 + .84147} = 1.6430$ , as  $\sin 1 = 0.84147$ , approximately, using a calculator. Thus  $f(0)$  and  $f(1)$  are of opposite signs. Therefore there exists a root of  $f(x) = 0$  lying between 0 and 1.

Now we randomly take some values between 0 and 1, say 0.3 and 0.5. Then

$$f(0.3) = .9 - 1.1381 = -0.23181 < 0$$

and

$$f(0.5) = 0.283683619 > 0.$$

Hence the root lies in  $]0.3, 0.5[$ .

Repeating the process once again with the values  $x = 0.35$  and  $0.41$  etc. we get,

$$f(0.35) < 0$$

and

$$f(0.41) > 0.$$

Therefore the root lies between  $0.35$  and  $0.41$ . This interval is small. If we stop the iteration here, we may either take  $0.41$ , since  $f(0.41)$  is closer to zero, or  $(0.35 + 0.41)/2 = 0.38$  as the required initial approximation.

- E2) Let  $f(x) = 2x - \tan x$ . Since we want a positive root of  $f(x) = 0$ , we evaluate  $f(x)$  for  $x > 0$ .

Let us consider  $x = 0, 1, 1.5$ . Then

$$f(0) = 0$$

$$f(1) = 0.443$$

and

$$f(1.5) = -11.1014$$

Therefore a root lies between  $1$  and  $1.5$ . Now if we consider values of  $f(x)$  for  $x = 1.1$  and  $1.2$ , we get

$$f(1) = 0.443 < 0$$

$$f(1.1) = -0.8648 < 0$$

and

$$f(1.2) = -0.1722 < 0$$

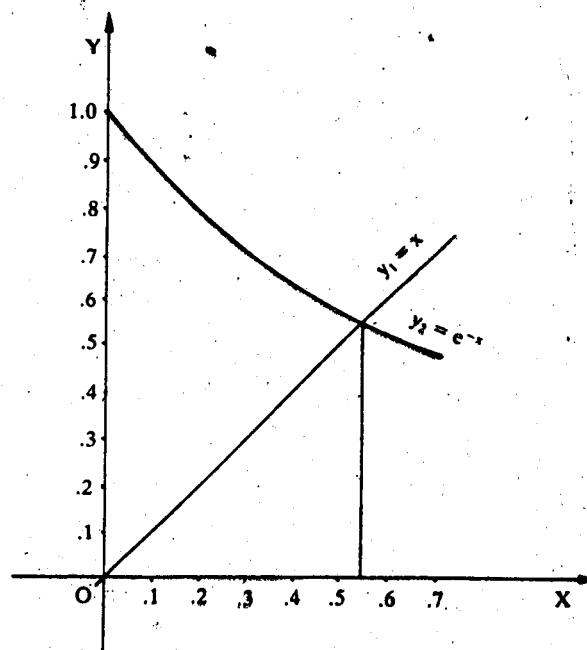
Therefore we get that a root lies in the interval  $]1, 1.1[$ . In fact the root lies more close to  $1$ . We may take  $(1 + 1.1)/2 = 1.05$  as an initial approximation.

- E3) a) Let  $f_1(x) = e^{-x}$

and

$$f_2(x) = x$$

The graphs of  $f_1$  and  $f_2$  are plotted in the following figure :

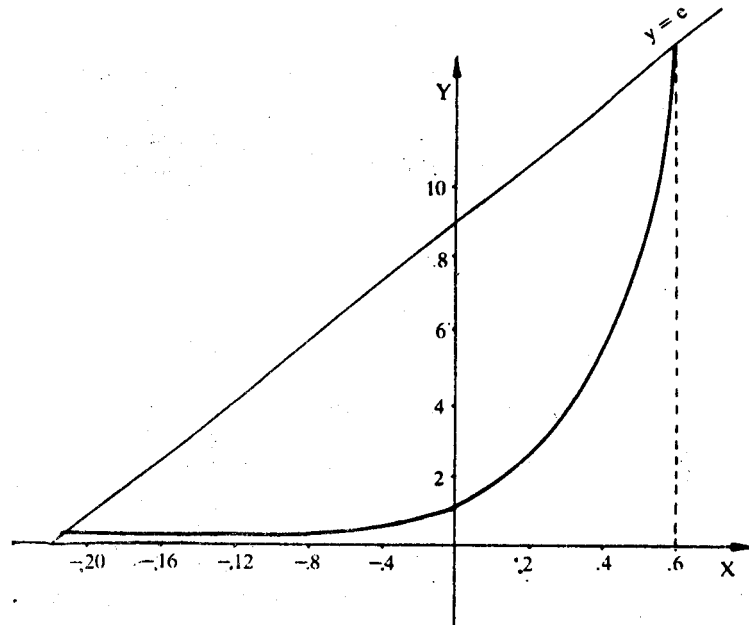


From the graph you can see that the x-coordinate of the point of intersection is approximately 0.55. Hence the root lies close to 0.55.

- b) The given equation can be written as

$$f(x) = e^{0.4x} - (0.4x + 9)$$

The graphs of  $e^{0.4x}$  and  $0.4x + 9$  are given in the following figure :



From the graph you can see that there are two points of intersections. x-coordinates of the points are approximately 6 and -22.5. Hence one root lies close to 6 and the other root lies close to -22.5.

- E4) a) We first note that the given function  $f(x) = e^x - 2 - x$  is continuous in the interval  $[1.0, 1.8]$ . Also

$$f(1) = e - 2 - 1 = e - 3 < 0$$

$$f(1.8) = e^{1.8} - 3 = 6.049647464 - 3 > 0$$

using a calculator.

Therefore the interval  $]1.0, 1.8[$  contains a root of the equation.

Middle point of the interval  $c = \frac{1 + 1.8}{2} = 1.4$ . Also,  $f(c) = e^{1.4} - 3 = 4.0552 - 3 > 0$ .

Therefore the root lies in the interval  $]1, 1.4[$ .

Repeating this process three times more, we get the intervals  $]1.0, 1.2[$ ,  $]1.1, 1.2[$  and  $]1.1, 1.15[$ . Therefore the improved interval after 4 bisections is  $]1.1, 1.15[$ . The width of this interval is 0.05. This shows that the required interval of width 0.05 which contains a root of the equation is  $]1.1, 1.15[$ .

- b) Using a calculator you can show that the intervals in each of the four bisections are given by  $]3.6, 4.0[$ ,  $]3.6, 3.8[$ ,  $]3.6, 3.7[$  and  $]3.65, 3.70[$ . The width of the last interval is  $3.70 - 3.65 = 0.05$ . Therefore the required interval is  $]3.65, 3.70[$ .

- E5) After 5 bisections the root lies in  $]1.7959, 1.7969[$ . Therefore the required root correct to two decimal places is 1.80.

E6)  $x = g(x) = \frac{2x + 8}{x}$ .



The iterations are given by

$$x_{i+1} = \frac{2x_i + 8}{x_i}$$

we have

$$x_1 = g(5) = 3.6$$

$$x_2 = g(3.6) = 4.2222$$

$$x_3 = 3.8947$$

$$x_4 = 4.0540$$

$$x_5 = 3.9733$$

$$x_6 = 4.0134$$

$$x_7 = 3.9933$$

$$x_8 = 4.0033$$

$$x_9 = 3.9983$$

$$x_{10} = 4.0008$$

$$x_{11} = 3.9996.$$

Since  $|x_{11} - x_{10}| = 0.001$ , we conclude that an approximate value of a root of  $f(x) = 0$  is 4.

E7) a) i) The iteration formula in fixed point iteration method is

$$x_{i+1} = g(x_i), i = 0, 1, 2, \dots$$

$$\text{i.e. } x_{i+1} = -45 + \frac{2}{x_i}, i = 0, 1, 2, \dots$$

Here  $x_0 = -20$ . Starting with  $x_0 = -20$ , the successive iterations are given by

$$x_1 = -45.1$$

$$x_2 = -45.04435$$

$$x_3 = -45.044401$$

$$x_4 = -45.044401$$

Since  $x_3$  and  $x_4$  are the same, we stop the iteration here. Hence the approximate root rounded off to four decimal places is  $-45.0444$ .

ii) The desired root is 1.4973.

b) The given equation is  $x^2 + 45x - 2 = 0$ . According to the quadratic formula, the two roots are

$$x_1 = \frac{-45 + \sqrt{(45)^2 + 8}}{2}, x_2 = \frac{-45 - \sqrt{(45)^2 + 8}}{2}$$

$$= 0.0444, \quad = -45.0444$$

Comparing with the result in part (a) (i), we find that the approximate root is the same as the exact root  $-45.0444$ .

---

## UNIT 3 CHORD METHODS FOR FINDING ROOTS

---

### Structure

- 3.1 Introduction
  - Objectives
- 3.2 Regula-Falsi Method
- 3.3 Newton—Raphson Method
- 3.4 Convergence Criterion
- 3.5 Summary
- 3.6 Solutions/Answers

---

### 3.1 INTRODUCTION

---

In the last unit we introduced you to two iteration methods for finding roots of an equation  $f(x) = 0$ . There we have shown that a root of the equation  $f(x) = 0$  can be obtained by writing the equation in the form  $x = g(x)$ . Using this form we generate a sequence of approximations  $x_{i+1} = g(x_i)$  for  $i = 0, 1, 2, \dots$ . We had also mentioned there that the success of the iteration methods depends upon the form of  $g(x)$  and the initial approximation  $x_0$ . In this unit, we shall discuss two iteration methods : regula-falsi and Newton—Raphson methods. These methods produce results faster than bisection method. The first two sections of this unit deal with derivations and the use of these two methods. You will be able to appreciate these iteration methods better if you can compare the efficiency of these methods. With this in view we introduce the concept of convergence criterion which helps us to check the efficiency of each method. Sec 3.4 is devoted to the study of rate of convergence of different iterative methods.

#### Objectives

After studying the unit you should be able to :

- apply regula-falsi and secant methods for finding roots
- apply Newton-Raphson method for finding roots
- define 'order of convergence' of an iterative scheme
- obtain the order of convergence of the following four methods :
  - i) bisection method
  - ii) fixed point iteration method
  - iii) secant method
  - iv) Newton—Raphson method

---

### 3.2 REGULA-FALSI METHOD (OR METHOD OF FALSE POSITION)

---

In this section we shall discuss the 'regula-falsi method'. The Latin word 'Regula Falsi' means rule of falsehood. It does not mean that the rule is a false statement. But it conveys that the roots that we get according to the rule are approximate roots and not necessarily exact roots. The method is also known as the method of false position. This method is similar to the bisection method you have learnt in Unit 3.

The bisection method for finding approximate roots has a drawback that it makes use of only the signs of  $f(a)$  and  $f(b)$ . It does not use the values  $f(a)$ ,  $f(b)$  in the computations. For example, if  $f(a) = 700$  and  $f(b) = -0.1$ , then by the bisection method the first approximate value of a root of  $f(x)$  is the mid value  $x_0$  of the interval  $[a, b]$ . But at  $x_0$ ,  $f(x_0)$  is nowhere

near 0. Therefore in this case it makes more sense to take a value near to  $-0.1$  than the middle value as the approximation to the root. This drawback is to some extent overcome by the regula-falsi method. We shall first describe the method geometrically.

Suppose we want to find a root of the equation  $f(x) = 0$  where  $f(x)$  is a continuous function. As in the bisection method, we first find an interval  $]a, b[$  such that  $f(a) f(b) < 0$ . Let us look at the graph of  $f(x)$  given in Fig. 1.

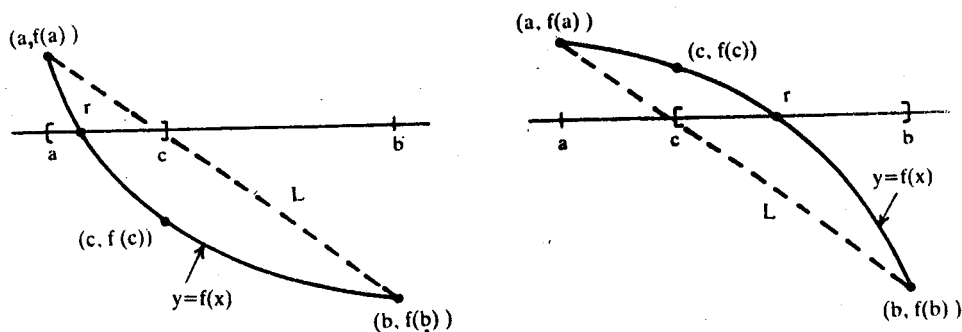


Fig. 1 : Regula-Falsi Method

The condition  $f(a) f(b) < 0$  means that the points  $(a, f(a))$  and  $(b, f(b))$  lie on the opposite sides of the x-axis. Let us consider the line joining  $(a, f(a))$  and  $(b, f(b))$ . This line crosses the x-axis at some point  $(c, 0)$  [see Fig. 1]. Then we take the x-coordinate of that point as the first approximation. If  $f(c) = 0$ , then  $x = c$  is the required root. If  $f(a) f(c) < 0$ , then the root lies in  $]a, c[$  (see Fig. 1 (a)). In this case the graph of  $y = f(x)$  is concave near the root  $r$ . Otherwise, if  $f(a) f(c) > 0$ , the root lies in  $]c, b[$  (see Fig. 1 (b)). In this case the graph of  $y = f(x)$  is convex near the root. Having fixed the interval in which the root lies, we repeat the above procedure.

Let us now write the above procedure in the mathematical form. Recall the formula for the line joining two points in the Cartesian plane [see MTE-05]. The line joining  $(a, f(a))$  and  $(b, f(b))$  is given by

$$y - f(a) = \frac{f(b) - f(a)}{b - a} (x - a)$$

We can rewrite this in the form

$$\frac{y - f(a)}{f(b) - f(a)} = \frac{x - a}{b - a} \quad \dots (1)$$

Since the straight line intersects the x-axis at  $(c, 0)$ , the point  $(c, 0)$  lies on the straight line. Putting  $x = c, y = 0$  in Eqn. (1), we get

$$\frac{-f(a)}{f(b) - f(a)} = \frac{c - a}{b - a}$$

$$\text{i.e. } \frac{c}{b - a} - \frac{a}{b - a} = \frac{-f(a)}{f(b) - f(a)}$$

$$\text{Thus } c = a - \frac{f(a)}{f(b) - f(a)} (b - a) \quad \dots (2)$$

This expression for  $c$  gives an approximate value of a root of  $f(x)$ . Simplifying (2), we can also write it as

$$c = \frac{a f(b) - b f(a)}{f(b) - f(a)}$$

Now, examine the sign of  $f(c)$  and decide in which interval  $]a, c[$  or  $]c, b[$ , the root lies. We thus obtain a new interval such that  $f(x)$  is of opposite signs at the end points of this interval. By repeating this process, we get a sequence of intervals  $]a, b[$ ,  $]a, a_1[$ ,  $]a, a_2[$ , ... as shown in Fig. 2.

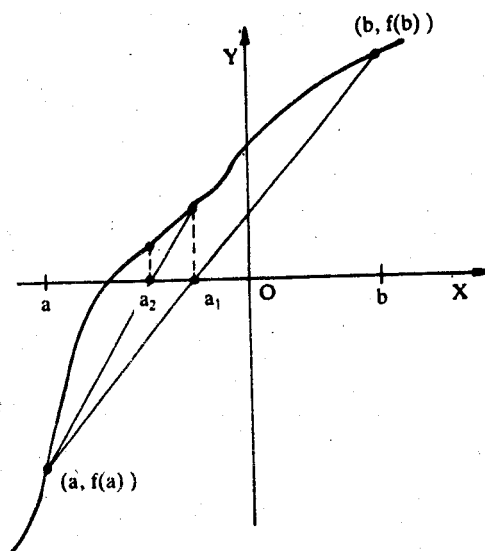


Fig. 2

We stop the process when either of the following holds.

- i) The interval containing the zero of  $f(x)$  is of sufficiently small length  
or
- ii) The difference between two successive approximations is negligible.

In the iteration format, the method is usually written as

$$x_2 = \frac{x_0 f(x_1) - x_1 f(x_0)}{f(x_1) - f(x_0)}$$

where  $]x_0, x_1[$  is the interval in which the root lies.

We now summarise this method in the algorithm form. This will enable you to solve problems easily.

**Step 1 :** Find numbers  $x_0$  and  $x_1$  such that  $f(x_0) f(x_1) < 0$ , using the tabulation method.

**Step 2 :** Set  $x_2 = \frac{x_0 f(x_1) - x_1 f(x_0)}{f(x_1) - f(x_0)}$ . This gives the first approximation.

**Step 3 :** If  $f(x_2) = 0$  then  $x_2$  is the required root. If  $f(x_2) \neq 0$  and  $f(x_0) f(x_2) < 0$ , then the next approximation lies in  $]x_0, x_2[$ . Otherwise it lies in  $]x_2, x_1[$ .

**Step 4 :** Repeat the process till the magnitude of the difference between two successive iterated values  $x_i$  and  $x_{i+1}$  is less than the accuracy required. (Note that  $|x_{i+1} - x_i|$  gives the error after  $i$ th iteration).

Let us now understand these steps through an example.

**Example 1 :** It is known that the equation  $x^3 + 7x^2 + 9 = 0$  has a root between  $-8$  and  $-7$ . Use the regula-falsi method to obtain the root rounded off to 3 decimal places. Stop the iteration when  $|x_{i+1} - x_i| < 10^{-4}$ .

**Solution :** For convenience we rewrite the given function  $f(x)$  as

$$\begin{aligned} f(x) &= x^3 + 7x^2 + 9 \\ &= x^2(x + 7) + 9 \end{aligned}$$

Since we are given that  $x_0 = -8$  and  $x_1 = -7$ , we do not have to use step 1. Now to get the first approximation, we apply the formula in Step 2.

Since,  $f(x_0) = f(-8) = -55$  and  $f(x_1) = f(-7) = 9$  we obtain

$$x_2 = \frac{(-8)9 - (-7)(-55)}{9 + 55} = -7.1406$$

Therefore our first approximation is  $-7.1406$ .

To find the next approximation we calculate  $f(x_2)$ . We have

$$\begin{aligned} f(x_2) &= f(-7.1406) = (-7.1406)^3 + 7(-7.1406)^2 + 9 \\ &= 1.862856 \end{aligned}$$

Now we compare the sign of  $f(x_2)$  with the signs of  $f(x_0)$  and  $f(x_1)$ . We can see that  $f(x_0)$  and  $f(x_2)$  are of opposite signs. Therefore a root lies in the interval  $]-8, -7.1406[$ . We apply the formula again by renaming the end points of the interval as  $x_1 = -8$ ,  $x_2 = -7.1406$ . Then we get the second approximation as

$$x_3 = \frac{-8 f(-7.1406) + 7.1406 f(-8)}{1.862856 + 55} = -7.168174.$$

We repeat this process using steps 2 and 3 given above. The iterated values are given in the following table.

Table 1

Number of iterations	Interval	Iterated Values $x_i$	The function value $f(x_i)$
1	$]-8, -7[$	$-7.1406$	$1.862856$
2	$]-8, -7.1406[$	$-7.168174$	$0.3587607$
3	$]-8, -7.168174[$	$-7.1735649$	$0.0683443$
4	$]-8, -7.1735649[$	$-7.1745906$	$0.012994$
5	$]-8, -7.1745906[$	$-7.1747855$	$0.00246959$
6	$]-8, -7.1747855[$	$-7.1748226$	$0.00046978$

From the table, we see that the absolute value of the difference between the 5th and 6th iterated values is  $|7.1748226 - 7.1747855| = .0000371$ . Therefore we stop the iteration here. Further, the values of  $f(x)$  at 6th iterated value is  $.00046978 = 4.6978 \times 10^{-4}$  which is close to zero. Hence we conclude that  $-7.175$  is an approximate root of  $x^3 + 7x^2 + 9 = 0$  rounded off to three decimal places.

Here is an exercise for you.

E1) Obtain an approximate root for the following equations rounded off to three decimal places, using regula-falsi method

a)  $x \log_{10} x - 1.2 = 0$

b)  $x \sin x - 1 = 0$

You note that in regula-falsi method, at each stage we find an interval  $]x_0, x_1[$  which contains a root and then apply iteration formula (3). This procedure has a disadvantage. To overcome this, regula-falsi method is modified. The modified method is known as secant method. In this method we choose  $x_0$  and  $x_1$  as any two approximations of the root. The Interval  $]x_0, x_1[$  need not contain the root. Then we apply formula (3) with  $x_0, x_1, f(x_0)$  and  $f(x_1)$ .

The iterations are now defined as :

$$x_2 = \frac{x_0 f(x_1) - x_1 f(x_0)}{f(x_1) - f(x_0)}$$

$$x_3 = \frac{x_1 f(x_2) - x_2 f(x_1)}{f(x_2) - f(x_1)}$$

$$x_{n+1} = \frac{x_{n-1} f(x_n) - x_n f(x_{n-1})}{f(x_n) - f(x_{n-1})} \quad \dots (4)$$

Note : Geometrically, in secant Method, we replace the graph of  $f(x)$  in the interval  $[x_n, x_{n+1}]$  by a straight line joining two points  $(x_n, f(x_n))$ ,  $(x_{n+1}, f(x_{n+1}))$  on the curve and take the point of intersection with x-axis as the approximate value of the root. Any line joining two points on the curve is called a secant line. That is why this method is known as secant method. (see Fig. 3).

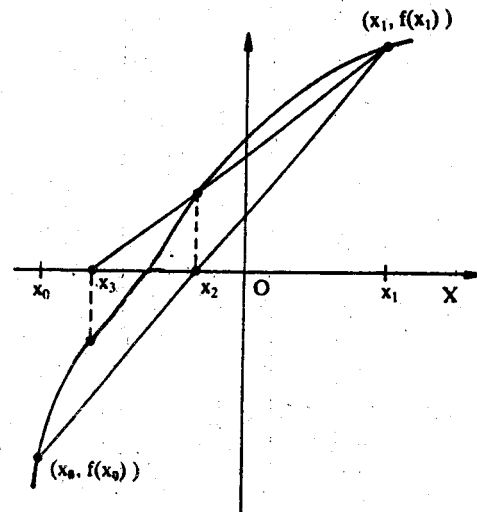


Fig. 3

Let us solve an example.

**Example 2 :** Determine an approximate root of the equation

$$\cos x - x e^x = 0$$

using

- i) secant method starting with the two initial approximations as  $x_0 = 1$  and  $x_1 = 1$  and
- ii) regula-falsi method.

(This example was considered in the book 'Numerical methods for scientific and engineering computation' by M.K. Jain, S.R.K. Iyengar and R.K. Jain).

**Solution :** Let  $f(x) = \cos x - x e^x$ .

Then  $f(0) = 1$  and  $f(1) = \cos 1 - e = -2.177979523$ . Now we apply formula (4) with  $x_0 = 0$  and  $x_1 = 1$ . Then

$$\begin{aligned} x_2 &= \frac{x_0 f(x_1) - x_1 f(x_0)}{f(x_1) - f(x_0)} = \frac{0(-2.177979523) + (-1)1}{-2.177979523 - 1} \\ &= \frac{-1}{-2.177979523 - 1} = \frac{1}{3.177979523} = 0.3146653378. \end{aligned}$$

Therefore the first iterated value is 0.3146653378. To get the 2nd iterated value, we apply Formula (4) with  $x_1 = 1$ ,  $x_2 = 0.3146653378$ . Now  $f(1) = -2.177979523$  and  $f(0.3146653378) = 0.519871175$ .

Therefore

$$\begin{aligned}
 x_3 &= \frac{x_1 f(x_2) - x_2 f(x_1)}{f(x_2) - f(x_1)} \\
 &= \frac{1(0.519871175) - 0.3146653378(-2.177979523)}{0.519871175 + 2.177979523} \\
 &= 0.4467281466
 \end{aligned}$$

We continue this process. The iterated values are tabulated in the following table.

**Table 2 : Secant Method**

Number of iterations	Iterated values $x_i$	$f(x_i)$
1	0.3146653378	0.519871
2	0.4467281466	0.203545
3	0.5317058606	-0.0429311
4	0.5169044676	.00259276
5	0.5177474653	0.00003011
6	0.5177573708	$-0.215132 \times 10^{-7}$
7	0.5177573637	$0.178663 \times 10^{-12}$
8	0.5177573637	$0.222045 \times 10^{-15}$

From the table we find that the iterated values for 7th and 8th iterations are the same. Also the value of the function at the 8th iteration is close to zero. Therefore we conclude that 0.5177573637 is an approximate root of the equation.

- ii) To apply regula-falsi method, let us first note that  $f(0) f(1) < 0$ . Therefore a root lies in the interval  $]0, 1[$ . Now we apply Formula (3) with  $x_0 = 0$  and  $x_1 = 1$ . Then the first approximation is

$$\begin{aligned}
 x_2 &= \frac{0(-2.177979523) + (-1)1}{-2.177979523 - 1} \\
 &= 0.3146653378
 \end{aligned}$$

You may have noticed that we have already calculated the expression on the right hand side of the above equation in part (i).

Now  $f(x_2) = 0.51987 > 0$ . This shows that the root lies in the interval  $]0.3146653378, 1[$ . To get the second approximation, we compute

$$x_3 = \frac{0.3146653378 f(1) - 1 f(0.3146653378)}{f(1) - f(0.3146653378)} = 0.4467281446$$

which is same as  $x_3$  obtained in (i). We find  $f(x_2) = 0.203545 > 0$ . Hence the root lies in  $]0.4467281446, 1[$ . To get the third approximation, we calculate

$$x_4 = \frac{0.4467281446 f(1) - 1 f(0.4467281446)}{f(1) - f(0.4467281446)}$$

The above expression on the right hand side is different from the expression for  $x_4$  in part (i). This is because when we use regula-falsi method, at each stage, we have to check the condition  $f(x_i) f(x_{i-1}) < 0$ .

The computed values of the rest of the approximations are given in Table 3.

**Table 3 : Regula-Falsi Method**

No.	Interval	Iterated value $x_i$	$f(x_i)$
1	]0, 1[	0.3146653378	0.519871
2	]0.04467281446, 1[	0.4467281446	0.203545
3	]0.4940153366, 1[	0.4940153366	$0.708023 \times 10^{-1}$
4	]0.5099461404, 1[	0.5099461404	$0.236077 \times 10^{-1}$
5	]0.5152010099, 1[	0.5152010099	$0.776011 \times 10^{-2}$
6	]0.5176683450, 1[	0.5177478783	$0.288554 \times 10^{-4}$
7	]0.5177478783, 1[	0.5177573636	$0.396288 \times 10^{-9}$

From the table, we observe that we have to perform 20 iterations using regula-falsi method to get the approximate value of the root 0.5177573637 which we obtained by secant method after 8 iterations. Note that the end point 1 is fixed in all iterations given in the table.

Here are some exercises for you.

- E2) Use secant method to find an approximate root of the equation  $x^2 - 2x + 1 = 0$ , rounded off to 5 decimal places, starting with  $x_0 = 2.6$  and  $x_1 = 2.5$ . Compare the result with the exact root  $1 + \sqrt{2}$ .
- E3) Find an approximate root of the cubic equation  $x^3 + x^2 - 3x - 3 = 0$  using
- i) regula-falsi method, correct to three decimal places.
  - ii) secant method starting with  $a = 1$ ,  $b = 2$ , rounded-off to three decimal places.
- b) compare the results obtained by (i) and (ii) in part (a).

Next we shall discuss another iteration method.

### 3.3 NEWTON—RAPHSON METHOD

This method is one of the most useful methods for finding roots of an algebraic equation.

Suppose that we want to find an approximate root of the equation  $f(x) = 0$ . If  $f(x)$  is continuous, then we can apply either bisection method or regula-falsi method to find approximate roots. Now if  $f(x)$  and  $f'(x)$  are continuous, then we can use a new iteration method called Newton—Raphson method. You will learn that this method gives the result more faster than the bisection or regula-falsi methods. The underlying idea of the method is due to mathematician Isac Newton. But the method as now used is due to the mathematician Raphson.

Let us begin with an equation  $f(x) = 0$  where  $f(x)$  and  $f'(x)$  are continuous. Let  $x_0$  be an initial approximation and assume that  $x_0$  is close to the exact root  $\alpha$  and  $f'(x_0) \neq 0$ . Let  $\alpha = x_0 + h$  where  $h$  is a small quantity in magnitude. Hence  $f(\alpha) = f(x_0 + h) = 0$

Now we expand  $f(x_0 + h)$  using Taylor's theorem. Note that  $f(x)$  satisfies all the requirements of Taylor's theorem. Therefore, we get

$$f(x_0 + h) = f(x_0) + hf'(x_0) + \dots = 0$$

Neglecting the terms containing  $h^2$  and higher powers we get

$$f(x_0) + hf'(x_0) = 0.$$

Then, 
$$h = \frac{-f(x_0)}{f'(x_0)}$$



This gives a new approximation to  $\alpha$  as

$$x_1 = x_0 + h = x_0 - \frac{f(x_0)}{f'(x_0)}$$

Now the iteration can be defined by

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

$$x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}$$

$$x_n = x_{n-1} - \frac{f(x_{n-1})}{f'(x_{n-1})} \quad \dots (5)$$

Eqn. (5) is called the **Newton—Raphson formula**. Before solving some examples we shall explain this method geometrically.

### Geometrical Interpretation of Newton—Raphson Method

Let the graph of the function  $y = f(x)$  be as shown in Fig. 4.

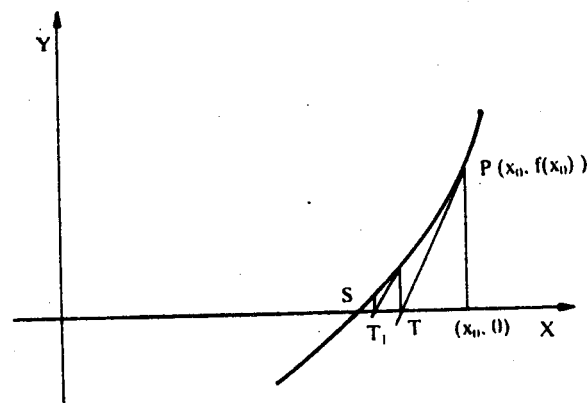


Fig. 4 : Newton—Raphson Method

If  $x_0$  is an initial approximation to the root, then the corresponding point on the graph is  $P(x_0, f(x_0))$ . We draw a tangent to the curve at  $P$ . Let it intersect the  $x$ -axis at  $T$  (see Fig. 4). Let  $x_1$  be the  $x$ -coordinate of  $T$ . Let  $S(\alpha, 0)$  denote the point on the  $x$ -axis where the curve cuts the  $x$ -axis. We know that  $\alpha$  is a root of the equation  $f(x) = 0$ . We take  $x_1$  as the new approximation which may be closer to  $\alpha$  than  $x_0$ . Now let us find the tangent at  $P(x_0, f(x_0))$ . The slope of the tangent at  $P(x_0, f(x_0))$  is given by  $f'(x_0)$ . Therefore by the point-slope form of the expression for a tangent to a curve (recall the expression from MTE-05), we can write

$$y - f(x_0) = f'(x_0) (x_1 - x_0)$$

This tangent passes through the point  $T(x_1, 0)$  (see Fig. 4). Therefore we get

$$0 - f(x_0) = f'(x_0) (x_1 - x_0)$$

$$\text{i.e. } x_1 f'(x_0) = x_0 f'(x_0) - f(x_0)$$

$$\text{i.e. } x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

This is the first iterated value. To get the second iterated value we again consider a tangent at the point  $P(x_1, f(x_1))$  on the curve (see Fig. 4) and repeat the process. Then we get a point

$T_1(x_2, 0)$  on the x-axis. From the figure, we observe that  $T_1$  is more closer to  $S(\alpha, 0)$  than  $T$ . Therefore after each iteration the approximation is coming closer and closer to the actual root. In practice we do not know the actual root of a given function.

Let us now take up some examples.

**Example 3 :** Find the smallest positive root of

$$2x - \tan x = 0$$

by Newton—Raphson method, correct to 5 decimal places.

**Solution :** Let  $f(x) = 2x - \tan x$ . Then  $f(x)$  is a continuous function and  $f'(x) = 2 - \sec^2 x$  is also a continuous function. Recall that the given equation has already appeared in an exercise in Unit 2 (see E2 in Unit 2). From that exercise we know that an initial approximation to the positive root of the equation is  $x = 1$ . Now we apply the Newton—Raphson iterated formula.

$$x_i = x_{i-1} - \frac{f(x_i)}{f'(x_i)}, \quad i = 1, 2, 3, \dots$$

Here  $x_0 = 1$ . Then  $f(x_0) = f(1) = 2 - \tan 1 = 0.4425922$

$$\begin{aligned} f'(x_0) &= f'(1) = 2 - \sec^2 1 = 2 - (1 + \tan^2 1) \\ &= 1 - \tan^2 1 \\ &= -1.425519 \end{aligned}$$

$$\begin{aligned} \text{Therefore } x_1 &= 1 - \frac{0.4425922}{-1.425519} \\ &= 1.31048 \end{aligned}$$

For  $i = 2$ , we get

$$\begin{aligned} x_2 &= 1.31048 - \frac{2 - \tan(1.31048)}{1 - \tan^2(1.31048)} \\ &= 1.22393 \end{aligned}$$

Similarly we get

$$\begin{aligned} x_3 &= 1.17605 \\ x_4 &= 1.165926 \\ x_5 &= 1.165562 \\ x_6 &= 1.165561 \end{aligned}$$

Now  $x_5$  and  $x_6$  are correct to five decimal places. Hence we stop the iteration process here. The root correct to 5 decimal places is 1.16556.

Next we shall consider an application of Newton—Raphson formula. We know that finding the square root of a number is not easy unless we use a calculator. Calculators use some algorithm to obtain this value. Now we shall illustrate how Newton—Raphson method enables us to obtain such an algorithm for calculating square roots. Let's consider an example.

**Example 4 :** Find an approximate value of  $\sqrt{2}$  using the Newton—Raphson formula.

**Solution :** Let  $x = \sqrt{2}$ . Then we have  $x^2 = 2$  i.e.  $x^2 - 2 = 0$ . Hence we need to find the positive root of the equation  $x^2 - 2 = 0$ . Let

$$f(x) = x^2 - 2.$$

Then  $f(x)$  satisfies all the conditions for applying Newton—Raphson method. We choose  $x_0 = 1$  as the initial approximation to the root. This is because we know that  $\sqrt{2}$  lies between  $\sqrt{1}$  and  $\sqrt{4}$  and therefore we can assume that the root will be close to 1.

Now we compute the iterated values.

The iteration formula is

$$\begin{aligned}x_i &= x_{i-1} - \frac{x_{i-1}^2 - 2}{2x_{i-1}} \\ &= \frac{1}{2} \left[ x_{i-1} + \frac{2}{x_{i-1}} \right]\end{aligned}$$

Putting  $i = 1, 2, 3, \dots$ , we get

$$\begin{aligned}x_1 &= \frac{1}{2} \left[ x_0 + \frac{2}{x_0} \right] = 1.5 \\ x_2 &= \frac{1}{2} \left[ 1.5 + \frac{2}{1.5} \right] = 1.4166667 \\ x_3 &= \frac{1}{2} \left[ 1.4166667 + \frac{2}{1.4166667} \right] \\ &= 1.41242157\end{aligned}$$

Similarly

$$\begin{aligned}x_4 &= 1.4142136 \\ x_5 &= 1.4142136\end{aligned}$$

Thus the value of  $\sqrt{2}$  correct to seven decimal places is 1.4142136. Now you can check this value with the calculator.

**Note 1 :** The method used in the above example is applicable for finding square root of any positive real number. For example suppose we want to find an approximate value of  $\sqrt{A}$  where  $A$  is a positive real number. Then we consider the equation  $x^2 - A = 0$ . The iterated formula in this case is

$$x_i = \frac{1}{2} \left[ x_{i-1} + \frac{A}{x_{i-1}} \right].$$

This formula involves only the basic arithmetic operations  $+$ ,  $-$ ,  $\times$  and  $\div$ .

**Note 2 :** From examples (3) and (4), we find that Newton—Raphson method gives the root very fast. One reason for this is that the derivative  $|f'(x)|$  is large compared to  $|f(x)|$  for any  $x = x_i$ . The quantity  $\left| \frac{f(x)}{f'(x)} \right|$  which is the difference between two iterated values is small in this case. In general we can say that if  $|f'(x_i)|$  is large compared to  $|f(x_i)|$ , then we can obtain the desired root very fast by this method.

The Newton—Raphson method has some limitations. In the following remarks we mention some of the difficulties.

**Remark 1 :** Suppose  $f'(x_i)$  is zero in a neighbourhood of the root, then it may happen that  $f'(x_i) = 0$  for some  $x_i$ . In this case we cannot apply Newton—Raphson formula, since division by zero is not allowed.

**Remark 2 :** Another difficulty is that it may happen that  $f'(x)$  is zero only at the roots. This happens in either of the situations.

- i)  $f(x)$  has multiple root at  $\alpha$ . Recall that a polynomial function  $f(x)$  has a multiple root  $\alpha$  of order  $N$  if we can write

$$f(x) = (x - \alpha)^N h(x)$$

where  $h(x)$  is a function such that  $h(\alpha) \neq 0$ . For a general function  $f(x)$ , this means  $f(\alpha) = 0 = f'(\alpha) = \dots = f^{N-1}(\alpha)$  and  $f^N(\alpha) \neq 0$ .

- ii)  $f(x)$  has a stationary point (point of maximum or minimum) point at the root [recall from your calculus course (MTE-01) that if  $f'(x) = 0$  at some point  $x$  then  $x$  is called a stationary point].

In such cases some modifications to the Newton—Raphson method are necessary to get an accurate result. We shall not discuss the modifications here as they are beyond the scope of this course.

You can try some exercises now. Wherever needed, you should use a calculator for computation.

- E4) Starting with  $x_0 = 0$  find an approximate root of the equation  $x^3 - 4x + 1 = 0$ , rounded off to five decimal places using Newton—Raphson method.
- E5) The motion of a planet in the orbit is governed by an equation of the form  $y = x - e \sin x$  where  $e$  stands for the eccentricity. Let  $y = 1$  and  $e = \frac{1}{2}$ . Then find an approximate root of  $2x - 2 - \sin x = 0$  in the interval  $[0, \pi]$  with error less than  $10^{-5}$ . Start with  $x_0 = 1.5$ .
- E6) Using Newton—Raphson square root algorithm, find the following roots within an accuracy of  $10^{-4}$ .
- $8^{1/2}$ , starting with  $x_0 = 3$
  - $9^{1/2}$ , starting with  $x_0 = 10$
- E7) Can Newton—Raphson iteration method be used to solve the equation  $x^{1/3} = 0$ ? Give reasons for your answer.

In the next section we shall discuss a criterion using which we can check the efficiency of an iteration process.

### 3.4 CONVERGENCE CRITERION

In this section we shall introduce a new concept called 'convergence criterion' related to an iteration process. This criterion gives us an idea of how many successive iterations have to be carried out to obtain the root to the desired accuracy. We begin with a definition.

**Definition 1 :** Let  $x_0, x_1, \dots, x_n, \dots$  be the successive approximations of an iteration process. We denote the sequence of these approximations as  $\{x_n\}_{n=0}^{\infty}$ . We say that

$\{x_n\}_{n=0}^{\infty}$  converges to a root  $\alpha$  with order  $p \geq 1$  if

$$|x_{n+1} - \alpha| \leq \lambda |x_n - \alpha|^p \quad \dots (6)$$

for some number  $\lambda > 0$ .  $p$  is called the **order of convergence** and  $\lambda$  is called the **asymptotic error constant**.

For each  $i$ , we denote by  $\epsilon_i = x_i - \alpha$ . Then the above inequality be written as

$$|\epsilon_{i+1}| \leq \lambda |\epsilon_i|^p \quad \dots (7)$$

This inequality shows the relationship between the error in successive approximations. For example, suppose  $p = 2$  and  $|\epsilon_i| \approx 10^{-2}$  for some  $i$ , then we can expect that

$|\epsilon_{i+1}| \approx \lambda 10^{-4}$ . Thus if  $p$  is large, the iteration converges rapidly. When  $p$  takes the integer values 1, 2, 3 then we say that the convergence is **linear, quadratic and cubic** respectively. In the case of linear convergence (i.e.  $p = 1$ ), then we require that  $\lambda < 1$ . In this case we can write (6) as

$$|x_{n+1} - \alpha| \leq \lambda |x_n - \alpha| \text{ for all } n \geq 0 \quad \dots (8)$$

If this condition is satisfied for an iteration process then we say that the iteration process converges linearly.

Setting  $n = 0$  in the inequality (8), we get

$$|x_1 - \alpha| \leq \lambda |x_0 - \alpha|$$

For  $n = 1$ , we get

$$|x_2 - \alpha| \leq \lambda |x_1 - \alpha| \leq \lambda^2 |x_0 - \alpha|$$

Similarly for  $n = 2$ , we get

$$|x_3 - \alpha| \leq \lambda |x_2 - \alpha| \leq \lambda^2 |x_1 - \alpha| \leq \lambda^3 |x_0 - \alpha|$$

Using induction on  $n$ , we get that

$$|x_n - \alpha| \leq \lambda^n |x_0 - \alpha| \text{ for } n \geq 0 \quad \dots (9)$$

If either of the inequalities (8) or (9) is satisfied, then we conclude that  $\{x_n\}_{n=0}^{\infty}$  converges to the root.

Now we shall find the order of convergence of the iteration methods which you have studied so far.

Let us first consider bisection method.

### Convergence of bisection method

Suppose that we apply the bisection method on the interval  $[a_0, b_0]$  for the equation  $f(x) = 0$ .

In this method you have seen that we construct intervals  $[a_0, b_0] \supset [a_1, b_1] \supset [a_2, b_2] \supset \dots$  each of which contains the required root of the given equation.

Recall that in each step the interval width is reduced by  $\frac{1}{2}$  i.e.

$$b_1 - a_1 = \frac{b_0 - a_0}{2}$$

$$b_2 - a_2 = \frac{b_1 - a_1}{2} = \frac{b_0 - a_0}{2^2}$$

and 
$$b_n - a_n = \frac{b_0 - a_0}{2^n} \quad \dots (10)$$

We know that the equation  $f(x) = 0$  has a root in  $[a_0, b_0]$ . Let  $\alpha$  be the root of the equation.

Then  $\alpha$  lies in all the intervals  $[a_i, b_i]$ ,  $i = 0, 1, 2, \dots$ . For any  $n$ , let  $c_n = \frac{a_n + b_n}{2}$  denote the middle point of the interval  $[a_n, b_n]$ . Then  $c_0, c_1, c_2, \dots$  are taken as successive approximations to the root  $\alpha$ . Let's check the inequality (8) for  $\{c_n\}_{n=0}^{\infty}$ .

For each  $n$ ,  $\alpha$  lies in the interval  $[a_n, b_n]$ . Therefore we have

$$|\alpha - c_{n+1}| \leq \frac{|\alpha - c_n|}{2}$$

Thus  $\{c_n\}_{n=0}^{\infty}$  converges to the root  $\alpha$ . Hence we can say that the bisection method always converges.

For practical purposes, we should be able to decide at what stage we can stop the iteration to have an acceptably good approximate value of  $\alpha$ . The number of iterations required to achieve a given accuracy for the bisection method can be obtained. Suppose that we want an approximate solution within an error bound of  $10^{-M}$  (Recall that you have studied error bounds in Unit 1, Sec 3.4). Taking logarithms on both sides of Eqn. (10), we find that the number of iterations required, say  $n$ , is approximately given by

$$n = \text{int} \left[ \frac{\ln(b_0 - a_0) - \ln 10^{-M}}{\ln 2} \right] \quad \dots (11)$$

where the symbol 'int' stands for the integral part of the number in the bracket and  $[a_0, b_0]$  is the initial interval in which a root lies.

Let us work out an example.

**Example 5 :** Suppose that the bisection method is used to find a zero of  $f(x)$  in the interval  $[0, 1]$ . How many times this interval be bisected to guarantee that we have an approximate root with absolute error less than or equal to  $10^{-5}$

**Solution :** Let  $n$  denote the required number. To calculate  $n$ , we apply the formula in Eqn. (11) with  $b_0 = 1, a_0 = 0$  and  $M = 5$ .

Then

$$n = \text{int} \left[ \frac{\ln 1 - \ln 10^{-5}}{\ln 2} \right]$$

Using a calculator, we find

$$\begin{aligned} n &= \text{int} \left[ \frac{11.51292547}{0.69314718} \right] \\ &= \text{int} [16.60964047] = 17 \end{aligned}$$

Similarly you can try the following exercise.

E8) For the problem given in Example 5, Unit 2, find the number  $n$  of bisections required to have an approximate root with absolute error less than or equal to  $10^{-7}$ .

The following table gives the minimum number of iterations required to find an approximate root in the interval  $[0, 1]$  for various acceptable errors.

E	$10^{-2}$	$10^{-3}$	$10^{-4}$	$10^{-5}$	$10^{-6}$	$10^{-7}$
n	7	10	14	17	20	24

This table shows that for getting an approximate value with an absolute error bounded by  $10^{-5}$ , we have to perform 17 iterations. Thus even though the bisection method is simple to use, it requires a large number of iterations to obtain a reasonably good approximate root. This is one of the disadvantages of the bisection method.

**Note :** The formula given in Eqn. (11) shows that, given an acceptable error, the number of iterations depends upon the initial interval and thereby depends upon the initial approximation of the root and not directly on the values of  $f(x)$  at these approximations.

Next we shall obtain the convergence criteria for the secant method.

#### Convergence criteria for Secant Method

Let  $f(x) = 0$  be the given equation. Let  $\alpha$  denote a simple root of the equation  $f(x) = 0$ . Then we have  $f(\alpha) = 0$ . The iteration scheme for the secant method is

$$x_{i+1} = x_i - \frac{x_i - x_{i-1}}{f(x_i) - f(x_{i-1})} f(x_i) \quad \dots (12)$$

For each  $i$ , set  $\epsilon_i = x_i - \alpha$ . Then  $x_i = \epsilon_i + \alpha$ . Substituting in Eqn. (12) we get

$$\varepsilon_{i+1} + \alpha = \varepsilon_i + \alpha - \frac{\varepsilon_i - \varepsilon_{i-1}}{f(\varepsilon_i + \alpha) - f(\varepsilon_{i-1} + \alpha)} f(\varepsilon_i + \alpha)$$

$$\varepsilon_{i+1} = \varepsilon_i - \frac{\varepsilon_i - \varepsilon_{i-1}}{f(\varepsilon_i + \alpha) - f(\varepsilon_{i-1} + \alpha)} f(\varepsilon_i + \alpha) \quad \dots (13)$$

Now we expand  $f(\varepsilon_i + \alpha)$  and  $f(\varepsilon_{i-1} + \alpha)$  using Taylor's theorem about the point  $x = \alpha$ .

We get  $f(\varepsilon_i + \alpha) = f(\alpha) + \frac{f'(\alpha)}{1} \varepsilon_i + \frac{f''(\alpha)}{2} \varepsilon_i^2 + \dots$

i.e.  $f(\varepsilon_i + \alpha) = f'(\alpha) \left[ \varepsilon_i + \frac{f''(\alpha)}{2 f'(\alpha)} \varepsilon_i^2 + \dots \right] \quad \dots (14)$

since  $f(\alpha) = 0$ .

Similarly,

$$f(\varepsilon_{i-1} + \alpha) = f'(\alpha) \left[ \varepsilon_{i-1} + \frac{f''(\alpha)}{2 f'(\alpha)} \varepsilon_{i-1}^2 + \dots \right] \quad \dots (15)$$

Therefore  $f(\varepsilon_i + \alpha) - f(\varepsilon_{i-1} + \alpha) = f'(\alpha) \left[ \varepsilon_i - \varepsilon_{i-1} + (\varepsilon_i^2 - \varepsilon_{i-1}^2) \frac{f''(\alpha)}{2 f'(\alpha)} + \dots \right]$

$$= f'(\alpha) (\varepsilon_i - \varepsilon_{i-1}) \left[ 1 + (\varepsilon_i + \varepsilon_{i-1}) \frac{f''(\alpha)}{2 f'(\alpha)} + \dots \right] \quad \dots (16)$$

Substituting Eqn. (14) and Eqn. (16) in Eqn. (13), we get

$$\varepsilon_{i+1} = \varepsilon_i - \left[ \varepsilon_i + \frac{1}{2} \varepsilon_i^2 \frac{f''(\alpha)}{f'(\alpha)} + \dots \right] \left[ 1 + \frac{1}{2} (\varepsilon_i + \varepsilon_{i-1}) \frac{f''(\alpha)}{f'(\alpha)} + \dots \right]^{-1}$$

$$= \varepsilon_i - \left[ \varepsilon_i + \frac{1}{2} \varepsilon_i^2 \frac{f''(\alpha)}{f'(\alpha)} + \dots \right] \left[ 1 - \frac{1}{2} (\varepsilon_i + \varepsilon_{i-1}) \frac{f''(\alpha)}{f'(\alpha)} + \dots \right]$$

$$= \varepsilon_i - \left[ \varepsilon_i + \frac{1}{2} \frac{f''(\alpha)}{f'(\alpha)} (\varepsilon_i^2 - \varepsilon_{i-1}^2 - \varepsilon_i \varepsilon_{i-1}) + \dots \right]$$

By neglecting the terms involving  $\varepsilon_i \varepsilon_{i-1}^2 + \varepsilon_{i-1}^2 \varepsilon_i'$  the above expression, we get

$$\varepsilon_{i+1} \approx \varepsilon_i \varepsilon_{i-1} \left[ \frac{f''(\alpha)}{2f'(\alpha)} \right] \quad \dots (17)$$

This relationship between the errors is called the error equation. Note that this relationship holds only if  $\alpha$  is a simple root. Now using Eqn. (17) we will find a numbers  $p$  and  $\lambda$  such that

$$\varepsilon_{i+1} = \lambda \varepsilon_i^p \quad i = 0, 1, 2, \dots \quad \dots (18)$$

Setting  $i = j - 1$ , we obtain

$$\varepsilon_j = \lambda \varepsilon_{j-1}^p$$

or

$$\varepsilon_i = \lambda \varepsilon_{i-1}^p$$

Taking  $p$ th root on both sides, we get

$$\varepsilon_i^{1/p} = \lambda^{1/p} \varepsilon_{i-1}$$

i.e.  $\varepsilon_{i-1} = \lambda^{-1/p} \varepsilon_i^{1/p} \quad \dots (19)$

Combining Eqns. (17) and (18), we get

$$\lambda \varepsilon_j^p = \varepsilon_i \varepsilon_{i-1} \frac{f''(\alpha)}{2f'(\alpha)}$$

Substituting the expression for  $\epsilon_{i-1}$  from Eqn. (19) in the above expression we get

$$\lambda \epsilon_i^p = \frac{f''(\alpha)}{2f'(\alpha)} \epsilon_i \lambda^{-1/p} \epsilon_i^{1/p}$$

i.e.  $\lambda \epsilon_i^p = \frac{f''(\alpha)}{2f'(\alpha)} \lambda^{-1/p} \epsilon_i^{1+1/p} \dots (20)$

Equating the powers of  $\epsilon_i$  on both sides of Eqn. (20) we get

$$p = 1 + \frac{1}{p} \text{ or } p^2 - p - 1 = 0.$$

This is a quadratic equation in p. The roots are given by

$$p = \frac{1 + \sqrt{5}}{2}$$

Since p cannot be negative we ignore the negative value. Hence we have,

$$p = \frac{1 + \sqrt{5}}{2} \approx 1.618.$$

Now, to get the number  $\lambda$ , we equate the constant terms on both sides of Eqn. (20). Then we get

$$\lambda = \left[ \frac{f''(\alpha)}{2f'(\alpha)} \right]^{p/1+p}$$

Hence the order of convergence of the secant method is  $p = 1.62$  and the asymptotic error

constant is  $\left[ \frac{f''(\alpha)}{2f'(\alpha)} \right]^{p/1+p}$

**Example 6 :** The following are the five successive iterations obtained by secant method to find the root  $\alpha = -2$  of the equation  $x^3 - 3x + 2 = 0$ .

$$x_1 = -2.6, x_2 = -2.4, x_3 = -2.106598985,$$

$$x_4 = -2.022641412, \text{ and } x_5 = -2.000022537.$$

Compute the asymptotic error constant and show that  $\epsilon_5 \approx \frac{2}{3} \epsilon_4$ .

**Solution :** Let  $f(x) = x^3 - 3x + 2$

Then

$$f'(x) = 3x^2 - 3, f'(-2) = 9$$

$$f''(x) = 6x, f''(-2) = -12$$

Therefore  $\lambda = \left[ -\frac{12}{18} \right]^{618}$   
 $= \left[ -\frac{2}{3} \right]^{618} = -0.778351205$

Now

$$\epsilon_5 = |x_5 - \alpha| = |-2.000022537 + 2|$$

$$= 0.000022537$$

and

$$\epsilon_4 = |-2.022641412 + 2| = 0.022641412.$$

Then  $\lambda \epsilon_4 = 0.778351205 \times 0.022641412$

$$= 0.000021246$$

$$\approx 0.00002253$$



Hence we get that  $\lambda \varepsilon_4 \approx \varepsilon_5$

### Convergence criterion for fixed point iteration method

Recall that in this method we write the equation in the form

$$x = g(x)$$

Let  $\alpha$  denote a root of the equation. Let  $x_0$  be an initial approximation to the root. The iteration formula is

$$x_{i+1} = g(x_i), i = 0, 1, 2, \dots \quad \dots (21)$$

We assume that  $g'(x)$  exists and is continuous and  $|g'(x)| < 1$  in an interval containing the root  $\alpha$ . We also assume that  $x_0, x_1, \dots$  lie in this interval.

Since  $g'(x)$  is continuous near the root and  $|g'(x)| < 1$ , there exists an interval  $]\alpha - h, \alpha + h[$ , where  $h > 0$ , such that  $|g'(x)| \leq k$  for some  $k$ , where  $0 < k < 1$ .

Since  $\alpha$  is a root of the equation, we have

$$\alpha = g(\alpha). \quad \dots (22)$$

Subtracting (22) from (21) we get

$$x_{i+1} - \alpha = g(x_i) - g(\alpha)$$

Now the function  $g(x)$  is continuous in the interval  $]x_i, \alpha[$  and  $g'(x)$  exists in this interval.

Hence  $g(x)$  satisfies all the conditions of the mean value theorem [see unit 1]. Then, by the mean value theorem there exists a  $\xi$  between  $x_i$  and  $\alpha$  such that

$$|x_{i+1} - \alpha| \leq |g(x_i) - g(\alpha)| \leq |g'(\xi)| |x_i - \alpha|$$

Note that  $\xi$  lies in  $]\alpha - h, \alpha + h[$  and therefore  $|g'(\xi)| < k$  and hence

$$|x_{i+1} - \alpha| \leq k |x_i - \alpha|$$

Setting  $i = 0, 1, 2, \dots, n$  we get

$$|x_1 - \alpha| \leq k |x_0 - \alpha|$$

$$|x_2 - \alpha| \leq k |x_1 - \alpha| \leq k^2 |x_0 - \alpha|$$

$$|x_n - \alpha| \leq k^n |x_0 - \alpha|$$

This shows that the sequence of approximations  $\{x_i\}$  converges to  $\alpha$  provided that the initial approximation is close to the root.

We summarise the result obtained for this iteration process in the following Theorem.

**Theorem :** If  $g(x)$  and  $g'(x)$  are continuous in an interval about a root  $\alpha$  of the equation  $x = g(x)$ , and if  $|g'(x)| < 1$  for all  $x$  in the interval, then the successive approximations  $x_1, x_2, \dots$  given by

$$x_i = g(x_{i-1}), i = 1, 2, 3, \dots$$

converges to the root  $\alpha$  provided that the initial approximation  $x_0$  is chosen in the above interval.

We shall now discuss the order of convergence of this method. From the previous discussions we have the result.

$$|x_{i+1} - \alpha| \leq |g'(\xi)| |x_i - \alpha|$$

Note that  $\xi$  is dependent on each  $x_i$ . Now we wish to determine the constants  $\lambda$  and  $p$  independent of  $x_i$  such that

$$|x_{i+1} - \alpha| \leq c |x_i - \alpha|^p$$

Note that as the approximations  $x_i$  get closer to the root  $\alpha$ ,  $g'(\xi)$  approaches a constant value  $g'(\alpha)$ . Therefore, in the limiting case, as  $i \rightarrow \infty$ , the approximations satisfy the relation

$$|x_{i+1} - \alpha| \leq |g'(\alpha)| |x_i - \alpha|$$

Therefore, we conclude that if  $g'(\alpha) \neq 0$ , then the convergence of the method is linear.

If  $g'(\alpha) = 0$ , then we have

$$\begin{aligned} x_{i+1} - \alpha &= g(x_i) - \alpha \\ &= g[(x_i - \alpha) + \alpha] - \alpha \\ &= g(\alpha) + (x_i - \alpha)g'(\alpha) + \frac{(x_i - \alpha)^2}{2} g''(\xi) - \alpha \\ &= \frac{(x_i - \alpha)^2}{2} g''(\xi) \end{aligned}$$

By applying Taylor's theorem to the function  $g(x)$  about  $\alpha$  and neglecting higher powers.

since  $g(\alpha) = \alpha$  and  $g'(\alpha) = 0$  and  $\xi$  lies between  $x_i$  and  $\alpha$ .

Therefore, in the limiting case we have,

$$|x_{i+1} - \alpha| \leq \frac{1}{2} |g''(\alpha)| |x_i - \alpha|^2$$

Hence, if  $g'(\alpha) = 0$  and  $g''(\alpha) \neq 0$ , then this iteration method is of order 2.

**Example 7 :** Suppose  $\alpha$  and  $\beta$  are the roots of the equation  $x^2 + ax + b = 0$ . Consider a rearrangement of this equation as

$$x = -\frac{(ax + b)}{x}$$

Show that the iteration  $x_{i+1} = -\frac{(ax_i + b)}{x_i}$  will converge near  $x = \alpha$  when  $|\alpha| > |\beta|$

**Solution :** The iterations are given by

$$x_{i+1} = g(x_i) = -\frac{(ax_i + b)}{x_i}, \quad i = 0, 1, 2, \dots$$

By Theorem 1, these iterations converge to  $\alpha$  if  $|g'(x)| < 1$  near  $\alpha$  i.e. if  $|g'(x)| = \left| -\frac{b}{x^2} \right| < 1$ . Note that  $g'(x)$  is continuous near  $\alpha$ . If the iterations converge to  $x = \alpha$ , then we require  $|g'(\alpha)| = \left| -\frac{b}{\alpha^2} \right| < 1$ .

Thus  $|b| < |\alpha|^2$

i.e.  $|\alpha|^2 > |b|$ . ... (23)

Now you recall from your elementary algebra course (MTE-04) that if  $\alpha$  and  $\beta$  are the roots, then

$$\alpha + \beta = -a \text{ and } \alpha\beta = b$$

Therefore  $|b| = |\alpha| |\beta|$ . Substituting in Eqn. (23), we get

$$|\alpha|^2 > |b| = |\alpha| |\beta|$$

Hence  $|\alpha| > |\beta|$ .

Similarly you can solve the following exercise.

E9) For the equation given in Example 7, show that the iteration  $x_{i+1} = \frac{b}{x_i + a}$  will converge to the root  $x = \alpha$ , when  $|\alpha| < |\beta|$ .

Finally we shall discuss the convergence of the Newton—Raphson method.

### Convergence of Newton—Raphson Method

Newton—Raphson iteration formula is given by

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} \quad \dots (24)$$

To obtain the order of the method we proceed as in the secant method. We assume that  $\alpha$  is a simple root of  $f(x) = 0$ . Let

$$x_i - \alpha = \varepsilon_i, \quad i = 0, 1, 2, \dots$$

Then we have

$$\varepsilon_{i+1} + \alpha = \varepsilon_i + \alpha - \frac{f(\varepsilon_i + \alpha)}{f'(\varepsilon_i + \alpha)}$$

i.e. 
$$\varepsilon_{i+1} = \frac{\varepsilon_i f'(\varepsilon_i + \alpha) - f(\varepsilon_i + \alpha)}{f'(\varepsilon_i + \alpha)}$$

Now we expand  $f(\varepsilon_i + \alpha)$  and  $f'(\varepsilon_i + \alpha)$ , using Taylor's theorem, about the point  $\alpha$ . We have

$$\varepsilon_{i+1} = \frac{\left[ \varepsilon_i \left\{ f'(\alpha) + \varepsilon_i f''(\alpha) + \frac{\varepsilon_i^2}{2} f'''(\alpha) + \dots \right\} - \left\{ f(\alpha) + \varepsilon_i f'(\alpha) + \frac{\varepsilon_i^2}{2} f''(\alpha) + \dots \right\} \right]}{f'(\alpha) + \varepsilon_i f''(\alpha) + \varepsilon_i^2 f'''(\alpha) + \dots}$$

But  $f(\alpha) = 0$  and  $f'(\alpha) \neq 0$ . Therefore

$$\begin{aligned} \varepsilon_{i+1} &= \left[ \frac{\varepsilon_i^2}{2} f''(\alpha) + \dots \right] \frac{1}{f'(\alpha)} \left[ 1 + \frac{\varepsilon_i f''(\alpha)}{f'(\alpha)} + \dots \right]^{-1} \\ &= \frac{1}{f'(\alpha)} \left[ \frac{\varepsilon_i^2}{2} f''(\alpha) + \dots \right] \left[ 1 - \varepsilon_i \frac{f''(\alpha)}{f'(\alpha)} + \dots \right] \end{aligned}$$

Hence, by neglecting higher powers of  $\varepsilon_i$ , we get

$$\varepsilon_{i+1} \approx \frac{f''(\alpha)}{2f'(\alpha)} \varepsilon_i^2$$

This shows that the errors satisfy Eqn. (6) with  $p = 2$  and  $\lambda = \frac{f''(\alpha)}{2f'(\alpha)}$ . Hence

Newton—Raphson method is of order 2. That is at each step, the error is proportional to the square of the previous error.

Now, we shall discuss an alternate method for showing that the order is 2. Note that we can write (24) in the form  $x = g(x)$  where

$$g(x) = x - \frac{f(x)}{f'(x)}$$

Then

$$g'(x) = \frac{d}{dx} \left[ x - \frac{f(x)}{f'(x)} \right] = 1 - \frac{[f'(x)]^2 - f(x) f''(x)}{[f'(x)]^2}$$

$$= \frac{f(x) f''(x)}{[f'(x)]^2}$$

Now,

$$g'(\alpha) = \frac{f(\alpha) f''(\alpha)}{[f'(\alpha)]^2} = 0, \text{ since } f(\alpha) = 0 \text{ and } f'(\alpha) \neq 0.$$

Hence by the conclusion drawn just above Example 7, the method is of order 2. Note that this is true only if  $\alpha$  is a simple root. If  $\alpha$  is a multiple root i.e. if  $f'(\alpha) = 0$ , then the convergence is not quadratic, but only linear. We shall not prove this result, but we shall illustrate this with an example.

Let us consider an example.

**Example 8 :** Let  $f(x) = (x - 2)^4 = 0$ . Starting with the initial approximation  $x_0 = 2.1$ , compute the iterations  $x_1, x_2, x_3$  and  $x_4$  using Newton—Raphson method. Is the sequence converging quadratically or linearly?

**Solution :** The given function has multiple roots at  $x = 2$  and is of order 4.

Newton—Raphson iteration formula for the given equation is

$$x_{i+1} = x_i - \frac{(x_i - 2)^4}{4(x_i - 2)^3} = x_i - \frac{1}{4} (x_i - 2)$$

$$= \frac{3}{4} (x_i - 2) \quad \dots (25)$$

Starting with  $x_0 = 2.1$ , the iterations are given by

$$x_1 = \frac{1}{4} (6.3 + 2) = \frac{8.3}{2} = 2.075$$

Similarly

$$x_2 = 2.05625$$

$$x_3 = 2.0421875$$

$$x_4 = 2.031640625$$

Now  $\epsilon_0 = x_0 - 2 = 0.1$ ,  $\epsilon_1 = x_1 - 2 = 0.075$ ,  $\epsilon_2 = 0.05625$ ,  $\epsilon_3 = 0.0421875$ ,  
 $\epsilon_4 = 0.031640625$ .

Then

$$\epsilon_1 = .075 = \frac{3}{4} \times 0.1 = \frac{3}{4} \epsilon_0$$

and

$$\epsilon_2 = \frac{3}{4} \epsilon_1$$

$$\epsilon_3 = \frac{3}{4} \epsilon_2$$

$$\epsilon_4 = \frac{3}{4} \epsilon_3$$

Thus the convergence is linear in this case. The error is reduced by a factor of  $\frac{3}{4}$  with each iteration. This result can also be obtained directly from Eqn. (25).

You can try this exercise now :

E10) The quadratic equation  $x^4 - 4x^2 + 4 = 0$  has a double root at  $x = \sqrt{2}$ . Starting with  $x_0 = 1.5$ , compute three successive approximations to the root by Newton—Raphson method. Does the result converge quadratically or linearly ?

We now end this unit by giving a summary of it.

### 3.5 SUMMARY

In this unit we have

- described the following methods for finding a root of an equation  $f(x) = 0$ 
  - i) **Regula-falsi method :**

The formula is

$$c = \frac{a f(b) - b f(a)}{f(b) - f(a)}$$

where  $[a, b]$  is an interval such that  $f(a) f(b) < 0$ .

- ii) **Secant method :**

The iteration formula is

$$x_{i+1} = \frac{x_{i-1} f(x_i) - x_i f(x_{i-1})}{f(x_i) - f(x_{i-1})}, i = 0, 1, 2, \dots$$

where  $x_0$  and  $x_1$  are any two given approximations of the root.

- iii) **Newton—Raphson method :**

The iteration formula is

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}, i = 0, 1, 2, \dots$$

where  $x_0$  is an initial approximation to the root.

- introduced the concept called convergence criterion of an iteration process
- discussed the convergence of the following iterative methods
  - i) Bisection method
  - ii) Fixed point iteration method
  - iii) Secant method
  - iv) Newton-Raphson method.

### 3.6 SOLUTIONS/ANSWERS

E1) i) Let  $f(x) = x \log_{10} x - 1.2 = 0$

We have to first find two numbers  $a$  and  $b$  such that  $f(a) f(b) < 0$ . Since the function  $\log_{10} x$  is defined only for positive values of  $x$ , we consider only positive numbers  $x$ . Let us take  $x = 1, 2, 3, \dots$ . Then, using a calculator,

$$f(1) = 1 (\log_{10} 1) - 1.2 = -1.2 < 0$$

$$f(2) = 2 (\log_{10} 2) - 1.2 = 2 (.30103) - 1.2 = -.59794 < 0$$

$$f(3) = 3 (\log_{10} 3) - 1.2 = 3 (.47712) - 1.2 = .23136 > 0$$

This shows that  $f(2)f(3) < 0$  and therefore a root lies in  $]2, 3[$ . Now put  $a = 2$  and  $b = 3$ . Then the first approximation of the root is

$$\begin{aligned} x_1 &= \frac{a f(b) - b f(a)}{f(b) - f(a)} \\ &= \frac{2(.23136) - 3(-.59794)}{.23136 + .59794} \\ &= 2.72102 \end{aligned}$$

Now  $f(2.72102) = 2.72102 (\log_{10} 2.72102) - 1.2 = 1.18291 - 1.2 < 0$ . Since  $f(2.72102)f(3) < 0$ , a root lies in the interval  $]2.72102, 3[$ . Hence the second approximation is

$$\begin{aligned} x_2 &= \frac{2.72102 f(3) - 3 f(2.72102)}{f(3) - f(2.72102)} \\ &= 2.7402 \end{aligned}$$

We find  $f(x_2) = -0.0004 < 0$ . Therefore the root lies in the interval  $]2.7402, 3[$ . The third approximation is obtained as

$$\begin{aligned} x_3 &= \frac{2.7402 f(3) - 3 f(2.7402)}{f(3) - f(2.7402)} \\ &= 2.7406 \end{aligned}$$

Since  $x_2$  and  $x_3$  rounded off to three decimal places are the same, we stop the process here. Hence the desired approximate value of the root rounded off to three decimal places is 2.740.

Let  $f(x) = x \sin x - 1$

Since  $f(0) = -1$  and  $f(2) = 0.818594854$ , a root lies in the interval  $]0, 2[$ . The first approximation is

$$x_1 = \frac{0 f(2) - 2 f(0)}{2} = 1.09975017$$

and  $f(x_1) = -0.02001921$

Since  $f(x_1) < 0$  and  $f(2) > 0$ , the root lies in  $] -0.02001921, 2[$ .

The second approximation is obtained as

$$x_2 = 1.2124074$$

and

$$f(x_2) = -0.00983461.$$

The root now lies in  $]1.2124074, 2[$ .

Similarly we can calculate the third and fourth approximations as

$$x_3 = 1.11416120$$

and

$$x_4 = 1.11415714$$

Since  $x_3$  and  $x_4$  rounded off to three decimal places are the same, we stop the process here. Hence the desired root is 1.114.

E2) Let  $f(x) = x^2 - 2x - 1$ . Starting with  $x_0 = 2.6$  and  $x_1 = 2.5$  the successive approximations are,

$$\begin{aligned} x_2 &= \frac{x_0 f(x_1) - x_1 f(x_0)}{f(x_1) - f(x_0)} \\ &= \frac{2.6 f(2.5) - 2.5 f(2.6)}{f(2.5) - f(2.6)} \\ &= \frac{2.6 (.25) - 2.5 (.56)}{.25 - .56} \end{aligned}$$

$$= 2.41935484$$

and  $f(x_2) = 0.0145682$ .

To find the next approximation we compute

$$\begin{aligned} x_3 &= \frac{x_1 f(x_2) - x_2 f(x_1)}{f(x_2) - f(x_1)} \\ &= \frac{2.5 (0.0145682) - (2.41935484) (.25)}{(0.0145682) - (.56)} \\ &= 2.41436464 \end{aligned}$$

Similarly you can calculate that

$$x_4 = 2.41421384$$

and

$$x_5 = 2.41421356$$

Since  $x_4$  and  $x_5$  rounded off to 5 decimal places are the same, we stop the process here. Therefore the required root rounded off to 5 decimal places is 2.41421.

Now we compare this root with the exact root  $1 + \sqrt{2}$ . Using a calculator we  $1 + \sqrt{2} = 2.41421$ , rounded off to five decimal places. Hence the computed root and exact root are the same when we round off to five decimal places.

Let  $f(x) = x^3 + x^2 - 3x - 3 = 0$

- i) We first note that  $f(1) < 0$  and  $f(2) > 0$ . Therefore a root lies in  $[1, 2]$ . The first approximation  $x_1$  is

$$x_1 = \frac{1 f(2) - 2 f(1)}{f(2) - f(1)} = \frac{11}{7} = 1.57142$$

and  $f(x_1) = -1.36449 < 0$

Therefore the root lies in  $]1.57142, 2[$ .

Proceeding similarly, we get the values as given in the following Table.

No.	Interval	Approximation $x_i$	$f(x_i)$
1	$]1, 2[$	1.57142	-1.36449
2	$]1.57142, 2[$	1.70540	-0.24784
3	$]1.70540, 2[$	1.72788	0.03936
4	$]1.72788, 2[$	1.73140	-0.00615
5	$]1.73140, 2[$	1.73194	

The table shows that  $x_5$  and  $x_4$  are correct to three decimal places. Therefore we stop the process here. Hence the root correct to three decimal places is 1.731.

- ii) In secant method we start with two approximations  $a = 1$  and  $b = 2$ . Then the first approximation is the same as in part (i), namely

$$x_1 = 1.57142$$

To calculate the next approximation  $x_2$  we take  $b$  and  $x_1$ . Here also we are getting the same value as in part (i), namely

$$x_2 = 1.70540$$

Then we take  $x_1 = 1.57142$  and  $x_2 = 1.70540$  to get the third approximation  $x_3$ .

We have

$$x_3 = \frac{1.57142 f(1.70540) - 1.70540 f(1.57142)}{f(1.70540) - f(1.57142)} = 1.73513.$$

The rest of the values are given by

$$x_4 = 1.73199$$

and

$$x_5 = 1.73205$$

Since  $x_4$  and  $x_5$  rounded off to three decimal places are the same, we stop here. Hence the root is 1.732, rounded off to three decimal places.

Let us now compare the two methods. We first note that  $|x_{i+1} - x_i|$  gives the error after  $i$ th iteration.

In regula-falsi method, the error after 5th iteration is

$$\begin{aligned} |x_5 - x_4| &= |1.73194 - 1.73140| \\ &= .00011 \end{aligned}$$

whereas in secant method, the error after 5th iteration is

$$\begin{aligned} |x_5 - x_4| &= |1.73205 - 1.73199| \\ &= .00006 \end{aligned}$$

This shows that the error in the case of secant method is smaller than that in regula-falsi method for the same number of iterations.

- E4) The given function  $f(x) = x^3 - 4x + 1$  and its derivative  $f'(x) = 3x^2 - 4$  are continuous everywhere.

The initial approximation is  $x_0 = 0$ .

The iteration formula is

$$x_i = x_{i-1} - \frac{f(x_{i-1})}{f'(x_{i-1})}, \quad i = 0, 1, 2, \dots$$

The first approximation is

$$x_1 = 0 - \frac{f(0)}{f'(0)} = -\frac{1}{(-4)} = 0.25$$

$$\text{and } f(x_1) = (0.25)^3 - 4(0.25) + 1 = 0.015625$$

$$f'(x_1) = 3(0.25)^2 - 4 = -3.8125$$

The second approximation is given by

$$\begin{aligned} x_2 &= x_1 - \frac{f(x_1)}{f'(x_1)} \\ &= .25 + \frac{0.015625}{3.8125} \\ &= 0.254098 \end{aligned}$$

Similarly we get

$$x_3 = 0.254101$$

Since  $x_2$  and  $x_3$  rounded off to four decimal places are the same, we stop the iteration here. Hence the root is 0.2541.

- E5) Let  $f(x) = 2x - 2 - \sin x$ . The  $f(x)$  and  $f'(x)$  are continuous everywhere. Starting with  $x_0 = 1.5$ , we compute the iterated values by the Newton—Raphson formula. The first iteration is

$$\begin{aligned} x_1 &= 1.5 - \frac{f(1.5)}{f'(1.5)} \\ &= 1.5 - \frac{1 - \sin(1.5)}{2 - \cos(1.5)} \end{aligned}$$



$$= 1.5 - \frac{0.002505}{1.929263} = 1.498702$$

Similarly,

$$x_2 = 1.498701$$

We find  $|x_2 - x_1| = |1.498701 - 1.498702| < 10^{-5}$

Therefore the required root is 1.498701

E6) i) Newton—Raphson iterated formula for computing the  $\sqrt{8}$  is

$$x_i = \frac{1}{2} \left[ x_{i-1} + \frac{8}{x_{i-1}} \right], i = 0, 1, 2, \dots$$

Starting with  $x_0 = 3$ , we obtain the iterated values as

$$x_1 = \frac{1}{2} \left[ 3 + \frac{8}{3} \right] = 2.833333$$

$$x_2 = \frac{1}{2} \left[ 2.833333 + \frac{8}{2.833333} \right] \\ = 2.828431$$

and  $x_3 = 2.828427$

Since  $|x_3 - x_2| < 10^{-4}$ , we stop the iteration. Therefore the approximate root is 2.8284.

ii) Here the Newton—Raphson formula is

$$x_i = \frac{1}{2} \left[ x_{i-1} + \frac{91}{x_{i-1}} \right], i = 0, 1, 2, \dots$$

and  $x_0 = 10$ . The iterated values are

$$x_1 = 9.55$$

$$x_2 = 9.539398$$

$$x_3 = 9.539392$$

Since  $|x_3 - x_2| < 10^{-4}$ , we get the approximate value as 9.5393.

E7) No, because  $f'(x) = \frac{1}{3x^{2/3}}$  is not continuous at the root  $x = 0$ .

$$E8) n = \text{int} \left[ \frac{\ln(0.01) - \ln 10^{-7}}{\ln 2} \right] = \text{int} \left[ \frac{11.512925}{0.693147} \right] = 17$$

E9) Here  $g(x) = -\frac{b}{x+a}$ . The iteration

$$x_{i+1} = g(x_i) = \frac{b}{x_i + a}$$

converges to  $\alpha$  if  $|g'(x)| = \left| \frac{b}{(x+a)^2} \right| < 1$  in an interval containing  $\alpha$ . In particular we require

$$|g'(\alpha)| = \left| \frac{b}{(\alpha+a)^2} \right| < 1$$

i.e.  $(\alpha+a)^2 < |b|$ .

But we have  $\alpha + \beta = -a$  and  $\alpha\beta = b$ . Therefore we get

$$\beta^2 > |b| = |\alpha| |\beta|$$

i.e.  $|\alpha| < |\beta|$ .

E10) The iterated formula is

$$x_{i+1} = x_i - \frac{(x_i^2 - 2)}{4x_i}$$

The three successive iterations are

$$x_1 = 1.458333333$$

$$x_2 = 1.436667143$$

$$x_3 = 1.425497619$$

Then we get  $\varepsilon_3 = \frac{1}{2} \varepsilon_2$  and  $\varepsilon_2 = \frac{1}{2} \varepsilon_1$ . This shows that the sequence is not quadratically convergent, it is linearly convergent.

# UNIT 4 APPROXIMATE ROOTS OF POLYNOMIAL EQUATIONS

## Structure

- 4.1 Introduction
  - Objectives
- 4.2 Some Results on Roots of Polynomial Equations
- 4.3 Birge-Vieta Method
- 4.4 Graeffe's Root Squaring Method
- 4.5 Summary
- 4.6 Solutions/Answers

## 4.1 INTRODUCTION

In the last two units we discussed methods for finding approximate roots of the equation  $f(x) = 0$ . In this unit we restrict our attention to polynomial equations. Recall that a polynomial equation is an equation of the form  $f(x) = 0$  where  $f(x)$  is a polynomial in  $x$ . Polynomial equations arise very frequently in all branches of science especially in physical applications. For example, the stability of electrical or mechanical systems is related to the real part of one of the complex roots of a certain polynomial equation. Thus there is a need to find all roots, real and complex, of a polynomial equation. The four iteration methods we have discussed so far, applies to polynomial equations also. But you have seen that all those methods are time consuming. Thus it is necessary to find some efficient methods for obtaining roots of polynomial equations.

The sixteenth century French mathematician Francois Vieta was the pioneer to develop methods for finding approximate roots of polynomial equations. Later, several other methods were developed for solving polynomial equations. In this unit we shall discuss two simple methods : Birge-Vieta's and Graeffe's root squaring methods. To apply these methods we should have some prior knowledge of location and nature of roots of a polynomial equation. You are already familiar with some results regarding location and nature of roots from the elementary algebra course MTE-04. We shall begin this unit by listing some of the important results about the roots of polynomial equations.

### Objectives

After reading this unit you should be able to :

- apply the following methods for finding approximate roots of polynomial equations
  - i) Birge-Vieta method
  - ii) Graeffe's root squaring method.
- list the advantages of the above methods over the methods discussed in the earlier units.

## 4.2 SOME RESULTS ON ROOTS OF POLYNOMIAL EQUATIONS

The main contribution in the study of polynomial equations is due to the French mathematician Rene Descartes. The results appeared in the third part of his famous paper 'La geometric' which means 'The geometry'.

Consider a polynomial equation of degree  $n$

$$p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0 \quad \dots (1)$$

where  $a_0, a_1, \dots, a_n$  are real numbers and  $a_n \neq 0$ . You know that the roots of a polynomial equation need not be real numbers, it can be complex numbers, that is numbers of the form  $z = a + ib$  where  $a$  and  $b$  are real numbers. The following results are basic to the study of roots of polynomial equations.

**Theorem 1 :** (Fundamental Theorem of Algebra) : Let  $p(x)$  be a polynomial of degree  $n \geq 1$  given by Eqn.(1). Then  $p(x) = 0$  has at least one root; that is there exists a number  $\alpha \in \mathbb{C}$  such that  $p(\alpha) = 0$ . In fact  $p(x)$  has  $n$  complex roots which may not be distinct.

**Theorem 2 :** Let  $p(x)$  be a polynomial of degree  $n$  and  $\alpha$  is a real number. Then

$$p(x) = (x - \alpha) q_0(x) + r_0 \quad \dots (2)$$

for some polynomial  $q_0(x)$  of degree  $n - 1$  and some constant number  $r_0$ .  $q_0(x)$  and  $r_0$  are called the **quotient polynomial** and the **remainder** respectively.

In particular, if  $\alpha$  is a root of the equation  $p(x) = 0$ , then  $r_0 = 0$ ; that is  $(x - \alpha)$  divides  $p(x)$ . Then we get

$$p(x) = (x - \alpha) q_0(x)$$

How do we determine  $q_0(x)$  and  $r_0$ ? We can find them by the method of synthetic division of a polynomial  $p(x)$ . Let us now discuss the synthetic division procedure.

Consider the polynomial  $p(x)$  as given in Eqn. (1)

$$p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$$

Dividing  $p(x)$  by  $x - \alpha$  we get

$$p(x) = q_0(x) (x - \alpha) + r_0 \quad \dots (3)$$

where  $q_0(x)$  is a polynomial of degree  $n - 1$  and  $r_0$  is a constant.

Let  $q_0(x)$  be represented as

$$q_0(x) = b_n x^{n-1} + b_{n-1} x^{n-2} + \dots + b_2 x + b_1$$

(Note that for convenience we are denoting the coefficients by  $b_1, \dots, b_n$  instead of  $b_0, b_1, \dots, b_{n-1}$ ). Set  $b_0 = r_0$ . Substituting the expressions for  $q_0(x)$  and  $r_0$  in Eqn. (3) we get

$$p(x) = (x - \alpha) (b_n x^{n-1} + b_{n-1} x^{n-2} + \dots + b_2 x + b_1) + b_0 \quad \dots (4)$$

Now, to find  $b_0, b_1, \dots, b_n$  we simplify the right hand side of Eqn. (4) and compare the coefficients of  $x^i$ ,  $i = 0, 1, \dots, n$  on both sides. Note that  $p(\alpha) = b_0$ . Comparing the coefficients we get

$$\text{Coefficient of } x^n : a_n = b_n, \quad b_n = a_n$$

$$\text{Coefficient of } x^{n-1} : a_{n-1} = b_{n-1} - \alpha b_n, \quad b_{n-1} = a_{n-1} + \alpha b_n$$

$$\text{Coefficient of } x^k : a_k = b_k - \alpha b_{k+1}, \quad b_k = a_k + \alpha b_{k+1}$$

$$\text{Coefficient of } x^0 : a_0 = b_0 - \alpha, \quad b_0 = a_0 + \alpha b_1$$

It is easy to perform the calculations if we write the coefficients of  $p(x)$  on a line and perform the calculations  $b_k = a_k + \alpha b_{k+1}$  below  $a_k$  as given in the table below.

**Table 1 : Horner's table for synthetic division procedure.**

$\alpha$	$a_n$	$a_{n-1}$	$a_{n-2}$	$\dots$	$a_k$	$\dots$	$a_2$	$a_1$	$a_0$
		$\alpha b_n$	$\alpha b_{n-1}$	$\dots$	$\alpha b_{k+1}$	$\dots$	$\alpha b_3$	$\alpha b_2$	$\alpha b_1$
	$b_n$	$b_{n-1}$	$b_{n-2}$	$\dots$	$b_k$	$\dots$	$b_2$	$b_1$	$b_0 = p_0(\alpha)$

We shall illustrate this procedure with an example.

**Example 1 :** Divide the polynomial

$$p(x) = x^5 - 6x^4 + 8x^3 + 8x^2 + 4x - 40$$

by  $x - 3$  by the synthetic division method and find the remainder.

**Solution :** Here  $p(x)$  is a polynomial of degree 5. If  $a_5, a_4, a_3, a_2, a_1, a_0$  are the coefficients of  $p(x)$ , then the Horner's table in this case is

**Table 2**

$a_5$	$a_4$	$a_3$	$a_2$	$a_1$	$a_0$
1	-6	8	8	4	-40
	3	-9	-3	15	57
1	-3	-1	5	19	17
$b_5$	$b_4$	$b_3$	$b_2$	$b_1$	$b_0$

Hence the quotient polynomial  $q_0(x)$  is

$$q_0(x) = x^4 - 3x^3 - x^2 + 5x + 19$$

and the remainder is  $r_0 = b_0 = 17$ . Thus we have  $p(3) = b_0 = 17$

Do the following exercises on the same lines.

- 
- E1) Find the quotient and the remainder when  $2x^3 - 5x^2 + 3x - 1$  is divided by  $x - 2$ .
  - E2) Using synthetic division check whether  $\alpha_0 = 3$  is a root of the polynomial equation  $x^4 + x^3 - 13x^2 - x + 12 = 0$  and find the quotient polynomial.
- 

**Theorem 3 :** Suppose that  $z = a + ib$  is a root of the polynomial equation  $p(x) = 0$ . Then the conjugate of  $z$ , namely  $\bar{z} = a - ib$  is also a root of the equation  $p(x) = 0$ , i.e. complex roots occur in pairs.

We denote by  $p(-x)$  the polynomial obtained by replacing  $x$  by  $-x$  in  $p(x)$ . We next give an important Theorem due to Rene Descarte.

**Theorem 4 :** (Descarte's Rule of signs) : A polynomial equation  $p(x) = 0$  cannot have more positive roots than the number of changes in sign of its coefficients. Similarly  $p(x) = 0$  cannot have more negative roots than the number of changes in sign of the coefficients of  $p(-x)$ .

For example, let us consider the polynomial equation

$$\begin{aligned} p(x) &= x^4 - 15x^2 + 7x - 11 = 0 \\ &= 1x^4 - 15x^2 + 7x - 11 = 0 \end{aligned}$$

we count the changes in the sign of the coefficients. Going from left to right there are changes between 1 and -15, between -15 and 7 and between 7 and -11. The total number

of changes is 3 and hence it can have at most 3 positive roots. Now we consider

$$p(-x) = (-x)^4 - 15(-x)^2 + 7(-x) - 11 = 0$$

$$= x^4 - 15x^2 - 7x - 11$$

Here there is only one change between 1 and -15 and hence the equation cannot have more than one negative root.

We now give another theorem which helps us in locating the real roots.

**Theorem 5 :** Let  $p(x) = 0$  be a polynomial equation of degree  $n \geq 1$ . Let  $a$  and  $b$  be two real numbers with  $a < b$ . Suppose further that  $p(a) \neq 0$  and  $p(b) \neq 0$ . Then,

- i) if  $p(a)$  and  $p(b)$  have opposite signs, the equation  $p(x) = 0$  has an odd number of roots between  $a$  and  $b$ .
- ii) if  $p(a)$  and  $p(b)$  have like signs, then  $p(x) = 0$  either has no root or an even number of roots between  $a$  and  $b$ .

**Note :** In this theorem multiplicity of the root is taken into consideration i.e. if  $a$  is a root of multiplicity  $k$  it has to be counted  $k$  times.

As a corollary of Theorem 5, we have the following results.

**Corollary 1 :** An equation of odd degree with real coefficients has at least one real root whose sign is opposite to that of the last term.

**Corollary 2 :** An equation of even degree whose constant term has the sign opposite to that of the leading coefficient, has at least two real roots one positive and the other negative.

**Corollary 3 :** The result given in Theorem 5(i) is the generalisation of the Intermediate value theorem.

The relationship between roots and coefficients of a polynomial equation is given below.

**Theorem 6 :** Let  $\alpha_1, \alpha_2, \dots, \alpha_n$  be  $n$  roots ( $n \geq 1$ ) of the polynomial equation

$$p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0 = 0.$$

$$\text{Then } \alpha_1 + \alpha_2 + \dots + \alpha_n = \frac{-a_{n-1}}{a_n}$$

$$\alpha_1 \alpha_2 + \alpha_2 \alpha_3 + \dots + \alpha_{n-1} \alpha_n = \frac{a_{n-2}}{a_n}$$

.....

.....

$$\alpha_1 \alpha_2 \dots \alpha_n = (-1)^n \frac{a_0}{a_n}$$

Now, you can try to solve some problems using the above theorems.

E3) How many negative roots does the equation  $3x^7 + x^5 + 4x^3 + 10x - 6 = 0$  have? Also determine the number of positive roots, if any.

E4) Show that the biquadratic equation

$p(x) = x^4 + x^3 - 2x^2 + 4x - 24 = 0$  has at least two real roots one positive and the other negative.

In the next section we shall discuss one of the simple methods for solving polynomial equations.

### 4.3 BIRGE-VIETA METHOD

We shall now discuss the Birge-Vieta method for finding the real roots of a polynomial equation. This method is based on an original method due to two English mathematicians Birge and Vieta. This method is a modified form of Newton - Raphson method.

Consider now, a polynomial equation of degree n, say

$$p_n(x) = a_n x^n + \dots + a_1 x + a_0 = 0. \quad \dots (5)$$

Let  $x_0$  be an initial approximation to the root  $\alpha$ . The Newton-Raphson iterated formula for improving this approximation is

$$x_i = x_{i-1} - \frac{p_n(x_{i-1})}{p'_n(x_{i-1})}, \quad i = 1, 2, \dots \quad \dots (6)$$

To apply this formula we should be able to evaluate both  $p_n(x)$  and  $p'_n(x_i)$  at any  $x_i$ . The most natural way is to evaluate

$$p_n(x_i) = a_n x_i^n + a_{n-1} x_i^{n-1} + \dots + a_2 x_i^2 + a_1 x_i + a_0$$

$$p'_n(x_i) = n a_n x_i^{n-1} + (n-1) a_{n-1} x_i^{n-2} + \dots + 2a_2 x_i + a_1$$

However, this is the most inefficient way of evaluating a polynomial, because of the amount of computations involved and also due to the possible growth of round off errors. Thus there is a need to look for some efficient method for evaluating  $p_n(x)$  and  $p'_n(x)$ .

Let us consider the evaluation of  $p_n(x)$  and  $p'_n(x)$  at  $x_0$  using Horner's method as discussed in the previous section.

We have

$$p_n(x) = (x - x_0) q_{n-1}(x) + r_0. \quad \dots (7)$$

where

$$q_{n-1}(x) = b_n x^{n-1} + b_{n-2} x^{n-2} + \dots + b_2 x + b_1$$

and  $b_0 = p_n(x_0) = r_0 \quad \dots (8)$

We have already discussed in the previous section how to find  $b_i, i = 1, 2, \dots, n$ .

Next we shall find the derivative  $p'_n(x_0)$  using Horner's method. We divide  $q_{n-1}(x)$  by  $(x - x_0)$  using Horner's method. That is, we write

$$q_{n-1}(x) = (x - x_0) q_{n-2}(x) + r_1$$

$$q_{n-1}(x) = c_n x^{n-2} + c_{n-1} x^{n-3} + \dots + c_3 x + c_2.$$

Comparing the coefficients, we get  $c_i$  as given in the following table

Table 3

	$b_n$	$b_{n-1}$	...	$b_k$	...	$b_2$	$b_1$
$x_0$		$x_0 c_n$	...	$x_0 c_{k+1}$	...	$x_0 c_3$	$x_0 c_2$
	$c_n = b_n$	$c_{n-1}$		$c_k$		$c_2$	$c_1$

As observed in Sec. 1, we have

$$c_1 = q_{n-1}(x_0). \quad \dots (9)$$

Now, from Eqns. (7) and (8), we have

$$p_n(x) = (x - x_0) q_{n-1}(x) + p_n(x_0). \quad \dots (10)$$



Francois Vieta (1540-1603)

Differentiating both sides of Eqn. (10) w.r.t.  $x$ , we get

$$p'_n(x) = q_{n-1}(x) + (x - x_0) q'_{n-1}(x). \quad \dots (11)$$

Putting  $x = x_0$  in Eqn. (11), we get

$$p'_n(x_0) = q_{n-1}(x_0). \quad \dots (12)$$

Comparing (9) and (12), we get

$$p'_n(x_0) = q_{n-1}(x_0) = c_1$$

Hence the Newton-Raphson method (Eqn. (6)) simplifies to

$$x_i = x_{i-1} - \frac{b_0}{c_1} \quad \dots (13)$$

We summarise the evaluation of  $b_i$  and  $c_i$  in the following table.

Table 4

	$a_n$	$a_{n-1}$	...	$a_k$	...	$a_2$	$a_1$	$a_0$
$x_0$		$x_0 b_n$	...	$x_0 b_{k+1}$	...	$x_0 b_3$	$x_0 b_2$	$x_0 b_1$
$x_0$	$a_n = b_n$	$b_{n-1}$	...	$b_k$	...	$b_2$	$b_1$	$b_0 = p_n(x_0)$
	$c_n = b_n$	$c_{n-1}$	...	$c_k$	...	$c_2$	$c_1$	$c_1 = p'_n(x_0)$

Let us consider an example.

**Example 2 :** Evaluate  $p'(3)$  for the polynomial

$$p(x) = x^5 - 6x^4 + 8x^3 + 8x^2 + 4x - 40.$$

**Solution :** Here the coefficients are  $a_0 = -40$ ,  $a_1 = 4$ ,  $a_2 = 8$ ,  $a_3 = 8$ ,  $a_4 = -6$  and  $a_5 = 1$ . To compute  $b_0$ , we form the following table.

Table 5

3	1	-6	8	8	4	-40
		3	-9	-3	15	57
3	1	-3	-1	5	19	$17 = p(3) = b_0$
		3	0	-3	6	
	1	0	-1	2		$25 = p'(3) = c_1$

Therefore  $p'(3) = 25$

To get some practice, why don't you try the following exercises.

E5) Using synthetic division, show that 2 is a simple root of the equation

$$p(x) = x^4 - 2x^3 - 7x^2 + 8x + 12 = 0.$$

E6) Evaluate  $p(0.5)$  and  $p'(0.5)$  for

$$p(x) = -8x^5 + 7x^4 - 6x^3 + 5x^2 - 4x + 3$$

Now we shall illustrate why this method is more efficient than the direct method. Let us consider an example. Suppose we want to evaluate the polynomial

$$p(x) = -8x^5 + 7x^4 - 6x^3 + 5x^2 - 4x + 3$$

for any given  $x$ .



When we evaluate by direct method, we compute each power of  $x$  by multiplying with  $x$  the preceding power of  $x$  as

$$x^3 = x(x^2), x^4 = x(x^3) \text{ etc.}$$

Thus each term  $cx^k$  takes two multiplications for  $k > 1$ . Then the total number of multiplications involved in the evaluation of  $p(x)$  is  $1 + 2 + 2 + 2 + 2 = 9$ .

When we use Horner's method the total number of multiplications is 5. The number of additions in both cases are the same. This shows that less computation is involved while using Horner's method and thereby reduces the error in computation.

Let us now solve some problems using Birge-Vieta method.

**Example 3 :** Use Birge-Vieta method to find all the positive real roots, rounded off to three decimal places, of the equation

$$x^4 + 7x^3 + 24x^2 + x - 15 = 0$$

Stop the iteration whenever  $|x_{i+1} - x_i| < 0.0001$

**Solution :** We first note that the given equation

$$p_4(x) = x^4 + 7x^3 + 24x^2 + x - 15 = 0$$

is of degree 4. Therefore, by Theorem 1, this equation has 4 roots. Since there is only one change of sign in the coefficients of this equation, Descartes's rule of signs (see Theorem 4), states that the equation can have at most one positive real root.

Now let us examine whether the equation has a positive real root.

Since  $p_4(0) = -15$  and  $p_4(1) = 19$ , by Intermediate value theorem, the equation has a root lying in  $]0, 1[$ .

We take  $x_0 = 0.5$  as the initial approximation to the root. The first iteration is given by

$$\begin{aligned} x_1 &= x_0 - \frac{p_4(x_0)}{p'_4(x_0)} \\ &= 0.5 - \frac{p_4(0.5)}{p'_4(0.5)} \end{aligned}$$

Now we evaluate  $p_4(0.5)$  and  $p'_4(0.5)$  using Horner's method. The results are given in the following table.

Table 6

	1	7	24	1	-15
0.5		0.5	3.75	13.875	7.4375 -
	1	7.5	27.75	14.875	-7.5625 = $p_4(0.5)$
0.5		0.5	4.00	15.875	
	1	8.0	31.75		30.750 = $p'_4(0.5)$

$$\text{Therefore } x_1 = 0.5 - \frac{-7.5625}{30.75} = 0.7459$$

The second iteration is given by

$$x_2 = x_1 - \frac{p_4(x_1)}{p'_4(x_1)} = 0.7459 - \frac{p_4(0.7459)}{p'_4(0.7459)}$$

Using synthetic division, we form the following table of values

Table 7

	1	7	24	1	-15
0.7459		0.7459	5.7777	22.2119	17.3138
	1	7.7459	29.7777	23.2119	2.3138
0.7459		0.7459	6.3340336		26.935717
	1	8.4918	36.111701		50.146879

$$\text{Therefore } x_2 = 0.7459 - \frac{2.3132}{50.1469} = 0.6998$$

Third iteration is given by

$$x_3 = x_2 - \frac{p_4(0.6998)}{p'_4(0.6998)}$$

Table 8

	1	7	24	1	-15
0.6998		0.6998	5.3881	20.5649	15.0905
	1	7.6998	29.3881	21.5649	0.0905
0.6998		.6998	5.8778	24.6780	
	1	8.3996	35.2659	46.2429	

$$x_3 = 0.6998 - \frac{0.0905}{46.2429} = 0.6978$$

For the fourth iteration we have

$$x_4 = x_3 - \frac{p_4(0.6978)}{p'_n(0.6978)}$$

Table 9

	1	7	24	1	-15
0.6978		0.6978	5.3715248	20.495459	14.999525
	1	7.6978	29.371525	21.495459	0.000475
0.6978		.6978	5.8584497	24.583476	
		8.3956	35.229975	46.078926	

$$x_4 = 0.6978 - \frac{0.0005}{46.0789} = 0.6978$$

Since  $x_3$  and  $x_4$  are the same, we get  $|x_4 - x_3| < 0.0001$  and therefore we stop the iteration here. Hence the approximate value of the root rounded off to three decimal places is 0.698.

Next we shall illustrate how Birge-Vieta's method helps us to find all real roots of a polynomial equation.

Consider Eqn. (4)

$$p(x) = (x - \alpha) (b_n x^{n-1} + b_{n-1} x^{n-2} + \dots + b_2 x + b_1) + b_0$$

If  $\alpha$  is a root of the equation  $p(x) = 0$ , then  $p(x)$  is exactly divisible by  $x - \alpha$ , that is,  $b_0 = 0$ . In finding the approximations to the root by the Birge-Vieta method, we find that  $b_0$  approaches zero ( $b_0 \rightarrow 0$ ) as  $x_i$  approaches  $\alpha$  ( $x_i \rightarrow \alpha$ ). Hence, if  $x_n$  is taken as the final approximation to the root satisfying the criterion  $|x_n - x_{n-1}| < \epsilon$ , then to this approximation, the required quotient is

$$q_{n-1}(x) = b_n x^{n-1} + b_{n-1} x^{n-2} + \dots + b_1$$

where  $b'_i$ 's are obtained by using  $x_n$  and the Horner's method. This polynomial is called the deflated polynomial or reduced polynomial. The next root is now obtained using  $q_{n-1}(x)$  and not  $p_n(x)$ . Continuing this process, we can successively reduce the degree of the polynomial and find one real root at a time.

Let us consider an example.

**Example 4 :** Find all the roots of the polynomial equation  $p_3(x) = x^3 + x - 3 = 0$  rounded off to three decimal places. Stop the iteration whenever  $|x_{i+1} - x_i| < 0.0001$ .

**Solution :** The equation  $p_3(x) = 0$  has three roots. Since there is only one change in the sign of the coefficients, by Descartes' rule of signs the equation can have at most one positive real root. The equation has no negative real root since  $p_3(-x) = 0$  has no change of sign of coefficients. Since  $p_3(x) = 0$  is of odd degree it has at least one real root. Hence the given equation  $x^3 + x - 3 = 0$  has one positive real root and a complex pair. Since  $p(1) = -1$  and  $p(2) = 7$ , by intermediate value theorem the equation has a real root lying in the interval  $]1, 2[$ . Let us find the real root using Birge-Vieta Method. Let the initial approximation be 1.1.

First iteration

**Table 10**

	1	0	14	-3
1.1		1.1	1.21	2.431
	1	1.1	2.21	-0.569
1.1		1.1	2.42	
	1	2.2	4.63	

Therefore  $x_1 = 1.1 - \frac{-0.569}{4.63} = 1.22289$

Similarly, we obtain

$$x_2 = 1.21347$$

$$x_3 = 1.21341$$

Since  $|x_2 - x_3| < 0.0001$ , we stop the iteration here. Hence the required value of the root is 1.213, rounded off to three decimal places. Next let us obtain the deflated polynomial of  $p_3(x)$ . To get the deflated polynomial, we have to find the polynomial  $q_2(x)$  by using the final approximation  $x_3 \approx 1.213$  (see Table 11).

**Table 11**

	1	0	1	-3
1.213		1.213	1.4714	2.9978
	1	1.213	2.4714	-0.0022

Note that  $p_3(1.213) = -0.0022$ . That is, the magnitude of the error in satisfying  $p_3(x_3) = 0$  is 0.0022.

We find  $q_2(x) = x^2 + 1.213x + 2.4714 = 0$

This is a quadratic equation and its roots are given by

$$\begin{aligned}
 x &= \frac{-1.213 \pm \sqrt{(1.213)^2 - 4 \times 2.4714}}{2} \\
 &= \frac{-1.213 \pm 2.9009 i}{2} \\
 &= 0.6065 \pm 1.4505 i
 \end{aligned}$$

Hence the three roots of the equation rounded off to three decimal places are 1.213, 0.6065 + 1.4505 i and -0.6065 - 1.4505 i.

**Remark :** We now know that we can determine all the real roots of a polynomial equation using deflated polynomials. This procedure reduces the amount of computations also. But this method has certain limitations. The computations using deflated polynomials can cause unexpected errors. If the roots are determined only approximately, the coefficients of the deflated polynomials will contain some errors due to rounding off. Therefore we can expect loss of accuracy in the remaining roots. There are some ways of minimizing this error. We shall not be going into the details of these refinements.

Before going into the next section, you can try these exercises.

E7) Find an approximation to one of the roots of the equation

$$p(x) = 2x^4 - 3x^2 + 3x - 4 = 0$$

using Birge-Vieta method starting with the initial approximation  $x_0 = -2$ . Stop the iteration whenever  $|x_{i+1} - x_i| < 0.4 \times 10^{-2}$ .

E8) Find all the roots of the equation  $x^3 - 2x - 5 = 0$  using Birge-Vieta method.

E9) Find the real root rounded off to two decimal places of the equation  $x^4 - 4x^3 - 3x + 23 = 0$  lying in the interval ]2, 3[ by Birge-Vieta method.

## 4.4 GRAEFFE'S ROOT SQUARING METHOD

In the last section we have discussed a method for finding real roots of polynomial equations. Here we shall discuss a direct method for solving polynomial equations. This method was developed independently by three mathematicians Dandelin, Lobachevsky and Graeffe. But Graeffe's name is usually associated with this method. The advantage of this method is that it finds all roots of a polynomial equation simultaneously; the roots may be real and distinct, real and equal (multiple) or complex roots.

The underlying idea of the method is based on the following fact : Suppose  $\beta_1, \beta_2, \dots, \beta_n$  are the  $n$  real and distinct roots of a polynomial equation of degree  $n$  such that they are widely separated, that is,

$$|\beta_1| \gg |\beta_2| \gg |\beta_3| \gg \dots \gg |\beta_n|$$

where  $\gg$  stands for 'much greater than'. Then we can obtain the roots approximately from the coefficients of the polynomial equation as follows :

Let the polynomial equation whose roots are  $\beta_1, \beta_2, \dots, \beta_n$  be

$$a_0 + a_1x + a_2x^2 + \dots + a_nx^n = 0, \quad a_n \neq 0.$$

Using the relations between the roots and the coefficients of the polynomial as given in Sec. 4.2, we get

$$\left. \begin{aligned} \beta_1 + \beta_2 + \dots + \beta_n &= -\frac{a_{n-1}}{a_n} \\ \beta_1\beta_2 + \beta_1\beta_3 + \dots + \beta_{n-1}\beta_n &= \frac{a_{n-2}}{a_n} \\ \beta_1\beta_2\beta_3 + \dots + \beta_{n-2}\beta_{n-1}\beta_n &= -\frac{a_{n-3}}{a_n} \\ \dots & \\ \beta_1\beta_2\dots\beta_n &= (-1)^n \frac{a_0}{a_n} \end{aligned} \right\} \dots (14)$$

Since  $|\beta_1| \gg |\beta_2| \gg |\beta_3| \gg \dots \gg |\beta_n|$ , we have from (14) the approximations

$$\left. \begin{aligned} \beta_1 &\approx -\frac{a_{n-1}}{a_n} \\ \beta_1 \beta_2 &\approx \frac{a_{n-2}}{a_n} \\ \beta_1 \beta_2 \beta_3 &\approx -\frac{a_{n-3}}{a_n} \\ &\dots \\ &\dots \\ \beta_1 \beta_2 \dots \beta_n &\approx (-1)^n \frac{a_0}{a_n} \end{aligned} \right\} \dots (15)$$

These approximations can be simplified as

$$\left. \begin{aligned} |\beta_1| &\approx \frac{a_{n-1}}{a_n} \\ |\beta_2| &\approx \frac{a_{n-2}}{a_n} \frac{a_n}{a_{n-1}} \approx \frac{a_{n-2}}{a_{n-1}} \\ |\beta_3| &\approx \frac{a_{n-3}}{a_n} \frac{a_{n-1}}{a_{n-2}} \frac{a_n}{a_{n-1}} = \frac{a_{n-3}}{a_{n-2}} \\ &\dots \\ &\dots \\ |\beta_n| &\approx \frac{a_0}{a_1} \end{aligned} \right\} \dots (16)$$

So the problem now is to find from the given polynomial equation, a polynomial equation whose roots are widely separated. This can be done by the method which we shall describe now.

In the present course we shall discuss the application of the method to a polynomial equation whose roots are real and distinct.

Let  $\alpha_1, \alpha_2, \dots, \alpha_n$  be the  $n$  real and distinct roots of the polynomial equation of degree  $n$  given by

$$a_0 + a_1x + a_2x^2 + \dots + a_nx^n = 0. \quad \dots (17)$$

where  $a_0, a_1, a_2, \dots, a_{n-1}, a_n$  are real numbers and  $a_n \neq 0$ . We rewrite Eqn. (17) by collecting all even terms on one side and all odd terms on the other side, i.e.

$$a_0 + a_2x^2 + a_4x^4 + \dots = -(a_1x + a_3x^3 + a_5x^5 + \dots) \quad \dots (18)$$

Squaring both sides of Eqn. (18), we get

$$(a_0 + a_2x^2 + a_4x^4 + \dots)^2 = (a_1x + a_3x^3 + a_5x^5 + \dots)^2$$

Now we expand both the right and left hand sides and simplify by collecting the coefficients. We get

$$\begin{aligned} a_0^2 - (a_1^2 - 2a_0a_2)x^2 + (a_2^2 - 2a_1a_3 + 2a_0a_4)x^4 - \\ (a_3^2 - 2a_2a_4 + 2a_1a_5 - 2a_0a_6)x^6 + \dots + (-1)^n a_n^2 x^{2n} = 0 \quad \dots (19) \end{aligned}$$

Putting  $x^2 = -y$  in Eqn. (19), we obtain a new equation given by

$$b_0 + b_1y + b_2y^2 + \dots + b_n = 0 \quad \dots (20)$$

where

$$b_0 = a_0^2$$

$$b_1 = a_1^2 - 2a_0a_2$$

$$b_2 = a_2^2 - 2a_1a_3 + 2a_0a_4$$

$$b_n = a_n^2$$

The following table helps us to compute the coefficients  $b_0, b_1, \dots, b_n$  of Eqn. (20) directly from Eqn. (17).

Table 12

$a_0$	$a_1$	$a_2$	$a_3 \dots$	$a_n$
$a_0^2$	$a_1^2$	$a_2^2$	$a_3^2$	$a_n^2$
0	$-2a_0a_2$	$-2a_1a_3$	$-2a_2a_4$	0
0	0	$2a_0a_4$	$-2a_1a_5$	0
0	0	0	$-2a_0a_6$	0
.	.	.	.	.
.	.	.	.	.
.	.	.	.	.
$b_0$	$b_1$	$b_2$	$b_3 \dots$	$b_n$

To form Table 12 we first write the coefficients  $a_0, a_1, a_2, \dots, a_n$  as the first row. Then we form  $(n + 1)$  columns as follows.

The terms in each column alternate in sign starting with a positive sign. The first term in each column is the square of the coefficients  $a_k, k = 0, 1, 2, \dots, n$ . The second term in each column is twice the product of the nearest neighbouring coefficients, if there are any, with a negative sign; otherwise put it as zero. For example, the second term in the first column is zero and second term in the second column is  $-2a_0a_2$ . Likewise the second term of the  $(k + 1)$ th column is  $-2a_{k-1}a_{k+1}$ . The third term in the  $(k + 1)$ th column is twice the product of the next neighbouring coefficients  $a_{k-2}$  and  $a_{k+2}$ , if there are any, otherwise put it as zero. This procedure is continued until there are no coefficients available to form the cross products. Then we add all the terms in each column. The sum gives the coefficient  $b_k$  for  $k = 0, 1, 2, \dots, n$  which are listed as the last term in each column. Since the substitution  $x^2 = -y$  is used, it is easy to see that if  $\alpha_1, \alpha_2, \dots, \alpha_n$  are the  $n$  roots of Eqn. (17), then  $-\alpha_1^2, -\alpha_2^2, \dots, -\alpha_n^2$  are the roots of Eqn. (20).

Thus, starting with a given polynomial equation, we obtained another polynomial equation whose roots are the squares of the roots of the original equation with negative sign.

We repeat the procedure for Eqn. (20) and obtain another equation

$$c_0 + c_1x + \dots + c_nx^n = 0.$$

whose roots are the squares of the roots of Eqn. (20) with a negative sign i.e., they are the powers of the roots of the original equation with a negative sign. Let this procedure be repeated  $n$  times. Then, we obtain an equation

$$q_0 + q_1x + \dots + q_nx^n = 0$$

whose roots  $\gamma_1, \gamma_2, \dots, \gamma_n$  are given by

$$\gamma_i = -\alpha_i^{2^m}, i = 0, 1, 2, \dots, n. \quad (22)$$

Now, since all the roots of Eqn. (17) are real and distinct, we have

$$|\alpha_1| > |\alpha_2| > \dots > |\alpha_n|$$

Hence  $|\gamma_1| \gg |\gamma_2| \gg \dots \gg |\gamma_n|$ .

We conclude that if the roots of Eqn. (17) are distinct then for large  $m$ , the  $2^m$ th powers of the roots are widely separated.

We stop this squaring process when the cross product terms become negligible in comparison to square terms.

Since roots of Eqn. (21) are widely separated we calculate the absolute values of the roots  $\gamma_1, \gamma_2, \dots, \gamma_n$  using Eqn. (16). We have

$$|\gamma_1| = |\alpha_1^{2^m}| = \left| \frac{q_{n-1}}{q_n} \right|$$

$$|\gamma_2| = |\alpha_2^{2^m}| = \left| \frac{q_{n-2}}{q_{n-1}} \right|$$

$$|\gamma_n| = |\alpha_n^{2^m}| = \left| \frac{q_0}{q_1} \right|$$

The magnitude of the roots of the original equation are therefore given by

$$|\alpha_1| = \sqrt[2^m]{\frac{q_{n-1}}{q_n}}$$

$$|\alpha_2| = \sqrt[2^m]{\frac{q_{n-2}}{q_{n-1}}}$$

$$|\alpha_n| = \sqrt[2^m]{\frac{q_0}{q_1}}$$

This gives the magnitudes of the roots. To determine the sign of the roots, we substitute these approximations in the original equation and verify whether positive or negative value satisfies it.

We shall now illustrate this method with an example.

**Example 5 :** Find all the roots of the cubic equation  $x^3 + 15x^2 + 62x + 72 = 0$  by Graeffe's method using three squarings.

**Solution :** Let  $P_3(x) = x^3 - 15x^2 + 62x - 72 = 0$ .

The equation has no negative real roots. Let us now apply the root squaring method successively. We get the following results :

**First Squaring**

**Table 13**

$a_0$	$a_1$	$a_2$	$a_3$
-72	62	-15	1
$a_0^2 = 5184$	$a_1^2 = 3844$	$a_2^2 = 225$	$a_3^2 = 1$
0	$-2a_0a_2 = -2160$	$-2a_1a_3 = -124$	0
5184	1684	101	1
$b_0$	$b_1$	$b_2$	$b_3$

Therefore the new equation is

$$x^3 + 101x^2 + 1684x + 5184 = 0.$$

Applying the squaring method to the new equation we get the following results.

**Second Squaring**

**Table 14**

5184	1684	101	1
26873856	2835856	10201	1
0	-1047168	-3368	0
26873856	1788688	6833	1

Thus the new equation is

$$x^3 + 6833x^2 + 1788688x + 26873856 = 0.$$

For the third squaring, we have the following results.

**Third Squaring**

**Table 15**

26873856	1788688	6833	1
$7.2220414 \times 10^{14}$	$3.1994048 \times 10^{12}$	46689889	1
0	$-.3672581 \times 10^{12}$	-3577376	0
$7.2220414 \times 10^{14}$	$2.83214 \times 10^{12}$	43112513	1
$q_0$	$q_1$	$q_2$	$q_3$

Hence the new equation is

$$x^3 + 43112513x^2 + (2.83214 \times 10^{12})x + (7.2220414 \times 10^{14}) = 0$$

After three squarings, the roots  $\gamma_1$ ,  $\gamma_2$  and  $\gamma_3$  of this equation are given by

$$|\gamma_1| = \left| \frac{q_2}{q_3} \right| = 43112513$$

$$|\gamma_2| = \left| \frac{q_1}{q_2} \right| = \frac{2.83214 \times 10^{12}}{43112513}$$

$$|\gamma_3| = \left| \frac{q_0}{q_1} \right| = \frac{7.22204 \times 10^{14}}{2.83214 \times 10^{12}}$$



Hence, the roots  $\alpha_1, \alpha_2, \alpha_3$  of the original equation are

$$|\alpha_1| = \sqrt[3]{443112513} = 9.0017$$

$$|\alpha_2| = \sqrt{\frac{2.83214 \times 10^{12}}{43112513}} = 4.0011$$

$$|\alpha_3| = \sqrt{\frac{7.22204 \times 10^{14}}{2.83214 \times 10^{12}}} = 1.9990$$

Since the equation has no negative real roots, all the roots are positive. Hence the roots can be taken as 9.0017, 4.0011 and 1.9990. If the approximations are rounded to 2 decimal places, we have the roots as 9, 4 and 2. Alternately, we can substitute the approximate roots in the given equation and find their sign.

You can try these exercises now.

E10) Determine all roots of the following equations by Graeffe's root squaring method, using three squarings.

i)  $x^3 + 6x^2 - 36x + 40 = 0$

ii)  $x^3 - 2x^2 - 5x + 6 = 0$

iii)  $x^3 - 5x^2 - 17x + 20 = 0$ .

We have seen that Graeffe's root squaring method obtains all real roots simultaneously. There is considerable saving in time also. The method can be extended to find multiple and complex roots also. However the method is not efficient to find these roots. We shall not discuss these extensions.

We shall end this block by summarising what we have covered in this unit.

## 4.5 SUMMARY

In this unit we have

- discussed the following methods for finding approximate roots of polynomial equations
  - i) Birge-Vieta method
  - ii) Graeffe's root squaring method
- mentioned the advantages and disadvantages of the above methods.

## 4.6 SOLUTIONS/ANSWERS

E1) Let  $p(x) = 2x^3 - 5x^2 + 3x - 1$

Here  $a_3 = 2, a_2 = -5, a_1 = 3, a_0 = -1$  and  $\alpha = 2$ . The Horner's table in this case is as follows:

2	2	-5	3	-1
		4	-2	2
	2	-1	1	1

Hence the quotient polynomial is  $q_0(x) = 2x^2 - x + 1$  and the remainder is  $r_0 = 1$ .

- E2) Form the Horner's table in this case. From the table you can see that the last term in the 3rd row is zero. Hence 3 is a root of the equation. The quotient polynomial is  $x^3 + 4x^2 - x - 4$ .
- E3) The equation  $p(x) = 3x^7 + x^5 + 4x^3 + 10x - 6 = 0$  has no negative real root, since there are no changes in the sign of coefficients of  $p(-x)$ .

Since there is one change in the sign of coefficients of  $p(x)$  the equation can have at most one positive real root. Since the equation is of odd degree it has at least one real root which is positive.

Since,  $f(0) = -6 < 0$

$f(1) = 12 > 0$ .

the equation has a positive root lying between 0 and 1 (by IV theorem).

- E4) The given equation  $p(x) = x^4 + x^3 - 2x^2 + 4x - 24 = 0$  is of degree 4 i.e. even degree. The sign of the constant term is negative whereas the sign of the leading coefficient is positive. Therefore by corollary 2, the equation has two real roots, one positive and the other negative.

- E5) The Horner's table is as follows :

2	1	-2	-7	8	12
		2	0	-14	-12
	1	0	-7	-6	0
2		2	4	-6	
	1	2	-3	-12	

Since  $p(2) = 0$  and  $p'(2) = -12$ , 2 is a simple root.

- E6)  $p(0.5) = 1.6875$ ,  $p'(0.5) = -3.875$

- E7) The given equation is  $p(x) = 2x^4 - 3x^2 + 3x - 4 = 0$ .

The initial approximation is  $x_0 = -2$ . Then the 1st iteration is

$$x_1 = x_0 - \frac{p(x_0)}{p'(x_0)} = -2 - \frac{p(-2)}{p'(-2)}$$

$p(-2)$  and  $p'(-2)$  are given by the following table.

-2	2	0	-3	3	-4
		-4	8	-10	14
-2	2	-4	5	-7	10 = p(-2)
		-4	6	-42	
	2	-8	21	-49	

Therefore  $x_1 = -2 - \frac{10}{-49} = 1.796$

Repeating the procedure to find  $x_2$ , we have

-1.796	2	0	-3	3	-4
		-3.592	6.451	-6.197	5.742
-1.796	2	-3.592	3.451	-3.197	1.742 = p(-1.796)
		-3.592	12.902	-29.368	
	2	-7.184	16.353	-32.565 = p'(-1.796)	

Therefore  $x_2 = -1.796 - \frac{1.742}{-32.565} = 1.7425$

To find  $x_3$ , we have

-1.7425	2	0	-3	3	-4
		-3.485	6.0726	-5.3540	3.8952
-1.7425	2	-3.485	3.0726	-2.3540	.1018 = p(-1.7425)
		-3.485	12.1452	-26.5229	
	2	-6.970	15.2212	-28.8770 = p'(-1.7425)	

Therefore  $x_3 = -1.7425 + \frac{.1018}{28.8770}$   
 $= 1.7390$

Since  $|x_3 - x_2| < 0.0035 < 0.4 \times 10^{-2}$ , we conclude that 1.7390 is the approximate root.

E8) Let  $p(x) = x^3 - 2x - 5$

Since there is only one change in the sign of the coefficients of  $p(x)$ , the equation has at most one real root. The equation has no negative real root since there is no change in the sign of the coefficients of  $p(-x)$ . Also

$p(2) = -1 < 0$

and

$p(3) = 16 > 0$

Therefore a root lies in  $]2, 3[$ . Using  $x_0 = 2.5$  as an initial approximation to the root, you can show that 2.0945 is an approximation to the real root.

The deflated polynomial is given by the following table

2.0945	1	0	-2	-5
		-2.0945	4.3869	4.9994
	1	2.0945	2.3869	.0006

Therefore we get the deflated polynomial as  $p_2(x) = x^2 + 2.0945x + 2.3869 = 0$ . The roots of this equation are given by

$$x = \frac{-2.0945 \pm \sqrt{(2.0945)^2 - 4 \times 2.3869}}{2}$$

$$= -1.0473 \pm i 1.1359$$

Hence the roots are given by 2.0945,  $-1.0473 + 1.1359i$ ,  $-1.0473 - 1.13359i$ .

E9) 2.05

E10) i) The given equation is  $x^3 + 6x^2 - 36x + 40 = 0$

First squaring

40	-36	6	1
1600	1296	36	1
	-480	72	
1600	816	108	1

**Second squaring**

1600	816	108	1
2560000	665856	11664	1
	- 345600	1632	
2560000	320256	10032	1

**Third squaring**

2560000	320256	10032	1
$.65536 \times 10^{13}$	$.10256 \times 10^{12}$	$.10064 \times 10^9$	1
	- $.513638 \times 10^{11}$	- $.640512 \times 10^{12}$	1
$.65536 \times 10^{13}$	$.05120 \times 10^{12}$	$10^8$	1

Hence the new equation is

$$x^3 + 10^8 x^2 + (.05120 \times 10^{12}) x + .65536 \times 10^{13} = 0.$$

The roots  $\gamma_1, \gamma_2$  and  $\gamma_3$  of this equation are given by

$$|\gamma_1| = 10^8$$

$$|\gamma_2| = \frac{.05120 \times 10^{12}}{10^8} = 0.5120 \times 10^4$$

$$|\gamma_3| = \frac{6.5536 \times 10^{13}}{.05120 \times 10^{12}} = 128$$

Hence the roots of the original equation are given by

$$|\alpha_1| = \sqrt[8]{10^8} = 10$$

$$|\alpha_2| = \sqrt[8]{.05120 \times 10^4} = \sqrt[8]{512} \approx 2.181015$$

$$|\alpha_3| = \sqrt[8]{128} \approx 1.83.$$

Substituting the computed values in the original equation, we get that the roots are approximately -10, 2.18 and 1.83. Therefore the roots are -10, 2 and 2.

ii) Computed values of the roots are 3.014443, 1.991424 and 0.9994937.

iii) Computed values of the roots are 7.017507, -2.974432, 0.9581706.



UTTAR PRADESH  
RAJARSHI TANDON OPEN UNIVERSITY

# UGMM - 10

## Numerical Analysis

Block

# 2

### SOLUTION OF LINEAR ALGEBRAIC EQUATIONS

---

#### UNIT 5

Direct Methods 5

---

#### UNIT 6

Inverse of a Square Matrix 31

---

#### UNIT 7

Iterative Methods 49

---

#### UNIT 8

Eigenvalues and Eigenvectors 67

---

---

## Course Design Committee

---

Prof. S.R.K. Iyengar (*Editor*)  
Dept. of Mathematics  
I.I.T., Delhi

Prof. J.P. Agarwal  
Dept. of Mathematics  
I.I.T., Kharagpur

Dr. U. Anantha Krishnaiah  
Dept. of Mathematics  
KREC, Surath Kal

Prof. R.K. Jain  
Dept. of Mathematics  
I.I.T., Delhi

Prof. C. Prabhakara Rao  
Dept. of Mathematics  
REC, Warangal

**Faculty Members**  
School of Sciences, IGNOU

Prof. R.K. Bose

Dr. V.D. Madan

Dr. Poornima Mital

Dr. Manik Patwardhan

Dr. Parvin Sinclair

Dr. Sujatha Varma

---

## Block Preparation Team

---

Prof. S.R.K. Iyengar (*Editor*)  
Dept. of Mathematics  
I.I.T., Delhi

Prof. R.K. Jain  
Dept. of Mathematics  
I.I.T. Delhi

Dr. Poornima Mital  
School of Sciences  
IGNOU

---

Course Coordinator

Dr. Poornima Mital

---

## Production

---

Mr. Balakrishna Selvaraj  
Registrar (PPD)  
IGNOU

Mr. M.P. Sharma  
Joint Registrar (PPD)  
IGNOU

---

June, 1993.

© Indira Gandhi National Open University, 1993

ISBN-81-7263-411-0

*All rights reserved. No part of this work may be reproduced in any form, by mimeograph or any other means, without permission in writing from the Indira Gandhi National Open University.*

*Further information on the Indira Gandhi National Open University courses may be obtained from the University's office at Maidan Garhi, New Delhi-110 068.*

---

Reproduced and reprinted with the permission of Indira Gandhi National Open University by Dr.A.K.Singh, Registrar, U.P.R.T.Open University, Allahabad (May, 2013)  
Reprinted by : Nitin Printers, 1 Old Katra, Manmohan Park, Allahabad.

---

## BLOCK 2 SOLUTION OF LINEAR ALGEBRAIC EQUATIONS

---

In Block 1, we have discussed various numerical methods for finding approximate roots of a non-linear equation in one unknown. But there are large number of physical, biological and engineering situations in which we need to find the solution or the inverse or the eigenvalues and eigenvectors of the system of linear equations. These systems arise, both directly in modelling physical situations and indirectly in the numerical solution of other mathematical models. For instance, problems such as determining the potential in certain electrical networks, stresses in a building frame, flow rates in a hydraulic systems etc. are all reduced to solving a set of algebraic equations simultaneously. Linear algebraic systems are also involved in the optimization theory, least squares fitting of data, numerical solution of boundary value problems for ordinary and partial differential equations, statistical inference etc. In this block we shall discuss both **direct** and **iterative** methods of solving linear algebraic system of equations.

This block consists of four units.

In **Unit 5**, which is the first unit of this block, we begin with a recall of a few definitions and properties of matrices and determinants which are necessary to understand the numerical methods of solving linear system of equations. We shall then discuss some direct methods i.e., the methods which, in the absence of round-off or other errors, yield the exact solution in a finite number of elementary arithmetic operations.

In **Unit 6**, we shall discuss the method of adjoints, the Gauss-Jordan reduction method and LU decomposition method for finding the inverse of a nonsingular square matrix.

In **Unit 7**, we shall discuss two iterative methods namely, the Jacobi iteration method and the Gauss-Seidel iteration method for solving the system of linear equations. These methods start with an initial approximation and by applying a suitably chosen algorithm, lead to successively better approximations.

In **Unit 8**, which is the last unit of this block, we shall deal with the problem of computation of the absolutely largest eigenvalue or smallest eigenvalue or even all the eigenvalues of a given square matrix along with the corresponding eigenvectors. More precisely, we shall discuss the power method and the inverse power method for solving the eigenvalue problems.

## Notations and Symbols

$$A = [a_{ik}]$$

$$\det A = |A|$$

$\infty$

$\rho$

$\nu$

$\mu$

$\lambda$

$\|A\|$

$i$

Matrix with the elements  $a_{ik}$

Determinant of a square matrix A

infinity

Rho

Nu

Mu

Lambda

Norm of a matrix A

Imaginary unit,  $i^2 = -1$ .

Also see the list given in Block 1.

---

### Acknowledgements

Prof. R.K. Bose for comments on the manuscript.

Mrs. Manju Sharma and Kiran for typing the manuscript.



# UNIT 5 DIRECT METHODS

## Structure

- 5.1 Introduction
- 5.2 Preliminaries
- 5.3 Cramer's Rule
- 5.4 Direct Methods for Special Matrices
- 5.5 Gauss Elimination Method
- 5.6 LU Decomposition Method
- 5.7 Summary
- 5.8 Solutions/Answers

## 5.1 INTRODUCTION

One of the commonly occurring problems in applied mathematics is finding one or more roots of an equation  $f(x)=0$ . In most cases explicit solutions are not available and we are satisfied with being able to find one or more roots to a specified degree of accuracy. In Block 1, we have discussed various numerical methods for finding the roots of an equation  $f(x)=0$ . There we have also discussed the convergence of these methods. Another important problem of applied mathematics is to find the solution of systems of linear equations. Systems of linear equations arise in a large number of areas, both directly in modelling physical situations and indirectly in the numerical solution of other mathematical models. These applications occur in all areas of the physical, biological and engineering sciences. For instance, in physics, the problem of steady state temperature in a plate is reduced to solving linear equations. Engineering problems such as determining the potential in certain electrical networks, stresses in a building frame, flow rates in a hydraulic system etc. are all reduced to solving a set of algebraic equations simultaneously. Linear algebraic systems are also involved in the optimization theory, least squares fitting of data, numerical solution of boundary value problems for ordinary and partial differential equations, statistical inference etc. Hence, the numerical solution of systems of linear algebraic equations play a very important role.

Numerical methods for solving linear algebraic systems may be divided into two types, **direct** and **iterative**. Direct methods are those which, in the absence of round-off or other errors, yield the exact solution in a finite number of elementary arithmetic operations. Iterative methods start with an initial approximation and by applying a suitably chosen algorithm, lead to successively better approximations.

To understand the numerical methods for solving linear system of equations it is necessary to have some knowledge of the properties of matrices. You might have already studied matrices, determinants and their properties in your linear algebra course (ref. MTE-02). However, we begin with a quick recall of few definitions here. In this unit, we have also discussed some direct methods for finding the solution of system of linear algebraic equations.

### Objectives

After studying this unit, you should be able to:

- state the difference between the direct and iterative methods of solving the system of linear algebraic equations;
- obtain the solution of system of linear algebraic equations by using the direct methods such as Cramer's rule, Gauss elimination method and LU decomposition method;
- use the pivoting technique while transforming the coefficient matrix to upper or lower triangular matrix.

## 5.2 PRELIMINARIES

As we have mentioned earlier, you might be already familiar with vectors, matrices, determinants and their properties (Ref. linear algebra MTE-02). A rectangular array of (real or complex) numbers of the form

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}$$

is called a **matrix**. The numbers  $a_{11}, a_{12}, \dots, a_{mn}$  are the **elements** of the matrix. The horizontal lines are called **rows** and the vertical lines are called **columns** of the matrix. A matrix with  $m$  rows and  $n$  columns is called an  $m \times n$  matrix (read as  $m$  by  $n$  matrix). We usually denote matrices by capital letters  $A, B$  etc., or by  $(a_{jk}), (b_{ik})$  etc.

If the matrix has the same number of rows and columns, we call it a **square matrix** and the number of rows or columns is called its order. If a matrix has only one column it is a **column matrix** or column vector and if it has only one row it is a **row matrix** or row vector.

The matrices  $A = \begin{bmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{n1} \end{bmatrix} = [a_{11} \ a_{21} \ \dots \ a_{n1}]^T$  and

$B = [a_{11} \ a_{12} \ \dots \ a_{1n}]$  are respectively the column and row matrices. We give below some special square matrices  $A = (a_{ij})$  of order  $n$ .

- 1) A matrix  $A = (a_{ij})$  in which  $a_{ij} = 0$  ( $i, j = 1, 2, \dots, n$ ) is called a **null matrix** and is denoted by  $0$ .

E.g.,

$$A = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \text{ is a } 2 \times 2 \text{ null matrix.}$$

- 2) A matrix  $A$  in which all the non-diagonal elements vanish i.e.,  $a_{ij} = 0$  for  $i \neq j$  is called a **diagonal matrix**.

$$\text{E.g., } A = \begin{bmatrix} a_{11} & 0 & 0 \\ 0 & a_{22} & 0 \\ 0 & 0 & a_{33} \end{bmatrix}$$

is a  $3 \times 3$  diagonal matrix.

- 3) The **identity matrix**  $I$  is a diagonal matrix in which all the diagonal elements are equal to one. The identity matrix of order 4 is

$$I = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

- 4) A square matrix is **lower triangular** if all the elements above the main diagonal vanish i.e.,  $a_{ij} = 0$  for  $j > i$ . A lower triangular matrix of order 3 has the form

$$A = \begin{bmatrix} a_{11} & 0 & 0 \\ a_{21} & a_{22} & 0 \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

Similarly **upper triangular** matrices are matrices in which,

$a_{ij} = 0$  for  $i > j$ .

$$\text{E.g., } A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ 0 & a_{22} & a_{23} \\ 0 & 0 & a_{33} \end{bmatrix}$$

Two matrices  $A = (a_{ij})$  and  $B = (b_{ij})$  are equal iff they have the same number of rows and columns and their corresponding elements are equal, that is,  $a_{ij} = b_{ij}$  for all  $i, j$ .

You must also be familiar with the addition and multiplication of matrices.

**Addition of matrices** is defined only for matrices of same order. The sum  $C = A + B$  of two matrices  $A$  and  $B$ , is obtained by adding the corresponding elements of  $A$  and  $B$ , i.e.,  $c_{ij} = a_{ij} + b_{ij}$ .

For example, if  $A = \begin{bmatrix} -4 & 6 & 3 \\ 0 & 1 & 2 \end{bmatrix}$  and  $B = \begin{bmatrix} 5 & -1 & 0 \\ 3 & 1 & 0 \end{bmatrix}$ , then

$$A+B = \begin{bmatrix} 1 & 5 & 3 \\ 3 & 2 & 2 \end{bmatrix}$$

**Product of an  $m \times n$  matrix  $A = (a_{ij})$  and an  $n \times p$  matrix  $B = (b_{jk})$  is an  $m \times p$  matrix  $C = AB$ , whose  $(i,k)$ th entry is**

$$c_{ik} = \sum_{j=1}^n a_{ij} b_{jk} = a_{i1} b_{1k} + a_{i2} b_{2k} + \dots + a_{in} b_{nk}$$

That is, to obtain the  $(i,k)$ th element of  $AB$ , take the  $i$ th row of  $A$  and  $k$ th column of  $B$ , multiply their corresponding elements and add up all these products. For example, if

$$A = \begin{bmatrix} 2 & 3 & -1 \\ 1 & 0 & 2 \end{bmatrix} \text{ and } B = \begin{bmatrix} 1 & 1 & 2 \\ 2 & 4 & 1 \\ 1 & 2 & 1 \end{bmatrix} \text{ then } (1,2)\text{th element}$$

of  $AB$  is

$$[2 \ 3 \ -1] \begin{bmatrix} 1 \\ 4 \\ 2 \end{bmatrix} = 2 \times 1 + 3 \times 4 + (-1) \times 2 = 12$$

**Note** that two matrices  $A$  and  $B$  can be multiplied only if the number of columns of  $A$  equals the number of rows of  $B$ . In the above example the product  $BA$  is not defined.

The matrix obtained by interchanging the rows and columns of  $A$  is called the transpose of  $A$  and is denoted by  $A^T$ .

$$\text{If } A = \begin{bmatrix} 2 & 3 \\ -1 & 1 \end{bmatrix} \text{ then } A^T = \begin{bmatrix} 2 & -1 \\ 3 & 1 \end{bmatrix}$$

**Determinant** is a number associated with square matrices.

$$\text{For a } 2 \times 2 \text{ matrix } A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$\det(A) = \det \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = a_{11}a_{22} - a_{12}a_{21}$$

$$\text{For a } 3 \times 3 \text{ matrix } A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

$$\det(A) = a_{11} \det \begin{bmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{bmatrix} - a_{12} \det \begin{bmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{bmatrix} + a_{13} \det \begin{bmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix}$$

A determinant can be expanded about any row or column. The determinant of an  $n \times n$  matrix  $A = (a_{ij})$  is given by  $\det(A) = (-1)^{i+1} a_{i1} \det(A_{i1}) + (-1)^{i+2} a_{i2} \det(A_{i2}) + \dots + (-1)^{i+n} a_{in} \det(A_{in})$ , where the determinant is expanded about the  $i$ th row and  $A_{ij}$  is the  $(n-1) \times (n-1)$  matrix obtained from  $A$  by deleting the  $i$ th row

and  $j$ th column and  $i \leq i \leq n$ . Obviously, computation is simple if  $\det(A)$  is expanded along a row or column that has maximum number of zeros. This reduces the number of terms to be computed.

The following example will help you to get used to calculating determinants.

**Example 1 :**

$$\text{If } A = \begin{bmatrix} 1 & 2 & 6 \\ 5 & 4 & 1 \\ 7 & 3 & 2 \end{bmatrix}, \text{ calculate } \det(A).$$

**Solution :** Let us expand by the first row. We have

We denote  $\det(A)$  by  $|A|$  also.

$$|A_{11}| = \begin{vmatrix} 4 & 1 \\ 3 & 2 \end{vmatrix} = 4 \times 2 - 1 \times 3 = 5, \quad |A_{12}| = \begin{vmatrix} 5 & 1 \\ 7 & 2 \end{vmatrix} = 5 \times 2 - 7 \times 1 = 3,$$

$$|A_{13}| = \begin{vmatrix} 5 & 4 \\ 7 & 3 \end{vmatrix} = 5 \times 3 - 4 \times 7 = -13.$$

Thus

$$|A| = (-1)^{1+1} \times 1 \times |A_{11}| + (-1)^{1+2} \times 2 \times |A_{12}| + (-1)^{1+3} \times 6 \times |A_{13}| = 5 - 6 - 78 = -79$$

You may now try this exercise.

E1) 
$$\text{If } A = \begin{bmatrix} -3 & -2 & 0 & 2 \\ 2 & 1 & 0 & -1 \\ 1 & 0 & 1 & 2 \\ 2 & 1 & -3 & 1 \end{bmatrix}, \text{ calculate } \det(A).$$

If the determinant of a square matrix  $A$  has the value zero, then the matrix  $A$  is called a **singular matrix**, otherwise,  $A$  is called a **nonsingular matrix**.

We shall now give some more definitions.

**Definition :** The inverse of an  $n \times n$  nonsingular matrix  $A$  is an  $n \times n$  matrix  $B$  having the property

$$A B = B A = I$$

where  $I$  is an identity matrix of order  $n \times n$ .

The inverse matrix  $B$  if it exists, is denoted by  $A^{-1}$  and is unique.

**Definition :** For a matrix  $A = (a_{ij})$ , the **cofactor**  $A_{ij}$  of the element  $a_{ij}$  is given by

$$A_{ij} = (-1)^{i+j} M_{ij}$$

where  $M_{ij}$  (minor) is the determinant of the matrix of order  $(n-1) \times (n-1)$  obtained from  $A$  after deleting its  $i$ th row and the  $j$ th column.

**Definition :** The **matrix of cofactors** associated with the  $n \times n$  matrix  $A$  is an  $n \times n$  matrix  $A^c$  obtained from  $A$  by replacing each element of  $A$  by its cofactor.

**Definition :** The transpose of the cofactor matrix  $A^c$  of  $A$  is called the **adjoint** of  $A$  and is written as  $\text{adj}(A)$ . Thus

$$\text{adj}(A) = (A^c)^T$$

Let us now consider a system of  $n$  linear algebraic equations in  $n$  unknowns

$$a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1$$

$$a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2$$

$$\vdots$$

$$\vdots$$

$$\vdots$$

$$a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n$$

where the coefficients  $a_{ij}$  and the constants  $b_i$  ( $i = 1, \dots, n$ ) are real and known. This system of equations in matrix form may be written as

... (1)

$$A x = b$$

(2)

where

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

A is called the coefficient matrix and has real elements.

Our problem is to find the values  $x_i, i=1,2,\dots,n$  if they exist, satisfying Eqn. (2). Before we discuss some methods of solving the system (2), we give the following definitions.

**Definition :** A system of linear Eqns. (2) is said to be **consistent** if it has at least one solution. If no solution exists, then the system is said to be **inconsistent**.

**Definition :** The system of Eqns. (2) is said to be **homogeneous** if  $b = 0$ , that is, all the elements  $b_1, b_2, \dots, b_n$  are zero, otherwise the system is called **nonhomogeneous**.

**In this unit, we shall consider only nonhomogeneous systems.**

You also know from your linear algebra that the nonhomogeneous system of Eqns. (2) has a unique solution, if the matrix A is nonsingular. You may recall the following basic theorem on the solvability of linear systems (Ref. Theorem 4, Sec. 9.5, Unit 9, Block 3, MTE-02).

**Theorem 1 :** A nonhomogeneous system of  $n$  linear equations in  $n$  unknowns has a unique solution if and only if the coefficient matrix A is nonsingular.

If A is nonsingular, then  $A^{-1}$  exists, and the solution of system (2) can be expressed as  $x = A^{-1}b$ .

In case the matrix A is singular, then the system (2) has no solution if  $b \neq 0$  or has an infinite number of solutions if  $b = 0$ . Here we assume that A is a nonsingular matrix.

As we have already mentioned in the introduction, the methods of solution of the system (2) may be classified into two types :

- i) **Direct Methods** : which in the absence of round-off errors give the exact solution in a finite number of steps.
- ii) **Iterative Methods** : Starting with an approximate solution vector  $x^{(0)}$ , these methods generate a sequence of approximate solution vectors  $\{x^{(k)}\}$ , which converge to the exact solution vector  $x$  as the number of iterations  $k \rightarrow \infty$ . Thus iterative methods are infinite processes. Since we perform only a finite number of iterations, these methods can only find some approximation to the solution vector  $x$ . We shall discuss iterative methods later in Units 7 and 8.

In this unit we shall discuss only the direct methods. You are familiar with one such method due to the mathematician Cramer and known as **Cramer's Rule**. Let us briefly review it.

---

### 5.3 CRAMER'S RULE

---

In the system (2), let  $d = \det(A) \neq 0$  and  $b \neq 0$ . Then the solution of the system is obtained as

$$x_i = d_i/d, \quad i = 1, 2, \dots, n \quad (3)$$

where  $d_i$  is the determinant of the matrix obtained from A by replacing the  $i$ th column of A by the column vector b. Let us illustrate the method through an example.

**Example 2 :** Solve the system of equations.

$$3x_1 + x_2 + 2x_3 = 3$$

$$2x_1 - 3x_2 - x_3 = -3$$

$$x_1 + 2x_2 + x_3 = 4$$

using Cramer's rule.

**Solution :** We have,

$$d = |A| = \begin{vmatrix} 3 & 1 & 2 \\ 2 & -3 & -1 \\ 1 & 2 & 1 \end{vmatrix} = 8$$

$$d_1 = \begin{vmatrix} 3 & 1 & 2 \\ -3 & -3 & -1 \\ 4 & 2 & 1 \end{vmatrix} = 8 \text{ (first column in A is replaced by the column vector b)}$$

$$d_2 = \begin{vmatrix} 3 & 3 & 2 \\ 2 & -3 & -1 \\ 1 & 4 & 1 \end{vmatrix} = 16 \text{ (second column in A is replaced by the column vector b)}$$

$$d_3 = \begin{vmatrix} 3 & 1 & 3 \\ 2 & -3 & -3 \\ 1 & 2 & 4 \end{vmatrix} = -8 \text{ (third column in A is replaced by the column vector b)}$$

Using (3), we get the solution

$$x_1 = d_1/d = 1; x_2 = d_2/d = 2; x_3 = d_3/d = -1$$

You may now try the following exercises.

E2) Solve the system of equations

$$3x_1 + 5x_2 = 8$$

$$-x_1 + 2x_2 - x_3 = 0$$

$$3x_1 - 6x_2 + 4x_3 = 1$$

using Cramer's rule.

E3) Solve the system of equations

$$x_1 + 2x_2 - 3x_3 + x_4 = -5$$

$$x_2 + 3x_3 + x_4 = 6$$

$$2x_1 + 3x_2 + x_3 + x_4 = 4$$

$$x_1 + x_3 + x_4 = 1$$

using Cramer's rule.

While going through the example and attempting the exercises you must have observed that in Cramer's method we need to evaluate  $n+1$  determinants each of order  $n$ , where  $n$  is the number of equations. If the number of operations required to evaluate a determinant is measured in terms of multiplications only, then to evaluate a determinant of second order, i.e.,

$$\begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11} a_{22} - a_{12} a_{21}$$

we need two multiplications or  $(2-1)!$  multiplications. To evaluate a determinant of third order

$$\begin{vmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{vmatrix} = a_{11}a_{22}a_{33} - a_{11}a_{23}a_{32} - a_{12}a_{21}a_{33} + a_{12}a_{23}a_{31} + a_{13}a_{21}a_{32} - a_{13}a_{22}a_{31}$$

we need 12 multiplications or  $(3-1)3!$  multiplications. In general, to evaluate a determinant of  $n$ th order we need  $(n-1)n!$  multiplications.

Also for a system of  $n$  equations, Cramer's rule requires  $n+1$  determinants each of order  $n$  and performs  $n$  divisions to obtain  $x_i, i = 1, 2, \dots, n$ . Thus the total number of multiplications and divisions needed to solve a system of  $n$  equations, using Cramer's rule becomes

$$M = \text{total number of multiplications} + \text{total number of divisions} \\ = (n+1)(n-1)n! + n$$

In Table 1, we have given the values of  $M$  for different values of  $n$ .

Table 1

Number of equations $n$	Number of operations $M$
2	8
3	51
4	364
5	2885
6	25206
7	241927
8	2540168
9	29030409
10	359251210

From the table, you will observe that as  $n$  increases, the number of operations required for Cramer's rule increases very rapidly. For this reason, Cramer's rule is not generally used for  $n > 4$ . Hence for solving large systems, we need more efficient methods. In the next section we describe some direct methods which depend on the form of the coefficient matrix.

## 5.4 DIRECT METHODS FOR SPECIAL MATRICES

We now discuss three special forms of matrix  $A$  in Eqn. (2) for which the solution vector  $x$  can be obtained directly.

**Case 1:**  $A = D$ , where  $D$  is a diagonal matrix. In this case the system of Eqns. (2) are of the form

$$\begin{aligned} a_{11}x_1 &= b_1 \\ a_{22}x_2 &= b_2 \\ &\vdots \\ a_{nn}x_n &= b_n \end{aligned}$$

and  $\det(A) = a_{11} a_{22} \dots a_{nn}$ .

Since the matrix  $A$  is nonsingular,  $a_{ii} \neq 0$  for  $i = 1, 2, \dots, n$  and we obtain the solution as

$$x_i = b_i/a_{ii}, i = 1, 2, \dots, n.$$

**Note** that in this case we need only  $n$  divisions to obtain the solution vector.

**Case 2:**  $A = L$ , where  $L$  is a lower triangular matrix ( $a_{ij} = 0, j > i$ ). The system of Eqns. (2) is now of the form

$$\begin{aligned} a_{11}x_1 &= b_1 \\ a_{21}x_1 + a_{22}x_2 &= b_2 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 &= b_3 \end{aligned} \tag{4}$$

$$a_{n1}x_1 + a_{n2}x_2 + a_{n3}x_3 + \dots + a_{nn}x_n = b_n$$

and  $\det(A) = a_{11}a_{22}\dots a_{nn}$ .

You may notice here that the first equation of the system (4) contains only  $x_1$ , the second equation contains only  $x_1$  and  $x_2$  and so on. Hence, we find  $x_1$  from the first equation,  $x_2$  from the second equation and proceed in that order till we get  $x_n$  from the last equation.

Since the coefficient matrix  $A$  is nonsingular,  $a_{ii} \neq 0, i = 1, 2, \dots, n$ . We thus obtain

$$x_1 = b_1/a_{11}$$

$$x_2 = (b_2 - a_{21}x_1)/a_{22}$$

$$x_3 = (b_3 - a_{31}x_1 - a_{32}x_2)/a_{33}$$

$$x_n = (b_n - \sum_{j=1}^{n-1} a_{nj} x_j)/a_{nn}$$

In general, we have for any  $i$

$$x_i = (b_i - \sum_{j=1}^{i-1} a_{ij} x_j)/a_{ii}, i = 1, 2, \dots, n \tag{5}$$

For example, consider the system of equations

$$5x_1 = 5$$

$$-x_1 - 2x_2 = -7$$

$$-x_1 + 3x_2 + 2x_3 = 5$$

From the first equation we have,

$$x_1 = 1$$

From the second equation we get,

$$x_2 = \frac{-7+x_1}{-2} = 3$$

and from the third equation we have,

$$x_3 = \frac{5+x_1-3x_2}{2} = -\frac{3}{2}$$

Since the unknowns in this method are obtained in the order  $x_1, x_2, \dots, x_n$ , this method is called the **forward substitution method**.

The total number of multiplications and divisions needed to obtain the complete solution vector  $x$ , using this method is

$$M = 1 + 2 + \dots + n = n(n+1)/2.$$

**Case 3:  $A = U$** , where  $U$  is an upper triangular matrix ( $a_{ij} = 0, j < i$ ). The system (2) is now of the form

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n &= b_1 \\ a_{22}x_2 + a_{23}x_3 + \dots + a_{2n}x_n &= b_2 \\ a_{33}x_3 + \dots + a_{3n}x_n &= b_3 \end{aligned} \tag{6}$$

$$\begin{aligned} a_{n-1,n-1}x_{n-1} + a_{n-1,n}x_n &= b_{n-1} \\ a_{nn} x_n &= b_n \end{aligned}$$

$$\text{and } \det(A) = a_{11}a_{22}\dots a_{nn}.$$

You may notice here that the  $n$ th (last) equation contains only  $x_n$ , the  $(n-1)$ th equation contains  $x_n$  and  $x_{n-1}$  and so on. We can obtain  $x_n$  from the  $n$ th equation,  $x_{n-1}$  from the  $(n-1)$ th equation and proceed in that order till we get  $x_1$  from the first equation. Since the coefficient matrix  $A$  is nonsingular,  $a_{ii} \neq 0, i = 1, 2, \dots, n$  and we obtain

$$x_n = b_n/a_{nn}$$

$$x_{n-1} = (b_{n-1} - a_{n-1,n}x_n)/a_{n-1,n-1}$$

$$x_1 = (b_1 - \sum_{j=2}^n a_{1j} x_j)/a_{11}$$



or in general

$$x_i = (b_i - \sum_{j=i+1}^n a_{ij} x_j) / a_{ii}, \quad i = 1, 2, \dots, n \quad (7)$$

Since the unknowns in this method are determined in the order  $x_n, x_{n-1}, \dots, x_1$ , this method is called the **back substitution method**. The total number of multiplications and divisions needed to obtain the complete solution vector  $x$  using this method is again  $n(n+1)/2$ .

Let us consider the following example.

**Example 3 :** Solve the linear system of equations

$$\begin{aligned} 2x_1 + 3x_2 - x_3 &= 5 \\ -2x_2 - x_3 &= -7 \\ -5x_3 &= -15 \end{aligned}$$

**Solution :** From the last equation, we have

$$x_3 = 3.$$

From the second equation, we have

$$x_2 = \frac{b_2 - a_{23}x_3}{a_{22}} = \frac{(-7+3)}{(-2)} = 2.$$

Hence, from the first equation, we get

$$x_1 = \frac{b_1 - a_{12}x_2 - a_{13}x_3}{a_{11}} = \frac{(5-3 \cdot 2+3)}{2} = 1$$

You may now try the following exercises :

**E4)** Solve the system of equations

$$\begin{aligned} x_1 &= 1 \\ 2x_1 - x_2 &= 1 \\ 3x_1 - x_2 - 2x_3 &= 0 \\ 4x_1 + x_2 - 3x_3 + x_4 &= 3 \\ 5x_1 - 2x_2 - x_3 - 2x_4 + x_5 &= 1 \end{aligned}$$

using forward substitution method.

**E5)** Solve the system of equations

$$\begin{aligned} x_1 - 2x_2 + 3x_3 - 4x_4 + 5x_5 &= 3 \\ x_2 - 2x_3 + 3x_4 - 4x_5 &= -2 \\ x_3 - 2x_4 + 3x_5 &= 2 \\ x_4 - 2x_5 &= -1 \\ x_5 &= 1 \end{aligned}$$

using backward substitution method.

In the above discussion you have observed that the system of Eqns. (2) can be easily solved if the coefficient matrix  $A$  in Eqns. (2) has one of the three forms  $D, L$  or  $U$  or if it can be transformed to one of these forms. Now, you would like to know how to reduce the given matrix  $A$  into one of these three forms? One such method which transforms the matrix  $A$  to the form  $U$  is the **Gauss elimination method** which we shall describe in the next section.

## 5.5 GAUSS ELIMINATION METHOD

Gauss elimination is one of the oldest and most frequently used methods for solving systems of algebraic equations. It is attributed to the famous German mathematician, Carl Friedrich Gauss (1777 – 1855). This method is the generalization of the familiar



Gauss (1777-1855)

method of eliminating one unknown between a pair of simultaneous linear equations. You must have learnt this method in your linear algebra course (Ref. : Sec 8.4, Unit 8, Block 2, MTE-02). In this method the matrix A is reduced to the form U by using the elementary row operations which include :

- i) interchanging any two rows
- ii) multiplying (or dividing) any row by a non-zero constant
- iii) adding (or subtracting) a constant multiple of one row to another row.

The operation  $R_i + mR_j$  is an elementary row operation, that means, add to the elements of the  $i$ th row  $m$  times the corresponding elements of the  $j$ th row. The elements in the  $j$ th row remain unchanged.

If any matrix A is transformed into another matrix B by a series of elementary row operations, we say that A and B are equivalent matrices. Formally, we have the following definition.

**Definition :** A matrix B is said to be row equivalent to a matrix A, if B can be obtained from A by using a finite number of elementary row operations.

Also two linear systems  $Ax = b$  and  $A'x = b'$  are equivalent provided any solution of one is a solution of the other. Thus, if a sequence of elementary operations on  $Ax = b$  produces the new system  $A^*x = b^*$  then the systems  $Ax = b$  and  $A^*x = b^*$  are equivalent.

To understand the Gauss elimination method let us consider a system of three equations :

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 &= b_1 \\ a_{21}x_1 + a_{22}x_2 + a_{23}x_3 &= b_2 \\ a_{31}x_1 + a_{32}x_2 + a_{33}x_3 &= b_3 \end{aligned} \tag{8}$$

Let  $a_{11} \neq 0$ . In the first stage of elimination we multiply the first equation in Eqns. (8) by  $m_{21} = (-a_{21}/a_{11})$  and add to the second equation. Then multiply the first equation by  $m_{31} = (-a_{31}/a_{11})$  and add to the third equation. This eliminates  $x_1$  from the second and third equations. The new system called the first derived system then becomes

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 &= b_1 \\ a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 &= b_2^{(1)} \\ a_{32}^{(1)}x_2 + a_{33}^{(1)}x_3 &= b_3^{(1)} \end{aligned} \tag{9}$$

where,

$$\begin{aligned} a_{22}^{(1)} &= a_{22} - \frac{a_{21}}{a_{11}} a_{12} \\ a_{23}^{(1)} &= a_{23} - \frac{a_{21}}{a_{11}} a_{13} \\ b_2^{(1)} &= b_2 - \frac{a_{21}}{a_{11}} b_1 \\ a_{32}^{(1)} &= a_{32} - \frac{a_{31}}{a_{11}} a_{12} \\ a_{33}^{(1)} &= a_{33} - \frac{a_{31}}{a_{11}} a_{13} \\ b_3^{(1)} &= b_3 - \frac{a_{31}}{a_{11}} b_1 \end{aligned} \tag{10}$$

In the second stage of elimination we multiply the second equation in (9) by  $m_{32} = (-a_{32}^{(1)}/a_{22}^{(1)})$ ,  $a_{22}^{(1)} \neq 0$  and add to the third equation. This eliminates  $x_2$  from the third equation. The new system called the second derived system becomes

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + a_{13}x_3 &= b_1 \\ a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 &= b_2^{(1)} \\ a_{33}^{(2)}x_3 &= b_3^{(2)} \end{aligned}$$

where

$$\begin{aligned} a_{33}^{(2)} &= a_{33}^{(1)} - \frac{a_{32}^{(1)}}{a_{22}^{(1)}} a_{23}^{(1)} \\ b_3^{(2)} &= b_3^{(1)} - \frac{a_{32}^{(1)}}{a_{22}^{(1)}} b_2^{(1)} \end{aligned} \quad (12)$$

You may note here that the system of Eqns. (11) is an upper triangular system of the form (6) and can be solved using the back substitution method provided  $a_{33}^{(2)} \neq 0$ .

Let us illustrate the method through an example.

**Example 4 :** Solve the following linear system

$$\begin{aligned} 2x_1 + 3x_2 - x_3 &= 5 \\ 4x_1 + 4x_2 - 3x_3 &= -3 \\ -2x_1 + 3x_2 - x_3 &= 1 \end{aligned} \quad (13)$$

using Gauss elimination method.

**Solution :** To eliminate  $x_1$  from the second and third equations of the system (13)

add  $\frac{-4}{2} = -2$  times the first equation to the second equation and add  $-(-2)/2 = 1$  times the first equation to the third equation. We obtain the new system as

$$\begin{aligned} 2x_1 + 3x_2 - x_3 &= 5 \\ -2x_2 - x_3 &= -7 \\ 6x_2 - 2x_3 &= 6 \end{aligned} \quad (14)$$

In the second stage, we eliminate  $x_2$  from the third equation of system (14). Adding  $-6/(-2) = 3$  times the second equation to the third equation, we get

$$\begin{aligned} 2x_1 + 3x_2 - x_3 &= 5 \\ -2x_2 - x_3 &= -7 \\ -5x_3 &= -15 \end{aligned} \quad (15)$$

System (15) is in upper triangular form and its solution is

$$x_3 = 3, x_2 = 2, x_1 = 1.$$

You may observe that we can write the above procedure more conveniently in matrix form. Since the arithmetic operations we have performed here affect only the elements of the matrix A and the vector b, we consider the **augmented matrix** i.e.  $[A|b]$  (the matrix A augmented by the vector b) and perform the elementary row operations on the augmented matrix.

$$[A|b] = \left[ \begin{array}{ccc|c} a_{11} & a_{12} & a_{13} & b_1 \\ a_{21} & a_{22} & a_{23} & b_2 \\ a_{31} & a_{32} & a_{33} & b_3 \end{array} \right] R_2 - \frac{a_{21}}{a_{11}} R_1, R_3 - \frac{a_{31}}{a_{11}} R_1$$

$$\approx \left[ \begin{array}{ccc|c} a_{11} & a_{12} & a_{13} & b_1 \\ & a_{22}^{(1)} & a_{23}^{(1)} & b_2^{(1)} \\ & a_{32}^{(1)} & a_{33}^{(1)} & b_3^{(1)} \end{array} \right] R_3 - \frac{a_{32}^{(1)}}{a_{22}^{(1)}} R_2$$

(symbol  $\approx$  means equivalent to)

$$\approx \left[ \begin{array}{ccc|c} a_{11} & a_{12} & a_{13} & b_1 \\ & a_{22}^{(1)} & a_{23}^{(1)} & b_2^{(1)} \\ & & a_{33}^{(2)} & b_3^{(2)} \end{array} \right]$$

which is in the desired form where,  $a_{22}^{(1)}, a_{23}^{(1)}, a_{32}^{(1)}, a_{33}^{(1)}, b_2^{(1)}, b_3^{(1)}, a_{33}^{(2)}, a_3^{(2)}$  are given by Eqns. (10) and (12).

**Definition :** The diagonal elements  $a_{11}, a_{22}^{(1)}$  and  $a_{33}^{(2)}$  which are used as divisors are called **pivots**.

You might have observed here that for a linear system of order 3, the elimination was performed in  $3-1=2$  stages. In general for a system of  $n$  equations given by Eqns. (2) the elimination is performed in  $(n-1)$  stages. At the  $i$ th stage of elimination, we eliminate  $x_i$ , starting from  $(i+1)$ th row upto the  $n$ th row. Sometimes, it may happen that the elimination process stops in less than  $(n-1)$  stages. But this is possible only when no equations containing the unknowns are left or when the coefficients of all the unknowns in remaining equations become zero. Thus if the process stops at the  $r$ th stage of elimination then we get a derived system of the form

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= b_1 \\ a_{22}^{(1)}x_2 + \dots + a_{2n}^{(1)}x_n &= b_2^{(1)} \\ &\vdots \\ a_{rr}^{(r-1)}x_r + \dots + a_{rn}^{(r-1)}x_n &= b_r^{(r-1)} \\ &0 = b_{r+1}^{(r-1)} \\ &\vdots \\ &0 = b_n^{(r-1)} \end{aligned} \tag{16}$$

where  $r \leq n$  and  $a_{11} \neq 0, a_{22}^{(1)} \neq 0, \dots, a_{rr}^{(r-1)} \neq 0$ .

In the solution of system of linear equations we can thus expect two different situations

- 1)  $r = n$       2)  $r < n$ .

Let us now illustrate these situations through examples.

**Example 5 :** Solve the system of equations

$$4x_1 + x_2 + x_3 = 4$$

$$x_1 + 4x_2 - 2x_3 = 4$$

$$-x_1 + 2x_2 - 4x_3 = 2$$

using Gauss elimination method.

**Solution :** Here we have

$$[A|b] = \left[ \begin{array}{ccc|c} 4 & 1 & 1 & 4 \\ 1 & 4 & -2 & 4 \\ -1 & 2 & -4 & 2 \end{array} \right] R_2 + \frac{1}{4} R_1, R_3 + \frac{1}{4} R_1$$

$$\approx \left[ \begin{array}{ccc|c} 4 & 1 & 1 & 4 \\ 0 & \frac{15}{4} & \frac{9}{4} & 3 \\ 0 & \frac{9}{4} & \frac{15}{4} & 3 \end{array} \right] R_3 - \frac{3}{5} R_2$$

$$\approx \left[ \begin{array}{ccc|c} 4 & 1 & 1 & 4 \\ 0 & \frac{15}{4} & -\frac{9}{4} & 3 \\ 0 & 0 & -\frac{12}{5} & \frac{6}{5} \end{array} \right]$$

using back substitution method, we get

$$x_3 = -1/2; x_2 = 1/2; x_1 = 1$$

$$\text{Also, } \det(A) = 4 \times \frac{15}{4} \times \frac{(-12)}{5} = -36$$

Thus in this case we observe that  $r = n = 3$  and the given system of equations has a unique solution. Also the coefficient matrix  $A$  in this case is nonsingular. Let us look at another example.

**Example 6 :** Solve the system of equations

$$3x_1 + 2x_2 + x_3 = 3$$

$$2x_1 + x_2 + x_3 = 0$$

$$6x_1 + 2x_2 + 4x_3 = 6$$

using Gauss elimination method. Does the solution exist?

**Solution :** We have.

$$[A|b] = \left[ \begin{array}{ccc|c} 3 & 2 & 1 & 3 \\ 2 & 1 & 1 & 0 \\ 6 & 2 & 4 & 6 \end{array} \right] R_2 - \frac{2}{3} R_1, R_3 - 2R_1$$

$$\approx \left[ \begin{array}{ccc|c} 3 & 2 & 1 & 3 \\ 0 & -\frac{1}{3} & \frac{1}{3} & -2 \\ 0 & -2 & 2 & 0 \end{array} \right] R_3 - 6R_2$$

$$\approx \left[ \begin{array}{ccc|c} 3 & 2 & 1 & 3 \\ 0 & -\frac{1}{3} & -\frac{1}{3} & -2 \\ 0 & 0 & \frac{0}{3} & 12 \end{array} \right]$$

In this case you can see that  $r < n$  and elements  $b_1$ ,  $b_2^{(1)}$  and  $b_3^{(2)}$  are all non-zero.

Since we cannot determine  $x_3$  from the last equation, the system has no solution. In such a situation we say that the equations are **inconsistent**. Also note that  $\det(A) = 0$  i.e., the coefficient matrix is singular.

A system of equations is called **inconsistent** if it does not have a solution.

We now consider a situation in which not all  $b$ 's are non-zero.

**Example 7 :** Solve the system of equations

$$16x_1 + 22x_2 + 4x_3 = -2$$

$$4x_1 - 3x_2 + 2x_3 = 9$$

$$12x_1 + 25x_2 + 2x_3 = -11$$

using Gauss elimination method.

**Solution :** In this case we have

$$[A|b] = \left[ \begin{array}{ccc|c} 16 & 22 & 4 & -2 \\ 4 & -3 & 2 & 9 \\ 12 & 25 & 2 & -11 \end{array} \right] R_2 - \frac{1}{4} R_1, R_3 - \frac{3}{4} R_1$$

$$\approx \left[ \begin{array}{ccc|c} 16 & 22 & 4 & -2 \\ 0 & -\frac{17}{2} & 1 & \frac{19}{4} \\ 0 & \frac{17}{2} & -1 & -\frac{19}{2} \end{array} \right] R_3 + R_2$$

$$\approx \left[ \begin{array}{ccc|c} 16 & 22 & 4 & -2 \\ 0 & -\frac{17}{2} & 1 & \frac{19}{2} \\ 0 & 0 & 0 & 0 \end{array} \right]$$

Now in this case  $r < n$  and elements  $b_1, b_2^{(1)}$  are non-zero, but  $b_3^{(2)}$  is zero. Also the last equation is satisfied for any value of  $x_3$ . Thus, we get

$$x_3 = \text{any value}$$

$$x_2 = -\frac{2}{17} \left( \frac{19}{2} - x_3 \right)$$

$$x_1 = \frac{1}{16} \left( -2 - 22x_2 - 4x_3 \right)$$

Hence the system of equations has infinitely many solutions.

Note that in this case also  $\det(A) = 0$ .

The conclusions derived from Examples 4, 5 and 6 are true for any system of linear equations. We now summarise these conclusions as follows :

- i) If  $r = n$ , then the system of Eqns. (2) has a unique solution which can be obtained using the back substitution method. Moreover, the coefficient matrix  $A$  in this case is nonsingular.
- ii) If  $r < n$  and all the elements  $b_{r+1}^{(r-1)}, b_{r+2}^{(r-1)}, \dots, b_n^{(r-1)}$  are not zero then the system has no solution. In this case we say that the system of equations is **inconsistent**.
- iii) If  $r < n$  and all the elements  $b_{r+1}^{(r-1)}, b_{r+2}^{(r-1)}, \dots, b_n^{(r-1)}$ , if present, are zero, then the system has infinite number of solutions. In this case the system has only  $r$  linearly independent rows.

In both the cases (ii) and (iii), the matrix  $A$  is singular.

Now we estimate the number of operations (multiplication and division) in the Gauss elimination method for a system of  $n$  linear equations in  $n$  unknowns as follows :

**No. of divisions**

1st step of elimination  $(n-1)$  divisions

2nd step of elimination  $(n-2)$  divisions

.....

$(n-1)$ th step of elimination 1 divisions

$$\begin{aligned} \therefore \text{Total number of divisions} &= (n-1) + (n-2) + \dots + 1 \\ &= \sum (n-1) = \frac{n(n-1)}{2} \end{aligned}$$

**No. of multiplications**

1st step of elimination  $n(n-1)$  multiplications

2nd step of elimination  $(n-1)(n-2)$  multiplications

.....

$(n-1)$ th step of elimination  $2 \cdot 1$  multiplications

$$\begin{aligned} \therefore \text{Total number of multiplications} &= n(n-1) + (n-1)(n-1) + \dots + 2 \cdot 1 \\ &= \sum n(n-1) \\ &= \sum n^2 - \sum n \\ &= \frac{n(n+1)(2n+1)}{6} - \frac{n(n+1)}{2} \\ &= \frac{1}{3} n(n+1)(n-1) \end{aligned}$$

Also the back substitution adds  $n$  divisions (one division at each step) and the number of multiplications added are

$(n-1)$ th equation 1 multiplication

$(n-2)$ th equation 2 multiplications

.....

1st equation  $(n-1)$  multiplication

$$\therefore \text{Total multiplications} = \sum (n-1) = \frac{n(n-1)}{2}$$

$$\text{Total operations added by back substitution} = \frac{n(n-1)}{2} + n = \frac{n(n+1)}{2}$$

The sum of first  $n$  natural numbers

$$\text{is } \sum_{i=1}^n i = \frac{n(n+1)}{2} \text{ and}$$

the sum of the squares of the first  $n$  natural numbers is

$$\sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{6}$$

You can verify these results for  $n = 3$  from Eqns (9) and (11).

Thus to find the solution vector  $x$  using the Gauss elimination method, we need

$$\begin{aligned} M &= \frac{n(n-1)}{2} + \frac{1}{3}n(n^2-1) + \frac{n}{2}(n+1) \\ &= \frac{n}{6} [2n^2 + 6n - 2] \\ &= \frac{n^3}{6} + n^2 - \frac{n}{3} \end{aligned}$$

operations. For large  $n$ , we may say that the total number of operations needed is  $\frac{1}{3}n^3$  (approximately). Thus, we find that Gauss elimination method needs much lesser number of operations compared to the Cramer's rule.

You may now try a few exercises.

E6) Use Gauss elimination method to solve the system of equations

$$\begin{aligned} x_1 + 2x_2 + x_3 &= 3 \\ 3x_1 - 2x_2 - 4x_3 &= -2 \\ 2x_1 + 3x_2 - x_3 &= -6 \end{aligned}$$

E7) Use Gauss elimination method to solve the system of equations

$$\begin{aligned} 3x_1 + 18x_2 + 9x_3 &= 18 \\ 2x_1 + 3x_2 + 3x_3 &= 117 \\ 4x_1 + x_2 + 2x_3 &= 283 \end{aligned}$$

E8) Solve the system of equations

$$\begin{bmatrix} 1 & 2 & -3 & 1 \\ 0 & 1 & 3 & 1 \\ 2 & 3 & 1 & 1 \\ 1 & 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} -5 \\ 6 \\ 4 \\ 1 \end{bmatrix}$$

using Gauss elimination method.

E9) Using the Gauss elimination method show that the system of equations

$$\begin{bmatrix} 3 & 2 & -1 & -4 \\ 1 & -1 & 3 & -1 \\ 2 & 1 & -3 & 0 \\ 0 & -1 & 8 & -5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 10 \\ -4 \\ 16 \\ 3 \end{bmatrix}$$

are inconsistent.

E10) Use Gauss elimination method to solve the system of equations

$$\begin{bmatrix} 2 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

It is clear from above that you can apply Gauss elimination method to a system of equations of any order. However, what happens if one of the diagonal elements i.e. the pivots in the triangularization process vanishes? Then the method will fail. In such situations we modify the Gauss elimination method and this procedure is called pivoting.

#### Pivoting

In the elimination procedure the pivots  $a_{11}, a_{22}^{(1)}, \dots, a_{nn}^{(n-1)}$  are used as divisors. If at any stage of the elimination one of these pivots say  $a_{ii}^{(i-1)}$ , ( $a_{11}^{(0)} = a_{11}$ ), vanishes then

the elimination procedure cannot be continued further (see Example 8). Also, it may happen that the pivot  $a_{ii}^{(i-1)}$ , though not zero, may be very small in magnitude compared to the remaining elements in the  $i$ th column. Using a small number as a divisor may lead to the growth of the round-off error. In such cases the multipliers

(e.g.  $\frac{-a_{i-1,i}^{(i-2)}}{a_{ii}^{(i-1)}}$ ,  $\frac{-a_{i-2,i}^{(i-3)}}{a_{ii}^{(i-1)}}$ ) will be larger than one in magnitude. The use of large

multipliers will lead to magnification of errors both during the elimination phase and during the back substitution phase of the solution. To avoid this we rearrange the remaining rows ( $i$ th row upto  $n$ th row) so as to obtain a non-vanishing pivot or to make it the largest element in magnitude in that column. The strategy is called pivoting (see Example 9). The pivoting is of the two types; partial pivoting and complete pivoting.

### Partial Pivoting

In the first stage of elimination, the first column is searched for the largest element in magnitude and this largest element is then brought at the position of the first pivot by interchanging the first row with the row having the largest element in magnitude in the first column. In the second stage of elimination, the second column is searched for the largest element in magnitude among the  $(n-1)$  elements leaving the first element and then this largest element in magnitude is brought at the position of the second pivot by interchanging the second row with the row having the largest element in the second column. This searching and interchanging of rows is repeated in all the  $n-1$  stages of the elimination. Thus we have the following algorithm to find the pivot.

For  $i = 1, 2, \dots, n$ , find  $j$  such that

$$|a_{ji}^{(i-1)}| = \max_k |a_{ki}^{(i-1)}|, \quad i \leq k \leq n,$$

and interchange rows  $i$  and  $j$ .

### Complete Pivoting

In the first stage of elimination, we search the entire matrix  $A$  for the largest element in magnitude and bring it at the position of the first pivot. In the second stage of elimination we search the square matrix of order  $n-1$  (leaving the first row and the first-column) for the largest element in magnitude and bring it to the position of second pivot and so on. This requires at every stage of elimination not only the interchanging of rows but also interchanging of columns. Complete pivoting is much more complicated and is not often used.

In this unit, by pivoting we shall mean only partial pivoting.

Let us now understand the pivoting procedure through examples.

**Example 8 :** Solve the system of equations

$$x_1 + x_2 + x_3 = 6$$

$$3x_1 + 3x_2 + 4x_3 = 20$$

$$2x_1 + x_2 + 3x_3 = 13$$

using Gauss elimination method with partial pivoting.

**Solution :** Let us first attempt to solve the system without pivoting. We have

$$[A|b] = \left[ \begin{array}{ccc|c} 1 & 1 & 1 & 6 \\ 3 & 3 & 4 & 20 \\ 2 & 1 & 3 & 13 \end{array} \right] \quad R_2 - 3R_1, R_3 - 2R_1$$

$$\approx \left[ \begin{array}{ccc|c} 1 & 1 & 1 & 6 \\ 0 & 0 & 1 & 2 \\ 0 & -1 & 1 & 1 \end{array} \right]$$



Note that in the above matrix the second pivot has the value zero and the elimination procedure cannot be continued further unless, pivoting is used.

Let us now use the partial pivoting. In the first column 3 is the largest element. Interchanging the rows 1 and 2, we have

$$[A|b] = \left[ \begin{array}{ccc|c} 3 & 3 & 4 & 20 \\ 1 & 1 & 1 & 6 \\ 2 & 1 & 3 & 13 \end{array} \right] R_2 - \frac{1}{3} R_1, R_3 - \frac{2}{3} R_1$$

$$\approx \left[ \begin{array}{ccc|c} 3 & 3 & 4 & 20 \\ 0 & 0 & -1/3 & -2/3 \\ 0 & -1 & 1/3 & -1/3 \end{array} \right]$$

In the second column, 1 is the largest element in magnitude leaving the first element. Interchanging the second and third rows we have

$$[A|b] \approx \left[ \begin{array}{ccc|c} 3 & 3 & 4 & 20 \\ 0 & -1 & 1/3 & -1/3 \\ 0 & 0 & -1/3 & -2/3 \end{array} \right]$$

You may observe here that the resultant matrix is in triangular form and no further elimination is required. Using back substitution method, we obtain the solution

$$x_3 = 2, x_2 = 1, x_1 = 3.$$

Let us consider another example.

**Example 9 :** Solve the system of equations

$$0.0003 x_1 + 1.566 x_2 = 1.569$$

$$0.3454 x_1 - 0.436 x_2 = 3.018 \quad (17)$$

Using Gauss elimination method with and without pivoting. Assume that the numbers in arithmetic calculations are rounded to four significant digits. The exact solution of the system (17) is  $x_1 = 10, x_2 = 1$ .

**Solution :** Without Pivoting

$$m_{21} = -\frac{a_{21}}{a_{11}} = -\frac{0.3454}{0.0003} = -1151.0 \quad (\text{rounded to four places})$$

$$a_{22}^{(1)} = -0.436 - 1.566 \times 1151$$

$$= -0.436 - 1802.0 - 1802.436$$

$$= -1802.0$$

$$b_2^{(1)} = 3.018 - 1.569 \times 1151.0$$

$$= 3.018 - 1806.0$$

$$= -1803.0$$

Thus, we get the system of equations

$$0.0003 x_1 + 1.566 x_2 = 1.569$$

$$-1802.0 x_2 = -1803.0$$

which gives

$$x_2 = \frac{1803.0}{1802.0} = 1.001$$

$$x_1 = \frac{1.569 - 1.566 \times 1.001}{0.0003} = \frac{1.569 - 1.568}{0.0003}$$

$$= 3.333$$

which is highly inaccurate compared to the exact solution.

We interchange the first and second equations in (17) and get

$$0.3454 x_1 - 0.436 x_2 = 3.018$$

$$0.0003 x_1 + 1.566 x_2 = 1.569$$

we obtain

$$m_{21} = -\frac{a_{21}}{a_{11}} = -0.0009$$

$$a_{22}^{(1)} = 1.566 - 0.0009 \times (0.436) \\ = 1.566 - 0.0004$$

$$= 1.566$$

$$b_2^{(1)} = 1.569 - 3.018 \times (0.0009)$$

$$= 1.569 - 0.0027$$

$$= 1.566$$

Thus, we get the system of equations

$$0.3454 x_1 - 0.436 x_2 = 3.018$$

$$1.566 x_2 = 1.566$$

which gives

$$x_2 = 1$$

$$x_1 = \frac{3.018 + 0.436}{0.3454} = \frac{3.454}{0.3454} = 10$$

which is the exact solution.

We now make the following two remarks about pivoting.

**Remark :** If the matrix  $A$  is diagonally dominant i.e.,

$$|a_{ii}| \geq \sum_{j=1, j \neq i}^n |a_{ij}|, \text{ then no pivoting is needed. See Example 5 in which } A \text{ is}$$

diagonally dominant.

**Remark :** If exact arithmetic is used throughout the computation, pivoting is not necessary unless the pivot vanishes. However, if computation is carried upto a fixed number of digits, we get accurate results if pivoting is used.

There is another convenient way of carrying out the pivoting procedure. Instead of physically interchanging the equations all the time, the  $n$  original equations and the various changes made in them can be recorded in a systematic way. Here we use an  $n \times (n+1)$  working array or matrix which we call  $W$  and is same as our augmented matrix  $[A|b]$ . Whenever some unknown is eliminated from an equation, the changed coefficients and right side for this equation are calculated and stored in the working array  $W$  in place of the previous coefficients and right side. Also, we use an  $n$ -vector which we call  $p = (p_i)$  to keep track of which equations have already been used as pivotal equation (and therefore should not be changed any further) and which equations are still to be modified. Initially, the  $i$ th entry  $p_i$  of  $p$  contains the integer  $i$ ,  $i=1, \dots, n$  and working array  $W$  is of the form

$$W = (w_{ij}) = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} & b_1 \\ a_{21} & a_{22} & \dots & a_{2n} & b_2 \\ \dots & \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} & b_n \end{bmatrix}$$

Further, one has to be careful in the selection of the pivotal equation for each step. For each step the pivotal equation must be selected on the basis of the current state of the system under consideration i.e. without foreknowledge of the effect of the selection on later steps. For this, we calculate initially the size  $d_i$  of row  $i$  of  $A$ , for  $i=1, \dots, n$ , where  $d_i$  is the number

$$d_i = \max |a_{ij}|$$

$$1 \leq j \leq n$$

At the beginning of say  $k$ th step of elimination, we pick as pivotal equation that one from the available  $n-k$ , which has the absolutely largest coefficient of  $x_k$  relative to the size of the equation. This means that the integer  $j$  is selected between  $k$  and  $n$  for which

$$\frac{|w_{p_j k}|}{d_{p_j}} \geq \frac{|w_{i k}|}{d_i}, \forall i = p_k, \dots, p_n.$$

We can also store the multipliers in the working array  $W$  instead of storing zeros. That is, if  $p_i$  is the first pivotal equation and we use the multipliers  $m_{p_i, 1}, i=2, \dots, n$  to eliminate  $x_1$  from the remaining  $(n-1)$  positions of the first column then in the first column we can store the multipliers  $m_{p_i, 1}, i=2, \dots, n$ , instead of storing zeros.

Let us now solve the following system of linear equations by scaled partial pivoting by storing the multipliers and maintaining pivotal vector.

**Example 10 :** Solve the following system of linear equations with pivoting:

$$x_1 - x_2 + 3x_3 = 3$$

$$2x_1 + x_2 + 4x_3 = 7$$

$$3x_1 + 5x_2 - 2x_3 = 6$$

**Solution :** Here the working matrix is

$$W = \begin{bmatrix} 1 & -1 & 3 & 3 \\ 2 & 1 & 4 & 7 \\ 3 & 5 & -2 & 6 \end{bmatrix}, \mathbf{p} = [p_1, p_2, p_3]^T = [1, 2, 3]^T$$

and  $d_1 = 3, d_2 = 4$  and  $d_3 = 5$ .

Note that  $d$ 's will not change in the successive steps.

$$\text{Step 1: Now } \frac{|w_{p_1, 1}|}{d_1} = \frac{1}{3}, \frac{|w_{p_2, 1}|}{d_2} = \frac{2}{4} = \frac{1}{2}, \frac{|w_{p_3, 1}|}{d_3} = \frac{3}{5}.$$

$$\text{Since } \frac{3}{5} > \frac{1}{2}, \frac{1}{3}.$$

$$\therefore p_1 = 3, p_2 = 2 \text{ and } p_3 = 1.$$

We use the third equation to eliminate  $x_1$  from first and second equations and store corresponding multipliers instead of storing zeros in the working matrix.

$$\text{The multipliers are } m_{p_i, 1} = \frac{w_{p_i, 1}}{w_{p_1, 1}}, i = 2, 3$$

$$\therefore m_{2, 1} = \frac{w_{p_2, 1}}{w_{p_1, 1}} = \frac{w_{2, 1}}{w_{3, 1}} = \frac{2}{3}$$

$$\text{and } m_{1, 1} = \frac{w_{p_3, 1}}{w_{p_1, 1}} = \frac{w_{1, 1}}{w_{3, 1}} = \frac{1}{3}$$

After the first step the working matrix is transformed to

$$W^{(1)} = \begin{bmatrix} \boxed{\frac{1}{3}} & -\frac{8}{3} & \frac{11}{3} & 1 \\ \boxed{\frac{2}{3}} & -\frac{7}{3} & \frac{16}{3} & 3 \\ \boxed{3} & 5 & -2 & 6 \end{bmatrix}, \mathbf{p} = (p_1, p_2, p_3)^T = (3, 2, 1)^T$$

$$\text{Step 2 : } \frac{|w_{p_{2,2}}|}{dp_2} = \frac{|w_{2,2}|}{d_2} = \frac{7/3}{4} = \frac{7}{12}$$

$$\frac{|w_{p_{3,2}}|}{dp_3} = \frac{|w_{1,2}|}{d_1} = \frac{8/3}{3} = \frac{8}{9}$$

Now  $\frac{8}{9} > \frac{7}{12}$  so that we have  $\mathbf{p} = (p_1, p_2, p_3)^T = (3, 1, 2)^T$ .

$$\text{Multiplier is } m_{p_{i,2}} = \frac{w_{p_{i,2}}}{w_{p_{2,2}}}, i = 3$$

$$\Rightarrow m_{p_{3,2}} = m_{2,2} = \frac{w_{2,2}}{w_{1,2}} = \frac{-7/3}{-8/3} = \frac{7}{8}$$

That is, we use the first equation as pivotal equation to eliminate  $x_2$  from second equation and also we store the multiplier. After the second step, we have the following working matrix.

$$W^{(2)} = \begin{bmatrix} \boxed{\frac{1}{3}} & \boxed{-\frac{8}{3}} & \frac{11}{3} & 1 \\ \boxed{\frac{2}{3}} & \boxed{\frac{7}{8}} & \frac{51}{24} & \frac{17}{8} \\ \boxed{3} & 5 & -2 & 6 \end{bmatrix}, \mathbf{p} = [3, 1, 2]^T$$

In the working matrix the circled numbers denote multipliers and squared ones denote pivotal elements. Rearranging the equations (i.e., 3rd equation becomes the first equation, 1st becomes the 2nd and 2nd becomes the third) we get the reduced upper triangular system which can be solved by back substitution.

$$\begin{aligned} 3x_1 + 5x_2 - 2x_3 &= 6 \\ -\frac{8}{3}x_2 + \frac{11}{3}x_3 &= 1 \\ \frac{51}{24}x_3 &= \frac{17}{8} \end{aligned}$$

By back substitution, we get  $x_1 = 1$ ,  $x_2 = 1$  and  $x_3 = 1$ .

We now make the following two remarks.

**Remark :** We do not interchange rows in Step 1 and 2, instead we maintain a pivotal vector and use it at the end to get upper triangular system.

**Remark :** We store multipliers in the working matrix so that we can easily solve  $Ax = c$ , once we have solved  $Ax = b$ . This will be explained to you in detail in Unit 6 when we discuss the method of obtaining inverse of a matrix A.

Here is now an exercise for you.

E11) Solve the system of equations

$$0.729x + 0.81y + 0.9z = 0.6867$$

$$x + y + z = 0.8338$$

$$1.331x + 1.21y + 1.1z = 1.000$$

using Gauss elimination method with and without pivoting. Round off the numbers in arithmetic calculations to four significant digits. The exact solution of the system rounded to four significant digit is

$$x = 0.2245, y = 0.2814, z = 0.3279$$

We shall now describe the triangularization method which is also a direct method for the solution of system of equations.

In this method the matrix of coefficients of the linear system being solved is factored into the product of two triangular matrices. This method is frequently used to solve a large system of equations. We shall discuss the method in the next section.

## 5.6 LU DECOMPOSITION METHOD

Let us consider the system of Eqns. (2), where  $A$  is a non-singular matrix. We first write the matrix  $A$  as the product of a lower triangular matrix  $L$  and an upper triangular matrix  $U$  in the form

$$A = LU \quad (18)$$

or in matrix form we write

$$\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} = \begin{bmatrix} l_{11} & 0 & \dots & 0 \\ l_{21} & l_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ l_{n1} & l_{n2} & \dots & l_{nn} \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & \dots & u_{1n} \\ 0 & u_{22} & \dots & u_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & u_{nn} \end{bmatrix} \quad (19)$$

The left side matrix  $A$  has  $n^2$  elements, whereas  $L$  and  $U$  have  $1+2+\dots+n = n(n+1)/2$  elements each. Thus, we have  $n^2+n$  unknowns in  $L$  and  $U$  which are to be determined. On comparing the corresponding elements on two sides in Eqn. (19), we get  $n^2$  equations in  $n^2+n$  unknowns and hence  $n$  unknowns are undetermined. Thus, we get a solution in terms of these  $n$  unknowns i.e., we get a  $n$  parameter family of solutions. In order to obtain a unique solution we either take all the diagonal elements of  $L$  as 1, or all the diagonal elements of  $U$  as 1.

For  $u_{ii} = 1, i = 1, 2, \dots, n$ , the method is called the **Crout LU decomposition method**. For  $l_{ii} = 1, i = 1, 2, \dots, n$  we have **Doolittle LU decomposition method**. Usually Crout's LU decomposition method is used unless it is specifically mentioned. We shall now explain the method for  $n = 3$  with  $u_{ii} = 1, i = 1, 2, 3$ . We have

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} \begin{bmatrix} 1 & u_{12} & u_{13} \\ 0 & 1 & u_{23} \\ 0 & 0 & 1 \end{bmatrix}$$

or

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} l_{11} & l_{11}u_{12} & l_{11}u_{13} \\ l_{21} & l_{21}u_{12}+l_{22} & l_{21}u_{13}+l_{22}u_{23} \\ l_{31} & l_{31}u_{12}+l_{32} & l_{31}u_{13}+l_{32}u_{23}+l_{33} \end{bmatrix}$$

On comparing the elements of the first column, we obtain

$$l_{11} = a_{11}, l_{21} = a_{21}, l_{31} = a_{31} \quad (20)$$

i.e., the first column of  $L$  is determined.

On comparing the remaining elements of the first row, we get

$$l_{11}u_{12} = a_{12}; l_{11}u_{13} = a_{13}$$

which gives

$$u_{12} = a_{12}/l_{11}; u_{13} = a_{13}/l_{11} \quad (21)$$

Hence the first row of  $U$  is determined.

On comparing the elements of the second column, we get

$$l_{21}u_{12} + l_{22} = a_{22}$$

$$l_{31}u_{12} + l_{32} = a_{32}$$

which gives

$$\begin{bmatrix} l_{22} = a_{22} - l_{21}u_{12} \\ l_{32} = a_{32} - l_{31}u_{12} \end{bmatrix} \quad (22)$$

Now the second column of  $L$  is determined.

On comparing the elements of the second row, we get

$$l_{21}u_{13} + l_{22}u_{23} = a_{23}$$

which gives  $u_{23} = (a_{23} - l_{21}u_{13})/l_{22}$  (23)

and the second row of  $U$  is determined.

On comparing the elements of the third column, we get

$$l_{31}u_{13} + l_{32}u_{23} + l_{33} = a_{33}$$

which gives  $l_{33} = a_{33} - l_{31}u_{13} - l_{32}u_{23}$  (24)

You must have observed that in this method, we alternate between getting a column of  $L$  and a row of  $U$  in that order. If instead of  $u_{ii} = 1, i = 1, 2, \dots, n$ , we take  $l_{ii} = 1, i = 1, 2, \dots, n$ , then we alternate between getting a row of  $U$  and a column of  $L$  in that order.

Thus, it is clear from Eqns. (20) – (24) that we can determine all the elements of  $L$  and  $U$  provided the nonsingular matrix  $A$  is such that

$$a_{11} \neq 0, \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \neq 0.$$

Similarly, for the general system of Eqns. (2), we obtain the elements of  $L$  and  $U$  using the relations

$$l_{ij} = a_{ij} - \sum_{k=1}^{j-1} l_{ik}u_{kj}, \quad i \geq j$$

$$u_{ij} = (a_{ij} - \sum_{k=1}^{j-1} l_{ik}u_{kj})/l_{ii}, \quad i \geq j$$

$$u_{ii} = 1$$

$$\text{Also, } \det(A) = l_{11}l_{22}\dots l_{nn}.$$

Thus we can say that every nonsingular matrix  $A$  can be written as the product of a lower triangular matrix and an upper triangular matrix if all the principal minors of  $A$  are nonsingular, i.e. if

$$a_{11} \neq 0, \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \neq 0, \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \neq 0, \dots, |A| \neq 0.$$

Once we have obtained the elements of the matrices  $L$  and  $U$ , we write the system of equations

$$A x = b \quad (25)$$

in the form

$$L U x = b \quad (26)$$

The system (26) may be further written as the following two systems

$$U x = y \quad (27)$$

$$L y = b \quad (28)$$

Now, we first solve the system (28), i.e.,

$$L y = b,$$

using the forward substitution method to obtain the solution vector  $y$ . Then using this  $y$ , we solve the system (27), i.e.,

$$U x = y,$$

using the backward substitution method to obtain the solution vector  $x$ .

The number of operations for this method remains the same as that in the Gauss-elimination method.

We now illustrate this method through an example.

**Example 11 :** Use the LU decomposition method to solve the system of equations

$$x_1 + x_2 + x_3 = 1$$

$$4x_1 + 3x_2 - x_3 = 6$$

$$3x_1 + 5x_2 + 3x_3 = 4$$

**Solution :** Using  $l_{ii} = 1, i = 1, 2, 3$ , we have

$$\begin{bmatrix} 1 & 1 & 1 \\ 4 & 3 & -1 \\ 3 & 5 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{bmatrix}$$

$$= \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ l_{21}u_{11} & l_{21}u_{12} + u_{22} & l_{21}u_{13} + u_{23} \\ l_{31}u_{11} & l_{31}u_{12} + l_{32}u_{22} & l_{31}u_{13} + l_{32}u_{23} + u_{33} \end{bmatrix}$$

On comparing the elements of row and column alternately, on both sides, we obtain

first row :  $u_{11} = 1, u_{12} = 1, u_{13} = 1$

first column :  $l_{21} = 4, l_{31} = 3$

second row :  $u_{22} = -1, u_{23} = -5$

second column :  $l_{32} = -2$

third row :  $u_{33} = -10$

Thus, we have

$$L = \begin{bmatrix} 1 & 0 & 0 \\ 4 & 1 & 0 \\ 3 & -2 & 1 \end{bmatrix}; U = \begin{bmatrix} 1 & 1 & 1 \\ 0 & -1 & -5 \\ 0 & 0 & -10 \end{bmatrix}$$

Now from the system

$$L y = b$$

or

$$\begin{bmatrix} 1 & 0 & 0 \\ 4 & 1 & 0 \\ 3 & -2 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 6 \\ 4 \end{bmatrix}$$

we get

$$y_1 = 1, y_2 = 2, y_3 = 5$$

and from the system

$$U x = y$$

or

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & -1 & -5 \\ 0 & 0 & -10 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 5 \end{bmatrix}$$

we get

$$x_3 = -1/2, x_2 = 1/2, x_1 = 1.$$

You may now try the following exercises :

E12) Use the LU decomposition method with  $u_{ii} = 1, i = 1, 2, 3$  to solve the system of equations given in Example 11.

E13) Use the LU decomposition method with  $l_{ii} = 1, i = 1, 2, 3$  to solve the system of equations given in E7.

E14) Use LU decomposition method to solve the system of equations given in E10.

## 5.7 SUMMARY

In this unit we have covered the following:

- 1) For a system of  $n$  equations

$$Ax = b$$

(see Eqn. (2))

in  $n$  unknowns, where  $A$  is an  $n \times n$  non-singular matrix, the methods of finding the solution vector  $x$  may be broadly classified into two types: (i) direct methods and (ii) iterative methods

- 2) **Direct methods** produces the exact solution in a finite number of steps provided there are no round-off errors. **Cramer's rule** is one such method. This method gives the solution vector as

$$x_i = \frac{d_i}{d} \quad i = 1, 2, \dots, n$$

where  $d = |A|$  and  $d_i$  is the determinant of the matrix obtained from  $A$  by replacing the  $i$ th column of  $A$  by the column vector  $b$ . Total number of operations required for Cramer's rule in solving a system of  $n$  equations are

$$M = (n+1)(n-1)n! + n$$

Since the number  $M$  increases very rapidly, Cramer's rule is not used for  $n > 4$ .

- 3) For larger systems, direct methods become more efficient if the coefficient matrix  $A$  is in one of the forms  $D$  (diagonal),  $L$  (lower triangular) or  $U$  (upper triangular).
- 4) **Gauss elimination** method is another direct method for solving large systems ( $n > 4$ ). In this method the coefficient matrix  $A$  is reduced to the form  $U$  by using the elementary row operations. The solution vector  $x$  is then obtained by using the back substitution method. For large  $n$ , the total number of operations required in Gauss elimination method are  $\frac{1}{3}n^3$  (approximately).
- 5) In Gauss elimination method if at any stage of the elimination any of the pivots vanishes or become small in magnitude, elimination procedure cannot be continued further. In such cases pivoting is used to obtain the solution vector  $x$ .
- 6) Every nonsingular matrix  $A$  can be written as the product of a lower triangular matrix and an upper triangular matrix, by the LU decomposition method, if all the principal minors of  $A$  are nonsingular. Thus, LU decomposition method, which is a modification of the Gauss elimination method can be used to obtain the solution vector  $x$ .

## 5.8 SOLUTIONS/ANSWERS

E1)  $\det(A) = 8$

E2)  $d = 11, d_1 = 11, d_2 = 11, d_3 = 11$

$$x_1 = x_2 = x_3 = 1$$

E3)  $d = 20, d_1 = 0, d_2 = 20, d_3 = 40, d_4 = -20$

$$x_1 = 0, x_2 = 1, x_3 = 2, x_4 = -1$$

E4)  $x_1 = x_2 = x_3 = x_4 = x_5 = 1$

E5)  $x_5 = x_4 = x_3 = x_2 = x_1 = 1$

E6) 
$$\left[ \begin{array}{ccc|c} 1 & 2 & 1 & 3 \\ 3 & -2 & 4 & -2 \\ 2 & 3 & -1 & -6 \end{array} \right] R_2 - 3R_1, R_3 - 2R_1$$

$$\approx \left[ \begin{array}{ccc|c} 1 & 2 & 1 & 3 \\ 0 & -8 & -7 & -11 \\ 0 & -1 & -3 & -12 \end{array} \right] R_3 - \frac{1}{8}R_2$$



$$\text{Final derived system: } \left[ \begin{array}{ccc|c} 1 & 2 & 1 & 3 \\ 0 & -8 & -7 & -11 \\ 0 & 0 & -\frac{17}{8} & -\frac{85}{8} \end{array} \right]$$

$$x_3 = 5, x_2 = -3, x_1 = 4$$

$$\text{E7) Final derived system: } \left[ \begin{array}{ccc|c} 3 & 18 & 9 & 18 \\ & -9 & -3 & 105 \\ & & -\frac{7}{3} & -\frac{28}{3} \end{array} \right]$$

$$x_3 = 4, x_2 = -13, x_1 = 72$$

E8) Final derived system :

$$\left[ \begin{array}{cccc|c} 1 & 2 & -3 & 1 & -5 \\ 0 & 1 & 3 & 1 & 6 \\ 0 & 0 & 10 & 0 & 20 \\ 0 & 0 & 0 & 2 & -2 \end{array} \right]$$

$$x_4 = -1, x_3 = 2, x_2 = 1, x_1 = 0$$

E9) Final derived system :

$$\left[ \begin{array}{cccc|c} 3 & 2 & -1 & -4 & 10 \\ 0 & -5/3 & 10/3 & 1/3 & -22/3 \\ 0 & 0 & -3 & 13/5 & 54/5 \\ 0 & 0 & 0 & 0 & 29 \end{array} \right]$$

We cannot determine  $x_4$  from the last equation.

E10) Final derived system :

$$\left[ \begin{array}{ccccc|c} 2 & -1 & 0 & 0 & 0 & 1 \\ 0 & 3/2 & -1 & 0 & 0 & 1/2 \\ 0 & 0 & 4/3 & -1 & 0 & 1/3 \\ 0 & 0 & 0 & 5/4 & -1 & 1/4 \\ 0 & 0 & 0 & 0 & 6/5 & 6/5 \end{array} \right]$$

$$x_5 = x_4 = x_3 = x_2 = x_1 = 1$$

E11) Solution without pivoting :

Using  $m_{21} = 1.372$

$$m_{31} = 1.826 \text{ and } m_{32} = 2.423$$

The final derived system is

$$\left[ \begin{array}{ccc|c} 0.7290 & 0.8100 & 0.9000 & 0.6867 \\ 0.0 & -0.1110 & -0.2350 & -0.1084 \\ 0.0 & 0.0 & 0.02640 & -0.0087 \end{array} \right]$$

The solution is

$$x = 0.2251, y = 0.2790, z = 0.3295$$

Solution with pivoting:

Interchanging first and the third row and using

$$m_{21} = 0.7513,$$

$$m_{31} = 0.5477$$

and  $m_{32} = 0.6171$ 

the final derived system is

$$\left[ \begin{array}{ccc|c} 1.331 & 1.210 & 1.100 & 1.000 \\ 0.0 & 0.1473 & 0.2975 & 0.1390 \\ 0.0 & 0.0 & -0.0100 & -0.003280 \end{array} \right]$$

The solution is  $x = 0.2246$ ,  $y = 0.2812$ ,  $z = 0.3280$ .

$$\text{E12) } L = \begin{bmatrix} 1 & 0 & 0 \\ 4 & -1 & 0 \\ 3 & 2 & -10 \end{bmatrix}; U = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 5 \\ 0 & 0 & 1 \end{bmatrix}$$

$$y = \left[ 1 \quad -2 \quad -\frac{1}{2} \right]^T; x = \left[ 1 \quad \frac{1}{2} \quad -\frac{1}{2} \right]^T$$

$$\text{E13) } L = \begin{bmatrix} 1 & 0 & 0 \\ 2/3 & 1 & 0 \\ 4/3 & 23/9 & 1 \end{bmatrix}; U = \begin{bmatrix} 3 & 18 & 9 \\ 0 & -9 & -3 \\ 0 & 0 & -7/3 \end{bmatrix}$$

$$y = \left[ 18 \quad 105 \quad -\frac{28}{3} \right]^T; x = \left[ 72 \quad -13 \quad 4 \right]^T$$

$$\text{E14) } L = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ -1 & 3/2 & 0 & 0 & 0 \\ 0 & -1 & 4/3 & 0 & 0 \\ 0 & 0 & -1 & 5/4 & 0 \\ 0 & 0 & 0 & -1 & 6/5 \end{bmatrix};$$

$$U = \begin{bmatrix} 1 & -1/2 & 0 & 0 & 0 \\ 0 & 1 & -2/3 & 0 & 0 \\ 0 & 0 & 1 & -3/4 & 0 \\ 0 & 0 & 0 & 1 & -4/5 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$y = \left[ \frac{1}{2} \quad \frac{1}{3} \quad \frac{1}{4} \quad \frac{1}{5} \quad 1 \right]^T; x = \left[ 1 \quad 1 \quad 1 \quad 1 \quad 1 \right]^T$$

# UNIT 6 INVERSE OF A SQUARE MATRIX

## Structure

- 6.1 Introduction
- 6.2 The Method of Adjoins
- 6.3 The Gauss-Jordan Reduction Method
- 6.4 LU Decomposition Method
- 6.5 Summary
- 6.6 Solutions/Answers

## 6.1 INTRODUCTION

In the previous unit, you have studied the Gauss elimination and LU decomposition methods for solving systems of algebraic equations  $Ax = b$ , when  $A$  is a  $n \times n$  nonsingular matrix. Matrix inversion is another problem associated with the problem of finding solutions of a linear system. If the inverse matrix  $A^{-1}$  of the coefficient matrix  $A$  is known then the solution vector  $x$  can be obtained from  $x = A^{-1}b$ . In general, inversion of matrices for solving system of equations should be avoided whenever possible. This is because, it involves greater amount of work and also it is difficult to obtain the inverse accurately in many problems. However, there are two cases in which the explicit computation of the inverse is desirable. **Firstly**, when several systems of equations, having the same coefficient matrix  $A$  but different right hand side  $b$ , have to be solved. Then computations are reduced if we first find the inverse matrix and then find the solution. **Secondly**, when the elements of  $A^{-1}$  themselves have some special physical significance. For instance, in the statistical treatment of the fitting of a function to observational data by the method of least squares, the elements of  $A^{-1}$  give information about the kind and magnitude of errors in the data.

In this unit, we shall study a few important methods for finding the inverse of a nonsingular square matrix.

### Objectives

After studying this unit, you should be able to :

- obtain the inverse by adjoint method for  $n < 4$ ;
- obtain the inverse by the Gauss-Jordan and LU decomposition methods;
- obtain the solution of a system of linear equations using the inverse method.

## 6.2 THE METHOD OF ADJOINTS

You already know that the transpose of the matrix of the cofactors of elements of  $A$  is called the **adjoint** matrix and is denoted by  $\text{adj}(A)$  (Ref. Unit 9, Block 3 of MTE-02, Linear Algebra).

Formally, we have the following definition.

**Definition :** The transpose of the cofactor matrix  $A^c$  of  $A$  is called the **adjoint** of  $A$  and is written as  $\text{adj}(A)$ .

Thus,

$$\text{adj}(A) = (A^c)^T$$

The inverse of a matrix can be calculated using the adjoint of a matrix.

We obtain the inverse matrix  $A^{-1}$  of  $A$  from

$$A^{-1} = \frac{1}{\det(A)} \text{adj}(A) \quad (1)$$

This method of finding the inverse of a matrix is called the **method of adjoints**.

Note that  $\det(A)$  in Eqn. (1) must not be zero and therefore the matrix  $A$  must be nonsingular.

We shall not be going into the details of the method here. We shall only illustrate it through examples.

**Example 1 :** Find  $A^{-1}$  for the matrix

$$A = \begin{bmatrix} 5 & 8 & 1 \\ 0 & 2 & 1 \\ 4 & 3 & -1 \end{bmatrix}$$

and solve the system of equations

$$A \mathbf{x} = \mathbf{b}$$

for

$$\text{i) } \mathbf{b} = \begin{bmatrix} 2 \\ -1 \\ 3 \end{bmatrix}, \quad \text{ii) } \mathbf{b} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \text{iii) } \mathbf{b} = \begin{bmatrix} 1 \\ -2 \\ 3 \end{bmatrix}$$

**Solution :** Since  $\det(A) = -1 \neq 0$ , the inverse of  $A$  exists. We obtain the cofactor matrix  $A^c$  from  $A$  by replacing each element of  $A$  by its cofactor as follows :

$$A^c = \begin{bmatrix} -5 & 4 & -8 \\ 11 & -9 & 17 \\ 6 & -5 & 10 \end{bmatrix}$$

$$\therefore \text{adj}(A) = (A^c)^T = \begin{bmatrix} -5 & 11 & 6 \\ 4 & -9 & -5 \\ -8 & 17 & 10 \end{bmatrix}$$

$$\text{Now } A^{-1} = \frac{1}{\det(A)} \text{adj}(A)$$

$$\therefore A^{-1} = - \begin{bmatrix} -5 & 11 & 6 \\ 4 & -9 & -5 \\ -8 & 17 & 10 \end{bmatrix} = \begin{bmatrix} 5 & -11 & -6 \\ -4 & 9 & 5 \\ 8 & -17 & -10 \end{bmatrix}$$

Also the solution of the given systems of equations are

$$\text{i) } \mathbf{x} = A^{-1}\mathbf{b} = \begin{bmatrix} 5 & -11 & -6 \\ -4 & 9 & 5 \\ 8 & -17 & -10 \end{bmatrix} \begin{bmatrix} 2 \\ -1 \\ 3 \end{bmatrix} = \begin{bmatrix} 3 \\ -2 \\ 3 \end{bmatrix}$$

$$\text{ii) } \mathbf{x} = A^{-1}\mathbf{b} = \begin{bmatrix} 5 & -11 & -6 \\ -4 & 9 & 5 \\ 8 & -17 & -10 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 5 \\ -4 \\ 8 \end{bmatrix}$$

$$\text{iii) } \mathbf{x} = A^{-1}\mathbf{b} = \begin{bmatrix} 5 & -11 & -6 \\ -4 & 9 & 5 \\ 8 & -17 & -10 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \\ 3 \end{bmatrix} = \begin{bmatrix} 9 \\ -7 \\ 12 \end{bmatrix}$$

We now take up an example in which the given matrix  $A$  is lower triangular and we shall show that its inverse is also a lower triangular matrix.

**Example 2 :** Find  $A^{-1}$  for the matrix

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 3 & 0 \\ 4 & 5 & 6 \end{bmatrix}$$

**Solution :** We have

$\det(A) = 18 \neq 0$ . Thus  $A^{-1}$  exists.

Now

$$A^c = \begin{bmatrix} 18 & -12 & -2 \\ 0 & 6 & -5 \\ 0 & 0 & 3 \end{bmatrix}$$

$$\therefore A^{-1} = \frac{(A^c)^T}{\det(A)} = \frac{1}{18} \begin{bmatrix} 18 & 0 & 0 \\ -12 & 6 & 0 \\ -2 & -5 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ -2/3 & 1/3 & 0 \\ -1/9 & -5/18 & 1/6 \end{bmatrix}$$

Thus,  $A^{-1}$  is again a lower triangular matrix. Similarly, we can illustrate that the inverse of an upper triangular matrix is again upper triangular.

**Example 3 :** Find  $A^{-1}$  for the matrix

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \end{bmatrix}$$

**Solution :** Since,  $\det(A) = 24 \neq 0$ ,  $A^{-1}$  exists.

We obtain

$$A^c = \begin{bmatrix} 24 & 0 & 0 \\ -12 & 6 & 0 \\ -2 & -5 & 4 \end{bmatrix}$$

$$\therefore A^{-1} = \frac{1}{24} \begin{bmatrix} 24 & -12 & -2 \\ 0 & 6 & -5 \\ 0 & 0 & 4 \end{bmatrix} = \begin{bmatrix} 1 & -1/2 & -1/12 \\ 0 & 1/4 & -5/24 \\ 0 & 0 & 1/6 \end{bmatrix}$$

which is again an upper triangular matrix.

You may now try the following exercises.

E1) Solve the system of equations

$$3x_1 + x_2 + 2x_3 = 3$$

$$2x_1 - x_2 - x_3 = 1$$

$$x_1 - 2x_2 + x_3 = -4$$

using the method of adjoints.

E2) Solve the system of equations

$$\begin{bmatrix} 2 & 3 & 4 & 1 \\ 1 & 2 & 0 & 1 \\ 2 & 3 & 1 & -1 \\ 1 & -2 & -1 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \\ 1 \\ 5 \end{bmatrix}$$

using the method of adjoints.

The method of adjoints provides a systematic procedure to obtain the inverse of a given matrix and for solving systems of linear equations. To obtain the inverse of an  $n \times n$  matrix, using this method, we need to evaluate one determinant of order  $n$ ,  $n$  determinants each of order  $n-1$  and perform  $n^2$  divisions. In addition, if this method

is used for solving a linear system we also need matrix multiplication. The number of operations (multiplications and divisions) needed, for using this method, increases very rapidly as  $n$  increases. For this reason, this method is not used when  $n > 4$ .

For large  $n$ , there are methods which are efficient and are frequently used for finding the inverse of a matrix and solving linear systems. We shall now discuss these methods.

### 6.3 THE GAUSS-JORDAN REDUCTION METHOD

This method is a variation of the Gauss elimination method. In the Gauss elimination method, using elementary row operations, we transform the matrix  $A$  to an upper triangular matrix  $U$  and obtain the solution by using back substitution method. In Gauss-Jordan reduction not only the elements below the diagonal but also the elements above the diagonal of  $A$  are made zero at the same time. In other words, we transform the matrix  $A$  to a diagonal matrix  $D$ . This diagonal matrix may then be reduced to an identity matrix by dividing each row by its pivot element. Alternately, the diagonal elements can also be made unity at the same time when the reduction is performed. This transforms the coefficient matrix into an identity matrix. Thus, on completion of the Gauss-Jordan method, we have

$$[A|b] \xrightarrow{\text{Gauss-Jordan}} [I|d] \quad (3)$$

The solution is then given by

$$x_i = d_i, \quad i = 1, 2, \dots, n \quad (4)$$

In this method also, we use elementary row operations that are used in the Gauss elimination method. We apply these operations both below and above the diagonal in order to reduce all the off-diagonal elements of the matrix to zero. Pivoting can be used to make the pivot non-zero or to make it the largest element in magnitude in that column as discussed in Unit 5. We illustrate the method through an example.

**Example 4 :** Solve the system of equations

$$x_1 + x_2 + x_3 = 1$$

$$4x_1 + 3x_2 - x_3 = 6$$

$$3x_1 + 5x_2 + 3x_3 = 4$$

using Gauss-Jordan method with pivoting.

**Solution :** We have

$$[A|b] = \left[ \begin{array}{ccc|c} 1 & 1 & 1 & 1 \\ 4 & 3 & -1 & 6 \\ 3 & 5 & 3 & 4 \end{array} \right] \quad (\text{interchanging first and second row})$$

$$\approx \left[ \begin{array}{ccc|c} 4 & 3 & -1 & 6 \\ 1 & 1 & 1 & 1 \\ 3 & 5 & 3 & 4 \end{array} \right] R_2 - \frac{1}{4} R_1, R_3 - \frac{3}{4} R_1$$

$$\approx \left[ \begin{array}{ccc|c} 4 & 3 & -1 & 6 \\ 0 & \frac{1}{4} & \frac{5}{4} & -\frac{1}{2} \\ 0 & \frac{11}{4} & \frac{15}{4} & -\frac{1}{2} \end{array} \right] \quad (\text{interchanging second and third row})$$

$$\approx \left[ \begin{array}{ccc|c} 4 & 3 & -1 & 6 \\ 0 & \frac{11}{4} & \frac{15}{4} & -\frac{1}{2} \\ 0 & \frac{1}{4} & \frac{5}{4} & -\frac{1}{2} \end{array} \right] R_3 - \frac{1}{11} R_2, R_1 - \frac{12}{11} R_2$$

$$\approx \left[ \begin{array}{ccc|c} 4 & 0 & -\frac{56}{11} & \frac{72}{11} \\ 0 & \frac{11}{4} & \frac{15}{4} & -\frac{1}{2} \\ 0 & 0 & \frac{10}{11} & -\frac{5}{11} \end{array} \right] R_1 + \frac{56}{10} R_3, R_2 - \frac{33}{8} R_3$$

$$\approx \left[ \begin{array}{ccc|c} 4 & 0 & 0 & 4 \\ 0 & \frac{11}{4} & 0 & \frac{11}{8} \\ 0 & 0 & \frac{10}{11} & \frac{5}{11} \end{array} \right] \begin{array}{l} \frac{1}{4} R_1 \text{ (divide first row by 4),} \\ \frac{4}{11} R_2 \text{ (divide second row by 11/4),} \\ \frac{11}{10} R_3 \text{ (divide third row by 10/11).} \end{array}$$

$$\approx \left[ \begin{array}{ccc|c} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & \frac{1}{2} \\ 0 & 0 & 1 & -\frac{1}{2} \end{array} \right]$$

which is the desired form.

Thus, we obtain

$$x_1 = 1, x_2 = \frac{1}{2}, x_3 = -\frac{1}{2}.$$

The method can be easily extended to a general system of  $n$  equations. Just as we calculated the number of operations needed for Gauss elimination method in Unit 5, in the same way you can verify that the total number of operations needed for this method is  $M = \frac{1}{2} n^3 + \frac{n^2}{2} + n$ .

E3) Verify that the total number of operations needed for Gauss Jordan reduction method is  $\frac{1}{2} n^3 + \frac{n^2}{2} + n$ .

Clearly this method requires more number of operations compared to the Gauss elimination method. We, therefore, do not use this method generally for solving system of equations but is very commonly used for finding the inverse matrix. This is done by augmenting the matrix  $A$  by the identity matrix  $I$  of the order same as that of  $A$ . Using elementary row operations on the augmented matrix  $[A|I]$  we reduce the matrix  $A$  to the form  $I$  and in the process the matrix  $I$  is transformed to  $A^{-1}$ .

That is

$$[A|I] \xrightarrow[\text{Jordan}]{\text{Gauss}} [I|A^{-1}] \quad (5)$$

We now illustrate the method through examples.

**Example 5 :** Find the inverse of the matrix

$$A = \begin{bmatrix} 3 & 1 & 2 \\ 2 & -3 & -1 \\ 1 & -2 & 1 \end{bmatrix}$$

using the Gauss-Jordan method.

**Solution :** We have :

$$[A|I] = \left[ \begin{array}{ccc|ccc} 3 & 1 & 2 & 1 & 0 & 0 \\ 2 & -3 & -1 & 0 & 1 & 0 \\ 1 & -2 & 1 & 0 & 0 & 1 \end{array} \right] \frac{1}{3} R_1$$

$$\approx \left[ \begin{array}{ccc|ccc} 1 & 1/3 & 2/3 & 1/3 & 0 & 0 \\ 2 & -3 & -1 & 0 & 1 & 0 \\ 1 & -2 & 1 & 0 & 0 & 1 \end{array} \right] R_2 - 2R_1, R_3 - R_1$$

$$\approx \left[ \begin{array}{ccc|ccc} 1 & 1/3 & 2/3 & 1/3 & 0 & 0 \\ 0 & -11/3 & -7/3 & -2/3 & 1 & 0 \\ 0 & -7/3 & 1/3 & -1/3 & 0 & 1 \end{array} \right] \left( -\frac{3}{11} R_2 \right)$$

$$\approx \left[ \begin{array}{ccc|ccc} 1 & 1/3 & 2/3 & 1/3 & 0 & 0 \\ 0 & 1 & 7/11 & 2/11 & -3/11 & 0 \\ 0 & -7/3 & 1/3 & -1/3 & 0 & 1 \end{array} \right] R_1 - \frac{1}{3} R_2, R_3 + \frac{7}{3} R_2$$

$$\approx \left[ \begin{array}{ccc|ccc} 1 & 0 & 5/11 & 3/11 & 1/11 & 0 \\ 0 & 1 & 7/11 & 2/11 & -3/11 & 0 \\ 0 & 0 & 20/11 & 1/11 & -7/11 & 1 \end{array} \right] \frac{11}{20} R_3$$

$$\approx \left[ \begin{array}{ccc|ccc} 1 & 0 & 5/11 & 3/11 & 1/11 & 0 \\ 0 & 1 & 7/11 & 2/11 & -3/11 & 0 \\ 0 & 0 & 1 & 1/20 & -7/20 & 11/20 \end{array} \right] R_1 - \frac{5}{11} R_3, R_2 - \frac{7}{11} R_3$$

$$\approx \left[ \begin{array}{ccc|ccc} 1 & 0 & 0 & 1/4 & 1/4 & -1/4 \\ 0 & 1 & 0 & 3/20 & -1/20 & -7/20 \\ 0 & 0 & 1 & 1/20 & -7/20 & 11/20 \end{array} \right]$$

Thus we obtain

$$A^{-1} = \begin{bmatrix} 1/4 & 1/4 & -1/4 \\ 3/20 & -1/20 & -7/20 \\ 1/20 & -7/20 & 11/20 \end{bmatrix}$$

**Example 6 :** Find the inverse of the matrix

$$A = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 1 & 1/2 & 0 & 0 \\ 2 & 0 & -3 & 0 \\ 1 & -7/2 & -17 & 55/3 \end{bmatrix}$$

using the Gauss-Jordan method

**Solution :** Here we have

$$[A|I] = \left[ \begin{array}{cccc|cccc} 2 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1/2 & 0 & 0 & 0 & 1 & 0 & 0 \\ 2 & 0 & -3 & 0 & 0 & 0 & 1 & 0 \\ 1 & -7/2 & -17 & 55/3 & 0 & 0 & 0 & 1 \end{array} \right] \frac{1}{2} R_1$$

$$\approx \left[ \begin{array}{cccc|cccc} 1 & 0 & 0 & 0 & 1/2 & 0 & 0 & 0 \\ 1 & 1/2 & 0 & 0 & 0 & 1 & 0 & 0 \\ 2 & 0 & -3 & 0 & 0 & 0 & 1 & 0 \\ 1 & -7/2 & -17 & 55/3 & 0 & 0 & 0 & 1 \end{array} \right]$$

$$R_2 - R_1, R_3 - 2R_1, R_4 - R_1$$



$$\approx \left[ \begin{array}{cccc|cccc} 1 & 0 & 0 & 0 & 1/2 & 0 & 0 & 0 \\ 0 & 1/2 & 0 & 0 & -1/2 & 1 & 0 & 0 \\ 0 & 0 & -3 & 0 & -1 & 0 & 1 & 0 \\ 0 & -7/2 & -17 & 55/3 & -1/2 & 0 & 0 & 1 \end{array} \right] 2R_2$$

$$\approx \left[ \begin{array}{cccc|cccc} 1 & 0 & 0 & 0 & 1/2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 2 & 0 & 0 \\ 0 & 0 & -3 & 0 & -1 & 0 & 1 & 0 \\ 0 & -7/2 & -17 & 55/3 & -1/2 & 0 & 0 & 1 \end{array} \right] R_4 + \frac{7}{2} R_2$$

$$\approx \left[ \begin{array}{cccc|cccc} 1 & 0 & 0 & 0 & 1/2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 2 & 0 & 0 \\ 0 & 0 & -3 & 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & -17 & 55/3 & -4 & 7 & 0 & 1 \end{array} \right] \left( -\frac{1}{3} R_3 \right)$$

$$\approx \left[ \begin{array}{cccc|cccc} 1 & 0 & 0 & 0 & 1/2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1/3 & 0 & -1/3 & 0 \\ 0 & 0 & -17 & 55/3 & -4 & 7 & 0 & 1 \end{array} \right] \left( -\frac{1}{17} R_4 \right)$$

$$\approx \left[ \begin{array}{cccc|cccc} 1 & 0 & 0 & 0 & 1/2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1/3 & 0 & -1/3 & 0 \\ 0 & 0 & -17 & 55/3 & 4/17 & -7/17 & 0 & -1/17 \end{array} \right] R_4 - R_3$$

$$\approx \left[ \begin{array}{cccc|cccc} 1 & 0 & 0 & 0 & 1/2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1/3 & 0 & -1/3 & 0 \\ 0 & 0 & 0 & -55/51 & -5/51 & -7/17 & 1/3 & -1/17 \end{array} \right] \left( -\frac{51}{55} R_4 \right)$$

$$\approx \left[ \begin{array}{cccc|cccc} 1 & 0 & 0 & 0 & 1/2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1/3 & 0 & -1/3 & 0 \\ 0 & 0 & 0 & 1 & 1/11 & 21/55 & -17/55 & 3/55 \end{array} \right]$$

Hence

$$A^{-1} = \begin{bmatrix} 1/2 & 0 & 0 & 0 \\ -1 & 2 & 0 & 0 \\ 1/3 & 0 & -1/3 & 0 \\ 1/11 & 21/55 & -17/55 & 3/55 \end{bmatrix}$$

is the inverse of the given lower triangular matrix.

Let us now consider the problem of finding the inverse of an upper triangular matrix.

**Example 7 :** Find the inverse of the matrix

$$A = \begin{bmatrix} 1 & 3/2 & 2 & 1/2 \\ 0 & 1 & -4 & 1 \\ 0 & 0 & 1 & 2/3 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

using the Gauss-Jordan method.

**Solution :** Here, we have

$$[A|I] = \left[ \begin{array}{cccc|cccc} 1 & 3/2 & 2 & 1/2 & 1 & 0 & 0 & 0 \\ 0 & 1 & -4 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 2/3 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \end{array} \right] R_1 - \frac{3}{2} R_2$$

$$\approx \left[ \begin{array}{cccc|cccc} 1 & 0 & 8 & -1 & 1 & -3/2 & 0 & 0 \\ 0 & 1 & -4 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 2/3 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \end{array} \right] R_1 - 8R_3, R_2 + 4R_3$$

$$\approx \left[ \begin{array}{cccc|cccc} 1 & 0 & 0 & -19/3 & 1 & -3/2 & -8 & 0 \\ 0 & 1 & 0 & 11/3 & 0 & 1 & 4 & 0 \\ 0 & 0 & 1 & 2/3 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \end{array} \right]$$

$$R_1 + \frac{19}{3} R_4, R_2 - \frac{11}{3} R_4, R_3 - \frac{2}{3} R_4$$

$$\approx \left[ \begin{array}{cccc|cccc} 1 & 0 & 0 & 0 & 1 & -3/2 & -8 & 19/3 \\ 0 & 1 & 0 & 0 & 0 & 1 & 4 & -11/3 \\ 0 & 0 & 1 & 0 & 0 & 0 & 1 & -2/3 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \end{array} \right]$$

Hence

$$A^{-1} = \begin{bmatrix} 1 & -3/2 & -8 & 19/3 \\ 0 & 1 & 4 & -11/3 \\ 0 & 0 & 1 & -2/3 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

which is the inverse of the given upper triangular matrix.

Note that in Examples 2,3,6 and 7, the inverse of a lower/upper triangular matrix is again a lower/upper triangular matrix. There is another method of finding the inverse of a matrix  $A$  which uses the pivoting strategy. Recall that in Sec. 5.5. of Unit 5, for the solution of system of linear algebraic equation  $Ax = b$ , we showed you how the multipliers  $m_{p,i,k}$ 's can be stored in working array  $W$  during the process of elimination. The main advantage of storing these multipliers is that if we have already solved the linear system of equations  $Ax = b$  or order  $n$ , by the elimination method and we want to solve the system  $Ax = c$  with the same coefficient matrix  $A$ , only the right side being different, then we do not have to go through the entire elimination process again. Since we have saved in the working matrix  $W$  all the multipliers used and also have saved the  $p$  vector, we have only to repeat the operations on the right hand side to obtain  $\tilde{c}$ , such that  $Ux = \tilde{c}$  is equivalent to  $Ax = c$ .

In order to understand the calculations necessary to derive  $\tilde{c}$ , from  $c$  consider the changes made in the right side  $b$  during the elimination process. Let  $k$  be an integer between 1 and  $n$ , and assume that the  $i$ th equation was used as pivotal equation during step  $k$  of the elimination process. Then  $i = p_k$ . Initially, the right side of equation  $i$  is just  $b_i$ .

If  $k > 1$ , then after Step 1, the right side is

$$b_i^{(1)} = b_i - m_{i1} b_{p_1}$$

If  $k > 2$ , then after Step 2, the right side is

$$\begin{aligned} b_i^{(2)} &= b_i^{(1)} - m_{i2} b_{p_2}^{(1)} \\ &= b_i - m_{i1} b_{p_1} - m_{i2} b_{p_2}^{(1)} \end{aligned}$$

In the same manner, we have the right side of equation  $i = p_k$  as

$$b_i^{(k-1)} = b_i - m_{i1} b_{p_1} - m_{i2} b_{p_2} - \dots - m_{i,k-1} b_{p_{k-1}} \quad (6)$$

Replacing  $i$  by  $p_k$  in Eqn. (6), we get

$$b_{p_k}^{(k-1)} = b_{p_k} - m_{p_k,1} b_{p_1} - m_{p_k,2} b_{p_2} - \dots - m_{p_k,k-1} b_{p_{k-1}} \quad (7)$$

$k = 1, 2, \dots, n.$

Also, since  $\tilde{b}_j = b_{p_j}^{(j-1)}$ ,  $j = 1, 2, \dots, n$ , we can rewrite Eqn. (7) as

$$\tilde{b}_k = b_{p_k} - m_{p_k,1} \tilde{b}_1 - m_{p_k,2} \tilde{b}_2 - \dots - m_{p_k,k-1} \tilde{b}_{k-1} \quad (8)$$

$k = 1, \dots, n.$

Eqn. (8) can then be used to calculate the entries of  $\tilde{\mathbf{b}}$ . But since the multipliers  $m'_{ij}$ 's are stored in entries  $w_{ij}$ 's of the working matrix  $W$ , we can also write Eqn. (8) in the form

$$\tilde{b}_k = b_{p_k} - \sum_{j=1}^{k-1} W_{p_k j} \tilde{b}_j, \quad k=1, \dots, n \quad (9)$$

Hence, if we just know the final content of the first  $n$  columns of  $W$  and the pivoting strategy  $\mathbf{p}$  then we can calculate the solution  $\mathbf{x}$  of  $A\mathbf{x} = \mathbf{b}$  by using the back substitution method and writing

$$x_k = \frac{\tilde{b}_k - \sum_{j=k+1}^n W_{p_k j} x_j}{W_{p_k k}}, \quad k=n, n-1, \dots, 1 \quad (10)$$

The vector  $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^T$  will then be the solution of  $A\mathbf{x} = \mathbf{b}$ .

For finding the inverse of an  $n \times n$  matrix  $A$ , we use the above algorithm. We first calculate the final contents of the  $n$  columns of the working matrix  $W$  and the pivoting vector  $\mathbf{p}$  and then solve each of the  $n$  systems

$$A\mathbf{x} = \mathbf{e}_j, \quad j=1, \dots, n \quad (11)$$

where  $\mathbf{e}_1 = [1 \ 0 \ \dots \ 0]^T$ ,  $\mathbf{e}_2 = [0 \ 1 \ 0 \ \dots \ 0]^T$ ,  $\dots$ ,  $\mathbf{e}_n = [0 \ 0 \ \dots \ 1]^T$ , with the help of Eqns (9) and (10). Then for each  $j=1, \dots, n$  the solution of system (11) will be the corresponding column of the inverse matrix  $A^{-1}$ . The following example will help you to understand the above procedure.

**Example 8 :** Find the inverse of the matrix

$$A = \begin{bmatrix} 1 & 2 & -1 \\ 2 & 1 & 0 \\ -1 & 1 & 2 \end{bmatrix}$$

using partial pivoting.

**Solution :** Initially  $\mathbf{p} = [p_1, p_2, p_3]^T = [1, 2, 3]^T$  and the working matrix is

$$W^{(0)} = \begin{bmatrix} 1 & 2 & -1 \\ 2 & 1 & 0 \\ -1 & 1 & 2 \end{bmatrix}$$

Now  $d_1 = 2$ ,  $d_2 = 2$ ,  $d_3 = 2$ .

$$\text{Step 1 : } \frac{|w_{p_1,1}|}{d_1} = \frac{1}{2}, \frac{|w_{p_2,1}|}{d_2} = \frac{2}{2} = 1, \frac{|w_{p_3,1}|}{d_3} = \frac{1}{2}.$$

$$1 > \frac{1}{2}, \frac{1}{2} \therefore p_1 = 2, p_2 = 1, p_3 = 3$$

We use the second equation to eliminate  $x_1$  from first and third equations and store corresponding multipliers instead of storing zeros in the working matrix. The multipliers are

$$m_{p_{i,1}} = \frac{w_{p_{i,1}}}{w_{p_{1,1}}}, \quad i = 2, 3$$

$$\therefore m_{p_{2,1}} = m_{11} = \frac{w_{p_{2,1}}}{w_{p_{1,1}}} = \frac{1}{2}$$

$$m_{p_{3,1}} = m_{31} = \frac{w_{p_{3,1}}}{w_{p_{1,1}}} = -\frac{1}{2}$$

we get the following working matrix

$$W^{(1)} = \begin{bmatrix} \boxed{\frac{1}{2}} & \frac{3}{2} & -1 \\ \boxed{2} & 1 & 0 \\ \boxed{-\frac{1}{2}} & \frac{3}{2} & 2 \end{bmatrix}, \quad p = (2, 1, 3)^T$$

$$\text{Step 2: } \frac{|w_{p_{2,2}}|}{dp_2} = \frac{|w_{1,2}|}{d_1} = \frac{3/2}{2} = \frac{3}{4}$$

$$\frac{|w_{p_{3,2}}|}{dp_3} = \frac{|w_{3,2}|}{d_3} = \frac{3/2}{2} = \frac{3}{4}$$

Since  $\frac{3}{4} = \frac{3}{4}$  so we take  $p = (2, 1, 3)^T$

$$\text{Now } m_{p_{i,2}} = \frac{w_{p_{i,2}}}{w_{p_{2,2}}}, \quad i = 3$$

$$\therefore m_{p_{3,2}} = m_{32} = \frac{w_{3,2}}{w_{1,2}} = \frac{3/2}{3/2} = 1$$

We use the first equation as pivotal equation to eliminate  $x_2$  from the third equation and also store the multipliers. After the second step we have the following working matrix

$$W^{(2)} = \begin{bmatrix} \boxed{\frac{1}{2}} & \boxed{\frac{3}{2}} & -1 \\ \boxed{2} & 1 & 0 \\ \boxed{-\frac{1}{2}} & \boxed{1} & 3 \end{bmatrix}, \quad p = (2, 1, 3)^T$$

Now in this case,  $W^{(2)}$  is our final working matrix with pivoting strategy  $p = (2, 1, 3)^T$

Note that circled ones denote multipliers and squared ones denote pivot elements in the working matrices.

To find the inverse of the given matrix  $A$ , we have to solve

$$Ax = e_1 = [b_1 \ b_2 \ b_3]^T$$

$$Ax = e_2 = [b_1 \ b_2 \ b_3]^T$$

$$Ax = e_3 = [b_1 \ b_2 \ b_3]^T$$

$$\text{where } e_1 = [1 \ 0 \ 0]^T, \quad e_2 = [0 \ 1 \ 0]^T, \quad e_3 = [0 \ 0 \ 1]^T$$

First we solve the system  $Ax = e_1$  and consider

$$\begin{bmatrix} \boxed{\frac{1}{2}} & \boxed{\frac{3}{2}} & -1 \\ \boxed{2} & 1 & 0 \\ \boxed{-\frac{1}{2}} & \boxed{1} & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \mathbf{p} = (2, 1, 3)^T \quad (12)$$

Using Eqn. (9), we get

$$\text{with } p_1 = 2, \tilde{b}_1 = b_2 = 0$$

$$\text{with } p_2 = 1, \tilde{b}_2 = b_1 - w_{11}\tilde{b}_1$$

$$= 1 - \left[\frac{1}{2}\right] \cdot 0$$

$$= 1$$

$$\text{with } p_3 = 3, \tilde{b}_3 = b_3 - w_{31}\tilde{b}_1 - w_{32}\tilde{b}_2$$

$$= 0 - \left[-\frac{1}{2}\right] \cdot 0 - 1 \cdot 1 = -1$$

Using Eqn. (10), we then get the following system of equations

$$3x_3 = -1$$

$$\frac{3}{2}x_2 - x_3 = 1$$

$$2x_1 + x_2 = 0$$

$$\text{which gives } x_3 = -\frac{1}{3}, x_2 = \frac{4}{9} \text{ and } x_1 = -\frac{2}{9}$$

$$\text{i.e., vector } \mathbf{x} = \left[-\frac{2}{9}, \frac{4}{9}, -\frac{1}{3}\right]^T \text{ is the solution of system (12).}$$

**Remember** that the solution of system (12) constitutes the first column of the inverse matrix  $A^{-1}$ .

In the same way we solve the system of equations  $Ax = e_2$  and  $Ax = e_3$ , or

$$\begin{bmatrix} \boxed{\frac{1}{2}} & \boxed{\frac{3}{2}} & -1 \\ \boxed{2} & \boxed{1} & 0 \\ \boxed{-\frac{1}{2}} & 1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \mathbf{p} = (2, 1, 3)^T \quad (13)$$

and

$$\begin{bmatrix} \frac{1}{2} & \frac{3}{2} & -1 \\ 2 & 1 & 0 \\ -\frac{1}{2} & 1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \mathbf{p} = (2, 1, 3)^T \quad (14)$$

Using Eqns (9) and (10), we obtain the solution of system (13) as

$$\mathbf{x} = \left[\frac{5}{9}, \frac{1}{9}, \frac{1}{3}\right]^T \text{ which is the second column of } A^{-1} \text{ and the solution of system}$$

$$(14), \text{ i.e. } \mathbf{x} = \left[-\frac{1}{9}, \frac{2}{9}, \frac{1}{3}\right]^T \text{ as the third column of } A^{-1}$$

$$\text{Hence } A^{-1} = \begin{bmatrix} -\frac{2}{9} & \frac{5}{9} & -\frac{1}{9} \\ \frac{4}{9} & -\frac{1}{9} & \frac{2}{9} \\ -\frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

You may now try the following exercises

E4) In Examples 6 and 7 verify that

$$A A^{-1} = A^{-1} A = I.$$

E5) Solve the system of equation

$$x_1 + 2x_2 + x_3 = 0$$

$$2x_1 + 2x_2 + 3x_3 = 3$$

$$-x_1 - 3x_2 = 2$$

using the Gauss-Jordan method with pivoting.

E6) Find the inverse of the matrix

$$A = \begin{bmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{bmatrix}$$

using the Gauss-Jordan method.

E7) Find the inverse of the matrix

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & \frac{1}{2} & 1 \end{bmatrix}$$

using the Gauss-Jordan method.

You may recall that in Sec. 5.6 of Unit 5 we discussed the LU decomposition method. Using this method we can factorize any nonsingular square matrix  $A$  into the product of a lower triangular matrix  $L$  and an upper triangular matrix  $U$ . That is, we can write  $A = LU$ . ... (15)

In the next section we shall discuss how form (15) can be used to find the inverse of nonsingular square matrices.

## 6.4 LU DECOMPOSITION METHOD

Let us consider Eqn. (15) and take the inverse on both the sides. If we use the fact that the inverse of the product of matrices is the product of their inverses taken in reverse order (ref. Theorem 6, Sec. 7.6 of Unit 7, Block 2, Linear Algebra MTE-02), then we obtain

$$A^{-1} = (LU)^{-1} = U^{-1}L^{-1} \quad (16)$$

We can now find the inverses of  $U$  and  $L$  separately and obtain the inverse matrix  $A^{-1}$  from Eqn. (16):

**Remark:** It may appear to you that finding an inverse of a matrix by this method is a lengthy process. But, in practice, this method is very useful because of the fact that here we deal with triangular matrices and triangular matrices are easily invertible. It involves only forward and backward substitutions.

Let us now consider an example to understand how the method works.

**Example 9 :** Find the inverse of the matrix

$$A = \begin{bmatrix} 3 & 1 & 2 \\ 2 & -3 & -1 \\ 1 & -2 & 1 \end{bmatrix}$$

using LU decomposition method.

**Solution :** We write,

$$A = \begin{bmatrix} 3 & 1 & 2 \\ 2 & -3 & -1 \\ 1 & -2 & 1 \end{bmatrix} = LU = \begin{bmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} \begin{bmatrix} 1 & u_{12} & u_{13} \\ 0 & 1 & u_{23} \\ 0 & 0 & 1 \end{bmatrix} \quad (17)$$

Comparing the coefficients on both sides of Eqn. (17), we obtain

$$l_{11} = 3, \quad l_{21} = 2, \quad l_{31} = 1 \quad (\text{multiplying the rows of L by the first column of U})$$

$$l_{11}u_{12} = 1, \quad u_{12} = \frac{1}{3} \quad (\text{multiplying first row of L by the}$$

$$l_{11}u_{13} = 2, \quad u_{13} = 2/3 \quad \text{second and third column of U})$$

The second column of L is obtained from

$$l_{21}u_{12} + l_{22} = a_{22}, \quad l_{22} = -3 - \frac{2}{3} = -\frac{11}{3}$$

$$l_{31}u_{12} + l_{32} = a_{32}, \quad l_{32} = -2 - \frac{1}{3} = -\frac{7}{3}$$

$u_{23}$  is obtained from

$$l_{21}u_{13} + l_{22}u_{23} = a_{23}, \quad u_{23} = \frac{-1 - 2(2/3)}{-11/3} = \frac{7}{11}$$

$l_{33}$  is obtained from

$$l_{31}u_{13} + l_{32}u_{23} + l_{33} = 1, \quad l_{33} = \frac{20}{11}$$

Thus we have

$$L = \begin{bmatrix} 3 & 0 & 0 \\ 2 & -\frac{11}{3} & 0 \\ 1 & -\frac{7}{3} & \frac{20}{11} \end{bmatrix} \quad \text{and} \quad U = \begin{bmatrix} 1 & 1/3 & 2/3 \\ 0 & 1 & 7/11 \\ 0 & 0 & 1 \end{bmatrix}$$

Now since L is a lower triangular matrix,  $L^{-1}$  is also a lower triangular matrix. Let us assume that

$$L^{-1} = \begin{bmatrix} l'_{11} & 0 & 0 \\ l'_{21} & l'_{22} & 0 \\ l'_{31} & l'_{32} & l'_{33} \end{bmatrix}$$

Using the identity  $LL^{-1} = I$ , we have

$$LL^{-1} = \begin{bmatrix} 3 & 0 & 0 \\ 2 & -\frac{11}{3} & 0 \\ 1 & -\frac{7}{3} & \frac{20}{11} \end{bmatrix} \begin{bmatrix} l'_{11} & 0 & 0 \\ l'_{21} & l'_{22} & 0 \\ l'_{31} & l'_{32} & l'_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

comparing the coefficients, we get

$$l'_{11} = \frac{1}{3}, \quad l'_{22} = -\frac{3}{11}, \quad l'_{33} = \frac{11}{20}$$

Also,

$$2l'_{11} - \frac{11}{3}l'_{21} = 0, \quad l'_{21} = \frac{6}{33} = \frac{2}{11}$$

$$l'_{11} - \frac{7}{3}l'_{21} + \frac{20}{11}l'_{31} = 0, \quad l'_{31} = \frac{1}{20}$$

$$-\frac{7}{3}l'_{22} + \frac{20}{11}l'_{32} = 0, \quad l'_{32} = -\frac{7}{20}$$

$$\therefore L^{-1} = \begin{bmatrix} 1/3 & 0 & 0 \\ 2/11 & -3/11 & 0 \\ 1/20 & -7/20 & 11/20 \end{bmatrix}$$

Similarly, since U is an upper triangular matrix,  $U^{-1}$  is also an upper triangular matrix. Using  $UU^{-1} = I$ , we obtain by backward substitution.

$$U = \begin{bmatrix} 1 & 1/3 & 2/3 \\ 0 & 1 & 7/11 \\ 0 & 0 & 1 \end{bmatrix} \quad \text{and} \quad U^{-1} = \begin{bmatrix} 1 & -1/3 & -5/11 \\ 0 & 1 & -7/11 \\ 0 & 0 & 1 \end{bmatrix}$$

Therefore, we have from Eqn. (16)

$$\begin{aligned} A^{-1} = U^{-1} L^{-1} &= \begin{bmatrix} 1 & -1/3 & -5/11 \\ 0 & 1 & -7/11 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1/3 & 0 & 0 \\ 2/11 & -3/11 & 0 \\ 1/20 & -7/20 & 11/20 \end{bmatrix} \\ &= \begin{bmatrix} 1/4 & -1/4 & -1/4 \\ 3/20 & -1/20 & -7/20 \\ 1/20 & -7/20 & 11/20 \end{bmatrix} \end{aligned}$$

And now a few exercises for you.

E8) Find the inverse of the matrix

$$A = \begin{bmatrix} 5 & 8 & 1 \\ 0 & 2 & 1 \\ 4 & 3 & -1 \end{bmatrix}$$

using the LU decomposition method.

E9) Find the inverse of the matrix

$$A = \begin{bmatrix} 3 & 1 & 2 \\ 2 & -1 & -1 \\ 1 & -2 & 1 \end{bmatrix}$$

using the LU decomposition method.

E10) Find the inverse of the matrix given in E6) using the LU decomposition method.

E11) Find the inverse of a matrix

$$A = \begin{bmatrix} 1 & 0 & 5 & 2 \\ -1 & 4 & 1 & 0 \\ 3 & 0 & 4 & 1 \\ -2 & 1 & 1 & 3 \end{bmatrix}$$

We now end this unit by giving a summary of what we have covered in it.

## 6.5 SUMMARY

In this unit we have covered the following :

- 1) Using the method of adjoints, the inverse of a given nonsingular matrix A can be obtained from

$$A^{-1} = \frac{1}{\det(A)} \text{adj}(A)$$

(see Eqn. (1))



Since the number of operations in the adjoint method to find the inverse of a  $n \times n$  nonsingular matrix  $A$  increases rapidly as  $n$  increases, the method is not generally used for  $n > 4$ .

- 2) For large  $n$ , the Gauss-Jordan reduction method, which is an extension of the Gauss elimination method can be used for finding the inverse matrix and solve the linear systems

$$Ax = b$$

(see Eqn. (2))

using the Gauss-Jordan method.

- a) the solution of system of Eqs (2) can be obtained by using elementary row operations.

$$[A|b] \xrightarrow{\text{reduced to}} [I|d]$$

- b) the inverse matrix  $A^{-1}$  can be obtained by using elementary row operations

$$[A|I] \xrightarrow{\text{reduced to}} [I|A^{-1}]$$

- 3) For large  $n$ , another useful method of finding the inverse matrix  $A^{-1}$  is LU decomposition method. Using this method any nonsingular matrix  $A$  is first decomposed into the product of a lower triangular matrix  $L$  and an upper triangular matrix  $U$ . That is

$$A = LU.$$

$U^{-1}$  and  $L^{-1}$  can be obtained by backward and forward substitutions. Then the inverse can be found from

$$A^{-1} = U^{-1}L^{-1}.$$

## 6.6 SOLUTIONS/ANSWERS

$$E1) A^c = \begin{bmatrix} -3 & -3 & -3 \\ -5 & 1 & 7 \\ 1 & 7 & -5 \end{bmatrix}; \det(A) = -18$$

$$A^{-1} = \begin{bmatrix} 1/6 & 5/18 & -1/18 \\ 1/6 & -1/18 & -7/18 \\ 1/6 & -7/18 & 5/18 \end{bmatrix}; x = \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix}$$

$$E2) A^c = \begin{bmatrix} 5 & 0 & -15 & -5 \\ 32 & -33 & 14 & -21 \\ -39 & 11 & 7 & 17 \\ -19 & 11 & 2 & -3 \end{bmatrix}; \det(A) = -55$$

$$A^{-1} = \begin{bmatrix} -1/11 & -32/55 & 39/55 & 19/55 \\ 0 & 3/5 & -1/5 & -1/5 \\ 3/11 & -14/55 & -7/55 & -2/55 \\ 1/11 & 21/55 & -17/55 & 3/55 \end{bmatrix}; x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

E3) No. of divisions

1st step of elimination :  $n$

2nd step of elimination :  $(n-1)$

.....

$n$ th step of elimination :  $1$

$\therefore$  Total number of divisions =  $n + (n-1) + \dots + 1$

$$= \sum n = \frac{n(n+1)}{2}$$

**No. of multiplication**

**1st step :** elimination in 2nd equation : n  
 elimination in 3rd equation : n  
 .....  
 elimination in nth equation : n  
 $\therefore$  Total of 1st step =  $n(n-1)$

**2nd step :** elimination in 1st equation :  $(n-1)$   
 elimination in 3rd equation :  $(n-1)$   
 .....  
 elimination in nth equation :  $(n-1)$   
 $\therefore$  Total of 2nd step =  $(n-1)(n-1)$

**3rd step :** elimination in 1st equation =  $(n-2)$   
 elimination in 2nd equation =  $(n-2)$   
 elimination in 4th equation =  $(n-2)$   
 .....  
 elimination in nth equation =  $(n-2)$   
 $\therefore$  Total of 3rd step =  $(n-1)(n-2)$

**(n-1)th step :** elimination in 1st equation = 1  
 elimination in 2nd equation = 1  
 .....  
 elimination in nth equation = 1  
 $\therefore$  Total of (n-1)th step =  $1 \cdot (n-1)$

$$\begin{aligned} \text{Total multiplications} &= n(n-1) + (n-1)(n-1) + (n-1)(n-2) + \dots + 1(n-1) \\ &= (n-1)[n + (n-1) + (n-2) + \dots + 1] \\ &= (n-1) \sum n = \frac{(n-1)n(n+2)}{2} \end{aligned}$$

We also need n divisions to find the solution vector

$$\begin{aligned} \therefore \text{Total operations} &= \frac{n(n+1)}{2} + \frac{(n-1)n(n+2)}{2} + n \\ &= \frac{n^3}{2} + \frac{n^2}{2} + n. \end{aligned}$$

E4) In Example 6

$$\begin{aligned} AA^{-1} &= \begin{bmatrix} 2 & 0 & 0 & 0 \\ 1 & 1/2 & 0 & 0 \\ 2 & 0 & -3 & 0 \\ 1 & -7/2 & -17 & 55/3 \end{bmatrix} \begin{bmatrix} -1/2 & 0 & 0 & 0 \\ -1 & 2 & 0 & 0 \\ 1/3 & 0 & -1/3 & 0 \\ 1/11 & 21/55 & -17/55 & 3/55 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = A^{-1}A. \end{aligned}$$

Similarly, check for Example 7.

E5) Elementary row operations required on the augmented matrix [A|I] in order are :

Interchange first and second row;  $R_2 - \frac{1}{2}R_1, R_3 + \frac{1}{2}R_1$ ;

interchanging second and third row;  $R_1 + \frac{1}{2}R_2, R_3 + \frac{1}{4}R_2$ ;

$R_2 - \frac{12}{7}R_3, R_1 + \frac{18}{7}R_3$

Final derived system is

$$\left[ \begin{array}{ccc|c} 2 & 0 & 0 & -2/7 \\ 0 & -4 & 0 & -20/7 \\ 0 & 0 & -\frac{7}{8} & -11/8 \end{array} \right]$$

Solution is  $x_1 = -\frac{1}{7}$ ,  $x_2 = -\frac{5}{7}$ ,  $x_3 = \frac{11}{7}$ .

E6) Elementary row operations required on the augmented matrix  $[A|I]$  in order are :

$$\frac{1}{2} R_1; R_2 + R_1; \frac{2}{3} R_2; R_1 + \frac{1}{2} R_2, R_3 + R_2; \frac{3}{4} R_3; R_1 + \frac{1}{3} R_3, R_2 + \frac{2}{3} R_3; R_4 + R_3; \frac{4}{5} R_4; R_1 + \frac{1}{4} R_4, R_2 + \frac{1}{2} R_4, R_3 + \frac{3}{4} R_4;$$

We find

$$A^{-1} = \begin{bmatrix} 4/5 & 3/5 & 2/5 & 1/5 \\ 3/5 & 6/5 & 4/5 & 2/5 \\ 2/5 & 4/5 & 6/5 & 3/5 \\ 1/5 & 2/5 & 3/5 & 4/5 \end{bmatrix}$$

E7) Elementary row operations required on the matrix  $[A|I]$  are

$$R_2 - 2R_1, R_3 + R_1; R_3 - \frac{1}{2} R_2$$

we obtain

$$A^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 2 & -\frac{1}{2} & 0 \end{bmatrix}$$

$$E8) L = \begin{bmatrix} 5 & 0 & 0 \\ 0 & 2 & 0 \\ 4 & -17/5 & -1/10 \end{bmatrix}; L^{-1} = \begin{bmatrix} 1/5 & 0 & 0 \\ 0 & 1/2 & 0 \\ 8 & -17 & -10 \end{bmatrix}$$

$$U = \begin{bmatrix} 1 & 8/5 & 1/5 \\ 0 & 1 & 1/2 \\ 0 & 0 & 1 \end{bmatrix}; U^{-1} = \begin{bmatrix} 1 & -8/5 & 3/5 \\ 0 & 1 & -1/2 \\ 0 & 0 & 1 \end{bmatrix}$$

$$A^{-1} = U^{-1}L^{-1} = \begin{bmatrix} 5 & -11 & -6 \\ -4 & 9 & 5 \\ 8 & -17 & -10 \end{bmatrix}$$

$$E9) L = \begin{bmatrix} 3 & 0 & 0 \\ 2 & -5/3 & 0 \\ 1 & -7/3 & 18/5 \end{bmatrix}; L^{-1} = \begin{bmatrix} 1/3 & 0 & 0 \\ 2/5 & -3/5 & 0 \\ 3/18 & -7/18 & 5/18 \end{bmatrix}$$

$$U = \begin{bmatrix} 1 & 1/3 & 2/3 \\ 0 & 1 & 7/5 \\ 0 & 0 & 1 \end{bmatrix}; U^{-1} = \begin{bmatrix} 1 & -1/3 & -1/5 \\ 0 & 1 & -7/5 \\ 0 & 0 & 1 \end{bmatrix}$$

$$A^{-1} = U^{-1}L^{-1} = \begin{bmatrix} 1/6 & 5/18 & -1/18 \\ 1/6 & -1/18 & -7/18 \\ 1/6 & -7/18 & 5/18 \end{bmatrix}$$

$$E10) L = \begin{bmatrix} 2 & 0 & 0 & 0 \\ -1 & 3/2 & 0 & 0 \\ 0 & -1 & 4/3 & 0 \\ 0 & 0 & -1 & 5/4 \end{bmatrix}$$

$$U = \begin{bmatrix} 1 & -1/2 & 0 & 0 \\ 0 & 1 & -2/3 & 0 \\ 0 & 0 & 1 & -3/4 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$L^{-1} = \begin{bmatrix} 1/2 & 0 & 0 & 0 \\ 1/3 & 2/3 & 0 & 0 \\ 1/4 & 2/4 & 3/4 & 0 \\ 1/5 & 2/5 & 3/5 & 4/5 \end{bmatrix}$$

$$U^{-1} = \begin{bmatrix} 1 & 1/2 & 1/3 & 1/4 \\ 0 & 1 & 2/3 & 2/4 \\ 0 & 0 & 1 & 3/4 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$A^{-1} = U^{-1}L^{-1} = \begin{bmatrix} 4/5 & 3/5 & 2/5 & 1/5 \\ 3/5 & 6/5 & 4/5 & 2/5 \\ 2/5 & 4/5 & 6/5 & 3/5 \\ 1/5 & 2/5 & 3/5 & 4/5 \end{bmatrix}$$

$$E11) L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ -1 & 4 & 0 & 0 \\ 3 & 0 & -11 & 0 \\ -2 & 1 & 19/2 & 24/11 \end{bmatrix}$$

$$U = \begin{bmatrix} 1 & 0 & 5 & 2 \\ 0 & 1 & 3/2 & 1/2 \\ 0 & 0 & 1 & 5/11 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$L^{-1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1/4 & 1/4 & 0 & 0 \\ 3/11 & 0 & -1/11 & 0 \\ -37/96 & -11/96 & 38/96 & 44/96 \end{bmatrix}$$

$$U^{-1} = \begin{bmatrix} 1 & 0 & -5 & 3/11 \\ 0 & 1 & -3/2 & 2/11 \\ 0 & 0 & 1 & -5/11 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$A^{-1} = U^{-1}L^{-1} = \begin{bmatrix} -15/32 & -1/32 & 9/16 & 1/8 \\ -11/48 & 11/48 & 5/24 & 1/12 \\ 43/96 & 5/96 & -13/48 & -5/24 \\ -37/96 & -11/96 & 19/48 & 11/24 \end{bmatrix}$$

---

# UNIT 7 ITERATIVE METHODS

---

## Structure

- 7.1 Introduction
- 7.2 The General Iteration Method
- 7.3 The Jacobi Iteration Method
- 7.4 The Gauss-Seidel Iteration Method
- 7.5 Summary
- 7.6 Solutions/Answers

---

## 7.1 INTRODUCTION

---

In the previous two units, you have studied direct methods for solving linear system of equations  $Ax = b$ ,  $A$  being  $n \times n$  non-singular matrix. Direct methods provide the exact solution in a finite number of steps provided exact arithmetic is used and there is no round-off error. Also, direct methods are generally used when the matrix  $A$  is **dense or filled**, that is, there are few zero elements, and the order of the matrix is not very large say  $n < 50$ .

Iterative methods, on the other hand, start with an initial approximation and by applying a suitably chosen algorithm, lead to successively better approximations. Even if the process converges, it would give only an approximate solution. These methods are generally used when the matrix  $A$  is **sparse** and the order of the matrix  $A$  is very large say  $n > 50$ . Sparse matrices have very few non-zero elements. In most cases these non-zero elements lie on or near the main diagonal giving rise to tri-diagonal, five diagonal or band matrix systems. It may be **noted** that there are no fixed rules to decide when to use direct methods and when to use iterative methods. However, when the coefficient matrix is sparse or large, the use of iterative methods is ideally suited to find the solution which take advantage of the sparse nature of the matrix involved.

In this unit we shall discuss two iterative methods, namely, Jacobi iteration and Gauss-Seidel iteration methods which are frequently used for solving linear system of equations.

### Objectives

After studying this unit, you should be able to:

- obtain the solution of system of linear equations,  $Ax = b$ , when the matrix  $A$  is large or sparse, by using the iterative method viz; Jacobi method or the Gauss-Seidel method;
- tell whether these iterative methods converge or not;
- obtain the rate of convergence and the approximate number of iterations needed for the required accuracy of these iterative methods.

---

## 7.2 THE GENERAL ITERATION METHOD

---

In iteration methods as we have already mentioned, we start with some initial approximate solution vector  $x^{(0)}$  and generate a sequence of approximants  $\{x^{(k)}\}$  which converge to the exact solution vector  $x$  as  $k \rightarrow \infty$ . If the method is convergent, each iteration produces a better approximation to the exact solution. We repeat the iterations till the required accuracy is obtained. Therefore, in an iterative method the amount of computation depends on the desired accuracy whereas in direct methods the amount of computation is fixed. The number of iterations needed to obtain the desired accuracy also depends on the initial approximation, closer the initial approximation to the exact solution, faster will be the convergence.

Consider the system of equations

$$Ax = b \quad \dots (1)$$

where  $A$  is an  $n \times n$  non-singular matrix.

Writing the system in expanded form, we get

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n &= b_2 \\ \dots & \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n &= b_n \end{aligned} \tag{2}$$

We assume that the diagonal coefficients  $a_{ii} \neq 0, (i=1, \dots, n)$ . If some of  $a_{ii} = 0$ , then we rearrange the equations so that this condition holds. We then rewrite system (2) as

$$\begin{aligned} x_1 &= -\frac{1}{a_{11}}(a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n) + \frac{b_1}{a_{11}} \\ x_2 &= -\frac{1}{a_{22}}(a_{21}x_1 + a_{23}x_3 + \dots + a_{2n}x_n) + \frac{b_2}{a_{22}} \\ \dots & \\ x_n &= -\frac{1}{a_{nn}}(a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn-1}x_{n-1}) + \frac{b_n}{a_{nn}} \end{aligned} \tag{3}$$

In matrix form, system (3) can be written as  $x = Hx + c$

where

$$H = \begin{bmatrix} 0 & -\frac{a_{12}}{a_{11}} & -\frac{a_{13}}{a_{11}} & \dots & -\frac{a_{1n}}{a_{11}} \\ -\frac{a_{21}}{a_{22}} & 0 & -\frac{a_{23}}{a_{22}} & \dots & -\frac{a_{2n}}{a_{22}} \\ \dots & \dots & \dots & \dots & \dots \\ -\frac{a_{n1}}{a_{nn}} & -\frac{a_{n2}}{a_{nn}} & \dots & -\frac{a_{n,n-1}}{a_{nn}} & 0 \end{bmatrix} \tag{4}$$

and the elements of  $c$  are  $c_i = \frac{b_i}{a_{ii}} (i = 1, 2, \dots, n)$ .

To solve system (3) we make an initial guess  $x^{(0)}$  of the solution vector and substitute into the r.h.s. of Eqn. (3). The solution of Eqn. (3) will then yield a vector  $x^{(1)}$ , which hopefully is a better approximation to the solution than  $x^{(0)}$ . We then substitute  $x^{(1)}$  into the r.h.s. of Eqn. (3) and get another approximation,  $x^{(2)}$ . We continue in this manner until the successive iterations  $x^{(k)}$  have converged to the required number of significant figures.

In general we can write the iteration method for solving the linear system of Eqns. (1) in the form

$$x^{(k+1)} = Hx^{(k)} + c, \quad k = 0, 1, \dots \tag{5}$$

where  $x^{(k)}$  and  $x^{(k+1)}$  are the approximations to the solution vector  $x$  at the  $k$ th and the  $(k+1)$ th iterations respectively.  $H$  is called the iteration matrix and depends on  $A$ .  $c$  is a column vector and depends on both  $A$  and  $b$ . The matrix  $H$  is generally a constant matrix.

When the method (5) is convergent, then

$$\lim_{k \rightarrow \infty} x^{(k)} = \lim_{k \rightarrow \infty} x^{(k+1)} = x$$

and we obtain from Eqn. (5)

$$x = Hx + c \tag{6}$$

If we define the error vector at the  $k$ th iteration as

$$e^{(k)} = x^{(k)} - x \tag{7}$$

then subtracting Eqn. (6) from Eqn. (5), we obtain

$$e^{(k+1)} = H e^{(k)} \tag{8}$$

Thus, we get from Eqn. (8)

$$e^{(k)} = H e^{(k-1)} = H^2 e^{(k-2)} = \dots = H^k e^{(0)} \tag{9}$$

where  $e^{(0)}$  is the error in the initial approximate vector. Thus, for the convergence of the iterative method, we must have

$$\lim_{k \rightarrow \infty} e^{(k)} = 0$$

independent of  $e^{(0)}$ .

Before we discuss the above convergence criteria, let us recall the following definitions from linear algebra, MTE-02.

**Definition :** For a square matrix  $A$  of order  $n$ , and a number  $\lambda$  the value of  $\lambda$  for which the vector equation  $Ax = \lambda x$  has a non-trivial solution  $x \neq 0$ , is called an **eigenvalue** or **characteristic value** of the matrix  $A$ .

**Definition :** The largest eigenvalue in magnitude of  $A$  is called the **spectral radius** of  $A$  and is denoted by  $\rho(A)$ .

The eigenvalues of the matrix  $A$  are obtained from the characteristic equation

$$\det (A - \lambda I) = 0$$

which is an  $n$ th degree polynomial in  $\lambda$ . The roots of this polynomial  $\lambda_1, \lambda_2, \dots, \lambda_n$  are the eigenvalues of  $A$ . Therefore, we have

$$\rho(A) = \max |\lambda_i| \tag{10}$$

We now state a theorem on the convergence of the iterative methods.

**Theorem 1 :** An iteration method of the form (5) is convergent for arbitrary initial approximate vector  $x^{(0)}$  **if and only if**  $\rho(H) < 1$ .

We shall not be proving this theorem here as its proof makes use of advanced concepts from linear algebra and is beyond the scope of this course.

We define the rate of convergence as follows:

**Definition :** The number  $\nu = -\log_{10} \rho(H)$  is called, **the rate of convergence** of an iteration method.

Obviously, smaller the value of  $\rho(H)$ , larger is the value of  $\nu$ .

**Definition :** The method is said to have converged to  $m$  significant digits if  $\max |\epsilon_i^{(k)}| \leq 10^{-m}$ , that is, largest element in magnitude, of the error vector

$\epsilon^{(k)} \leq 10^{-m}$ . Also the number of iterations  $k$  that will be needed to make

$$\max |\epsilon_i^{(k)}| \leq 10^{-m}$$

is given by

$$k = \frac{m}{\nu} \tag{11}$$

Therefore, the number of iterations that are required to achieve the desired accuracy depends on  $\nu$ . For a method having higher rate of convergence, lesser number of iterations will be needed for a fixed accuracy and fixed initial approximation.

There is another convergence criterion for iterative methods which is based on the norm of a matrix.

The **norm** of a square matrix  $A$  of order  $n$  can be defined in the same way as we define the norm of an  $n$ -vector by comparing the size of  $Ax$  with the size of  $x$  (an  $n$ -vector) as follows:

i)  $\|A\|_2 = \max \frac{\|Ax\|_2}{\|x\|_2}$

based on the euclidean vector norm,  $\|x\|_2 = \sqrt{|x_1|^2 + |x_2|^2 + \dots + |x_n|^2}$

and

ii)  $\|A\|_\infty = \max \frac{\|Ax\|_\infty}{\|x\|_\infty}$ , based on the maximum vector norm,  $\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$

$\|A\|$  denotes the norm of  $A$ .

In (i) and (ii) above the maximum is taken over all (non zero)  $n$ -vectors. The most commonly used norms is the maximum norm  $\|A\|_\infty$ , as it is easier to calculate. It can be calculated in any of the following two ways:

$$\|A\|_\infty = \max_k \sum_i |a_{ik}| \text{ (maximum absolute column-sum)}$$

or

$$\|A\|_\infty = \max_i \sum_k |a_{ik}| \text{ (maximum absolute row sum)}$$

The norm of a matrix is a non-negative number which in addition to the property

$$\|AB\| \leq \|A\| \|B\|$$

satisfies all the properties of a vector norm, viz.,

- a)  $\|A\| \geq 0$  and  $\|A\| = 0$  iff  $A = 0$
- b)  $\|\alpha A\| = |\alpha| \|A\|$ , for all numbers  $\alpha$ .
- c)  $\|A+B\| \leq \|A\| + \|B\|$

where  $A$  and  $B$  are square matrices of order  $n$ .

We now state a theorem which gives the convergence criterion for iterative methods in terms of the norm of a matrix.

**Theorem 2 :** The iteration method of the form (5) for the solution of system (1) converges to the exact solution for any initial vector, if  $\|H\| < 1$ .

Also note that

$$\|H\| \geq \rho(H).$$

This can be easily proved by considering the eigenvalue problem  $Ax = \lambda x$ .

$$\text{Then } \|Ax\| = \|\lambda x\| = |\lambda| \|x\|$$

$$\text{or } |\lambda| \|x\| = \|Ax\| \leq \|A\| \|x\|$$

$$\text{i.e., } |\lambda| \leq \|A\| \text{ since } \|x\| \neq 0$$

Since this result is true for all eigenvalues, we have

$$\rho(A) \leq \|A\|.$$

The criterion given in Theorem 2 is **only a sufficient condition**, it is not necessary. Therefore, for a system of equations for which the matrix  $H$  is such that either

$$\max_k \sum_{i=1}^n |h_{ik}| < 1 \text{ or } \max_i \sum_{k=1}^n |h_{ik}| < 1, \text{ the iteration always converges, but if the}$$

condition is violated it is not necessary that the iteration diverges.

There is another **sufficient** condition for convergence as follows:

**Theorem 3 :** If the matrix  $A$  is strictly diagonally dominant that is,

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i=1,2,\dots,n.$$

then the iteration method (5) converges for any initial approximation  $x_{(0)}$ .

If no better initial approximation is known, we generally take  $x^{(0)} = 0$ .

We shall mostly use the criterion given in Theorem 1, which is both **necessary and sufficient**.

For using the iteration method (5), we need the matrix  $H$  and the vector  $c$  which depend on the matrix  $A$  and the vector  $b$ . The well-known iteration methods are based on the splitting of the matrix  $A$  in the form

$$A = D + L + U \tag{12}$$

where  $D$  is the diagonal matrix,  $L$  and  $U$  are respectively the lower and upper triangular matrices with zero diagonal elements. Based on the splitting (12), we now discuss two iteration methods of the form (5).

### 7.3 THE JACOBI ITERATION METHOD

We write the system of Eqn. (1) in the form (2), viz.,

$$a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1$$

$$a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2$$

$$\vdots$$

$$a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n$$





**Example 1 :** Perform four iterations of the Jacobi method for solving the system of equations

$$\begin{bmatrix} -8 & 1 & 1 \\ 1 & -5 & 1 \\ 1 & 1 & -4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 16 \\ 7 \end{bmatrix} \quad (15)$$

with  $x^{(0)} = 0$ , the exact solution is  $x = [-1 \ -4 \ -3]^T$ .

**Solution :** The Jacobi method when applied to the system of Eqns. (15) becomes

$$\begin{aligned} x_1^{(k+1)} &= \frac{1}{8} [x_2^{(k)} + x_3^{(k)} - 1] \\ x_2^{(k+1)} &= \frac{1}{5} [x_1^{(k)} + x_3^{(k)} - 16] \\ x_3^{(k+1)} &= \frac{1}{4} [x_1^{(k)} + x_2^{(k)} - 7], \quad k = 0, 1, \dots \end{aligned} \quad (16)$$

Starting with  $x^{(0)} = [0 \ 0 \ 0]^T$ , we obtain from Eqns (16), the following results:

**k = 0**

$$x_1^{(1)} = \frac{1}{8} [0 + 0 - 1] = -0.125$$

$$x_2^{(1)} = \frac{1}{5} [0 + 0 - 16] = -3.2$$

$$x_3^{(1)} = \frac{1}{4} [0 + 0 - 7] = -1.75$$

**k = 1**

$$x_1^{(2)} = \frac{1}{8} [-3.2 - 1.75 - 1] = -0.7438$$

$$x_2^{(2)} = \frac{1}{5} [-0.125 - 1.75 - 16] = -3.5750$$

$$x_3^{(2)} = \frac{1}{4} [-0.125 - 3.2 - 7] = -2.5813$$

**k = 2**

$$x_1^{(3)} = \frac{1}{8} [-3.5750 - 2.5813 - 1] = -0.8945$$

$$x_2^{(3)} = \frac{1}{5} [-0.7438 - 2.5813 - 16] = -3.8650$$

$$x_3^{(3)} = \frac{1}{4} [-0.7438 - 3.5750 - 7] = -2.8297$$

**k = 3**

$$x_1^{(4)} = \frac{1}{8} [-3.8650 - 2.8297 - 1] = -0.9618$$

$$x_2^{(4)} = \frac{1}{5} [-0.8945 - 2.8297 - 16] = -3.9448 \quad (17)$$

$$x_3^{(4)} = \frac{1}{4} [-0.8945 - 3.8650 - 7] = -2.9399$$

Thus, after four iterations we get the solution as given in Eqns (17). We find that after each iteration, we get better approximation to the exact solution.

**Example 2 :** Jacobi method is used to solve the system of equations

$$\begin{bmatrix} 4 & -1 & 1 \\ 4 & -8 & 1 \\ -2 & 1 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 7 \\ -21 \\ 15 \end{bmatrix} \quad (18)$$

Determine the rate of convergence of the method and the number of iterations needed to make  $\max |\epsilon_i^{(k)}| \leq 10^{-2}$

Perform these number of iterations starting with initial approximation  $\mathbf{x}^{(0)} = [1 \ 2 \ 2]^T$  and compare the result with the exact solution  $[2 \ 4 \ 3]^T$

**Solution :** The Jacobi method when applied to the system of Eqns. (18), gives the iteration matrix

$$\begin{aligned} H &= - \begin{bmatrix} \frac{1}{a_{11}} & 0 & 0 \\ 0 & \frac{1}{a_{22}} & 0 \\ 0 & 0 & \frac{1}{a_{33}} \end{bmatrix} \begin{bmatrix} 0 & a_{12} & a_{13} \\ a_{21} & 0 & a_{23} \\ a_{31} & a_{32} & 0 \end{bmatrix} \\ &= - \begin{bmatrix} \frac{1}{4} & 0 & 0 \\ 0 & \frac{1}{8} & 0 \\ 0 & 0 & \frac{1}{5} \end{bmatrix} \begin{bmatrix} 0 & -1 & 1 \\ 4 & 0 & 1 \\ -2 & 1 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 0 & \frac{1}{4} & -\frac{1}{4} \\ \frac{1}{2} & 0 & \frac{1}{8} \\ \frac{2}{5} & -\frac{1}{5} & 0 \end{bmatrix} \end{aligned}$$

The eigenvalues of the matrix H are the roots of the characteristic equation.

$$\det (H - \lambda I) = 0$$

Now

$$\det (H - \lambda I) = \begin{vmatrix} -\lambda & \frac{1}{4} & -\frac{1}{4} \\ \frac{1}{2} & -\lambda & \frac{1}{8} \\ \frac{2}{5} & -\frac{1}{5} & -\lambda \end{vmatrix} = \lambda^3 - \frac{3}{80} = 0$$

All the three eigenvalues of the matrix H are equal and they are equal to

$$\lambda = 0.3347$$

The spectral radius is

$$\rho(H) = 0.3347 \quad (19)$$

We obtain the rate of convergence as

$$\nu = -\log_{10}(0.3347) = 0.4753$$

The number of iterations needed for the required accuracy is given by

$$k = \frac{2}{\nu} \approx 5 \quad (20)$$

The Jacobi method when applied to the system of Eqns. (18) becomes

$$\mathbf{x}^{(k+1)} = \begin{bmatrix} 0 & \frac{1}{4} & -\frac{1}{4} \\ \frac{1}{2} & 0 & \frac{1}{8} \\ \frac{2}{5} & -\frac{1}{5} & 0 \end{bmatrix} \mathbf{x}^{(k)} + \begin{bmatrix} \frac{7}{4} \\ -\frac{21}{8} \\ 3 \end{bmatrix}, \quad k=0,1,\dots \quad (21)$$

starting with the initial approximation  $\mathbf{x}^{(0)} = [1 \ 2 \ 2]^T$ , we get from Eqn. (21)

$$\mathbf{x}^{(1)} = [1.75 \ 3.375 \ 3.0]^T$$

$$\mathbf{x}^{(2)} = [1.8437 \ 3.875 \ 3.025]^T$$

$$\mathbf{x}^{(3)} = [1.9625 \ 3.925 \ 2.9625]^T$$

$$\mathbf{x}^{(4)} = [1.9906 \ 3.9766 \ 3.0000]^T$$

$$\mathbf{x}^{(5)} = [1.9941 \ 3.9953 \ 3.0009]^T$$

which is the result after five iterations. Thus, you can see that result obtained after five iterations is quite close to the exact solution  $[2 \ 4 \ 3]^T$

**Example 3 :** Perform four iterations of the Jacobi method for solving the system of equations

$$\begin{bmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} \tag{22}$$

with  $x^{(0)} = [0.5 \ 0.5 \ 0.5 \ 0.5]^T$ . What can you say about the solution obtained if the exact solution is  $x = [1 \ 1 \ 1 \ 1]^T$ ?

**Solution :** The Jacobi method when applied to the system of Eqns. (22) becomes

$$\begin{aligned} x_1^{(k+1)} &= \frac{1}{2} [1 + x_2^{(k)}] \\ x_2^{(k+1)} &= \frac{1}{2} [x_1^{(k)} + x_3^{(k)}] \\ x_3^{(k+1)} &= \frac{1}{2} [x_2^{(k)} + x_4^{(k)}] \\ x_4^{(k+1)} &= \frac{1}{2} [1 + x_3^{(k)}], \quad k = 0, 1, \dots \end{aligned} \tag{23}$$

Using  $x^{(0)} = [0.5 \ 0.5 \ 0.5 \ 0.5]^T$ , we obtain

$$\begin{aligned} x^{(1)} &= [0.75 \ 0.5 \ 0.5 \ 0.75]^T \\ x^{(2)} &= [0.75 \ 0.625 \ 0.625 \ 0.75]^T \\ x^{(3)} &= [0.8125 \ 0.6875 \ 0.6875 \ 0.8125]^T \\ x^{(4)} &= [0.8438 \ 0.75 \ 0.75 \ 0.8438]^T \end{aligned}$$

You may notice here that the solution is improving after each iteration. Also, the solution obtained after four iterations is not a good approximation to the exact solution  $x = [1 \ 1 \ 1 \ 1]^T$ . This shows that we require a few more iterations to get a good approximation.

**Example 4 :** Find the spectral radius of the iteration matrix when the Jacobi method, is applied to the system of equations

$$\begin{bmatrix} 1 & 0 & 2 \\ 0 & 1 & -2 \\ 1 & -1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -1 \\ 5 \\ -3 \end{bmatrix}$$

Verify that the iterations do not converge to the exact solution  $x = [1 \ 3 \ -1]^T$ .

**Solution :** The iteration matrix H in this case becomes

$$\begin{aligned} H &= - \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 2 \\ 0 & 0 & -2 \\ 1 & -1 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 0 & 0 & -2 \\ 0 & 0 & 2 \\ -1 & 1 & 0 \end{bmatrix} \end{aligned}$$

and  $c = [-1 \ 5 \ -3]^T$ .

The eigenvalue of H are the roots of the characteristic equation

$\det (H - \lambda I) = 0$ . This gives us

$$-\lambda(\lambda^2 - 4) = 0$$

i.e.,  $\lambda = 0, \pm 2$

$$\therefore \rho(H) = 2 > 1.$$

Thus, the condition in Theorem 1 is violated. The iteration method does not converge.

We now perform few iterations and see what happens actually. Taking  $\mathbf{x}^{(0)} = \mathbf{0}$  and using the Jacobi method

$$\mathbf{x}^{(k+1)} = \begin{bmatrix} 0 & 0 & -2 \\ 0 & 0 & 2 \\ -1 & 1 & 0 \end{bmatrix} \mathbf{x}^{(k)} + \begin{bmatrix} -1 \\ 5 \\ -3 \end{bmatrix},$$

we obtain

$$\mathbf{x}^{(1)} = (-1 \ 5 \ -3)^T$$

$$\mathbf{x}^{(2)} = (5 \ -1 \ 3)^T$$

$$\mathbf{x}^{(3)} = (-7 \ 11 \ -9)^T$$

$$\mathbf{x}^{(4)} = (17 \ -13 \ 15)^T$$

$$\mathbf{x}^{(5)} = (-31 \ 35 \ -33)^T$$

and so on, which shows that the iterations are diverging fast. You may also try to obtain the solution with other initial approximations.

---

E1) Perform five iterations of the Jacobi method for solving the system of equations given in Example 4 with  $\mathbf{x}^{(0)} = [1 \ 1 \ 1]^T$ .

---

Let us now consider an example to show that the convergence criterion given in Theorem 3 is only a sufficient condition. That is, there are system of equations which are not diagonally dominant but, the Jacobi iteration method converges.

**Example 5 :** Perform iterations of the Jacobi method for solving the system of equations

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 2 & 0 \\ 0 & 3 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}$$

with  $\mathbf{x}^{(0)} = [0 \ 1 \ 1]^T$ . What can you say about the solution obtained if the exact solution is  $\mathbf{x} = [0 \ 1 \ 2]^T$ ?

**Solution :** The Jacobi method when applied to the given system of equations becomes

$$x_1^{(k+1)} = [3 - x_2^{(k)} - x_3^{(k)}]$$

$$x_2^{(k+1)} = 1$$

$$x_3^{(k+1)} = [-1 + 3x_2^{(k)}], \quad k=0,1,\dots$$

Using  $\mathbf{x}^{(0)} = [0 \ 1 \ 1]^T$ , we obtain

$$\mathbf{x}^{(1)} = [1 \ 1 \ 2]^T$$

$$\mathbf{x}^{(2)} = [0 \ 1 \ 2]^T$$

$$\mathbf{x}^{(3)} = [0 \ 1 \ 2]^T$$

You may notice here that the coefficient matrix is not diagonally dominant but the iterations converge to the exact solution after only two iterations.

And now a few exercises for you.

---

E2) Perform four iterations of the Jacobi method for solving the system of equations

$$\begin{bmatrix} 5 & 2 & 2 \\ 2 & 5 & 3 \\ 2 & 1 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ -6 \\ -4 \end{bmatrix}$$

with  $\mathbf{x}^{(0)} = \mathbf{0}$ . Exact solution is  $\mathbf{x} = (1 \ -1 \ -1)^T$

E3) Perform four iterations of the Jacobi method for solving the system of equations

$$\begin{bmatrix} 2 & -\frac{1}{2} & 0 \\ -\frac{3}{2} & 2 & -\frac{1}{2} \\ 0 & -\frac{3}{2} & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} \frac{3}{2} \\ 0 \\ \frac{1}{2} \end{bmatrix}$$

with  $x^{(0)} = 0$ . The exact solution is  $x = (1 \ 1 \ 1)^T$

E4) Perform four iterations of the Jacobi method for solving the system of equations

$$\begin{bmatrix} 5 & -1 & -1 & -1 \\ -1 & 10 & -1 & -1 \\ -1 & -1 & 5 & -1 \\ -1 & -1 & -1 & 10 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} -4 \\ 12 \\ 8 \\ 34 \end{bmatrix}$$

with  $x^{(0)} = 0$ . The exact solution is  $x = [1 \ 2 \ 3 \ 4]^T$

E5) Set up the Jacobi method in matrix form for solving the system of equations

$$\begin{bmatrix} 1 & 0 & -\frac{1}{4} & -\frac{1}{4} \\ 0 & 1 & -\frac{1}{4} & -\frac{1}{4} \\ -\frac{1}{4} & -\frac{1}{4} & 1 & 0 \\ -\frac{1}{4} & -\frac{1}{4} & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

and perform four iterations. Exact solution is  $x = (1 \ 1 \ 1 \ 1)^T$ . Take  $x^{(0)} = 0$ .

E6) Jacobi method is used to solve the system of equations

$$\begin{bmatrix} 2 & 0 & 1 \\ 0 & 3 & 2 \\ 1 & 2 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}$$

Determine the rate of convergence of the method and the number of iterations needed to make  $\max |\epsilon_i^{(k)}| \leq 10^{-2}$ . Perform four iterations and compare the result with the exact solution  $(1 \ 1 \ -1)^T$ .

We have already mentioned that iterative methods are usually applied to large linear systems with a sparse coefficient matrix. For sparse matrices, the number of non-zero entries is small, and hence the number of arithmetic operations to be performed per step is small. However, iterative methods may not always converge, and even when they converge, they may require a large number of iterations.

We shall now discuss the Gauss-Seidel method which is a simple modification of the method of simultaneous displacements and has improved rate of convergence.

## 7.4 THE GAUSS-SEIDEL ITERATION METHOD

Consider the system of Eqns. (2) written in form (3). For this system of equations, we define the Gauss-Seidel method as:

$$\begin{aligned}
 x_1^{(k+1)} &= -\frac{1}{a_{11}} (a_{12}x_2^{(k)} + a_{13}x_3^{(k)} + \dots + a_{1n}x_n^{(k)} - b_1) \\
 x_2^{(k+1)} &= -\frac{1}{a_{22}} (a_{21}x_1^{(k+1)} + a_{23}x_3^{(k)} + \dots + a_{2n}x_n^{(k)} - b_2) \\
 &\vdots \\
 x_n^{(k+1)} &= -\frac{1}{a_{nn}} (a_{n1}x_1^{(k+1)} + a_{n2}x_2^{(k+1)} + \dots + a_{n,n-1}x_{n-1}^{(k+1)} - b_n)
 \end{aligned}
 \tag{24}$$

or

$$x_i^{(k+1)} = -\frac{1}{a_{ii}} \left[ \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} + \sum_{j=i+1}^n a_{ij} x_j^{(k)} - b_i \right], i = 1, 2, \dots, n$$

You may notice here that in the first equation of system (24), we substitute the initial approximation  $(x_2^{(0)}, x_3^{(0)}, \dots, x_n^{(0)})$  on the right hand side. In the second equation, we substitute  $(x_1^{(1)}, x_3^{(0)}, \dots, x_n^{(0)})$  on the right hand side. In the third equation, we substitute  $(x_1^{(1)}, x_2^{(1)}, x_4^{(0)}, \dots, x_n^{(0)})$  on the right hand side. We continue in this manner until all the components have been improved. At the end of this first iteration, we will have an improved vector  $(x_1^{(1)}, x_2^{(1)}, \dots, x_n^{(1)})$ . The entire process is then repeated. In other words, the method uses an improved component as soon as it becomes available. It is for this reason the method is also called the method of successive displacements.

We can also write the system of Eqns. (24) as follows:

$$\begin{aligned}
 a_{11} x_1^{(k+1)} &= -a_{12} x_2^{(k)} - a_{13} x_3^{(k)} - \dots - a_{1n} x_n^{(k)} + b_1 \\
 a_{21} x_1^{(k+1)} + a_{22} x_2^{(k+1)} &= -a_{23} x_3^{(k)} - \dots - a_{2n} x_n^{(k)} + b_2 \\
 &\vdots \\
 a_{n1} x_1^{(k+1)} + a_{n2} x_2^{(k+1)} + \dots + a_{nn} x_n^{(k+1)} &= b_n
 \end{aligned}$$

In matrix form, this system can be written as

$$(D+L) x^{(k+1)} = -U x^{(k)} + b \tag{25}$$

where D is the diagonal matrix

$$D = \begin{bmatrix} a_{11} & & 0 \\ 0 & a_{22} & \\ \vdots & & \vdots \\ 0 & & a_{nn} \end{bmatrix}$$

and L and U are respectively the lower and upper triangular matrices with the zeros along the diagonal and are of the form

$$L = \begin{bmatrix} 0 & 0 & \dots & 0 \\ a_{21} & 0 & \dots & 0 \\ a_{31} & a_{32} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \quad U = \begin{bmatrix} 0 & a_{12} & a_{13} & \dots & a_{1n} \\ 0 & 0 & a_{23} & \dots & a_{2n} \\ 0 & 0 & 0 & \dots & a_{3n} \\ \vdots & & & & a_{n-1, n} \\ 0 & \dots & \dots & \dots & 0 \end{bmatrix}$$

From Eqn. (25), we obtain

$$x^{(k+1)} = -(D+L)^{-1} U x^{(k)} + (D+L)^{-1} b \tag{26}$$

which is of the form (5) with

$$H = -(D+L)^{-1} U \text{ and } c = (D+L)^{-1} b.$$

It may again be noted here, that if A is diagonally dominant then the iteration always converges.

Gauss-Seidel method will generally converge if the Jacobi method converges, and will converge at a faster rate. For symmetric A, it can be shown that

$$\rho(\text{Gauss-Seidel iteration method}) = [\rho(\text{Jacobi iteration method})]^2$$

Hence the rate of convergence of the Gauss-Seidel method is twice the rate of convergence of the Jacobi method. This result is usually true even when  $A$  is not symmetric.

We shall illustrate this fact through examples.

**Example 6 :** Perform four iterations (rounded to four decimal places) using the Gauss-Seidel method for solving the system of equations

$$\begin{bmatrix} -8 & 1 & 1 \\ 1 & -5 & 1 \\ 1 & 1 & -4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 16 \\ 7 \end{bmatrix} \quad (27)$$

with  $x^{(0)} = 0$ . The exact solution is  $x = (-1 \ -4 \ -3)^T$ .

**Solution :** The Gauss-Seidel method, for the system (25) is

$$\begin{aligned} x_1^{(k+1)} &= \frac{1}{8} [x_2^{(k)} + x_3^{(k)} - 1] \\ x_2^{(k+1)} &= \frac{1}{5} [x_1^{(k+1)} + x_3^{(k+1)} - 16] \\ x_3^{(k+1)} &= \frac{1}{4} [x_1^{(k+1)} + x_2^{(k+1)} - 7], \quad k = 0, 1, \dots \end{aligned} \quad (28)$$

Taking  $x^{(0)} = 0$ , we obtain the following iterations.

$k = 0$

$$x_1^{(1)} = \frac{1}{8} [0 + 0 - 1] = -0.125$$

$$x_2^{(1)} = \frac{1}{5} [-0.125 + 0 - 16] = -3.225$$

$$x_3^{(1)} = \frac{1}{4} [-0.125 - 3.225 - 7] = -2.5875$$

$k = 1$

$$x_1^{(2)} = \frac{1}{8} [-3.225 - 2.5875 - 1] = -0.8516$$

$$x_2^{(2)} = \frac{1}{5} [-0.8516 - 2.5875 - 16] = -3.8878$$

$$x_3^{(2)} = \frac{1}{4} [-0.8516 - 3.8878 - 7] = -2.9349$$

$k = 2$

$$x_1^{(3)} = \frac{1}{8} [-3.8878 - 2.9349 - 1] = -0.9778$$

$$x_2^{(3)} = \frac{1}{5} [-0.9778 - 2.9349 - 16] = -3.9825$$

$$x_3^{(3)} = \frac{1}{4} [-0.9778 - 3.9825 - 7] = -2.9901$$

$k = 3$

$$x_1^{(4)} = \frac{1}{8} [-3.9825 - 2.9901 - 1] = -0.9966$$

$$x_2^{(4)} = \frac{1}{5} [-0.9966 - 2.9901 - 16] = -3.9973$$

$$x_3^{(4)} = \frac{1}{4} [-0.9966 - 3.9973 - 7] = -2.9985$$

which is a good approximation to the exact solution  $x = (-1 \ -4 \ -3)^T$  with maximum absolute error 0.0034. Comparing with the results obtained in Example 1, we find that the values of  $x_i$ ,  $i=1,2,3$  obtained here are better approximates to the exact solution than the one obtained in Example 1.



**Example 7 :** Gauss-Seidel Method is used to solve the system of equations

$$\begin{bmatrix} 4 & -1 & 1 \\ 4 & -8 & 1 \\ -2 & 1 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 7 \\ -21 \\ 15 \end{bmatrix} \quad (29)$$

Determine the rate of convergence of the method and the number of iterations needed to make  $\max |\epsilon_i^{(k)}| \leq 10^{-2}$ . Perform these number of iterations with  $\mathbf{x}^{(0)} = [1 \ 2 \ 2]^T$  and compare the results with the exact solution  $\mathbf{x} = [2 \ 4 \ 3]^T$ .

**Solution :** The Gauss-Seidel method (26) when applied to the system of Eqns. (29) gives the iteration matrix.

$$H = - \begin{bmatrix} 4 & 0 & 0 \\ 4 & -8 & 0 \\ -2 & 1 & 5 \end{bmatrix}^{-1} \begin{bmatrix} 0 & -1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

Since the inverse of a lower triangular matrix is also a lower triangular matrix let

$$L = \begin{bmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} = \begin{bmatrix} 4 & 0 & 0 \\ 4 & -8 & 0 \\ -2 & 1 & 5 \end{bmatrix}^{-1}$$

then

$$\begin{bmatrix} 4 & 0 & 0 \\ 4 & -8 & 0 \\ -2 & 1 & 5 \end{bmatrix} \begin{bmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\therefore 4l_{11} = 1, \quad l_{11} = \frac{1}{4}$$

$$4l_{11} - 8l_{21} = 0, \quad l_{21} = \frac{1}{8}$$

$$-8l_{22} = 1, \quad l_{22} = -\frac{1}{8}$$

$$-2l_{11} + l_{21} + 5l_{31} = 0, \quad l_{31} = \frac{3}{40}$$

$$-l_{22} + 5l_{32} = 0, \quad l_{32} = \frac{1}{40}$$

$$5l_{33} = 1, \quad l_{33} = \frac{1}{5}$$

$$\therefore L = \begin{bmatrix} \frac{1}{4} & 0 & 0 \\ \frac{1}{8} & -\frac{1}{8} & 0 \\ \frac{3}{40} & \frac{1}{40} & -\frac{1}{5} \end{bmatrix}$$

Hence

$$H = \begin{bmatrix} -\frac{1}{4} & 0 & 0 \\ -\frac{1}{8} & \frac{1}{8} & 0 \\ -\frac{3}{40} & -\frac{1}{40} & \frac{1}{5} \end{bmatrix} \begin{bmatrix} 0 & -1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & \frac{1}{4} & \frac{1}{4} \\ 0 & \frac{1}{8} & 0 \\ 0 & \frac{3}{40} & -\frac{1}{10} \end{bmatrix}$$

The eigenvalues of the matrix H are the roots of the characteristic equation

$$\det(H - \lambda I) = \begin{vmatrix} -\lambda & \frac{1}{4} & \frac{1}{4} \\ 0 & \frac{1}{8} - \lambda & 0 \\ 0 & \frac{3}{40} & -(\frac{1}{10} + \lambda) \end{vmatrix} = 0$$

We have

$$\lambda(80\lambda^2 - 2\lambda - 1) = 0$$

which gives

$$\lambda = 0, 0.125, -0.1$$

Therefore, we have

$$\rho(H) = 0.125$$

The rate of convergence of the method is given by

$$\nu = -\log_{10}(0.125) = 0.9031$$

The number of iterations needed for obtaining the desired accuracy is given by

$$k = \frac{2}{\nu} \approx \frac{2}{0.9031} \approx 3$$

The Gauss-Seidel method when applied to the system of Eqns. (29) becomes

$$x_1^{(k+1)} = \frac{1}{4} [7 - x_3^{(k)} + x_2^{(k)}]$$

$$x_2^{(k+1)} = -\frac{1}{8} [-21 - 4x_1^{(k+1)} - x_3^{(k)}] \tag{30}$$

$$x_3^{(k+1)} = \frac{1}{5} [15 + 2x_1^{(k+1)} - x_2^{(k+1)}]$$

The successive iterations are obtained as

$$x^{(1)} = [1.75 \quad 3.75 \quad 2.95]^T$$

$$x^{(2)} = [1.95 \quad 3.9688 \quad 2.9863]^T$$

$$x^{(3)} = [1.9956 \quad 3.9961 \quad 2.9990]^T$$

which is an approximation to the exact solution after three iterations. Comparing the results obtained in Example 2, we conclude that the Gauss-Seidel method converges faster than the Jacobi method.

**Example 8 :** Use the Gauss-Seidel method for solving the following system of equations.

$$\begin{bmatrix} 2 & -1 & 0 & 1 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} \tag{31}$$

with  $x^{(0)} = [0.5 \quad 0.5 \quad 0.5 \quad 0.5]^T$ . Compare the results with those obtained in Example 3 after four iterations. The exact solution is  $x = [1 \quad 1 \quad 1 \quad 1]^T$ .

**Solution :** The Gauss-Seidel method, when applied to the system of Eqns. (31) becomes

$$x_1^{(k+1)} = \frac{1}{2} [1 + x_2^{(k)}]$$

$$x_2^{(k+1)} = \frac{1}{2} [x_1^{(k+1)} + x_3^{(k)}] \tag{32}$$

$$x_3^{(k+1)} = \frac{1}{2} [x_2^{(k+1)} + x_4^{(k)}]$$

$$x_4^{(k+1)} = \frac{1}{2} [1 + x_3^{(k+1)}], \quad k = 0, 1, \dots$$

Starting with the initial approximation  $\mathbf{x}^{(0)} = [0.5 \ 0.5 \ 0.5 \ 0.5]^T$ , we obtain the following iterates

$$\mathbf{x}^{(1)} = [0.75 \ 0.625 \ 0.5625 \ 0.7813]^T$$

$$\mathbf{x}^{(2)} = [0.8125 \ 0.6875 \ 0.7344 \ 0.8672]^T$$

$$\mathbf{x}^{(3)} = [0.8438 \ 0.7891 \ 0.8282 \ 0.9141]^T$$

$$\mathbf{x}^{(4)} = [0.8946 \ 0.8614 \ 0.8878 \ 0.9439]^T$$

In Example 3, the result obtained after four iterations by the Jacobi method was

$$\mathbf{x}^{(4)} = [0.8438 \ 0.75 \ 0.75 \ 0.8438]^T$$

**Remark :** The matrix formulations of the Jacobi and Gauss-Seidel methods are used whenever we want to check whether the iterations converges or to find the rate of convergence. If we wish to iterate and find solutions of the systems, we shall use the equation form of the methods.

And now a few exercises for you.

You may now attempt the following exercises.

- 
- E7) Perform four iterations of the Gauss-Seidel method for solving the system of equations given in E2).
- E8) Perform four iterations of the Gauss-Seidel method for solving the system of equations given in E3).
- E9) Perform four iterations of the Gauss-Seidel method for solving the system of equations given in E4).
- E10) Set up the matrix formulation of the Gauss-Seidel method for solving the system of equations given in E5). Perform four iterations of the method.
- E11) Gauss-Seidel method is used to solve the system of equations given in E6). Determine the rate of convergence and the number of iterations needed to make  $\max |e_i^{(k)}| \leq 10^{-2}$ . Perform four iterations and compare the results with the exact solution.
- 

We now end this unit by giving a summary of what we have covered in it.

---

## 7.5 SUMMARY

---

In this unit, we have covered the following:

- 1) Iterative methods for solving linear system of equations

$$A\mathbf{x} = \mathbf{b}$$

(see Eqn. (1))

where  $A$  is an  $n \times n$ , non-singular matrix. Iterative methods are generally used when the system is large and the matrix  $A$  is sparse. The process is started using an initial approximation and lead to successively better approximations.

- 2) General iterative method for solving the linear system of Eqn. (1) can be written in the form

$$\mathbf{x}^{(k+1)} = H\mathbf{x}^{(k)} + \mathbf{c}, \quad k = 0, 1, \dots \text{(see Eqn. (5))}$$

where  $\mathbf{x}^{(k)}$  and  $\mathbf{x}^{(k+1)}$  are the approximations to the solution vector  $\mathbf{x}$  at the  $k$ th and the  $(k+1)$ th iterations respectively.  $H$  is the iteration matrix which depends on  $A$  and is generally a constant matrix.  $\mathbf{c}$  is a column vector and depends on both  $A$  and  $\mathbf{b}$ .

- 3) Iterative method of the form given in 2) above converges for any initial vector, if  $\|H\| < 1$ , which is a sufficient condition for convergence. The necessary and sufficient condition for convergence is  $\rho(H) < 1$ , where  $\rho(H)$  is the spectral radius of  $H$ .
- 4) In the Jacobi iteration method or the method of simultaneous displacements.

$$H = -D^{-1}(L+U); \quad \mathbf{c} = D^{-1}\mathbf{b}$$

where  $D$  is a diagonal matrix,  $L$  and  $U$  are respectively the lower and upper triangular matrices with zero diagonal elements.

- 5) In the Gauss-Seidel iteration method or the method of successive displacements  
 $H = -(D + L)^{-1}U$  and  $c = (D + L)^{-1}b$ .
- 6) If the matrix  $A$  in Eqn. (1) is strictly diagonally dominant then the Jacobi and Gauss-Seidel methods converge. Gauss-Seidel method converges faster than the Jacobi method.

## 7.6 SOLUTIONS/ANSWERS

$$\begin{aligned} \text{E1) } \mathbf{x}^{(1)} &= (-3 \ 7 \ -3)^T \\ \mathbf{x}^{(2)} &= (5 \ -1 \ 7)^T \\ \mathbf{x}^{(3)} &= (-15 \ 19 \ -9)^T \\ \mathbf{x}^{(4)} &= (17 \ -13 \ 31)^T \\ \mathbf{x}^{(5)} &= (-63 \ 67 \ -33)^T \end{aligned}$$

Iterations do not converge.

$$\begin{aligned} \text{E2) } \mathbf{x}^{(1)} &= [0.2 \ -1.2 \ -0.8]^T \\ \mathbf{x}^{(2)} &= [1.0 \ -0.8 \ -0.64]^T \\ \mathbf{x}^{(3)} &= [0.776 \ -1.216 \ -1.04]^T \\ \mathbf{x}^{(4)} &= [1.1024 \ -0.8864 \ -0.8672]^T \end{aligned}$$

$$\begin{aligned} \text{E3) } \mathbf{x}^{(1)} &= [0.75 \ 0.0 \ 0.25]^T \\ \mathbf{x}^{(2)} &= [0.75 \ 0.625 \ 0.4375]^T \\ \mathbf{x}^{(3)} &= [0.9063 \ 0.6719 \ 0.7188]^T \\ \mathbf{x}^{(4)} &= [0.9180 \ 0.8594 \ 0.7539]^T \end{aligned}$$

$$\begin{aligned} \text{E4) } \mathbf{x}^{(1)} &= [-0.8 \ 1.2 \ 1.6 \ 3.4]^T \\ \mathbf{x}^{(2)} &= [0.44 \ 1.62 \ 2.36 \ 3.6]^T \\ \mathbf{x}^{(3)} &= [0.716 \ 1.84 \ 2.732 \ 3.842]^T \\ \mathbf{x}^{(4)} &= [0.8828 \ 1.9290 \ 2.8796 \ 3.9288]^T \end{aligned}$$

$$\text{E5) } \mathbf{x}^{(k+1)} = \begin{bmatrix} 0 & 0 & \frac{1}{4} & \frac{1}{4} \\ 0 & 0 & \frac{1}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{4} & 0 & 0 \\ \frac{1}{4} & \frac{1}{4} & 0 & 0 \end{bmatrix} \mathbf{x}^{(k)} + \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

$$\begin{aligned} \mathbf{x}^{(1)} &= [0.5 \ 0.5 \ 0.5 \ 0.5]^T \\ \mathbf{x}^{(2)} &= [0.75 \ 0.75 \ 0.75 \ 0.75]^T \\ \mathbf{x}^{(3)} &= [0.875 \ 0.875 \ 0.875 \ 0.875]^T \\ \mathbf{x}^{(4)} &= [0.9375 \ 0.9375 \ 0.9375 \ 0.9375]^T \end{aligned}$$

$$\text{E6) } H = \begin{bmatrix} 0 & 0 & -\frac{1}{2} \\ 0 & 0 & -\frac{2}{3} \\ -\frac{1}{4} & -\frac{1}{2} & 0 \end{bmatrix}$$

$$\rho(H) = \sqrt{\frac{11}{24}} = 0.6770 \text{ (spectral radius)}$$

$$\nu = 0.1694 \text{ (rate of convergence)}$$

$$k \approx 12 \text{ (number of iterations)}$$

$$\mathbf{x}^{(1)} = [0.5 \quad 0.3333 \quad -0.25]^T$$

$$\mathbf{x}^{(2)} = [0.625 \quad 0.500 \quad -0.5417]^T$$

$$\mathbf{x}^{(3)} = [0.7709 \quad 0.6945 \quad -0.6563]^T$$

$$\mathbf{x}^{(4)} = [0.8282 \quad 0.7709 \quad -0.7900]^T$$

$$\text{E7) } \mathbf{x}^{(1)} = [0.2 \quad -1.28 \quad -0.624]^T$$

$$\mathbf{x}^{(2)} = [0.9616 \quad -1.2102 \quad -0.9426]^T$$

$$\mathbf{x}^{(3)} = [1.0611 \quad -1.0589 \quad -1.0127]^T$$

$$\mathbf{x}^{(4)} = [1.0286 \quad -1.0038 \quad -1.0107]^T$$

$$\text{E8) } \mathbf{x}^{(1)} = [0.75 \quad 0.5625 \quad 0.6719]^T$$

$$\mathbf{x}^{(2)} = [0.8906 \quad 0.8359 \quad 0.8769]^T$$

$$\mathbf{x}^{(3)} = [0.9590 \quad 0.9385 \quad 0.9539]^T$$

$$\mathbf{x}^{(4)} = [0.9846 \quad 0.9769 \quad 0.9827]^T$$

$$\text{E9) } \mathbf{x}^{(1)} = [0.8 \quad 1.12 \quad 1.664 \quad 3.5984]^T$$

$$\mathbf{x}^{(2)} = [0.4765 \quad 1.7739 \quad 2.7698 \quad 3.9020]^T$$

$$\mathbf{x}^{(3)} = [0.8891 \quad 1.9561 \quad 2.9494 \quad 3.9795]^T$$

$$\mathbf{x}^{(4)} = [0.9770 \quad 1.9906 \quad 2.9894 \quad 3.9957]^T$$

$$\text{E10) } (\mathbf{D}+\mathbf{L})^{-1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ \frac{1}{4} & \frac{1}{4} & 1 & 0 \\ \frac{1}{4} & \frac{1}{4} & 0 & 1 \end{bmatrix}$$

$$\mathbf{H} = -(\mathbf{D}+\mathbf{L})^{-1} \mathbf{U} = \begin{bmatrix} 0 & 0 & \frac{1}{4} & \frac{1}{4} \\ 0 & 0 & \frac{1}{4} & \frac{1}{4} \\ 0 & 0 & \frac{1}{8} & \frac{1}{8} \\ 0 & 0 & \frac{1}{8} & \frac{1}{8} \end{bmatrix}$$

$$\mathbf{c} = (\mathbf{D}+\mathbf{L})^{-1} \mathbf{b} = \left[ \frac{1}{2} \quad \frac{1}{2} \quad \frac{3}{4} \quad \frac{3}{4} \right]^T$$

$$\mathbf{x}^{(k+1)} = \begin{bmatrix} 0 & 0 & \frac{1}{4} & \frac{1}{4} \\ 0 & 0 & \frac{1}{4} & \frac{1}{4} \\ 0 & 0 & \frac{1}{8} & \frac{1}{8} \\ 0 & 0 & \frac{1}{8} & \frac{1}{8} \end{bmatrix} \mathbf{x}^{(k)} + \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{3}{4} \\ \frac{3}{4} \end{bmatrix}$$

$$\mathbf{x}^{(1)} = [0.5 \quad 0.5 \quad 0.75 \quad 0.75]^T$$

$$\mathbf{x}^{(2)} = [0.875 \quad 0.875 \quad 0.9375 \quad 0.9375]^T$$

$$\mathbf{x}^{(3)} = [0.9688 \quad 0.9688 \quad 0.9844 \quad 0.9844]^T$$

$$\mathbf{x}^{(4)} = [0.9922 \quad 0.9922 \quad 0.9961 \quad 0.9961]^T$$

$$E11) (D+L)^{-1} = \begin{bmatrix} \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & 0 \\ -\frac{1}{8} & -\frac{1}{6} & \frac{1}{4} \end{bmatrix}$$

$$H = \begin{bmatrix} 0 & 0 & -\frac{1}{2} \\ 0 & 0 & -\frac{2}{3} \\ 0 & 0 & \frac{11}{24} \end{bmatrix}$$

$$\rho(H) = \frac{11}{24}$$

$$\rho = 0.3388$$

$$k = \frac{2}{\rho} \approx 6$$

$$x^{(1)} = [0.5 \quad 0.3333 \quad -0.5417]^T$$

$$x^{(2)} = [0.7709 \quad 0.6945 \quad -0.7900]^T$$

$$x^{(3)} = [0.8950 \quad 0.8600 \quad -0.9038]^T$$

$$x^{(4)} = [0.9519 \quad 0.9359 \quad -0.9559]^T$$

---

# UNIT 8 EIGENVALUES AND EIGENVECTORS

---

## Structure

- 8.1 Introduction
- 8.2 The Eigenvalue Problem
- 8.3 The Power Method
- 8.4 The Inverse Power Method
- 8.5 Summary
- 8.6 Solutions/Answers

---

## 8.1 INTRODUCTION

---

In Unit 7, you have seen that eigenvalues of the iteration matrix play a major role in the study of convergence of iterative methods for solving linear system of equations. Eigenvalues are also of great importance in many physical problems. The stability of an aircraft is determined by the location of the eigenvalues of a certain matrix in the complex plane. The natural frequencies of the vibrations of a beam are actually eigenvalues of a matrix. Thus the computation of the absolutely largest eigenvalue or smallest eigenvalue, or even all the eigenvalues of a given matrix is an important problem.

For a given system of equations of the form

$$Ax = \lambda x \quad (1)$$

or

$$(A - \lambda I)x = 0 \quad (2)$$

the values of the parameter  $\lambda$ , for which the system of Eqn. (2) has a nonzero solution, are called the **eigenvalues** of  $A$ . Corresponding to these eigenvalues, the nonzero solutions of Eqn. (2) i.e. the vectors  $x$ , are called the **eigenvectors** of  $A$ . The problem of finding the eigenvalues and the corresponding eigenvectors of a square matrix  $A$  is known as the **eigenvalue problem**. In this unit, we shall discuss the eigenvalue problem. To begin with, we shall give you some definitions and properties related to eigenvalues.

### Objectives

After studying this unit, you should be able to:

- Solve simple eigenvalue problems;
- Obtain the largest eigenvalue in magnitude and the corresponding eigenvector of a given matrix by using the power method;
- Obtain the smallest eigenvalue in magnitude and an eigenvalue closest to any chosen number along with the corresponding eigenvector of a given matrix by using the inverse power method.

---

## 8.2 THE EIGENVALUE PROBLEM

---

In the previous three units, we were concerned with the nonhomogeneous system of linear equations,  $Ax = b$ . We know that this system has a unique solution iff the matrix  $A$  is nonsingular. But, if the vector  $b = 0$ , then the system reduces to the homogeneous system

$$Ax = 0 \quad (3)$$

If the coefficient matrix  $A$ , in Eqn. (3) is nonsingular, then system has only the zero solution,  $x = 0$ . For the homogeneous system (3) to have a nonzero solution, the matrix  $A$  must be singular and in this case the solution is not unique (ref. Theorem 5, Unit 9, Block 3, MTE-02).

The homogeneous system of Eqn. (2) will have a nonzero solution only when the coefficient matrix  $(A - \lambda I)$  is singular, that is,

$$\det(A - \lambda I) = 0 \quad (4)$$

If the matrix  $A$  is an  $n \times n$  matrix then Eqn. (4) gives a polynomial of degree  $n$  in  $\lambda$ . This polynomial is called the characteristic equation of  $A$ . The  $n$  roots  $\lambda_1, \lambda_2, \dots, \lambda_n$  of this polynomial are the eigenvalues of  $A$ . For each eigenvalue  $\lambda_i$ , there exists a vector  $x_i$  (the eigenvector) which is the nonzero solution of the system of equations

$$(A - \lambda_i) x_i = 0 \quad (5)$$

The eigenvalues have a number of interesting properties. We shall now state and prove a few of these properties which we shall be using frequently.

**P1 :** A matrix  $A$  is singular if and only if it has a zero eigenvalue.

**Proof :** If  $A$  has a zero eigenvalue then

$$\det(A - 0 I) = 0$$

$$\Rightarrow \det(A) = 0$$

$$\Rightarrow A \text{ is singular.}$$

Conversely, if  $A$  is singular then

$$\det(A) = 0$$

$$\Rightarrow \det(A - 0 I) = 0$$

$$\Rightarrow 0 \text{ is an eigenvalue of the matrix } A.$$

**P2 :**  $A$  and  $A^T$  have the same eigenvalues.

**Proof :** If  $\lambda$  is an eigenvalue of  $A$  then

$$\det(A - \lambda I) = 0$$

$$\Rightarrow \det(A - \lambda I)^T = 0 \text{ (ref. P6 Sec. 9.3, Unit 9, Block 3, MTE-02)}$$

$$\Rightarrow \det(A^T - \lambda I^T) = 0 \text{ (Ref. Theorem 3, Sec. 7.3, Unit 7, Block 2, MTE-02)}$$

$$\Rightarrow \det(A^T - \lambda I) = 0$$

$$\Rightarrow \lambda \text{ is an eigenvalue of } A^T.$$

Hence the result.

However, the eigenvectors of  $A$  and  $A^T$  are not the same.

**P3 :** If the eigenvalues of a matrix  $A$  are  $\lambda_1, \lambda_2, \dots, \lambda_n$  then the eigenvalues of  $A^m$ ,  $m$  any positive integer, are  $\lambda_1^m, \lambda_2^m, \dots, \lambda_n^m$ . Also both the matrices  $A$  and  $A^m$  have the same set of eigenvectors.

**Proof :** Since  $\lambda_i$  ( $i = 1, 2, \dots, n$ ) are the eigenvalues of  $A$ , we have

$$Ax = \lambda_i x, \quad i = 1, 2, \dots, n \quad (6)$$

Premultiplying Eqn. (6) by  $A$  on both sides, we get

$$A^2x = A \lambda_i x = \lambda_i (Ax) = \lambda_i^2 x \quad (7)$$

which implies that  $\lambda_1^2, \lambda_2^2, \dots, \lambda_n^2$  are the eigenvalues of  $A^2$ . Further,  $A$  and  $A^2$  have the same eigenvectors. Premultiplying Eqn. (7)  $(m-1)$  times by  $A$  on both sides the general result follows.

**P4 :** If  $\lambda_1, \lambda_2, \dots, \lambda_n$  are the eigenvalues of  $A$ , then  $1/\lambda_1, 1/\lambda_2, \dots, 1/\lambda_n$  are the eigenvalues of  $A^{-1}$ . Also both the matrices  $A$  and  $A^{-1}$  have the same set of eigenvectors.

**Proof :** Since  $\lambda_i$  ( $i=1, 2, \dots, n$ ), are the eigenvalues of  $A$ , we have

$$Ax = \lambda_i x, \quad i = 1, 2, \dots, n \quad (8)$$

Premultiplying Eqn. (8) on both sides by  $A^{-1}$ , we get

$$A^{-1}Ax = \lambda_i A^{-1}x$$

which gives

$$x = \lambda_i A^{-1}x$$

$$\text{or } A^{-1}x = \frac{1}{\lambda_i} x$$

and hence the result.



**P5 :** If  $\lambda_1, \lambda_2, \dots, \lambda_n$  are the eigenvalues of  $A$ , then  $\lambda_i - q, i=1,2,\dots,n$  are the eigenvalues of  $A - qI$  for any real number  $q$ . Both the matrices  $A$  and  $A - qI$  have the same set of eigenvectors.

**Proof :** Since  $\lambda_i$  is an eigenvalue of  $A$ , we have

$$Ax = \lambda_i x, i = 1, 2, \dots, n \quad (9)$$

Subtracting  $q x$  from both sides of Eqn. (9), we get

$$Ax - qx = \lambda_i x - qx$$

which gives

$$(A - qI)x = (\lambda_i - q)x$$

and the result follows.

**P6 :** If  $\lambda_i, i = 1, 2, \dots, n$  are the eigenvalues of  $A$  then  $\frac{1}{\lambda_i - q}, i=1, 2, \dots, n$  are the eigenvalues of  $(A - qI)^{-1}$  for any real number  $q$ . Both the matrices  $A$  and  $(A - qI)^{-1}$  have the same set of eigenvectors.

P6 can be proved by combining P4 and P5. We leave the proof to you.

**E1) Prove P6**

We now give you a direct method of calculating the eigenvalues and eigenvectors of a matrix.

**Example 1 :** Find the eigenvalues of the matrix

$$\text{a) } A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}; \text{ b) } A = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 3 & 0 \\ 4 & 5 & 6 \end{bmatrix}$$

$$\text{c) } A = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \end{bmatrix}$$

**Solution :** a) Using Eqns. (4), we obtain the characteristic equations as

$$\det(A - \lambda I) = \begin{vmatrix} 1-\lambda & 0 & 0 \\ 0 & 2-\lambda & 0 \\ 0 & 0 & 3-\lambda \end{vmatrix} = 0$$

which gives  $(1-\lambda)(2-\lambda)(3-\lambda) = 0$ .

and hence the eigenvalues of  $A$  are  $\lambda_1=1, \lambda_2=2, \lambda_3=3$ .

$$\text{b) } \det(A - \lambda I) = \begin{vmatrix} 1-\lambda & 0 & 0 \\ 2 & 3-\lambda & 0 \\ 4 & 5 & 6-\lambda \end{vmatrix} = 0$$

which gives  $(1-\lambda)(3-\lambda)(6-\lambda) = 0$ .

Eigenvalues of  $A$  are  $\lambda_1=1, \lambda_2=3, \lambda_3=6$ .

$$\text{c) } \det(A - \lambda I) = \begin{vmatrix} 1-\lambda & 2 & 3 \\ 0 & 4-\lambda & 5 \\ 0 & 0 & 6-\lambda \end{vmatrix} = 0$$

Therefore,  $(1-\lambda)(4-\lambda)(6-\lambda) = 0$ .

Eigenvalues of  $A$  are  $\lambda_1=1, \lambda_2=4, \lambda_3=6$ .

**Remark :** Observe that in Example 1 (a), the matrix  $A$  is diagonal and in parts (b) and (c), it is lower and upper triangular respectively. In these cases the eigenvalues of  $A$  are the diagonal elements. This is true for any diagonal, lower triangular or upper triangular matrix. Formally, we give the result in the following theorem.

**Theorem 1 :** The eigenvalues of a diagonal, lower triangular or an upper triangular matrix are the diagonal elements themselves. Let us consider another example.

**Example 2 :** Find the eigenvalues and the corresponding eigenvectors of the matrices.

a)  $\begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix}$ ; b)  $A = \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix}$  and c)  $\begin{bmatrix} 1 & -2 \\ 2 & 1 \end{bmatrix}$

**Solution :** a) Using Eqns. (4), we obtain the characteristic equation as

$$|A - \lambda I| = \begin{vmatrix} 2 - \lambda & 2 \\ 1 & 3 - \lambda \end{vmatrix} = 0,$$

which gives the polynomial

$$\lambda^2 - 5\lambda + 4 = 0$$

i.e.,  $(\lambda - 1)(\lambda - 4) = 0$

The matrix A has two **distinct real eigenvalues**  $\lambda_1 = 1$ ,  $\lambda_2 = 4$ . To obtain the corresponding eigenvectors we solve the system of Eqns. (5) for each value of  $\lambda$ .

For  $\lambda = 1$ , we obtain the system of equations

$$x_1 + 2x_2 = 0$$

$$x_1 + 2x_2 = 0$$

which reduces to a single equation

$$x_1 + 2x_2 = 0$$

Taking  $x_2 = k$ , we get  $x_1 = -2k$ ,  $k$  being arbitrary nonzero constant. Thus, the eigenvector is of the form

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -2k \\ k \end{bmatrix} = k \begin{bmatrix} -2 \\ 1 \end{bmatrix}$$

For  $\lambda = 4$ , we obtain the system of equations

$$-2x_1 + 2x_2 = 0$$

$$x_1 - x_2 = 0$$

which reduces to a single equation

$$x_1 - x_2 = 0$$

Taking  $x_2 = k$ , we get  $x_1 = k$  and the corresponding eigenvector is

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = k \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

**Note:** In practice we usually omit  $k$  and say that  $[-2 \ 1]^T$  and  $[1 \ 1]^T$  are the eigenvectors of A corresponding to the eigenvalues  $\lambda = 1$  and  $\lambda = 4$  respectively. Moreover, the eigenvectors in this case are **linearly independent**.

b) The characteristic equation in this case becomes

$$(\lambda - 1)^2 = 0.$$

Therefore, the matrix A has a **repeated real eigenvalue**. The eigenvector corresponding to  $\lambda = 1$  is the solution of the system of Eqns. (5), which reduces to a single equation

$$x_2 = 0.$$

Taking  $x_1 = k$ , we obtain the eigenvector as

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = k \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

**Note:** that, in this case of repeated eigenvalues, we got linearly dependent eigenvectors.

c) The characteristic equation in this case becomes

$$\lambda^2 - 2\lambda + 5 = 0$$

which gives two complex eigenvalues  $\lambda = 1 \pm 2i$ .

The eigenvector corresponding to  $\lambda = 1+2i$  is the solution of the system of Eqns. (5). In this case we obtain the following equations

$$ix_1 + x_2 = 0$$

$$x_1 - ix_2 = 0$$

which reduces to the single equation

$$x_1 - ix_2 = 0$$

Taking  $x_2 = k$ , we get the eigenvector

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = k \begin{bmatrix} i \\ 1 \end{bmatrix}$$

Similarly, for  $\lambda = 1 - 2i$ , we obtain the eigenvector

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = k \begin{bmatrix} -i \\ 1 \end{bmatrix}$$

In the above problem you may note that corresponding to complex eigenvalues, we got complex eigenvectors. Let us now consider an example of  $3 \times 3$  matrix.

**Example 3 :** Determine the eigenvalues and the corresponding eigenvectors for the matrices

$$\text{a) } A = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}; \text{ b) } A = \begin{bmatrix} 6 & -2 & 2 \\ -2 & 3 & -1 \\ 2 & -1 & 3 \end{bmatrix}$$

**Solution :** a) The characteristic equation in this case becomes

$$\begin{vmatrix} 2-\lambda & -1 & 0 \\ -1 & 2-\lambda & -1 \\ 0 & -1 & 2-\lambda \end{vmatrix} = 0$$

which gives the polynomial

$$(2-\lambda)(\lambda^2 - 4\lambda + 2) = 0$$

Therefore, the eigenvalues of A are  $= 2, 2 + \sqrt{2}$  and  $2 - \sqrt{2}$ .

The eigenvector of A corresponding to  $\lambda = 2$  is the solution of the system of Eqns. (5), which reduces to

$$\begin{aligned} x_2 &= 0 \\ x_1 + x_3 &= 0 \end{aligned}$$

Taking  $x_3 = k$ , we obtain the eigenvector

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = k \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$$

The eigenvector of A corresponding to  $\lambda = 2 + \sqrt{2}$  is the solution of the system of equations

$$\begin{bmatrix} -\sqrt{2} & -1 & 0 \\ -1 & -\sqrt{2} & -1 \\ 0 & -1 & -\sqrt{2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad (10)$$

To find the solution of system of Eqns. (10), we use Gauss elimination method.

Performing  $R_2 - \frac{1}{\sqrt{2}}R_1$ , we get

$$\begin{bmatrix} -\sqrt{2} & -1 & 0 \\ 0 & -1/\sqrt{2} & -1 \\ 0 & -1 & -\sqrt{2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Again performing  $R_3 - \sqrt{2} R_2$ , we get

$$\begin{bmatrix} -\sqrt{2} & -1 & 0 \\ 0 & -1/\sqrt{2} & -1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

which give the equations

$$-\sqrt{2} x_1 - x_2 = 0$$

$$-x_2 - \sqrt{2} x_3 = 0$$

Taking  $x_3 = k$ , we obtain the eigenvector

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = k \begin{bmatrix} 1 \\ \sqrt{2} \\ 1 \end{bmatrix}$$

Similarly, corresponding to the eigenvalue  $\lambda = 2 - \sqrt{2}$ , the eigenvector is the solution of system of equations

$$\begin{bmatrix} -\sqrt{2} & -1 & 0 \\ -1 & \sqrt{2} & -1 \\ 0 & -1 & \sqrt{2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Using the Gauss elimination method, the system reduces to the equations

$$\sqrt{2} x_1 - x_2 = 0$$

$$x_2 - \sqrt{2} x_3 = 0$$

Taking  $x_3 = k$ , we obtain the eigenvector

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = k \begin{bmatrix} 1 \\ \sqrt{2} \\ 1 \end{bmatrix}$$

b) The characteristic equation in this case becomes

$$(\lambda - 8)(\lambda - 2)^2 = 0$$

Therefore, the matrix A has the real eigenvalues 8, 2 and 2. The eigenvalue 2 is repeated two times.

The eigenvector corresponding to  $\lambda = 8$  is solution of system of Eqns. (5), which reduces to

$$x_1 + x_2 - x_3 = 0$$

$$2x_1 + 5x_2 + x_3 = 0$$

$$2x_1 - x_2 - 5x_3 = 0$$

(11)

Subtracting the last equation of system (11) from the second equation we obtain the system of equations

$$x_1 + x_2 - x_3 = 0$$

$$x_2 + x_3 = 0$$

Taking  $x_3 = k$ , the eigenvector is

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = k \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix}$$

The eigenvector corresponding to  $\lambda = 2$  is the solution of system of Eqns. (5), which reduces to a single equation.

$$2x_1 - x_2 + x_3 = 0$$

(12)

We can take any values for  $x_1$  and  $x_2$  which need not be related to each other. The two linearly independent solutions can be written as:

$$k \begin{bmatrix} 1 \\ 0 \\ -2 \end{bmatrix} \text{ or } k \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

Note that in Eqn. (12), it is not necessary that we always assign values to  $x_1$  and  $x_2$ . We can assign values to any of the two variables and obtain the corresponding value of the third variable.

On the basis of Examples 2 and 3, we can make in general, the following observations:

For a given  $n \times n$  matrix  $A$ , the characteristic Eqn. (4) is a polynomial of degree  $n$  in  $\lambda$ . The  $n$  roots of this polynomial  $\lambda_1, \dots, \lambda_n$ , called the eigenvalues of  $A$  may be real or complex, distinct or repeated. Then,

- i) For distinct, real eigenvalues, we obtain linearly independent eigenvectors. (Examples 2(a), 3(a))
- ii) For a repeated eigenvalue, there may or may not be linearly independent eigenvectors. (Examples 2(b) and 3(b))
- iii) For a complex eigenvalue, we obtain a complex eigenvector.
- iv) An eigenvector is not unique. Any non-zero multiple of it is again an eigenvector.

How about trying some exercises now?

---

Determine the Eigenvalues and the corresponding eigenvectors of the following matrices.

$$E1) A = \begin{bmatrix} 1 & \sqrt{2} & 2 \\ \sqrt{2} & 3 & \sqrt{2} \\ 2 & \sqrt{2} & 1 \end{bmatrix}$$

$$E2) A = \begin{bmatrix} -15 & 4 & 3 \\ 10 & -12 & 6 \\ 20 & -4 & 2 \end{bmatrix}$$

$$E3) A = \begin{bmatrix} -2 & 2 & -3 \\ 2 & 1 & -6 \\ 1 & -2 & 0 \end{bmatrix}$$

$$E4) A = \begin{bmatrix} 2 & -1 & -1 \\ 3 & -2 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$


---

In the examples considered so far, it was possible for us to find all the roots of the characteristic equation exactly. But this may not always be possible. This is particularly true for  $n > 3$ . In such cases some iterative method like Newton-Raphson method may have to be used to find a particular eigenvalue or all the eigenvalues from the characteristic equation. However, in many practical problems, we do not require all the eigenvalues but need only a selected eigenvalue. For example, when we use iterative methods for solving a nonhomogeneous system of linear equations  $Ax = b$ , we need to know only the largest eigenvalue in magnitude of the iteration matrix  $H$ , to find out whether the method converges or not. One iterative method, which is frequently used to determine the largest eigenvalue in magnitude (also called the dominant eigenvalue) and the corresponding eigenvector for a given square matrix  $A$  is the power method. In this method we do not find the characteristic equation. This method is applicable only when all the eigenvalues are real and distinct. If the magnitude of two or more eigenvalues is the same then the method converges slowly.

---

### 8.3 THE POWER METHOD

---

Let us consider the eigenvalue problem

$$Ax = \lambda x.$$

Let  $\lambda_1, \lambda_2, \dots, \lambda_n$  be the  $n$  real and distinct eigenvalues of  $A$  such that

$$|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|$$

Therefore,  $\lambda_1$  is the dominant eigenvalue of  $A$ .

In this method, we start with an arbitrary nonzero vector  $y^{(0)}$  (not an eigenvector), and form a sequence of vectors ( $y^{(k)}$ )

$$y^{(k+1)} = Ay^{(k)}, k = 0, 1, \dots \quad (13)$$

In the limit as  $k \rightarrow \infty$ ,  $y^{(k)}$  converges to the eigenvector corresponding to the dominant eigenvalue of the matrix  $A$ . We can stop the iteration when the largest element in magnitude in  $y^{(k+1)} - y^{(k)}$  is less than the predefined error tolerance. For simplicity, we usually take the initial vector  $y^{(0)}$  with all its elements equal to one.

Vector for which scaling has been done is called a **scaled** vector otherwise, it is **unscaled**.

Note that in the process of multiplying the matrix  $A$  with the vector  $y^{(k)}$ , the elements of the vector  $y^{(k+1)}$  may become very large. To avoid this, we **normalize (or scale)** the vector  $y^{(k)}$  at each step by dividing  $y^{(k)}$  by its largest element in magnitude. This will make the largest element in magnitude in the vector  $y^{(k+1)}$  as one and the remaining elements less than one.

If  $y^{(k)}$  represents the unscaled vector and  $v^{(k)}$  the scaled vector then, we have the power method.

$$y^{(k+1)} = Av^{(k)} \quad (14)$$

$$v^{(k+1)} = \frac{1}{m_{k+1}} y^{(k+1)}, k = 0, 1, \dots \quad (15)$$

with,  $v^{(0)} = y^{(0)}$  and  $m_{k+1}$  being the largest element in magnitude of  $y^{(k+1)}$ . We then obtain the dominant eigenvalue by taking the limit

$$\lambda_1 = \lim_{k \rightarrow \infty} \frac{(y^{(k+1)})_r}{(v^{(k)})_r} \quad (16)$$

where  $r$  represents the  $r$ th component of that vector. Obviously, there are  $n$  ratios of numbers. As  $k \rightarrow \infty$  all these ratios tend to the same value, which is the largest eigenvalue in magnitude i.e.,  $\lambda_1$ . The iteration is stopped when the magnitude of the difference of any two ratios is less than the prescribed tolerance.

The corresponding eigenvector is then  $v^{(k+1)}$  obtained at the end of the last iteration performed.

We now illustrate the method through an example.

**Example 4 :** Find the dominant eigenvalue and the corresponding eigenvector correct to two decimal places of the matrix

$$A = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$

using the power method.

**Solution :** We take

$$y^{(0)} = v^{(0)} = (1 \ 1 \ 1)^T$$

Using Eqn. (14), we obtain

$$y^{(1)} = Av^{(0)} = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

$$\text{Now } m_1 = 1 \text{ and } v^{(1)} = \frac{1}{m_1} y^{(1)} = (1 \ 0 \ 1)^T.$$

Again,

$$y^{(2)} = Av^{(1)} = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ -2 \\ 2 \end{bmatrix}$$

$$m_2 = 2 \text{ and } v^{(2)} = \frac{1}{m_2} y^{(2)} = \frac{1}{2} y^{(2)} = (1 \ -1 \ 1)^T.$$

Proceeding in this manner, we have

$$y^{(3)} = Av^{(2)} = [3 \quad -4 \quad 3]^T$$

$$m_3 = 4$$

$$v^{(3)} = \frac{1}{4}y^{(3)} = [0.75 \quad -1 \quad 0.75]^T$$

$$y^{(4)} = Av^{(3)} = [2.5 \quad -3.5 \quad 2.5]^T$$

$$m_4 = 3.5$$

$$v^{(4)} = \frac{1}{3.5}y^{(4)} = [0.7143 \quad -1 \quad 0.7143]^T$$

$$y^{(5)} = Av^{(4)} = [2.4286 \quad -3.4286 \quad 2.4286]^T$$

$$m_5 = 3.4286$$

$$v^{(5)} = \frac{1}{3.4286}y^{(5)} = [0.7083 \quad -1 \quad 0.7083]^T$$

$$y^{(6)} = Av^{(5)} = [2.4166 \quad -3.4166 \quad 2.4166]^T$$

$$m_6 = 3.4166$$

$$v^{(6)} = \frac{1}{3.4166}y^{(6)} = [0.7073 \quad -1 \quad 0.7073]^T$$

$$y^{(7)} = Av^{(6)} = [2.4146 \quad -3.4146 \quad 2.4146]^T$$

$$m_7 = 3.4146$$

$$v^{(7)} = \frac{1}{3.4146}y^{(7)} = [0.7071 \quad -1 \quad 0.7071]^T$$

After 7 iterations, the ratios  $\frac{(y^{(7)})_r}{(v^{(6)})_r}$  are given as 3.4138, 3.4146 and 3.4138. The

maximum error in these ratios is 0.0008. Hence the dominant eigenvalue can be taken as 3.414 and the corresponding eigenvector is  $[0.7071 \quad -1 \quad 0.7071]^T$

Note that the exact dominant eigenvalue of A as obtained in Example 3 was  $2 + \sqrt{2} = 3.4142$  and the corresponding eigenvector was  $[1 \quad -\sqrt{2} \quad 1]^T$  which can also be written as  $[\frac{1}{\sqrt{2}} \quad -1 \quad \frac{1}{\sqrt{2}}]^T = [0.7071 \quad -1 \quad 0.7071]^T$

You may now try the following exercises.

Using four iterations of the power method and taking the initial vector  $y^{(0)}$  with all its elements equal to one, find the dominant eigenvalue and the corresponding eigenvector for the following matrices.

$$E5) A = \begin{bmatrix} 1 & \sqrt{2} & 2 \\ \sqrt{2} & 3 & \sqrt{2} \\ 2 & \sqrt{2} & 1 \end{bmatrix}$$

$$E6) A = \begin{bmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{bmatrix}$$

You must have realised that an advantage of the power method is that the eigenvector corresponding to the dominant eigenvalue is also generated at the same time. Usually, for most of the methods of determining eigenvalues, we need to do separate computations to obtain the eigenvector.

In some problems, the most important eigenvalue is the eigenvalue of least magnitude. We shall discuss now the inverse power method which gives the least eigenvalue in magnitude.

## 8.4 THE INVERSE POWER METHOD

We first note that if  $\lambda$  is the smallest eigenvalue in magnitude of  $A$ , then  $\frac{1}{\lambda}$  is the largest eigenvalue in magnitude of  $A^{-1}$ . The corresponding eigenvectors are same. If we apply the power method to  $A^{-1}$ , we obtain its largest eigenvalue and the corresponding eigenvector. This eigenvalue is then the smallest eigenvalue in magnitude of  $A$  and the eigenvector is same. Since power method is applied to  $A^{-1}$ , it is called the **inverse power method**.

Consider the method

$$y^{(k+1)} = A^{-1}v^{(k)}, k=0,1,2,\dots \quad (17)$$

$$v^{(k+1)} = \frac{1}{m_{k+1}}y^{(k+1)} \text{ with } v^{(0)} = y^{(0)}$$

where  $y^{(0)}$  is an arbitrary nonzero vector different from the eigenvector of  $A$ .

However, algorithm (17) is not in suitable form, as one has to find  $A^{-1}$ . Alternately, we write Eqn. (17) as

$$Ay^{(k+1)} = v^{(k)}$$

$$v^{(k+1)} = \frac{1}{m_{k+1}}y^{(k+1)}, k=0,1,2,\dots \quad (18)$$

We now need to solve a system of equations for  $y^{(k+1)}$ , which can be obtained using any of the method discussed in the previous units. The largest eigenvalue of  $A^{-1}$  is again given by

$$\mu = \lim_{k \rightarrow \infty} \frac{(y^{(k+1)})_r}{(v^{(k)})_r}$$

The corresponding eigenvector is  $v^{(k+1)}$ .

We now illustrate the method through an example.

**Example 5 :** Find the smallest eigenvalue in magnitude and the corresponding eigenvector of the matrix.

$$A = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$

using four iterations of the inverse power method.

**Solution :** Taking  $v^{(0)} = [1 \ 1 \ 1]^T$ , we write

**First iteration**

$$Ay^{(1)} = v^{(0)}$$

or

$$\begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad (19)$$

For solving the system of Eqns. (19), we use the LU decomposition method. We write

$$A = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix} = LU = \begin{bmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} \begin{bmatrix} 1 & u_{12} & u_{13} \\ 0 & 1 & u_{23} \\ 0 & 0 & 1 \end{bmatrix} \quad (20)$$

comparing the coefficients on both sides of Eqns. (20), we obtain

$$A = LU = \begin{bmatrix} 2 & 0 & 0 \\ -1 & \frac{3}{2} & 0 \\ 0 & -1 & \frac{4}{3} \end{bmatrix} \begin{bmatrix} 1 & -\frac{1}{2} & 0 \\ 0 & 1 & -\frac{2}{3} \\ 0 & 0 & 1 \end{bmatrix}$$



Solving  $Lz = v^{(0)}$

and then  $Uy^{(1)} = z$

we obtain

$$y^{(1)} = \left[ \frac{3}{2} \quad 2 \quad \frac{3}{2} \right]^T = \left[ 1.5 \quad 2.0 \quad 1.5 \right]^T$$

$$m_1 = 2.0$$

$$\therefore v^{(1)} = \frac{1}{m_1} y^{(1)} = \left[ 0.75 \quad 1.0 \quad 0.75 \right]^T$$

**Second iteration**

$$Ay^{(2)} = v^{(1)}$$

Solving  $Lz = v^{(1)}$

and  $Uy^{(2)} = z$

we obtain

$$y^{(2)} = \left[ 1.25 \quad 1.75 \quad 1.25 \right]^T$$

$$m_2 = 1.75$$

$$v^{(2)} = \frac{1}{m_2} y^{(2)} = \left[ 0.7143 \quad 1 \quad 0.7143 \right]^T$$

**Third iteration**

$$Ay^{(3)} = v^{(2)}$$

$$y^{(3)} = \left[ 1.2143 \quad 1.7143 \quad 1.2143 \right]^T$$

$$m_3 = 1.7143$$

$$v^{(3)} = \frac{1}{m_3} y^{(3)} = \left[ 0.7083 \quad 1 \quad 0.7083 \right]^T$$

**Fourth iteration**

$$Ay^{(4)} = v^{(3)}$$

$$y^{(4)} = \left[ 1.2083 \quad 1.7083 \quad 1.2083 \right]^T$$

$$m_4 = 1.7083$$

$$v^{(4)} = \frac{1}{m_4} y^{(4)} = \left[ 0.7073 \quad 1 \quad 0.7073 \right]^T$$

After 4 iterations, the ratios  $\frac{(y^{(4)})_r}{(v^{(3)})_r}$  are given as 1.7059, 1.7083, 1.7059. The maximum error in these ratios is 0.0024. Hence the dominant eigenvalue of  $A^{-1}$  can be taken as 1.70. Therefore,  $\frac{1}{1.70} = 0.5882$  is the smallest eigenvalue of  $A$  in magnitude and the corresponding eigenvector is given by  $[0.7073 \quad 1 \quad 0.7073]^T$ .

Note that the smallest eigenvalue in magnitude of  $A$  as calculated in Example 3 was  $2 - \sqrt{2} = 0.5858$  and the corresponding eigenvector was  $[1 \quad \sqrt{2} \quad 1]^T$  or  $[0.7071 \quad 1 \quad 0.7071]^T$ .

You may now try the following exercise :

E7) Find the smallest eigenvalue in magnitude and the corresponding eigenvector of the matrix

$$A) \begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix}$$

with  $v^{(0)} = [-1 \quad 1]^T$ , using four iterations of the inverse power method.

The inverse power method can be further generalized to find some other selected eigenvalues of  $A$ . For instance, one may be interested to find the eigenvalue of  $A$  which is nearest to some chosen number  $q$ . You know from P6 of Sec. 8.2 that the matrices  $A$  and  $A - qI$  have the same set of eigenvectors. Further, for each eigenvalue  $\lambda_i$  of  $A$ ,  $\lambda_i - q$  is the eigenvalue of  $A - qI$ .

We can therefore use the iteration

$$y^{(k+1)} = (A - qI)^{-1} v^{(k)} \quad (21)$$

with scaling as described in Eqns. (14) – (16). We determine the dominant eigenvalue  $\mu$  of  $(A - qI)^{-1}$  using the procedure given in Eqn. (18), i.e.

$$\begin{aligned} (A - qI) y^{(k+1)} &= v^{(k)} \\ v^{(k+1)} &= \frac{1}{m_{k+1}} y^{(k+1)} \end{aligned} \quad (22)$$

Using P6, we have the relation

$$\mu = \frac{1}{\lambda - q}, \text{ where } \lambda \text{ is an eigenvalue of } A.$$

$$\text{i.e., } \lambda = \frac{1}{\mu} + q \quad (23)$$

Now since  $\mu$  is the largest eigenvalue in magnitude of  $(A - qI)^{-1}$ ,  $\frac{1}{\mu}$  must be the smallest eigenvalue in magnitude of  $A - qI$ . Hence, the eigenvalue  $\frac{1}{\mu} + q$  of  $A$  is closest to  $q$ .

**Example 6 :** Find the eigenvalue of the matrix  $A$ , nearest to 3 and also the corresponding eigenvector using four iterations of the inverse power method where,

$$A = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$

**Solution :** In this case  $q = 3$ . Thus we have

$$A - 3I = \begin{bmatrix} -1 & -1 & 0 \\ -1 & -1 & -1 \\ 0 & -1 & -1 \end{bmatrix}$$

To find  $y^{(k+1)}$ , we need to solve the system

$$\begin{bmatrix} -1 & -1 & 0 \\ -1 & -1 & -1 \\ 0 & -1 & -1 \end{bmatrix} y^{(k+1)} = v^{(k)} \quad (24)$$

and normalise  $y^{(k+1)}$  as given in Eqn. (22).

**First iteration**

Starting with  $v^{(0)} = [1 \ 1 \ 1]^T$  and using the Gauss elimination method to solve the system (24), we obtain

$$y^{(1)} = [0 \ -1 \ 0]^T$$

$$m_1 = 1$$

$$v^{(1)} = \frac{1}{m_1} y^{(1)} = [0 \ -1 \ 0]^T$$

**Second iteration**

$$Ay^{(2)} = v^{(1)}$$

$$y^{(2)} = [1 \ -1 \ 1]^T$$

$$m_2 = 1$$

$$v^{(2)} = \frac{1}{m_2} y^{(2)} = [1 \ -1 \ 1]^T$$

**Third iteration**

$$Ay^{(3)} = v^{(2)}$$

$$y^{(3)} = [2 \ -3 \ 2]^T$$

$$m_3 = 3$$

$$v^{(3)} = \frac{1}{m_3} y^{(3)} = \left[ \frac{2}{3} \ -1 \ \frac{2}{3} \right]^T$$

## Fourth iteration

$$Ay^{(4)} = v^{(3)}$$

$$y^{(4)} = \begin{bmatrix} \frac{5}{3} & -\frac{7}{3} & \frac{5}{3} \end{bmatrix}^T$$

$$m_4 = \frac{7}{3} = 2.333$$

$$v^{(4)} = \frac{1}{m_4} y^{(4)} = \begin{bmatrix} \frac{5}{7} & -1 & \frac{5}{7} \end{bmatrix}^T$$

After four iterations, the ratios  $\frac{(y^{(4)})_r}{(v^{(3)})_r}$  are given as 2.5, 2.333, 2.5. The maximum error in these ratios is 0.1667. Hence the dominant eigenvalue of  $(A-3I)^{-1}$  can be taken as 2. Thus the eigenvalue  $\lambda$  of  $A$  closest to 3 as given by Eqn. (23) is

$$\begin{aligned} \lambda &= \frac{1}{\mu} + 3 \\ &= \frac{1}{2} + 3 = \frac{7}{2} = 3.5 \end{aligned}$$

and the corresponding eigenvector is  $v^{(4)} = \begin{bmatrix} \frac{5}{7} & -1 & \frac{5}{7} \end{bmatrix}^T = \begin{bmatrix} 0.7143 & -1 & 0.7143 \end{bmatrix}^T$

Note that the eigenvalue of  $A$  closest to 3 as obtained in Example 3 was  $2 + \sqrt{2} = 3.4142$ .

The eigenvector corresponding to this eigenvalue was  $\begin{bmatrix} 0.7071 & -1 & 0.7071 \end{bmatrix}^T$

And now a few exercises for you.

E8) Find the eigenvalue which is nearest to  $-1$  and the corresponding eigenvector for the matrix

$$A = \begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix}$$

with  $v^{(0)} = \begin{bmatrix} -1 & 1 \end{bmatrix}^T$ , using four iterations of the inverse power method.

E9) Using four iterations of the inverse power method, find the eigenvalue which is nearest to 5 and the corresponding eigenvector for the matrix

$$A = \begin{bmatrix} 3 & 2 \\ 3 & 4 \end{bmatrix} \text{ (exact eigenvalues are } = 1 \text{ and } 6)$$

$$\text{with } v^{(0)} = \begin{bmatrix} 1 & 1 \end{bmatrix}^T$$

The eigenvalues of a given matrix can also be estimated. That is, for a given matrix  $A$ , we can find the region in which all its eigenvalues lie. This can be done as follows:

Let  $\lambda_i$  be an eigenvalue of  $A$  and  $x_i$  be the corresponding eigenvector, i.e.,

$$Ax_i = \lambda_i x_i \quad (25)$$

or

$$\begin{aligned} a_{11}x_{i,1} + a_{12}x_{i,2} + \dots + a_{1n}x_{i,n} &= \lambda_i x_{i,1} \\ a_{21}x_{i,1} + a_{22}x_{i,2} + \dots + a_{2n}x_{i,n} &= \lambda_i x_{i,2} \\ &\vdots \\ a_{k1}x_{i,1} + a_{k2}x_{i,2} + \dots + a_{kn}x_{i,n} &= \lambda_i x_{i,k} \\ &\vdots \\ a_{n1}x_{i,1} + a_{n2}x_{i,2} + \dots + a_{nn}x_{i,n} &= \lambda_i x_{i,n} \end{aligned} \quad (26)$$

Let  $|x_{i,k}|$  be the largest element in magnitude of the vector  $[x_{i,1}, x_{i,2}, \dots, x_{i,n}]^T$ . Consider the  $k$ th equation of the system (26) and divide it by  $x_{i,k}$ . We then have

$$a_{k1} \left( \frac{x_{i,1}}{x_{i,k}} \right) + a_{k2} \left( \frac{x_{i,2}}{x_{i,k}} \right) + \dots + a_{kk} + \dots + a_{kn} \left( \frac{x_{i,n}}{x_{i,k}} \right) = \lambda_i \quad (27)$$

Taking the magnitudes on both sides of Eqn. (27), we get

$$\begin{aligned} |\lambda_i| &\leq |a_{k1}| \left| \frac{x_{i,1}}{x_{i,k}} \right| + |a_{k2}| \left| \frac{x_{i,2}}{x_{i,k}} \right| + \dots + |a_{kk}| + \dots + |a_{kn}| \left| \frac{x_{i,n}}{x_{i,k}} \right| \\ &\leq |a_{k1}| + |a_{k2}| + \dots + |a_{kk}| + \dots + |a_{kn}| \end{aligned} \quad (28)$$

since  $\left| \frac{x_{i,j}}{x_{i,k}} \right| \leq 1$  for  $j = 1, 2, \dots, n$ .

Since eigenvalues of  $A$  and  $A^T$  are same (Ref. P2), Eqn.(28) can also be written as

$$|\lambda_i| \leq |a_{1k}| + |a_{2k}| + \dots + |a_{kk}| + \dots + |a_{nk}| \quad (29)$$

Since  $|x_{i,k}|$ , the largest element in magnitude, is unknown, we approximate Eqns.(28) and (29) by

$$|\lambda| \leq \max_i \left[ \sum_{j=1}^n |a_{ij}| \right] \text{ (maximum absolute row sum)} \quad (30)$$

and

$$|\lambda| \leq \max_j \left[ \sum_{i=1}^n |a_{ij}| \right] \text{ (maximum absolute column sum)} \quad (31)$$

We can also rewrite Eqn. (27) in the form

$$|\lambda_i - a_{kk}| = a_{k1} \left( \frac{x_{i,1}}{x_{i,k}} \right) + a_{k2} \left( \frac{x_{i,2}}{x_{i,k}} \right) + \dots + a_{kn} \left( \frac{x_{i,n}}{x_{i,k}} \right)$$

and taking magnitude on both sides, we get

$$|\lambda_i - a_{kk}| \leq \sum_{\substack{j=1 \\ i \neq j}}^n |a_{ij}| \quad (32)$$

Again, since  $A$  and  $A^T$  have the same eigenvalues Eqn.(32) can be written as

$$|\lambda_i - a_{kk}| \leq \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}| \quad (33)$$

A square matrix  $A$  is symmetric if  $A=A^T$ .

Note that since the eigenvalues can be complex, the bounds (30), (31), (32) and (33) represents circles in the complex plane. If the eigenvalues are real, then they represent intervals. For example, when  $A$  is symmetric then the eigenvalues of  $A$  are real.

Again in Eqn. (32), since  $k$  is not known, we replace the circle by the union of the  $n$  circles

$$|\lambda - a_{ii}| \leq \sum_{\substack{j=1 \\ i \neq j}}^n |a_{ij}|, \quad i = 1, 2, \dots, n. \quad (34)$$

Similarly from Eqn. (33), we have that eigenvalues of  $A$  lie in the union of circles

$$|\lambda - a_{jj}| \leq \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}|, \quad j = 1, 2, \dots, n. \quad (35)$$

The bounds derived in Eqns. (30), (31), (34) and (35) for eigenvalues are all independent bounds. Hence the eigenvalues must lie in the intersection of these bounds. The circles derived above are called the Gerschgorin circles and the bounds are called the Gerschgorin bounds.

Let us now consider the following examples:

**Example 7 :** Estimate the eigenvalues of the matrix

$$A = \begin{bmatrix} 1 & -1 & 2 \\ 2 & 1 & 3 \\ -1 & 3 & 2 \end{bmatrix}$$

using the Gerschgorin bounds.

**Solution :** The eigenvalues of A lie in the following regions:

i) absolute row sums are 4, 6 and 6. Hence

$$|\lambda| \leq \max [4, 6, 6] = 6 \quad (36)$$

ii) absolute column sums are 4, 5 and 7. Hence

$$|\lambda| \leq 7 \quad (37)$$

iii) union of the circles [using (34)]

$$|\lambda - 1| \leq 3$$

$$|\lambda - 1| \leq 5$$

$$|\lambda - 2| \leq 4$$

iv) union of the circles [using (35)]

$$|\lambda - 1| \leq 3$$

$$|\lambda - 1| \leq 4$$

$$|\lambda - 2| \leq 5$$

union of circles in (iii) is  $|\lambda - 1| \leq 5$  (38)

union of circles in (iv) as  $|\lambda - 2| \leq 5$  (39)

The eigenvalues lie in all the circles (36), (37), (38) and (39) i.e., in the intersection of these circles as shown by shaded region in Fig. 1.

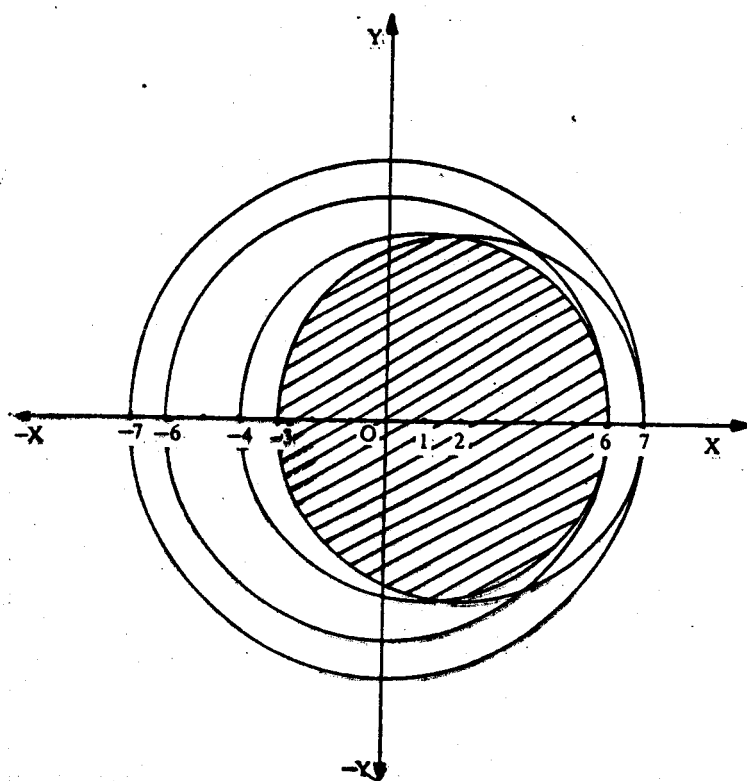


Fig. 1

**Example 8 :** Estimate the eigenvalues of the symmetric matrix

$$A = \begin{bmatrix} 1 & -1 & 2 \\ -2 & 1 & 2 \\ 2 & 2 & -2 \end{bmatrix}$$

by the Gerschgorin bounds.

**Solution** The eigenvalues lie in the following regions:

- i)  $|\lambda| \leq \max [4, 4, 6] = 6$
- ii) union of the circles
  - a)  $|\lambda - 1| \leq 3$
  - b)  $|\lambda - 1| \leq 3$
  - c)  $|\lambda + 2| \leq 4$

Since  $A$  is symmetric, it has real eigenvalues. Therefore, the eigenvalues lie in the intervals

- i)  $-6 \leq \lambda \leq 6$
  - ii) union of
    - a)  $-3 \leq \lambda - 1 \leq 3$ , i.e.  $-2 \leq \lambda \leq 4$
    - b)  $-4 \leq \lambda + 2 \leq 4$ , i.e.  $-6 \leq \lambda \leq 2$
- union of (a) and (c) is  $-6 \leq \lambda \leq 4$ .

Intersection of (i) and (ii) is  $-6 \leq \lambda \leq 4$ . Hence the eigenvalues of  $A$  lie in the interval  $-6 \leq \lambda \leq 4$ .

**Note** that in Example 8, since the matrix  $A$  is symmetric, the bounds (30) and (31) are same and also the bounds (34) and (35) are same.

You may now try the following exercise.

---

E10) Estimate the eigenvalues of the matrix  $A$  given in Example 3(a) and 3(b), using the Gerschgorin bounds.

---

We now end this unit by giving a summary of what we have covered in it.

## 8.5 SUMMARY

In this unit we have covered the following:

- 1) For a given system of equations of the form
 
$$Ax = \lambda x \quad (\text{see Eqn. (1)})$$
 the values of  $\lambda$  for which Eqn. (1) has a nonzero solution are called the eigenvalues and the corresponding nonzero solutions (which are not unique) are called the eigenvectors of the matrix  $A$ .
- 2) The following are the steps involved in solving an eigenvalue problem
  - i) Find the  $n$ th degree polynomial (called the characteristic equation) in  $\lambda$  from  $\det (A - \lambda I) = 0$ .
  - ii) Find the  $n$  roots  $\lambda_i$ ,  $i = 1, 2, \dots, n$  of the characteristic equation.
  - iii) Find the eigenvectors corresponding to each  $\lambda_i$ .
- 3) For  $n \geq 3$ , it may not be possible to find the roots of the characteristic equation exactly. In such cases, we use some iterative method like Newton Raphson method to find these roots. However,

- i) when only the largest eigenvalue in magnitude is to be obtained, we use the power method. In this method we obtain a sequence of vectors  $\{y^{(k)}\}$ , using the iterative scheme

$$y^{(k+1)} = A y^{(k)}, k=0, 1, \dots \text{ (see Eqn. (13))}$$

which in the limit as  $k \rightarrow \infty$ , converges to the eigenvector corresponding to the dominant eigenvalue of the matrix  $A$ . The vector  $y^{(0)}$  is an arbitrary non-zero vector (different from the eigenvector of  $A$ ).

- ii) we use the inverse power method with the iteration scheme

$$y^{(k+1)} = (A - qI)^{-1} v^{(k)},$$

$$\text{i.e., } (A - qI) y^{(k+1)} = v^{(k)}, k = 0, 1, 2, \dots$$

where  $y^{(0)} = v^{(0)}$  is an arbitrary non-zero vector (not an eigenvector)

- a) with  $q = 0$ , if only the least eigenvalue of  $A$  in magnitude and the corresponding eigenvector are to be obtained and  
 b) with any  $q$ , if the eigenvalue of  $A$ , nearest to some chosen number  $q$  and the corresponding eigenvector are to be obtained.

## 8.6 SOLUTIONS/ANSWERS

E1) Characteristic equation :  $\lambda^3 - 5\lambda^2 - \lambda + 5 = 0$

eigenvalues:  $-1, 1, 5$

eigenvectors:  $[-1, 0, 1]^T; [1, -\sqrt{2}, 1]^T; [1, \sqrt{2}, 1]^T$

E2) Characteristic equation:  $\lambda^3 + 25\lambda^2 + 50\lambda - 1000 = 0$

eigenvalues :  $-20, -10, 5$

eigenvectors :  $[-1, 1/2, 1]^T; [-1, -2, 1]^T; [1/4, 1/2, 1]^T$

E3) Characteristic equation:  $\lambda^3 + \lambda^2 - 21\lambda - 45 = 0$

eigenvalues:  $-3, -3, 5$

eigenvectors:  $[1, 0, 1/3]^T; [0, 1, 2/3]^T; [-1, -2, 1]^T$

E4) Characteristic equation:  $\lambda^3 - \lambda^2 - \lambda + 1 = 0$

eigenvalues:  $-1, 1, 1$

eigenvectors:  $[1/3, 1, 0]^T; [1, 1, 0]^T; [1, 1, 0]^T$

E5)  $y^{(1)} = [4.4142, 5.8284, 4.4142]^T; m_1 = 5.8284$

$$v^{(1)} = [0.7574, 1, 0.7574]^T$$

$y^{(2)} = [3.6864, 5.1422, 3.6864]^T; m_2 = 5.1422$

$$v^{(2)} = [0.7169, 1, 0.7169]^T$$

$y^{(3)} = [3.5649, 5.0276, 3.5649]^T; m_3 = 5.0276$

$$v^{(3)} = [0.7090, 1, 0.7090]^T$$

$y^{(4)} = [3.5412, 5.0054, 3.5412]^T; m_4 = 5.0053$

$$v^{(4)} = [0.7075, 1, 0.7075]^T$$

After 4 iterations the ratios  $\frac{(y^{(4)})_r}{(v^{(3)})_r}$  are given by 4.9946, 5.0054, 4.9946. The

maximum error in these ratios is 0.0108. Thus the dominant eigenvalue of  $A$  can be taken as 5.00 and the corresponding eigenvector is  $[0.7075, 1, 0.7075]^T$

E6)  $y^{(1)} = [1, 0, 0, 1]^T; m_1 = 1$

$$v^{(1)} = [1, 0, 0, 1]^T$$

$y^{(2)} = [2, -1, -1, 2]^T; m_2 = 2$

$$v^{(2)} = [1, -0.5, -0.5, 1]^T$$

$y^{(3)} = [2.5, -1.5, -1.5, 2.5]^T; m_3 = 2.5$

$$v^{(3)} = [1, -0.6, -0.6, 1]^T$$

$y^{(4)} = [2.6, -1.6, -1.6, 2.6]^T; m_4 = 2.6$

$$v^{(4)} = [1, -0.6154, -0.6154, 1]^T$$

After 4 iterations the ratios  $\frac{(y^{(4)})_r}{(v^{(3)})_r}$  are given by 2.6, 2.6667, 2.6667, 2.6. The maximum error in these ratios is 0.0667. Thus the dominant eigenvalue of A can be taken as 2 and the corresponding eigenvector is  $[1 \ -0.6154 \ -0.6154 \ 1]^T$

E7) Starting with  $v^{(0)} = [-1 \ 1]^T$  and solving  $Ay^{(1)} = v^{(0)}$ , we get

$$y^{(1)} = \begin{bmatrix} -\frac{5}{4} & \frac{3}{4} \end{bmatrix}^T; m_1 = \frac{5}{4} = 1.25$$

$$v^{(1)} = \begin{bmatrix} -1 & \frac{3}{5} \end{bmatrix}^T$$

$$y^{(2)} = \begin{bmatrix} -\frac{21}{20} & \frac{11}{20} \end{bmatrix}^T; m_2 = \frac{21}{20} = 1.05$$

$$v^{(2)} = \begin{bmatrix} -1 & \frac{11}{21} \end{bmatrix}^T$$

$$y^{(3)} = \begin{bmatrix} -\frac{37}{42} & \frac{16}{42} \end{bmatrix}^T; m_3 = \frac{37}{42} = 0.8810$$

$$v^{(3)} = \begin{bmatrix} -1 & \frac{16}{37} \end{bmatrix}^T$$

$$y^{(4)} = \begin{bmatrix} -\frac{143}{148} & \frac{69}{148} \end{bmatrix}^T; m_4 = \frac{143}{148} = 0.9662$$

$$v^{(4)} = \begin{bmatrix} -1 & \frac{69}{143} \end{bmatrix}^T$$

After 4 iterations, the ratios  $\frac{(y^{(4)})_r}{(v^{(3)})_r}$  are 0.9662, 1.0781. The maximum error in these ratios is 0.1119. Hence the dominant eigenvalue of  $A^{-1}$  can be taken as 0.9. The smallest eigenvalue of A is therefore  $\frac{1}{0.9} = 1.1111$ .

The corresponding eigenvector is  $\begin{bmatrix} -1 & \frac{69}{143} \end{bmatrix}^T$

$$E8) [A + I] = \begin{bmatrix} 6 & 4 \\ 1 & 4 \end{bmatrix};$$

Starting with  $v^{(0)} = [-1 \ 1]^T$  and solving  $\begin{bmatrix} 3 & 2 \\ 1 & 4 \end{bmatrix} y^{(1)} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$

we get

$$y^{(1)} = \begin{bmatrix} -\frac{6}{10} & \frac{4}{10} \end{bmatrix}^T; m_1 = \frac{3}{5} = 0.6$$

$$v^{(1)} = \begin{bmatrix} -1 & \frac{2}{3} \end{bmatrix}^T$$

Similarly,

$$y^{(2)} = \begin{bmatrix} -\frac{8}{15} & \frac{3}{10} \end{bmatrix}^T; m_2 = \frac{8}{15} = 0.5333$$

$$v^{(2)} = \begin{bmatrix} -1 & \frac{9}{16} \end{bmatrix}^T$$

$$y^{(3)} = \begin{bmatrix} -\frac{41}{80} & \frac{43}{160} \end{bmatrix}^T; m_3 = \frac{41}{80} = 0.5125$$

$$v^{(3)} = \begin{bmatrix} -1 & \frac{43}{82} \end{bmatrix}^T$$



$$y^{(4)} = \left[ -\frac{207}{410} \quad \frac{211}{820} \right]^T; m_4 = \frac{207}{410} = 0.5049$$

$$v^{(4)} = \left[ -1 \quad \frac{211}{414} \right]^T$$

After 4 iterations, the ratios  $\frac{(y^{(4)})_r}{(v^{(3)})_r}$  are 0.5049, 0.4907. The maximum error in these ratios is 0.0142. Hence the dominant eigenvalue of  $(A+I)^{-1}$  can be taken as  $\mu = 0.5$ . The eigenvalue of  $A$  which is nearest to  $-1$  is obtained from

$$\lambda = \frac{1}{\mu} - 1 = \frac{1}{0.5} - 1 = 1$$

The corresponding eigenvector is  $\left[ -1 \quad \frac{211}{414} \right]^T$

$$E9) [A-5I] = \begin{bmatrix} -2 & 2 \\ 3 & -1 \end{bmatrix}$$

$$\text{Starting with } v^{(0)} = [1 \quad 1]^T \text{ and solving } \begin{bmatrix} -2 & 2 \\ 3 & -1 \end{bmatrix} y^{(1)} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

we get

$$y^{(1)} = \left[ \frac{3}{4} \quad \frac{5}{4} \right]^T; m_1 = \frac{5}{4} = 1.25$$

$$v^{(1)} = \left[ \frac{3}{5} \quad 1 \right]^T$$

Similarly,

$$y^{(2)} = \left[ \frac{13}{20} \quad \frac{19}{20} \right]^T; m_2 = \frac{19}{20} = 0.95$$

$$v^{(2)} = \left[ \frac{13}{19} \quad 1 \right]^T$$

$$y^{(3)} = \left[ \frac{51}{76} \quad \frac{77}{76} \right]^T; m_3 = \frac{77}{76} = 1.0132$$

$$v^{(3)} = \left[ \frac{51}{77} \quad 1 \right]^T$$

$$y^{(4)} = \left[ \frac{205}{308} \quad \frac{307}{308} \right]^T; m_4 = \frac{307}{308} = 0.9968$$

$$v^{(4)} = \left[ \frac{205}{307} \quad 1 \right]^T$$

After 4 iterations, the ratios  $\frac{(y^{(4)})_r}{(v^{(3)})_r}$  are 1.005, 0.9968. The maximum error in these ratios is 0.0082. Hence the dominant eigenvalue of  $(A-5I)^{-1}$  can be taken as  $\mu = 0.99$ .

The eigenvalue of  $A$  which is nearest to 5 is obtained from

$$\lambda = \frac{1}{\mu} + 5$$

$$= \frac{1}{0.99} + 5 = 6.0101$$

The corresponding eigenvector is  $\left[ \frac{205}{307} \quad 1 \right]^T$

Solution of Linear Algebraic Equations E10) a) The eigenvector of A lie in the following regions:

i)  $|\lambda| \leq \max [3, 4, 3] = 4$

ii) union of the circles

a)  $|\lambda - 2| \leq 1$

b)  $|\lambda - 2| \leq 2$

c)  $|\lambda - 2| \leq 1$

Since A is symmetric, it has real eigenvalues.

Therefore, the eigenvalues lie in the intervals

i)  $-4 \leq \lambda \leq 4$ .

ii) union of

a)  $-1 \leq \lambda - 2 \leq 1, \quad 1 \leq \lambda \leq 3$

b)  $-2 \leq \lambda - 2 \leq 2, \quad 0 \leq \lambda \leq 4$

union of (a) and (b) is  $1 \leq \lambda \leq 4$ .

Intersection of (i), (ii) is  $1 \leq \lambda \leq 4$ . Hence the eigenvalues of A lie in the interval  $1 \leq \lambda \leq 4$ .

b) The eigenvalues of A lie in the intervals

i)  $-10 \leq \lambda \leq 10$

ii) union of

a)  $-4 \leq \lambda - 6 \leq 4, \quad 2 \leq \lambda \leq 10$

b)  $-3 \leq \lambda - 3 \leq 3, \quad 0 \leq \lambda \leq 6$

union of (a) and (b) is  $2 \leq \lambda \leq 10$ .

Intersection of (i) and (ii) is  $2 \leq \lambda \leq 10$ . Hence the eigenvalues of A lie in the interval  $2 \leq \lambda \leq 10$ .

**NOTES**

**NOTES**



Uttar Pradesh  
Rajarshi Tandon Open University

UGMM - 10

## NUMERICAL ANALYSIS

Block

# 3

### INTERPOLATION

---

#### UNIT 9

Lagrange's Form

5

---

#### UNIT 10

Newton Form of the Interpolating Polynomial

16

---

#### UNIT 11

Interpolation at Equally Spaced Points

30

---

---

## Course Design Committee

---

Prof. S.R.K. Iyengar (*Editor*)  
Dept. of Mathematics  
I.I.T., Delhi.

Prof. J.P. Agarwal  
Dept. of Mathematics  
I.I.T., Kharagpur

Dr. U. Anantha Krishnaiah  
Dept. of Mathematics  
KREC, Surath Kal

Prof. R.K. Jain  
Dept. of Mathematics  
I.I.T., Delhi

Prof. C. Prabhakara Rao  
Dept. of Mathematics  
REC, Warangal

Faculty Members  
School of Sciences, IGNOU

Prof. R.K. Bose  
Dr. V.D. Madan  
Dr. Poornima Mital  
Dr. Manik Patwardhan  
Dr. Parvin Sinclair  
Dr. Sujatha Varma

---

## Block Preparation Team

---

Prof. S.R.K. Iyengar (*Editor*)  
Dept. of Mathematics  
I.I.T., Delhi

Prof. C. Prabhakara Rao  
Dept. of Mathematics  
REC, Warangal

Prof. R.K. Bose  
School of Sciences  
IGNOU

Course Coordinator : Dr. Poornima Mital

---

## Production

---

Mr. Balakrishna Selvaraj  
Registrar (PPD)  
IGNOU

---

## Acknowledgements

Dr. Poornima Mital and Dr. Parvin Sinclair for comments on the manuscript.  
Ms. Kiran for typing the manuscript.

October, 1993

© Indira Gandhi National Open University, 1993

ISBN-81-7263-482-X

*All rights reserved. No part of this work may be reproduced in any form, by mimeograph or any other means, without permission in writing from the Indira Gandhi National Open University.*

*Further information on the Indira Gandhi National Open University courses may be obtained from the University's office at Maidan Garhi, New Delhi-110 068.*

Reproduced and reprinted with the permission of Indira Gandhi National Open University  
by Dr.A.K.Singh, Registrar, U.P.R.T.Open University, Allahabad (May, 2013)  
Reprinted by : Nitin Printers, 1 Old Katra, Manmohan Park, Allahabad.

---

## BLOCK 3 INTERPOLATION

---

The interpolation has been defined as the art of reading between the lines of a table, and in elementary mathematics the term usually denotes the process of computing intermediate values of a function from a set of given values of that function. For example, consider the table that lists the population of Delhi. The population census is taken every ten years and the table gives population for the years 1901, 1911, 1961, 1971, 1981 and 1991 in Delhi. We would like to know whether this table could be used to estimate the population of Delhi in 1936 say or even in 1996. Such estimates of population can be made using a function that fits the given data.

The general problem of interpolation, however, is much more complex than this. In higher mathematics we often have to deal with functions whose analytical form is either totally unknown or else is of such a nature (complicated or otherwise) that the function cannot easily be subjected to certain operations like differentiation and integration etc. In either case, it is desirable to replace the given function by another which can be more readily handled.

In this block, we study the polynomial interpolation in detail. We derive various forms of the interpolating polynomial. Polynomials are used as the basic means of approximation in nearly all areas of numerical analysis. One major reason for their importance is that they (uniformly) approximate continuous functions, that is, given any continuous function defined on a closed and bounded interval  $[a, b]$  there exists a polynomial that is "close" to the given function. Another important reason for considering the class of polynomials in the approximation of functions is that the derivative and indefinite integral of any polynomial are easy to calculate and the results are again polynomials. You must have encountered one application of polynomial approximation in Taylor polynomials, in your calculus course. But Taylor polynomials, have the property that all information used in the approximation is concentrated at one point. For general computational purpose, it is more efficient to use methods that uses information at various points. In the sequel, it is the construction of this type of polynomials which we are going to consider.

We have divided our discussion on polynomial interpolation into 3 units. In Unit 9, we discuss the Lagrange form of interpolating polynomial to prove the existence and uniqueness of the interpolating polynomial for unequally spaced nodes. Also the general expression for the error of polynomial interpolation is proved which gives the estimates of the error in polynomial approximation.

In the next unit, Unit 10, we deal with another very useful form of interpolating polynomial called the Newton form of interpolating polynomial. Newton's Form is derived using the concept of divided differences. We also obtain another expression for the error term for the interpolating polynomial in terms of divided differences.

In Unit 11, we deal with some useful forms of interpolating polynomials for equally spaced nodes like Newton's backward and Newton's forward difference forms, and Gauss-Stirling central difference interpolating polynomial after introducing the concepts of forward, backward and central differences. In last block (Block 4), you will come across the use of polynomial interpolation in numerical differentiation, integration etc.

## NOTATIONS AND SYMBOLS

$f[x_i]$	Zeroeth divided difference
$f[x_i, x_{i+1}]$	Divided difference of order 1
$f[x_i, x_{i+1}, \dots, x_{i+k}]$	Divided difference of order k
$\delta$	Delta
$E_2(\bar{x})$	Interpolation error of f at $\bar{x}$

Also see the list given in Blocks 1 and 2.



---

# UNIT 9 LAGRANGE'S FORM

---

## Structure

- 9.1 Introduction
  - Objectives
- 9.2 Lagrange's Form
- 9.3 Inverse Interpolation
- 9.4 General Error Term
- 9.5 Summary
- 9.6 Solutions/Answers

---

## 9.1 INTRODUCTION

---

Let  $f$  be a real-valued function defined on the interval  $[a, b]$  and we denote  $f(x_k)$  by  $f_k$ . Suppose that the values of the function  $f(x)$  are given to be  $f_0, f_1, f_2, \dots, f_n$  when  $x = x_0, x_1, x_2, \dots, x_n$  respectively where  $x_0 < x_1 < x_2 \dots < x_n$  lying in the interval  $[a, b]$ . The function  $f(x)$  may not be known to us. The technique of determining an approximate value of  $f(x)$  for a non-tabular value of  $x$  which lies in the interval  $[a, b]$  is called interpolation. The process of determining the value of  $f(x)$  for a value of  $x$  lying outside the interval  $[a, b]$  is called extrapolation. In this unit, we derive a polynomial  $P(x)$  of degree  $\leq n$  which agrees with the values of  $f(x)$  at the given  $(n + 1)$  distinct points, called nodes or abscissas. In other words, we can find a polynomial  $P(x)$  such that  $P(x_j) = f_j, j = 0, 1, 2, \dots, n$ . Such a polynomial  $P(x)$  is called the interpolating polynomial of  $f(x)$ .

In Section 9.2 we prove the existence of an interpolating polynomial by actually constructing one such polynomial having the desired property. The uniqueness is proved by invoking the corollary of the fundamental theorem of Algebra. In Section 9.3 we derive general expression for error in approximating the function by the interpolating polynomial at a point and this allows us to calculate a bound on the error over an interval. In proving this we make use of the general Rolle's theorem.

### Objectives

After reading this unit, you should be able to:

- find the Lagrange's form of interpolating polynomial interpolating  $f(x)$  at  $n + 1$  distinct nodal points;
- compute the approximate value of  $f$  at a non-tabular point;
- compute the value of  $\bar{x}$  (approximately) given a number  $\bar{y}$  such that  $f(\bar{x}) = (\bar{y})$  (inverse interpolation);
- compute the error committed in interpolation, if the function is known, at a non-tabular point of interest;
- find an upper bound in the magnitude of the error.

---

## 9.2 LAGRANGE'S FORM

---

Let us recall the fundamental theorem of algebra and its useful corollaries.

**Theorem 1:** If  $P(x)$  is a polynomial of degree  $n \geq 1$ , that is,  $P(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$ , with  $a_0, \dots, a_n$  real or complex numbers and  $a_n \neq 0$ , then  $P(x)$  has at least one zero, that is, there exists a real or complex number  $\xi$  such that  $P(\xi) = 0$ .

**Lemma 1:** If  $z_1, z_2, \dots, z_k$  are distinct zeros of the polynomial  $P(x)$ , then

$$P(x) = (x - z_1)(x - z_2) \dots (x - z_k)R(x)$$

for some polynomial  $R(x)$ .

**Corollary:** If  $P_k(x)$  and  $Q_k(x)$  are two polynomials of degree  $\leq k$  which agree at the  $k + 1$  distinct points  $z_0, z_1, \dots, z_k$  then  $P_k(x) = Q_k(x)$  identically.

You have come across Rolle's theorem in Section 1.2. We need a generalized version of this theorem in the Section 9.4 (General Error Term). This is stated below.

**Theorem 2: (Generalized Rolle's Theorem).** Let  $f$  be a real-valued function defined on  $[a, b]$  which is  $n$  times differentiable on  $]a, b[$ . If  $f$  vanishes at the  $n + 1$  distinct points  $x_0, \dots, x_n$  in  $[a, b]$ , then a number  $c$  in  $]a, b[$  exists such that  $f^{(n)}(c) = 0$ .

We now show the existence of an interpolating polynomial and also show that it is unique. The form of the interpolating polynomial that we are going to discuss in this section is called the Lagrange form of the interpolating polynomial. We start with a relevant theorem.

**Theorem 3:** Let  $x_0, x_1, \dots, x_n$  be  $n + 1$  distinct points on the real line and let  $f(x)$  be a real-valued function defined on some interval  $I = [a, b]$  containing these points. Then, there exists exactly one polynomial  $P_n(x)$  of degree  $\leq n$ , which interpolates  $f(x)$  at  $x_0, \dots, x_n$ , that is,  $P_n(x_i) = f(x_i)$ ,  $i = 0, 1, 2, \dots, n$ .

**Proof:** First we discuss the uniqueness of the interpolating polynomial, and then exhibit one explicit construction of an interpolating polynomial (Lagrange's Form).

Let  $P_n(x)$  and  $Q_n(x)$  be two distinct interpolating polynomials of degree  $\leq n$ , which interpolate  $f(x)$  at  $(n + 1)$  distinct points  $x_0, x_1, \dots, x_n$ . Let  $h(x) = P_n(x) - Q_n(x)$ . Note that  $h(x)$  is also a polynomial of degree  $\leq n$ . Also

$$h(x_i) = P_n(x_i) - Q_n(x_i) = f(x_i) - f(x_i) = 0, \quad i = 0, 1, 2, \dots, n.$$

That is,  $h(x)$  has  $(n + 1)$  distinct zeros. But  $h(x)$  is of degree  $\leq n$  and from the Corollary to Lemma 1, we have  $h(x) \equiv 0$ . That is  $P_n(x) \equiv Q_n(x)$ . This proves the uniqueness of the polynomial.

Since the data is given at the points  $(x_0, f_0), (x_1, f_1), \dots, (x_n, f_n)$  let the required polynomial be written as

$$P_n(x) = L_0(x)f_0 + L_1(x)f_1 + \dots + L_n(x)f_n = \sum_{i=0}^n L_i(x) f_i \quad (1)$$

Setting  $x = x_j$  in (1), we get

$$P_n(x_j) = \sum_{i=0}^n L_i(x_j) f_i \quad (2)$$

Since this polynomial fits the data exactly, we must have

$$L_j(x_j) = 1$$

and  $L_i(x_j) = 0, \quad i \neq j$

or  $L_i(x_j) = \delta_{ij} \quad (3)$

The polynomials  $L_i(x)$  which are of degree  $\leq n$  are called the Lagrange fundamental

polynomials. It is easily verified that these polynomials are given by

$$L_i(x) = \frac{(x - x_0)(x - x_1) \dots (x - x_{i-1})(x - x_{i+1}) \dots (x - x_n)}{(x_i - x_0)(x_i - x_1) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_i - x_n)}$$

$$= \left[ \prod_{\substack{j=0 \\ j \neq i}}^n (x - x_j) \right] / \left[ \prod_{\substack{j=0 \\ j \neq i}}^n (x_i - x_j) \right] \quad (4)$$

Substitution of (4) in (1) gives the required Lagrange form of the interpolating polynomial.

**Remark:** The Lagrange form (Eqn. (1)) of interpolating polynomial makes it easy to show the existence of an interpolating polynomial. But its evaluation at a point  $x_i$  involves a lot of computation.

A more serious drawback of the Lagrange form arises in practice due to the following: One calculates a linear polynomial  $P_1(x)$ , a quadratic polynomial  $P_2(x)$  etc., by increasing the number of interpolation points, until a satisfactory approximation  $P_k(x)$  to  $f(x)$  has been found. In such a situation Lagrange form does not take any advantage of the availability of  $P_{k-1}(x)$  in calculating  $P_k(x)$ . Later on, we shall see how in this respect, Newton form, discussed in the next unit, is more useful.

Let us consider some examples to construct this form of interpolation polynomials.

**Example 1:** If  $f(1) = -3$ ,  $f(3) = 9$ ,  $f(4) = 30$  and  $f(6) = 132$ , find the Lagrange's interpolation polynomial of  $f(x)$ .

**Solution:** We have  $x_0 = 1$ ,  $x_1 = 3$ ,  $x_2 = 4$ ,  $x_3 = 6$  and  $f_0 = -3$ ,  $f_1 = 9$ ,  $f_2 = 30$ ,  $f_3 = 132$ .

The Lagrange's interpolating polynomial  $P(x)$  is given by

$$P(x) = L_0(x) f_0 + L_1(x) f_1 + L_2(x) f_2 + L_3(x) f_3 \quad (5)$$

where

$$L_0(x) = \frac{(x - x_1)(x - x_2)(x - x_3)}{(x_0 - x_1)(x_0 - x_2)(x_0 - x_3)}$$

$$= \frac{(x - 3)(x - 4)(x - 6)}{(1 - 3)(1 - 4)(1 - 6)}$$

$$= -\frac{1}{30} (x^3 - 13x^2 + 54x - 72)$$

$$L_1(x) = \frac{(x - x_0)(x - x_2)(x - x_3)}{(x_1 - x_0)(x_1 - x_2)(x_1 - x_3)}$$

$$= \frac{(x - 1)(x - 4)(x - 6)}{(3 - 1)(3 - 4)(3 - 6)}$$

$$= \frac{1}{6} (x^3 - 11x^2 + 34x - 24)$$

$$L_2(x) = \frac{(x - x_0)(x - x_1)(x - x_3)}{(x_2 - x_0)(x_2 - x_1)(x_2 - x_3)}$$

$$= \frac{(x - 1)(x - 4)(x - 6)}{(4 - 1)(4 - 3)(4 - 6)}$$

$$= -\frac{1}{6} (x^3 - 10x^2 + 27x - 18)$$

$$L_3(x) = \frac{(x - x_0)(x - x_1)(x - x_2)}{(x_3 - x_0)(x_3 - x_1)(x_3 - x_2)}$$

$$= \frac{(x - 1)(x - 3)(x - 4)}{(6 - 1)(6 - 3)(6 - 4)}$$

$$= \frac{1}{30} (x^3 - 8x^2 + 19x - 12)$$

Substituting  $L_j(x)$  and  $f_j$ ,  $j = 0, 1, 2, 3$  in Eqn. (5), we get

$$\begin{aligned} P(x) &= -\frac{1}{30} [x^3 - 13x^2 + 54x - 72] (-3) + \frac{1}{6} [x^3 - 11x^2 + 34x - 24] (9) \\ &\quad - \frac{1}{6} [x^3 - 10x^2 + 27x - 18] (30) + \frac{1}{30} [x^3 - 8x^2 + 19x - 12] (132) \\ &= \frac{1}{10} [x^3 - 13x^2 + 54x - 72] + \frac{3}{2} [x^3 - 11x^2 + 34x - 24] \\ &\quad - 5 [x^3 - 10x^2 + 27x - 18] + \frac{22}{5} [x^3 - 8x^2 + 19x - 12] \end{aligned}$$

which gives on simplification

$$P(x) = x^3 - 3x^2 + 5x - 6$$

which is the Lagrange's interpolating polynomial of  $f(x)$ .

**Example 2:** Using Lagrange's interpolation formula, find the value of  $f$  when  $x = 1.4$  from the following table.

x	1.2	1.7	1.8	2.0
f	3.3201	5.4739	6.0496	7.3891

**Solution:** The Lagrange's interpolating formula with 4 points is

$$\begin{aligned} P(x) &= \frac{(x - x_1)(x - x_2)(x - x_3)}{(x_0 - x_1)(x_0 - x_2)(x_0 - x_3)} f_0 + \frac{(x - x_0)(x - x_2)(x - x_3)}{(x_1 - x_0)(x_1 - x_2)(x_1 - x_3)} f_1 + \\ &\quad \frac{(x - x_0)(x - x_1)(x - x_3)}{(x_2 - x_0)(x_2 - x_1)(x_2 - x_3)} f_2 + \frac{(x - x_0)(x - x_1)(x - x_2)}{(x_3 - x_0)(x_3 - x_1)(x_3 - x_2)} f_3 \end{aligned} \quad (6)$$

Substituting

$$x_0 = 1.2, x_1 = 1.7, x_2 = 1.8, x_3 = 2.0 \text{ and}$$

$$f_0 = 3.3201, f_1 = 5.4739, f_2 = 6.0496, f_3 = 7.3891$$

in (6), we get

$$\begin{aligned} P(x) &= \frac{(x - 1.7)(x - 1.8)(x - 2.0)}{(1.2 - 1.7)(1.2 - 1.8)(1.2 - 2.0)} \times 3.3201 + \\ &\quad \frac{(x - 1.2)(x - 1.8)(x - 2.0)}{(1.7 - 1.2)(1.7 - 1.8)(1.7 - 2.0)} \times 5.4739 + \\ &\quad \frac{(x - 1.2)(x - 1.7)(x - 2.0)}{(1.8 - 1.2)(1.8 - 1.7)(1.8 - 2.0)} \times 6.0496 + \\ &\quad \frac{(x - 1.2)(x - 1.7)(x - 1.8)}{(2.0 - 1.2)(2.0 - 1.7)(2.0 - 1.8)} \times 7.3891 \end{aligned} \quad (7)$$

Putting  $x = 1.4$  on both sides of (7), we get

$$\begin{aligned} f(1.4) = P(1.4) &= \frac{(1.4 - 1.7)(1.4 - 1.8)(1.4 - 2.0)}{(-0.5)(-0.6)(-0.8)} \times 3.3201 + \\ &\quad \frac{(1.4 - 1.2)(1.4 - 1.8)(1.4 - 2.0)}{(0.5)(-0.1)(0.3)} \times 5.4739 + \\ &\quad \frac{(1.4 - 1.2)(1.4 - 1.7)(1.4 - 2.0)}{(0.6)(0.1)(-0.2)} \times 6.0496 + \\ &\quad \frac{(1.4 - 1.2)(1.4 - 1.7)(1.4 - 1.8)}{(0.8)(0.3)(0.2)} \times 7.3891 \end{aligned}$$

$$\begin{aligned}
&= \frac{(-0.3)(-0.4)(-0.6)}{(-0.5)(-0.6)(-0.8)} \times 3.3201 + \\
&\quad \frac{(0.2)(-0.4)(-0.6)}{(0.5)(-0.1)(-0.3)} \times 5.4739 + \\
&\quad \frac{(0.2)(-0.3)(-0.6)}{(0.6)(0.1)(-0.2)} \times 6.0496 + \\
&\quad \frac{(0.2)(-0.3)(-0.4)}{(0.8)(0.3)(0.2)} \times 7.3891 \\
&= 0.99603 + 17.51648 - 18.1488 + 3.69455 \\
&= 4.05826
\end{aligned}$$

$$\therefore f(x) = 4.05826.$$

Now you can try some exercises.

E1) Show that

$$(i) \sum_{i=0}^n L_i(x) = 1$$

$$(ii) \sum_{i=0}^n L_i(x) x_i^k = x^k, \quad k \leq n$$

where  $L_i(x)$  are Lagrange fundamental polynomials

E2) Let  $w(x) = \prod_{k=0}^n (x - x_k)$ . Show that the interpolating polynomial of degree  $\leq n$  with the nodes  $x_0, x_1, \dots, x_n$  can be written as

$$P_n(x) = w(x) \sum_{k=0}^n \frac{f(x_k)}{(x - x_k) w'(x_k)}$$

### 9.3 INVERSE INTERPOLATION

In inverse interpolation in a table of values of  $x$  and  $y = f(x)$ , one is given a number  $\bar{y}$  and wishes to find the point  $\bar{x}$  so that  $f(\bar{x}) = \bar{y}$ , where  $f(x)$  is the tabulated function. This problem can always be solved if  $f(x)$  is (continuous/and) strictly increasing or decreasing (that is, the inverse of  $f$  exists). This is done by considering the table of values  $x_i, f(x_i), i = 0, 1, \dots, n$  to be a table of values  $y_i, g(y_i), i = 0, 1, 2, \dots, n$  for the inverse function  $g(y) = f^{-1}(y) = x$  by taking  $y_i = f(x_i), g(y_i) = x_i, i = 0, 1, 2, \dots, n$ . Then we can interpolate for the unknown value  $g(\bar{y})$  in this table.

$$P_n(y) = \left[ \sum_{i=0}^n x_i \prod_{\substack{j=0 \\ j \neq i}}^n \frac{(y - y_j)}{(y_i - y_j)} \right]$$

and  $\bar{x} = P_n(\bar{y})$ . This process is called inverse interpolation.

Let us consider some examples.

**Example 3:** From the following table, find the Lagrange's interpolating polynomial which agrees with the values of  $x$  at the given values of  $y$ . Hence find the value of  $x$  when  $y = 2$ .

Interpolation

x	1	19	49	101
y	1	3	4	5

Solution: Let  $x = g(y)$ . The Lagrange's interpolating polynomial  $P(y)$  of  $g(y)$  is given by

$$\begin{aligned}
 P(y) &= \frac{(y-3)(y-4)(y-5)}{(1-3)(1-4)(1-5)} \times 1 + \frac{(y-1)(y-4)(y-5)}{(3-1)(3-4)(3-5)} \times 19 \\
 &+ \frac{(y-1)(y-3)(y-5)}{(4-1)(4-3)(4-5)} \times 49 + \frac{(y-1)(y-3)(y-4)}{(5-1)(5-3)(5-4)} \times 101 \\
 &= -\frac{1}{24} [y^3 - 12y^2 + 47y - 60] + \frac{19}{4} [y^3 - 10y^2 + 29y - 20] \\
 &- \frac{49}{3} [y^3 - 9y^2 + 23y - 15] + \frac{101}{8} [y^3 - 8y^2 + 19y - 12]
 \end{aligned}$$

which, on simplification, gives

$$P(y) = y^3 - y^2 + 1.$$

The Lagrange's interpolating polynomial of  $x$  is given by  $P(y)$ .

$$\therefore x = P(y) = y^3 - y^2 + 1$$

$$\therefore \text{when } y = 2, x = P(2) = 5.$$

Example 4: Find the value of  $x$  when  $y = 3$  from the following table of values.

x	4	7	10	12
y	-1	1	2	4

Solution: The Lagrange's interpolation polynomial of  $x$  is given by

$$\begin{aligned}
 P(y) &= \frac{(y-1)(y-2)(y-4)}{(-2)(-3)(-5)} (4) + \frac{(y+1)(y-2)(y-4)}{2(1)(-3)} (7) \\
 &+ \frac{(y+1)(y-1)(y-4)}{(3)(1)(-2)} (10) + \frac{(y+1)(y-1)(y-2)}{(5)(3)(2)} (12) \\
 \therefore P(3) &= \frac{(2)(1)(-1)}{-2(3)(5)} (4) + \frac{(4)(1)(-1)}{(2)(3)} (7) \\
 &+ \frac{(4)(2)(-1)}{-3(2)} (10) + \frac{(4)(2)(1)}{(5)(3)(2)} (12) \\
 &= \frac{4}{15} - \frac{14}{3} + \frac{40}{3} + \frac{48}{15} \\
 &= \frac{182}{15} = 12.1333
 \end{aligned}$$

$$\therefore x(3) = P(3) = 12.1333.$$

Now you try some exercises.

E3) Find the Lagrange's interpolation polynomial of  $f(x)$  from the following data. Hence obtain  $f(2)$ .

x	0	1	4	5
f(x)	8	11	68	123

- E4) Using the Lagrange's interpolation formula, find the value of  $f(x)$  when  $x = 0$  from the following table:

x	3	2	1	-1
f(x)	3	12	15	-21

- E5) Find the value of  $y$  when  $x = 6$  from the following table:

x	1	2	7	8
y	4	5	5	4

- E6) From the following table of values, find the value of  $y$  when  $x = 2.5$

x	0	1	2	4
y	5	14	41	98

- E7) Find the value of  $f(5)$  from the following table:

x	0	1	3	4	7
f(x)	4	1	43	112	655

- E8) Using the Lagrange's interpolation formula, find the value of  $y$  when  $x = 10$ .

x	5	6	9	11
y	12	13	14	16

- E9) In the following table,  $h$  is the height above the sea level and  $p$  is the barometric pressure. Calculate  $p$  when  $h = 5280$ .

h	0	4763	6942	10594
p	27	25	23	20

- E10) In the following table,  $y$  represents the percentage of the number of workers in a factory whose age is less than  $x$  years. Find what percentage of workers have their age less than 35 years.

x	25	30	40	50
y	52	67.3	84.1	94.4

Now we are going to find the error committed in approximating the value of the function by  $P_n(x)$ .

## 9.4 ERROR

Let  $E_n(x) = f(x) - P_n(x)$  be the error involved in approximating the function  $f(x)$  by an interpolating polynomial. We derive an expression for  $E_n(x)$  in the following theorem. This result helps us in estimating a useful bound on the error as explained in an example.

**Theorem 4:** Let  $x_0, x_1, \dots, x_n$  be distinct numbers in the interval  $[a, b]$  and  $f$  has (continuous) derivatives upto order  $(n + 1)$  in the open interval  $]a, b[$ . If  $P_n(x)$  is the interpolating polynomial of degree  $\leq n$ , which interpolates  $f(x)$  at the points  $x_0, \dots, x_n$ , then for each  $x \in ]a, b[$ , a number  $\xi(x)$  in  $]a, b[$  exists such that

$$E_n(x) = f(x) - P_n(x) = \frac{f^{(n+1)}(\xi(x))}{(n+1)!} (x - x_0)(x - x_1) \dots (x - x_n). \quad (8)$$

**Proof:** If  $x \neq x_k$  for any  $k = 0, 1, 2, \dots, n$ , define the function  $g$  for  $t$  in  $[a, b]$  by

$$g(t) = f(t) - P_n(t) - [f(x) - P_n(x)] \prod_{j=0}^n \frac{(t - x_j)}{(x - x_j)}$$

Since  $f(t)$  has continuous derivatives upto order  $(n + 1)$  and  $P(t)$  has derivatives of all orders,  $g(t)$  has continuous derivatives upto  $(n + 1)$  order. Now, for  $k = 0, 1, 2, \dots, n$ , we have

$$\begin{aligned} g(x_k) &= f(x_k) - P_n(x_k) - [f(x) - P_n(x)] \prod_{i=0}^n \frac{(x_k - x_i)}{(x - x_i)} \\ &= 0 - [f(x) - P_n(x)] \cdot 0 = 0 \end{aligned}$$

$$\begin{aligned} \text{Furthermore, } g(x) &= f(x) - P_n(x) - [f(x) - P_n(x)] \prod_{i=0}^n \frac{(x - x_i)}{(x - x_i)} \\ &= f(x) - P_n(x) - [f(x) - P_n(x)] \cdot 1 = 0 \end{aligned}$$

Thus  $g$  has continuous derivatives upto order  $(n + 1)$  and  $g$  vanishes at the  $(n + 2)$  distinct points  $x, x_0, \dots, x_n$ . By the generalized Rolle's Theorem (Theorem 2) there exists  $\xi(x)$  in  $]a, b[$  for which  $g^{(n+1)}(\xi) = 0$ . Differentiating  $g(t)$ ,  $(n + 1)$  times (with respect to  $t$ ) and evaluating at  $\xi$ , we get

$$0 = g^{(n+1)}(\xi) = f^{(n+1)}(\xi) - (n + 1)! \frac{[f(x) - P_n(x)]}{\prod_{i=0}^n (x - x_i)}$$

Simplifying we get (error at  $x = \bar{x}$ )

$$E_n(\bar{x}) = f(\bar{x}) - P_n(\bar{x}) = \frac{f^{(n+1)}(\xi(\bar{x}))}{(n+1)!} \prod_{i=0}^n (\bar{x} - x_i) \quad (9)$$

The error formula (Eqn. (9)) derived above, is an important theoretical result because Lagrange interpolating polynomials are extensively used in deriving important formulae for numerical differentiation and numerical integration.

It is to be noted that  $\xi = \xi(\bar{x})$  depends on the point  $\bar{x}$  at which the error estimate is required. This dependence need not even be continuous. This error formula is of limited utility since  $f^{(n+1)}(x)$  is not known (when we are given a set of data at specific nodes) and the point  $\xi$  is hardly known. But the formula can be used to obtain a bound on the error of interpolating polynomial. Let us see how, by an example.

**Example 5:** The following table gives the values of  $f(x) = e^x$ . If we fit an interpolating polynomial of degree four to the data, find the magnitude of the maximum possible error in the computed value of  $f(x)$  when  $x = 1.25$ .

$x$	1.2	1.3	1.4	1.5	1.6
$f(x)$	3.3201	3.6692	4.0552	4.4817	4.9530

**Solution:** From Eqn. (9), the magnitude of the error associated with the 4th degree polynomial approximation is given by

$$\begin{aligned} |E_4(x)| &= \left| (x - x_0)(x - x_1)(x - x_2)(x - x_3)(x - x_4) \frac{f^{(5)}(\xi)}{5!} \right| \\ &= \left| (x - 1.2)(x - 1.3)(x - 1.4)(x - 1.5)(x - 1.6) \frac{f^{(5)}(\xi)}{5!} \right| \quad (10) \end{aligned}$$



Since  $f(x) = e^x$ ,  $f^{(5)}(x) = e^x$ .

When  $x$  lies in the interval  $[1.2, 1.6]$ ,

$$\text{Max } |f^{(5)}(x)| = e^{1.6} = 4.9530 \quad (11)$$

Substituting (11) in (10), and putting  $x = 1.25$ , the upper bound on the magnitude of the error

$$= |(0.05)(-0.05)(-0.15)(-0.25)(-0.35)| \times \frac{4.9530}{120}$$

$$= 0.00000135.$$

You may now try the following exercises.

E11) For the data of Example 5 with last one omitted, i.e., considering only first four nodes, if we fit a polynomial of degree 3, find an estimate of the magnitude of the error in the computed value of  $f(x)$  when  $x = 1.25$ . Also find an upper bound in the magnitude of the error.

E12) The following table gives the values of  $x$  and  $f(x) = \text{Sinh}x$ . If the value of  $f(x)$  when  $x = 0.53$  is computed from the second degree interpolation polynomial, find the estimate of the magnitude of the error.

$x$	0.50	0.55	0.60	0.65	0.70
$f(x)$	0.52110	0.57815	0.63665	0.69675	0.75858

E13) Find the value of  $x$  when  $y = 3$  from the following table:

$x$	12	18	24	42
$y$	-2	1	2	4

E14) Find the value of  $x$  when  $y = 4$  from the table given below:

$x$	8	16	20	72
$y$	-1	1	3	5

E15) Find the interpolating polynomial which fits the following data taking  $x$  as the independent variable.

$x$	-1	0	1	2
$f(x)$	1	1	1	-5

E16) Using Lagrange's interpolation formula, find the value of  $f(4)$  from the following data:

$x$	1	3	7	13
$f(x)$	2	5	12	20

Let us take a brief look at what you have studied in this unit.

## 9.5 SUMMARY

In this unit, we have seen how to derive the Lagrange's form of interpolating polynomial for a given data. It has been shown that the interpolating polynomial for a given data is unique. Moreover the Lagrange form of interpolating polynomial can be determined for equally spaced or unequally spaced nodes. We have also seen how the Lagrange's

interpolation formula can be applied with  $y$  as the independent variable and  $x$  as the dependent variable so that the value of  $x$  corresponding to a given value of  $y$  can be calculated approximately when some conditions are satisfied. Finally, we have derived the general error formula and its use has been illustrated to judge the accuracy of our calculation. The mathematical formulae derived in this unit are listed below for your easy reference.

### 1) Lagrange's Form

$$P_n(x) = \sum_{i=0}^n f(x_i) L_i(x) \text{ where}$$

$$L_i(x) = \left[ \prod_{\substack{j=0 \\ j \neq i}}^n (x - x_j) \right] / \left[ \prod_{\substack{j=0 \\ j \neq i}}^n (x_i - x_j) \right]$$

### 2) Inverse Interpolation

$$P_n(y) = \sum_{i=0}^n x_i \left[ \prod_{\substack{j=0 \\ j \neq i}}^n \frac{(y - y_j)}{(y_i - y_j)} \right]$$

### 3) Interpolation Error

$$E_n(x) = f(x) - P_n(x) = \frac{f^{(n+1)}(\xi(x))}{(n+1)!} \prod_{i=0}^n (x - x_i)$$

## 9.6 SOLUTIONS/ANSWERS

E1) If  $f(x)$  is a polynomial of degree  $\leq n$ , then

$$f(x) = P_n(x) = \sum_{i=0}^n L_i(x) f(x_i) \text{ by uniqueness of interpolating polynomial.}$$

When  $f(x) \equiv 1$ , we get (i).

When  $f(x) = x^k$ ,  $k \leq n$ , we get (ii) by the same argument.

$$E2) P_n(x) = \sum_{i=0}^n f(x_i) L_i(x)$$

$$\text{where } L_i(x) = \left[ \prod_{\substack{j=0 \\ j \neq i}}^n (x - x_j) \right] / \left[ \prod_{\substack{j=0 \\ j \neq i}}^n (x_i - x_j) \right], i = 0, \dots, n.$$

$$\text{Since } w(x) = \prod_{j=0}^n (x - x_j)$$

$$w'(x_i) = \prod_{\substack{j=0 \\ j \neq i}}^n (x_i - x_j)$$

$$\text{Also } \prod_{\substack{j=0 \\ j \neq i}}^n (x - x_j) = \frac{\prod_{j=0}^n (x - x_j)}{(x - x_i)} = \frac{w(x)}{(x - x_i)}$$

$$\text{Hence } L_i(x) = \frac{w(x)}{(x - x_i) \cdot w'(x_i)}$$

$$\text{Thus } P_n(x) = \sum_{i=0}^n \frac{f(x_i) w(x)}{(x - x_i) w'(x_i)}$$

E3)  $x^3 - x^2 + 3x + 8, 18$

E4) 6

E5)  $\frac{17}{3}$

E6) 57.265625

E7) 229

E8) 14.6667

E9) 24.5493

E10) 77.405

E11) From Eqn. (9), the magnitude of the error associated with the 3rd degree polynomial approximation is given by

$$|E(x)| = \left| (x - x_0)(x - x_1)(x - x_2)(x - x_3) \frac{f^{(4)}(\xi)}{4!} \right|$$

$$|E(1.25)| = \left| (1.25 - 1.2)(1.25 - 1.3)(1.25 - 1.4)(1.25 - 1.5) \frac{e^{\xi}}{4!} \right|$$

Since  $f^{(4)}(x) = e^x$ , when  $x$  lies in the interval  $[1.2, 1.5]$ ,

$$\text{Max } |f^{(4)}(x)| = e^{1.5}$$

$$\text{Hence } |E(1.25)| \leq \frac{(0.5)(-.05)(-.15)(-.25) \cdot e^{1.5}}{25}$$

E12)  $8.4 \times 10^{-6}$

E13) 32

E14) 36.75

E15)  $-x^3 + x + 1$

E16) 6.6875.

---

# UNIT 10 NEWTON FORM OF THE INTERPOLATING POLYNOMIAL

---

## Structure

- 10.1 Introduction
  - Objectives
- 10.2 Divided Differences
- 10.3 Newton's General Form of Interpolating Polynomial
- 10.4 The Error of the Interpolating Polynomial
- 10.5 Divided Differences and Derivatives
- 10.6 Further Results on Interpolation Error
- 10.7 Summary
- 10.8 Solutions/Answers

---

## 10.1 INTRODUCTION

---

The Lagrange's form of the interpolating polynomial derived in Unit 9 has some drawbacks compared to Newton form of interpolating polynomial that we are going to consider now.

In practice, one is often not sure as to how many interpolation points to use. One often calculates  $P_1(x)$ ,  $P_2(x)$ , .... increasing the number of interpolation points, and hence the degrees of the interpolating polynomials till one gets a satisfactory approximation  $P_k(x)$  to  $f(x)$ . In such an exercise, Lagrange form seems to be wasteful as in calculating  $P_k(x)$ , no advantage is taken of the fact that one has already constructed  $P_{k-1}(x)$ , whereas in Newton form it is not so.

Before deriving Newton's general form of interpolating polynomial, we introduce the concept of divided difference and the tabular representation of divided differences. Also the error of the interpolating polynomial in this case is derived in terms of divided differences. Using the two different expressions for the error term we get a relationship between  $n$ th order divided difference and  $n$ th order derivative.

### Objectives

After studying this unit, you should be able to :

- obtain a divided difference in terms of function values;
- form a table of divided differences and find divided differences with a given set of arguments from the table;
- show that divided difference is independent of the order of its arguments;
- obtain the Newton's divided differences interpolating polynomial for a given data;
- find an estimate of  $f(x)$  for a given non – tabular value of  $x$  from a table of values of  $x$  and  $y [ f(x) ]$ ;
- relate the  $k$ th order derivative of  $f(x)$  with the  $k$ th order divided difference from the expression for the error term.

## 10.2 DIVIDED DIFFERENCES

Suppose that we have determined a polynomial  $P_{k-1}(x)$  of degree  $\leq k-1$  which interpolates  $f(x)$  at the points  $x_0, x_1, \dots, x_{k-1}$ . In order to make use of  $P_{k-1}(x)$  in calculating  $P_k(x)$  we consider the following problem: What function  $g(x)$  should be added to  $P_{k-1}(x)$  to get  $P_k(x)$ ? Let  $g(x) = P_k(x) - P_{k-1}(x)$ . Now,  $g(x)$  is a polynomial of degree  $\leq k$  and  $g(x_i) = P_k(x_i) - P_{k-1}(x_i) = f(x_i) - f(x_i) = 0$  for  $i = 0, 1, \dots, k-1$ .

Suppose that  $P_n(x)$  is the Lagrange polynomial of degree at most  $n$  that agrees with the function  $f$  at the distinct numbers  $x_0, x_1, \dots, x_n$ .  $P_n(x)$  can have the following representation, called Newton form.

$$P_n(x) = A_0 + A_1(x - x_0) + A_2(x - x_0)(x - x_1) + \dots + A_n(x - x_0) \dots (x - x_{n-1}) \quad (1)$$

for appropriate constants  $A_0, A_1, \dots, A_n$ .

Evaluating  $P_n(x)$  (Eqn. (1)) at  $x_0$  we get  $A_0 = P_n(x_0) = f(x_0)$ . Similarly when  $P_n(x)$  is evaluated at  $x_1$ , we get  $A_1 = \frac{f(x_1) - f(x_0)}{x_1 - x_0}$ . Let us introduce the notation for divided differences and define it at this stage: The zeroth divided difference of the function  $f$ , with respect to  $x_i$ , is denoted by  $f[x_i]$  and is simply the evaluation of  $f$  at  $x_i$ , that is,  $f[x_i] = f(x_i)$ . The first divided difference of  $f$  with respect to  $x_i$  and  $x_{i+1}$  is denoted by  $f[x_i, x_{i+1}]$  and defined as

$$f[x_i, x_{i+1}] = \frac{f[x_{i+1}] - f[x_i]}{x_{i+1} - x_i}$$

The remaining divided differences of higher orders are defined inductively as follows. The  $k$ th divided differences relative to  $x_i, x_{i+1}, \dots, x_{i+k}$  is defined as

$$f[x_i, x_{i+1}, \dots, x_{i+k}] = \frac{f[x_{i+1}, \dots, x_{i+k}] - f[x_i, \dots, x_{i+k-1}]}{x_{i+k} - x_i}$$

where the  $(k-1)$ st divided differences  $f[x_i, \dots, x_{i+k-1}]$  and  $f[x_{i+1}, \dots, x_{i+k}]$  have been determined. This shows that the  $k$ th divided difference is the divided differences of  $(k-1)$ st divided differences justifying the name. The divided difference  $f[x_1, x_2, \dots, x_k]$  is invariant under all permutations of the arguments  $x_1, x_2, \dots, x_k$ . To show this we proceed as follows giving another expression for the divided difference.

For any integer  $k$  between 0 and  $n$ , let  $Q_k(x)$  be the sum of the first  $k+1$  terms in form (1), i.e.,

$$Q_k(x) = A_0 + A_1(x - x_0) + \dots + A_k(x - x_0) \dots (x - x_{k-1}).$$

Since each of the remaining terms in Eqn. (1) has the factor  $(x - x_0)(x - x_1) \dots (x - x_k)$ , Eqn. (1) can be rewritten as

$P_n(x) = Q_k(x) + (x - x_0) \dots (x - x_k) R(x)$  for some polynomial  $R(x)$ . As the term  $(x - x_0)(x - x_1) \dots (x - x_k) R(x)$  vanishes at each of the points  $x_0 \dots x_k$ , we have  $f(x_i) = P_n(x_i) = Q_k(x_i)$ ,  $i = 0, 1, 2, \dots, k$ . Since  $Q_k(x)$  is a polynomial of degree  $\leq k$ , by uniqueness of interpolating polynomial  $Q_k(x) = P_k(x)$ .

This shows that  $P_n(x)$  can be constructed step by step with the addition of the next term in Eqn. (1), as one constructs the sequence  $P_0(x), P_1(x) \dots$  with  $P_k(x)$  obtained from  $P_{k-1}(x)$  in the form

$$P_k(x) = P_{k-1}(x) + A_k(x - x_0) \dots (x - x_{k-1}) \quad (2)$$

That is,  $g(x)$  is a polynomial of degree  $\leq k$  having (at least) the  $k$  distinct zeros  $x_0, \dots, x_{k-1}$ .

$\therefore P_k(x) - P_{k-1}(x) = g(x) = A_k(x - x_0) \dots (x - x_{k-1})$ , for some constant  $A_k$ . This constant  $A_k$  is called the  $k$ th divided difference of  $f(x)$  at the points  $x_0, \dots, x_k$  for reasons discussed below and is denoted by  $f[x_0, x_1, \dots, x_k]$ . This coefficient depends only on the values of

$f(x)$  at the points  $x_0, \dots, x_k$ . Thus Eqn. (2) can be rewritten as

$$P_k(x) = P_{k-1}(x) + f[x_0, \dots, x_k] (x - x_0) (x - x_1) \dots (x - x_{k-1}) \quad (3)$$

To get an explicit expression for  $f[x_0, \dots, x_k]$  we make use of Lagrange form of interpolating polynomial and the uniqueness of interpolating polynomial.

From Eqn. (3) we have

$$P_k(x) = P_{k-1}(x) + f[x_0, \dots, x_k] (x - x_0) \dots (x - x_{k-1}),$$

since  $(x - x_0) (x - x_1) \dots (x - x_{k-1}) = x^k +$  a polynomial of degree  $< k$ , we can rewrite  $P_k(x)$  as  $P_k(x) = f[x_0, \dots, x_k] x^k +$  a polynomial of degree  $< k$  (4)

(as  $P_{k-1}(x)$  is a polynomial of degree  $< k$ ).

But considering the Lagrange form of interpolating polynomial we have

$$\begin{aligned} P_k(x) &= \sum_{i=0}^k f(x_i) \prod_{\substack{j=0 \\ j \neq i}}^k \frac{(x - x_j)}{(x_i - x_j)} \\ &= \sum_{i=0}^k \left[ \frac{f(x_i)}{\prod_{\substack{j=0 \\ j \neq i}}^k (x_i - x_j)} \right] x^k + \text{a polynomial of degree } < k. \end{aligned}$$

Therefore, on comparison with Eqn. (4) we have

$$f[x_0, \dots, x_k] = \sum_{i=0}^k \frac{f(x_i)}{(x_i - x_0) \dots (x_i - x_{i-1}) \dots (x_i - x_{i+1}) \dots (x_i - x_k)} \quad (5)$$

This shows that

$$f[y_0, \dots, y_k] = f[x_0, \dots, x_k]$$

if  $y_0, \dots, y_k$  is a reordering of the sequence  $x_0, \dots, x_k$ . We have defined the zeroth divided difference of  $f(x)$  at  $x_0$  by  $f[x_0] = f(x_0)$  which is consistent with Eqn. (5).

For  $k = 1$ , we have from Eqn. (5)

$$f[x_0, x_1] = \frac{f(x_0)}{x_0 - x_1} + \frac{f(x_1)}{x_1 - x_0} = \frac{f(x_0) - f(x_1)}{x_0 - x_1} = \frac{f[x_1] - f[x_0]}{x_1 - x_0}$$

This shows that the first divided difference is really a divided difference.

For  $k = 2$ , it can be shown (using Eqn. 5) that

$$f[x_0, x_1, x_2] = \frac{f[x_0, x_1] - f[x_1, x_2]}{x_0 - x_2}$$

This shows that the second divided difference is a divided difference of divided differences.

We show below in Theorem 1 that for  $k > 2$

$$f[x_0, \dots, x_k] = \frac{f[x_1, \dots, x_k] - f[x_0, \dots, x_{k-1}]}{x_k - x_0} \quad (6)$$

This shows that the  $k$ th divided difference is the divided difference of  $(k - 1)$ st divided differences justifying the name. If  $M = (x_0, \dots, x_n)$  and  $N$  denotes any  $n - 1$  elements of  $M$  and the remaining two elements are denoted by  $\alpha$  and  $\beta$ , then

$$f[x_0, \dots, x_n] = \frac{[(n - 1)\text{st divided difference on } N \text{ and } \alpha] - (n - 1)\text{st divided difference on } N \text{ and } \beta}{\alpha - \beta} \quad (7)$$

**Theorem 1:**

$$f[x_0, \dots, x_j] = \frac{f[x_1, \dots, x_j] - f[x_0, x_1, \dots, x_{j-1}]}{x_j - x_0} \quad (8)$$

**Proof:** Let  $P_{i-1}(x)$  be the polynomial of degree  $\leq i-1$  which interpolates  $f(x)$  at  $x_0, \dots, x_{i-1}$  and let  $Q_{j-1}(x)$  be the polynomial of degree  $\leq j-1$  which interpolates  $f(x)$  at the points  $x_1, \dots, x_j$ . Let us define  $P(x)$  as

$$P(x) = \frac{x - x_0}{x_j - x_0} Q_{j-1}(x) + \frac{x_j - x}{x_j - x_0} P_{j-1}(x).$$

This is a polynomial of degree  $\leq j$ , and  $P(x_i) = f(x_i)$  for  $i = 0, 1, \dots, j$ . By uniqueness of the interpolating polynomial we have  $P(x) = P_j(x)$ . Therefore

$$P_j(x) = \frac{x - x_0}{x_j - x_0} Q_{j-1}(x) + \frac{x_j - x}{x_j - x_0} P_{j-1}(x).$$

Equating the coefficient of  $x^j$  from both sides of Eqn. (8), we obtain (leading) coefficient of

$$x^j \text{ in } P_j(x) = \frac{\text{leading coefficient of } Q_{j-1}(x)}{x_j - x_0} - \frac{\text{leading coefficient of } P_{j-1}(x)}{x_j - x_0}.$$

$$\text{That is } f[x_0, \dots, x_j] = \frac{f[x_1, \dots, x_j] - f[x_0, \dots, x_{j-1}]}{x_j - x_0}.$$

We now illustrate this theorem with the help of a few examples but before that we give the table of divided differences of various orders.

**Table of divided differences**

Suppose we denote, for convenience, a first order divided difference of  $f(x)$  with any two arguments by  $f[...]$ , a second order divided difference with any three arguments by  $f[...]$  and so on. Then the table of divided differences can be written as follows

Table 1

$x$	$f[.]$	$f[.,.]$	$f[.,.,.]$	$f[.,.,.,.]$	$f[.,.,.,.,.]$
$x_0$	$f_0$				
		$f[x_0, x_1]$			
$x_1$	$f_1$		$f[x_0, x_1, x_2]$		
		$f[x_1, x_2]$		$f[x_0, x_1, x_2, x_3]$	
$x_2$	$f_2$		$f[x_1, x_2, x_3]$		$f[x_0, x_1, x_2, x_3, x_4]$
		$f[x_2, x_3]$		$f[x_1, x_2, x_3, x_4]$	
$x_3$	$f_3$		$f[x_2, x_3, x_4]$		
		$f[x_3, x_4]$			
$x_4$	$f_4$				

**Example 1:** If  $f(x) = x^3$ , find the value of  $f[a, b, c]$ .

**Solution:**  $f[a, b] = \frac{f(b) - f(a)}{b - a} = \frac{b^3 - a^3}{b - a}$

$$= b^2 + ba + a^2 = a^2 + ab + b^2$$

Similarly,

$$\begin{aligned}
 f[b,c] &= c^2 + cb + b^2 = b^2 + bc + c^2 \\
 \therefore f[a,b,c] &= \frac{f[b,c] - f[a,b]}{c - a} \\
 &= \frac{(b^2 + bc + c^2) - (a^2 + ab + b^2)}{c - a} \\
 &= \frac{(c^2 - a^2) + b(c - a)}{c - a} \\
 &= \frac{(c - a)(c + a + b)}{(c - a)} \\
 &= a + b + c \\
 \therefore f[a,b,c] &= a + b + c.
 \end{aligned}$$

**Example 2:** If  $f(x) = \frac{1}{x}$ , show that

$$f[a,b,c,d] = -\frac{1}{abcd}$$

**Solution:**  $f[a,b] = \frac{\frac{1}{b} - \frac{1}{a}}{b - a} = \frac{a - b}{ab(b - a)} = -\frac{1}{ab}$

Similarly,

$$\begin{aligned}
 f[b,c] &= -\frac{1}{bc}, \quad f[c,d] = -\frac{1}{cd} \\
 \therefore f[a,b,c] &= \frac{-\frac{1}{bc} + \frac{1}{ab}}{c - a} = \frac{\frac{1}{ab} - \frac{1}{bc}}{c - a} \\
 &= \frac{\left[ \frac{c - a}{abc} \right]}{c - a} = \frac{1}{abc}
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 f[b,c,d] &= \frac{1}{bcd} \\
 \therefore f[a,b,c,d] &= \frac{\left[ \frac{c - a}{abc} \right]}{c - a} = \frac{1}{abc} \\
 &= \frac{\left[ \frac{a - d}{abcd} \right]}{d - a} \\
 &= -\frac{1}{abcd} \\
 \therefore f[a,b,c,d] &= -\frac{1}{abcd}.
 \end{aligned}$$

In next section we shall make use of the divided differences to derive Newton's general form of interpolating polynomial.

---

### 10.3 NEWTON'S GENERAL FORM OF INTERPOLATING POLYNOMIAL

---

In Sec.10.2 we have shown how  $P_n(x)$  can be constructed step by step as one constructs the sequence  $F_0(x), P_1(x), P_2(x), \dots$ , with  $P_k(x)$  obtained from  $P_{k-1}(x)$  with the addition of the



next term in Eqn.(3), that is,

$$P_k(x) = P_{k-1}(x) + (x - x_0)(x - x_1) \dots (x - x_{k-1}) f[x_0, \dots, x_k]$$

Using this Eqn. (1) can be rewritten as

$$P_n(x) = f[x_0] + (x - x_0) f[x_0, x_1] + (x - x_0)(x - x_1) f[x_0, x_1, x_2] + \dots + (x - x_0)(x - x_1) \dots (x - x_{n-1}) f[x_0, x_1, \dots, x_n]. \quad (9)$$

This can be written compactly as follows :

$$P_n(x) = \sum_{i=0}^n f[x_0, \dots, x_i] \prod_{j=0}^{i-1} (x - x_j). \quad (10)$$

This is the Newton's form of interpolating polynomial.

**Example 3:** From the following table of values, find the Newton's form of interpolating polynomial approximating  $f(x)$ .

x	-1	0	3	6	7
f(x)	3	-6	39	822	1611

**Solution:** We notice that the values of  $x$  are not equally spaced. We are required to find a polynomial which approximates  $f(x)$ . We form the table of divided differences of  $f(x)$ .

Table 2

x	f[.]	f[...]	f[.....]	f[.....]	f[.....]
-1	3				
		9			
0	-6				
		15			
3	39		41		
		261		13	
6	822		132		
		789			
7	1611				

Since the divided difference upto order 4 are available, the Newton's interpolating polynomial  $P_4(x)$  is given by

$$P_4(x) = f(x_0) + (x - x_0) f[x_0, x_1] + (x - x_0)(x - x_1) f[x_0, x_1, x_2] + (x - x_0)(x - x_1)(x - x_2) f[x_0, x_1, x_2, x_3] + (x - x_0)(x - x_1)(x - x_2)(x - x_3) f[x_0, x_1, x_2, x_3, x_4] \quad (11)$$

where  $x_0 = -1, x_1 = 0, x_2 = 3, x_3 = 6$  and  $x_4 = 7$ .

The divided differences  $f(x_0), f[x_0, x_1], f[x_0, x_1, x_2], f[x_0, x_1, x_2, x_3]$  and  $f[x_0, x_1, x_2, x_3, x_4]$  are those which lie along the diagonal at  $f(x_0)$  as shown by the dotted line. Substituting the values of  $x_i$  and the values of the divided differences in Eqn. (11), we get

$$P_4(x) = 3 + (x + 1)(-9) + (x + 1)x(6) + (x + 1)x(x - 3)(5) + (x + 1)x(x - 3)(x - 6)(1)$$

which on simplification gives

$$P_4(x) = x^4 - 3x^3 + 5x^2 - 6$$

$$\therefore f(x) \simeq P_4(x) = x^4 - 3x^3 + 5x^2 - 6$$

We now consider an example to show how Newton's interpolating polynomial can be used to obtain the approximate value of the function  $f(x)$  at any non-tabular point.

**Example 4:** Find the approximate values of  $f(x)$  at  $x = 2$  and  $x = 5$  in Example 3.

**Solution:** Since  $f(x) \simeq P_4(x)$ , from Example 3, we get

$$f(2) \simeq P_4(2) = 16 - 24 + 20 - 6 = 6$$

and

$$f(5) \simeq P(5) = 625 - 375 + 125 - 6 = 369$$

**Note 1:** When the values of  $f(x)$  for given values of  $x$  are required to be found, it is not necessary to find the interpolating polynomial  $P_4(x)$  in its simplified form given above. We can obtain the required values by substituting the values of  $x$  in Eqn.(11) itself. Thus,

$$P_4(2) = 3 + (3)(-9) + (3)(2)(6) + (3)(2)(-1)(5) + (3)(2)(-1)(-4)1$$

$$\therefore P_4(2) = 3 - 27 + 36 - 30 + 24 = 6.$$

Similarly,

$$\begin{aligned} P_4(5) &= 3 + (6)(-9) + (6)(5)(6) + (6)(5)(2)(5) + (6)(5)(2)(-1)(1) \\ &= 3 - 54 + 180 + 300 - 60 = 369. \end{aligned}$$

Then  $f(2) \simeq P_4(2) = 6$

and  $f(5) \simeq P(5) = 369.$

**Example 5:** Obtain the divided differences interpolation polynomial and the Lagrange's interpolating polynomial of  $f(x)$  from the following data and show that they are same.

x	0	2	3	4
f(x)	-4	6	26	64

**Solution:** (a) Divided differences interpolation polynomial:

Table 3

x	f[x]	f[...]	f[...]	f[...]
0	-4			
2	6	5		
3	26	20	5	
4	64	38	9	1

$$\begin{aligned} P(x) &= -4 + x(5) + x(x-2)(5) + x(x-2)(x-3)(1) \\ &= x^3 + x - 4 \end{aligned}$$

(b) Lagrange's interpolation polynomial:

$$\begin{aligned}
 P(x) &= \frac{(x-2)(x-3)(x-4)}{(-2)(-3)(-4)}(-4) + \frac{x(x-3)(x-4)}{(2)(-1)(-2)} \quad (6) \\
 &+ \frac{x(x-2)(x-4)}{(3)(1)(-1)} \quad (26) + \frac{x(x-2)(x-3)}{(4)(2)(1)} \quad (64) \\
 &= \frac{1}{6}(x^3 - 9x^2 + 26x - 24) + \frac{3}{2}(x^3 - 7x^2 + 12x) \\
 &\quad - \frac{26}{3}(x^3 - 6x^2 + 8x) + 8(x^3 - 5x^2 + 6x).
 \end{aligned}$$

On simplifying, we get

$$P(x) = x^3 + x - 4.$$

Thus, we find that both polynomials are the same.

You may now try the following exercises:

E1) Find the Lagrange's interpolating polynomial of  $f(x)$  from the table of values given below and show that it is the same as the Newton's divided differences interpolating polynomial.

x	0	1	4	5
f(x)	8	11	68	123

E2) From the table of values given below, obtain the value of  $y$  when  $x = 1.5$  using

(a) divided differences interpolation formula.

(b) Lagrange's interpolation formula.

x	0	1	2	4	5
f(x)	5	14	41	98	122

E3) Using Newton's divided differences interpolation formula, find the values of  $f(8)$  and  $f(15)$  from the following table.

x	4	5	7	10	11	13
f(x)	48	100	294	900	1210	2028

In Unit 9 we have derived the general error term i.e. the error committed in approximating  $f(x)$  by  $P_n(x)$ . In the next section we derive another expression for the error term in term of divided difference.

## 10.4 THE ERROR OF THE INTERPOLATING POLYNOMIAL

Let  $P_n(x)$  be the Newton form of interpolating polynomial of degree  $\leq n$  which interpolates  $f(x)$  at  $x_0, \dots, x_n$ . The interpolating error  $E_n(x)$  of  $P_n(x)$  is given by

$$E_n(x) = f(x) - P_n(x). \quad (12)$$

Let  $\bar{x}$  be any point different from  $x_0, \dots, x_n$ . If  $P_{n+1}(x)$  is the Newton form of interpolating polynomial which interpolates  $f(x)$  at  $x_0, \dots, x_n$  and  $\bar{x}$ , then  $P_{n+1}(\bar{x}) = f(\bar{x})$ . Then by (10) we have

$$P_{n+1}(x) = P_n(x) + f[x_0, \dots, x_n, \bar{x}] \prod_{j=0}^n (x - x_j)$$

Putting  $x = \bar{x}$  in the above, we have

$$f(\bar{x}) = P_{n+1}(\bar{x}) = P_n(\bar{x}) + f[x_0, \dots, x_n, \bar{x}] \prod_{j=0}^n (\bar{x} - x_j)$$

$$\text{i.e. } E_n(\bar{x}) = f(\bar{x}) - P_n(\bar{x}) = f[x_0, \dots, x_n, \bar{x}] \prod_{j=0}^n (\bar{x} - x_j). \quad (13)$$

This shows that the error is like the next term in the Newton form.

## 10.5 DIVIDED DIFFERENCE AND DERIVATIVE OF THE FUNCTION

Comparing Eqn.(13) with the error formula derived in Unit 9 Eqn. (9), we can establish a relationship between divided differences and the derivatives of the function

$$\begin{aligned} E_n(\bar{x}) &= \frac{f^{(n+1)}[\xi(\bar{x})]}{(n+1)!} \prod_{i=0}^n (\bar{x} - x_i) \\ &= f[x_0, x_1, \dots, x_n, \bar{x}] \prod_{i=0}^n (\bar{x} - x_i). \end{aligned}$$

Comparing, we have  $f[x_0, x_1, \dots, x_{n+1}] = \frac{f^{(n+1)}(\xi)}{(n+1)!}$

(considering  $\bar{x} = x_{n+1}$ )

Further it can be shown that  $\xi \in ]\min x_i, \max x_i[$ .

We state these results in the following theorem.

**Theorem 2:** Let  $f(x)$  be a real-valued function, defined on  $[a, b]$  and  $n$  times differentiable in  $]a, b[$ . If  $x_0, \dots, x_n$  are  $n+1$  distinct points in  $[a, b]$ , then there exists  $\xi \in ]a, b[$  such that

$$f[x_0, \dots, x_n] = \frac{f^{(n)}(\xi)}{n!}$$

**Corollary 1:**

If  $f(x) = x^n$ , then

$$f[x_0, \dots, x_n] = \frac{n!}{n!} = 1.$$

**Corollary 2:**

If  $f(x) = x^k$ ,  $k < n$ , then

$$f[x_0, \dots, x_k] = 0$$

since  $n$ th derivative of  $x^k$ ,  $k < n$ , is zero.

For example, consider the first divided difference

$$f[x_0, x_1] = \frac{f(x_1) - f(x_0)}{x_1 - x_0}.$$

By Mean value theorem  $f(x_1) = f(x_0) + (x_1 - x_0) f'(\xi)$ ,  $x_0 < \xi < x_1$ .

Substituting, we get

$$f[x_0, x_1] = f'(\xi), \quad x_0 < \xi < x_1.$$

Example 6: If  $f(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$ , then find  $f[x_0, x_1, \dots, x_n]$ .

Solution: From Corollaries 1 and 2 we have  $f[x_0, x_1, \dots, x_n] = a_n \cdot \frac{n!}{n!} + 0 = a_n$ .

Let us consider another example.

Example 7: If  $f(x) = 2x^3 + 3x^2 - x + 1$ , find

$$f[1, -1, 2, 3], f[a, b, c, d], f[4, 6, 7, 8].$$

Solution: Since  $f(x)$  is a cubic polynomial, the 3rd order divided differences of  $f(x)$  with any set of arguments are constant and equal to 2, the coefficient of  $x^3$  in  $f(x)$ .

Thus, it follows that  $f[1, -1, 2, 3]$ ,  $f[a, b, c, d]$ , and  $f[4, 6, 7, 8]$  are each equal to 2.

You may now try the following exercises:

---

E4) If  $f(x) = 2x^3 - 3x^2 + 7x + 1$ , what is the value of  $f[1, 2, 3, 4]$ ?

E5) If  $f(x) = 3x^2 - 2x + 5$ , find  $f[1, 2]$ ,  $f[2, 3]$  and  $f[1, 2, 3]$ .

---

In the next section, we are going to discuss about bounds on the interpolation error.

## 10.5 FURTHER RESULTS ON INTERPOLATION ERROR

We have derived the error formula

$$E_n(x) = f(x) - P_n(x) = \frac{f^{(n+1)}(\xi(x))}{(n+1)!} \prod_{i=0}^n (x - x_i).$$

We assume that  $f(x)$  is  $(n+1)$  times continuously differentiable in the interval of interest  $[a, b] = I$  that contains  $x_0, \dots, x_n$ , and  $x$ . Since  $\xi(x)$  is unknown we may replace  $f^{(n+1)}(\xi(x))$  by  $\max_{x \in I} |f^{(n+1)}(x)|$ . If we denote  $(x - x_0)(x - x_1)\dots(x - x_n)$  by  $\psi_n(x)$  then we have

$$|E_n(x)| = |f(x) - P_n(x)| \leq \frac{\max_{t \in I} |f^{(n+1)}(t)|}{(n+1)!} \max_{x \in I} |\psi_n(x)| \quad (14)$$

Consider now the case when the nodes are equally spaced, that is,  $x_j = x_0 + jh$ ,  $j = 0, \dots, N$ , and  $h$  is the spacing between consecutive nodes. For the case  $n=1$  we have linear interpolation. If  $x \in [x_{i-1}, x_i]$ , then we approximate  $f(x)$  by  $P_1(x)$  which interpolates at

$x_{i-1}$ , and  $x_i$ . From Eqn. (14) we have  $|E_1(x)| \leq \frac{1}{2} \max_{t \in I} |f''(t)| \max_{t \in I} |\psi_1(x)|$

where  $\psi_1(x) = (x - x_{i-1})(x - x_i)$ .

Now,

$$\frac{d\psi_1}{dx} = x - x_{i-1} + x - x_i = 0$$

gives  $x = (x_{i-1} + x_i)/2$ .

Hence, the maximum value of  $|(x - x_{i-1})(x - x_i)|$  occurs at  $x = x^* = (x_{i-1} + x_i)/2$ .

The maximum value is given by

$$|\Psi_1(x^*)| = \frac{(x_i - x_{i-1})^2}{4} = \frac{h^2}{4}.$$

Thus, we have for linear interpolation, for any  $x \in I$

$$\begin{aligned} |E_1(x)| = |f(x) - P_1(x)| &\leq \frac{(x_i - x_{i-1})^2}{4} \frac{1}{2} \max_{x \in I} |f''(x)| \\ &= \frac{h^2}{8} M. \end{aligned} \quad (15)$$

where  $|f''(x)| \leq M$  on  $I$ .

For the case  $n=2$ , it can be shown that for any  $x \in [x_{i-1}, x_{i+1}]$ ,

$$|E_2(x)| \leq \frac{h^3 M}{9\sqrt{3}} \text{ where } |f'''(x)| \leq M \text{ on } I. \quad (16)$$

**Example 8:** Determine the spacing  $h$  in a table of equally spaced values of the function of  $f(x) = \sqrt{x}$  between 1 and 2, so that interpolation with a first degree polynomial in this table will yield seven place accuracy.

**Solution:** Here

$$f''(x) = -\frac{1}{4}x^{-3/2}$$

$$\max_{1 \leq x \leq 2} |f''(x)| = \frac{1}{4}.$$

and  $|E_1(x)| \leq \frac{h^2}{32}.$

For seven place accuracy,  $h$  is to be chosen such that

$$\frac{h^2}{32} < 5 \cdot 10^{-8}.$$

or  $h^2 < (160)10^{-8}$  that is  $h < .0013$ .

E6) If  $f(x)$  takes the values  $-21, 15, 12$  and  $3$  respectively when  $x$  assumes the values  $-1, 1, 2$  and  $3$ , find the polynomial which approximates  $f(x)$ .

E7) Using the following table of values, find the polynomial which approximates  $f(x)$ . Hence obtain the value of  $f(5)$ .

$x$	$-1$	$0$	$2$	$3$	$7$	$10$
$f(x)$	$-11$	$1$	$1$	$1$	$141$	$561$

E8) Find the polynomial which approximates  $f(x)$ , tabulated below

$x$	$-4$	$-1$	$0$	$2$	$5$
$f(x)$	$1245$	$33$	$5$	$9$	$1335$

Also find an approximate value of  $f(x)$  at  $x = 1$  and  $x = -2$ .

E9) If  $f(3) = 168, f(7) = 120, f(9) = 72$  and  $f(10) = 63$ , find an approximate value of  $f(6)$ .

E10) The following table gives steam pressures  $P$  at different temperatures  $T$ , measured in degrees. Find the pressure at temperature  $372.1$  degrees.

T	361	367	378	387	399
P	154.9	167.0	191.0	212.5	244.2

E11) From the following table, find the value of y when x = 102

x	93.0	96.2	100.0	104.2	108.7
y	11.38	12.80	14.70	17.07	19.91

E12) From the following table of values, obtain the value of y at x = 3

x	0	1	2	4	5
y	0	16	48	88	0

E13) Obtain the polynomial which agrees with the values of f(x) as shown below

x	0	1	2	5
f(x)	2	3	12	147

E14) Determine the spacing h in a table of equally spaced values of the function  $f(x) = \sqrt{x}$  between 1 and 2, so that interpolation with a second-degree polynomial in this table yields seven-place accuracy.

We now end this unit by giving a summary of what we have covered in it.

## 10.6 SUMMARY

In this unit we have derived a form of interpolating polynomial called Newton's general form, which has some advantages over the Lagrange's form discussed in Unit 9. This form is useful in deriving some other interpolating formulas. We have introduced the concept of divided differences and discussed some of its important properties before deriving Newton's general form. The error term has also been derived and utilizing the error term we have established a relationship between the divided difference and the derivative of the function f(x) for which the interpolating polynomial has been obtained. The main formulas derived are listed below:

$$1. \quad f[x_0, \dots, x_j] = \frac{f[x_1, \dots, x_j] - f[x_0, \dots, x_{j-1}]}{x_j - x_0}$$

$$2. \quad P_n(x) = \sum_{i=0}^n f[x_0, \dots, x_i] \prod_{j=0}^{i-1} (x - x_j)$$

$$3. \quad E_n(x) = f[x_0, \dots, x_n, x] \prod_{j=0}^n (x - x_j)$$

$$4. \quad f[x_0, \dots, x_n] = \frac{f^{(n)}(\xi)}{n!}, \quad \xi \in ]\min x_i, \max x_i[$$

## 10.7 SOLUTIONS AND ANSWERS

E1)  $x^3 - x^2 + 3x + 8$

E2) 26.35156

E3) We form the divided differences table of f(x) below

Table 4

x	f[.]	f[...]	f[.....]	f[.....]
4	48			
		52		
5	100		15	
		97		1
7	294		21	
		202		1
10	900		27	
		310		1
11	1210		33	
		409		
13	2028			

From the Newton's divided difference interpolation formula, we have

$$f(x) \approx 48 + (x-4)(52) + (x-4)(x-5)(15) + (x-4)(x-5)(x-7)(1).$$

Substituting  $x = 8$  in the above get

$$\begin{aligned} f(8) &\approx 48 + 4 \times 52 + 4 \times 3 \times 15 + 4 \times 3 \times 1 \\ &= 48 + 208 + 180 + 12 = 448. \end{aligned}$$

substituting  $x = 15$ , we get

$$f(15) \approx 3150$$

E4) 3

E5)  $6\xi - 2, \xi \in ]1, 2[; 6\eta - 2, \eta \in ]2, 3[, \text{ and } 6$

E6)  $x^3 - 9x^2 + 17x + 6$

E7)  $x^3 - 5x^2 + 6x + 1, 31$

E8)  $3x^4 - 5x^3 + 6x^2 - 14x + 5, -5, 145$

E9) 147

E10) 177.4

E11) 15.79

E12) 84

E13)  $x^3 + x^2 - x + 2$

E14)  $f'''(x) = \frac{3}{8}x^{-5/2}$ ; hence  $\max_{1 \leq x \leq 2} |f'''(x)| = \frac{3}{8}$ .



$$|E_2(x)| \leq \frac{2h^3}{3\sqrt{3}} \left[ \frac{3}{8} \right] \frac{1}{6} = \frac{h^3}{24\sqrt{3}}$$

For seven place accuracy, h has to be chosen such that

$$\frac{h^3}{24\sqrt{3}} < 5 \cdot 10^{-8}. \text{ This gives } h \approx 0.0128.$$

$$\text{The number of interval is } N = \frac{2 - 1}{h} \approx 79.$$

---

# UNIT 11 INTERPOLATION AT EQUALLY SPACED POINTS

---

## Structure

- 11.1 Introduction
  - Objectives
- 11.2 Differences
  - 11.2.1 Forward Differences
  - 11.2.2 Backward Differences
  - 11.2.3 Central Differences
- 11.3 Difference Formulas
  - 11.3.1 Newton's Forward Difference Formula
  - 11.3.2 Newton's Backward Difference Formula
  - 11.3.3 Stirling's Central Difference Formula
- 11.4 Summary
- 11.5 Solutions/Answers

---

## 11.1 INTRODUCTION

---

Suppose that  $y$  is a function of  $x$ . The exact functional relation  $y = f(x)$  between  $x$  and  $y$  may or may not be known. But, the values of  $y$  at  $(n + 1)$  equally spaced values of  $x$  are supposed to be known, i.e.,  $(x_i, y_i); i = 0, \dots, n$  are known where  $x_i - x_{i-1} = h$  (fixed),  $i = 1, 2, \dots, n$ . Suppose that we are required to determine an approximate value of  $f(x)$  or its derivative  $f'(x)$  for some values of  $x$  in the interval of interest. The methods for solving such problems are based on the concept of finite differences. We have introduced the concept of forward, backward and central differences and discussed their interrelationship in Sec. 11.2.

We have already introduced two important forms of the interpolating polynomial in Units 9 and 10. These forms simplify when the nodes are equidistant. For the case of equidistant nodes, we have derived the Newton's forward, backward difference forms and Stirling's central difference form of interpolating polynomial, each suitable for use under a specific situation. We have derived these methods in Sec. 11.3, and also given the corresponding error term.

### Objectives

After reading this unit, you should be able to

- write a forward difference in terms of function values from a table of forward differences and locate a difference of given order at a given point;
- write a backward difference in terms of function values from a table of backward differences and identify differences of various orders at any given point from the table;
- expand a central difference in terms of function values and form a table of central differences;
- establish relations between  $\Delta$ ,  $\nabla$ ,  $\delta$  and divided difference;
- obtain the interpolating polynomial of  $f(x)$  for a given data by applying any one of the interpolating formulas;

- compute  $f(x)$  approximately when  $x$  lies near the beginning of the table and estimate the error;
- compute  $f(x)$  approximately when  $x$  lies near the end of the table and estimate the error;
- estimate the value of  $f(x)$  when  $x$  lies near the middle of the table and estimate the error.

## 11.2 DIFFERENCES

Suppose that we are given a table of values  $(x_i, y_i)$ ,  $i = 0, 1, 2, \dots, N$  where  $y_i = f(x_i) = f_i$ . Let the nodal points be equidistant. That is

$$x_i = a + ih, i = 0, \dots, N, \text{ with } N = (b - a)/h \quad (1)$$

For simplicity we introduce a linear change of variables

$$s = s(x) = \frac{x - x_0}{h}, \text{ so that } x = x(s) = x_0 + sh \quad (2)$$

and introduce the notation

$$f(x) = f(x_0 + sh) = f_s \quad (3)$$

The linear change of variables in Eqn. (2) transforms polynomials of degree  $n$  in  $x$  into polynomials of degree  $n$  in  $s$ . We have already introduced the divided-difference table to calculate a polynomial of degree  $\leq n$  which interpolates  $f(x)$  at  $x_0, x_1, \dots, x_n$ . For equally spaced nodes, we shall deal with three types of differences, namely, forward, backward and central and discuss their representation in the form of a table. We shall also derive the relationship of these differences with divided differences and their interrelationship.

### 11.2.1 Forward Differences

We denote the forward differences of  $f(x)$  of  $i$ th order at  $x = x_0 + sh$  by  $\Delta^i f_s$  and define it as follows:

$$\Delta^i f_s = \begin{cases} f_s & i = 0 \\ \Delta(\Delta^{i-1} f_s) = \Delta^{i-1} f_{s+1} - \Delta^{i-1} f_s, & i > 0. \end{cases}$$

Where  $\Delta$  denotes forward difference operator.

When  $s = k$ , that is,  $x = x_k$ , we have

for  $i = 1$   $\Delta f_k = f_{k+1} - f_k$

for  $i = 2$   $\Delta^2 f_k = \Delta f_{k+1} - \Delta f_k$

$$= f_{k+2} - f_{k+1} - [f_{k+1} - f_k]$$

$$= f_{k+2} - 2f_{k+1} + f_k$$

Similarly  $\Delta^3 f_k = f_{k+3} - 3f_{k+2} + 3f_{k+1} - f_k$

We recall the binomial theorem

$$(a + b)^s = \sum_{j=0}^s \binom{s}{j} a^j b^{s-j} \quad (4)$$

where  $s$  is a real and non-negative integer.

We give below in Lemma 1 the relationship between the forward and divided differences. This relation will be utilized to derive the Newton's forward-difference formula which interpolates  $f(x)$  at  $x_k + ih$ ,  $i = 0, 1, \dots, n$ .

**Lemma 1:** For all  $i \geq 0$

$$f[x_k, \dots, x_{k+i}] = \frac{1}{i! h^i} \Delta^i f_k \quad (5)$$

**Proof:** We prove the result by induction.

For  $i = 0$ , both sides of relation (5) are same by convention, that is,

$$f[x_k] = f(x_k) = f_k = \Delta^0 f_k.$$

Assuming that relation (5) holds for  $i = n \geq 0$ , we have for  $i = n + 1$

$$\begin{aligned} f[x_k, x_{k+1}, \dots, x_{k+n+1}] &= \frac{f[x_{k+1}, \dots, x_{k+n+1}] - f[x_k, \dots, x_{k+n}]}{x_{k+n+1} - x_k} \\ &= \frac{[\Delta^n f_{k+1} / n! h^n] - [\Delta^n f_k / n! h^n]}{x_0 + (k+n+1)h - x_0 - kh} \\ &= \frac{\Delta^n f_{k+1} - \Delta^n f_k}{(n+1)! h^{n+1}} = \frac{\Delta^{n+1} f_k}{(n+1)! h^{n+1}} \end{aligned}$$

This shows that relation (5) holds for  $i = n + 1$  also. Hence (5) is proved. We now give a result which immediately follows from this theorem in the following corollary.

**Corollary:** If  $P_n(x)$  is a polynomial of degree  $n$  with leading coefficient  $a_n$ , and  $x_0$  is an arbitrary point, then

$$\Delta^n P_n(x_0) = a_n n! h^n$$

and  $\Delta^{n+1} P_n(x_0) = 0$ , i.e., all higher differences are zero.

**Proof:** Taking  $k = 0$  in relation (5) we have

$$f[x_0, \dots, x_i] = \frac{1}{i! h^i} \Delta^i f_0. \quad (6)$$

Let us recall that

$$f[x_0, \dots, x_i] = \frac{f^{(i)}(\xi)}{i!} \quad (7)$$

where  $f(x)$  is a real-valued function defined on  $[a, b]$  and  $i$  times differentiable in  $]a, b[$  and  $\xi \in ]a, b[$ .

Taking  $i = n$  and  $f(x) = P_n(x)$  in Eqns. (6) and (7), we get

$$\begin{aligned} \Delta^n P_n(x_0) &= n! h^n P_n[x_0, \dots, x_n] = n! h^n \frac{P_n^{(n)}(\xi)}{n!} \\ &= h^n n! a_n. \end{aligned}$$

Since  $\Delta^{n+1} P_n(x_0) = \Delta^n P_n(x_1) - \Delta^n P_n(x_0)$

$$= h^n n! a_n - h^n n! a_n = 0.$$

This completes the proof

The shift operator  $E$  is defined as

$$E f_i = f_{i+1} \quad (8)$$

In general  $E f(x) = f(x + h)$ .

We have  $E^s f_i = f_{i+s}$ .

For example,

$$E^3 f_i = f_{i+3}, E^{1/2} f_i = f_{i+1/2} \text{ and } E^{-1/2} f_i = f_{i-1/2}$$

Now,

$$\Delta f_i = f_{i+1} - f_i = E f_i - f_i = (E - 1) f_i$$

Hence the shift and forward difference operations are related by

$$\Delta = E - 1$$

or  $E = 1 + \Delta$ .

Operating  $s$  times, we get

$$\Delta^s = (E - 1)^s = \sum_{j=0}^s \binom{s}{j} E^j (-1)^{s-j} \tag{9}$$

Making use of relation (8) in Eqn. (9), we get

$$\Delta^s f_i = \sum_{j=0}^s (-1)^{s-j} \binom{s}{j} f_{j+i}$$

We now give in Table 1, the forward differences of various orders using 5 values.

Table 1 : Forward Difference Table

$x$	$f(x)$	$\Delta f$	$\Delta^2 f$	$\Delta^3 f$	$\Delta^4 f$
$x_0$	$f_0$	$\Delta f_0$	$\Delta^2 f_0$	$\Delta^3 f_0$	$\Delta^4 f_0$
$x_1$	$f_1$	$\Delta f_1$	$\Delta^2 f_1$	$\Delta^3 f_1$	
$x_2$	$f_2$	$\Delta f_2$	$\Delta^2 f_2$		
$x_3$	$f_3$	$\Delta f_3$			
$x_4$	$f_4$				

Note that the forward difference  $\Delta^k f_0$  lie on a straight line sloping downward to the right.

### 11.2.2 Backward Differences

Let  $f$  be a real-valued function of  $x$ . Let the values of  $f(x)$  at  $n + 1$  equally spaced points  $x_0, x_1, \dots, x_n$  be  $f_0, f_1, \dots, f_n$  respectively.

The backward differences of  $f(x)$  of  $i$ th order at  $x_k = x_0 + kh$  are denoted by  $\nabla^i f_k$ . They are defined as follows:

$$\nabla^i f_k = \begin{cases} f_k, & i = 0 \\ \nabla^{i-1}(\nabla f_k) = \nabla^{i-1}[f_k - f_{k-1}], & i \geq 1 \end{cases} \tag{10}$$

where  $\nabla$  denotes backward difference operator.

Using (10), we have for

$$i = 1; \quad \nabla f_k = f_k - f_{k-1}$$

$$\begin{aligned}
 i = 2; \quad \nabla^2 f_k &= \nabla(f_k - f_{k-1}) \\
 &= \nabla f_k - \nabla f_{k-1} \\
 &= f_k - 2f_{k-1} + f_{k-2}
 \end{aligned}$$

$$\begin{aligned}
 i = 3; \quad \nabla^3 f_k &= \nabla^2[f_k - f_{k-1}] = \nabla^2 f_k - \nabla^2 f_{k-1} = \nabla[\nabla f_k] - \nabla[\nabla f_{k-1}] \\
 &= \nabla[f_k - f_{k-1}] - \nabla[f_{k-1} - f_{k-2}] \\
 &= \nabla f_k - \nabla f_{k-1} - \nabla f_{k-1} + f_{k-2} \\
 &= f_k - f_{k-1} - 2[f_{k-1} + f_{k-2}] + f_{k-2} - f_{k-3} \\
 &= f_k - 3f_{k-2} + 3f_{k-2} - f_{k-3}
 \end{aligned}$$

By induction we can prove the following lemma which connects the divided difference with the backward difference.

**Lemma 2:** The following relation holds

$$f[x_{n-k}, \dots, x_n] = \frac{1}{k!h^k} \nabla^k f(x_n) \tag{11}$$

The relation between the backward difference operator  $\nabla$  and the shift operator  $E$  is given by

$$\nabla = 1 - E^{-1} \text{ or } E = (1 - \nabla)^{-1}.$$

Since  $\nabla f_k = f_k - f_{k-1} = f_k - E^{-1}f_k = [1 - E^{-1}]f_k.$

Operating  $s$  times, we get

$$\begin{aligned}
 \nabla^s f_k &= [1 - E^{-1}]^s f_k = \left[ \sum_{m=0}^s (-1)^m \binom{s}{m} E^{-m} \right] f_k \\
 &= \sum_{m=0}^s (-1)^m \binom{s}{m} f_{k-m}.
 \end{aligned} \tag{12}$$

We can extend the binomial coefficient notation to include negative numbers, by letting

$$\binom{-s}{i} = \frac{-s(-s-1)(-s-2)\dots(-s-i+1)}{i!} = (-1)^i \frac{s(s+1)\dots(s+i-1)}{i!}$$

The backward differences of various orders with 5 nodes are given in Table 2.

**Table 2 : Backward Difference Table**

$x$	$f(x)$	$\nabla f$	$\nabla^2 f$	$\nabla^3 f$	$\nabla^4 f$
$x_0$	$f_0$				
$x_1$	$f_1$	$\nabla f_1$			
$x_2$	$f_2$	$\nabla f_2$	$\nabla^2 f_2$		
$x_3$	$f_3$	$\nabla f_3$	$\nabla^2 f_3$	$\nabla^3 f_3$	
$x_4$	$f_4$	$\nabla f_4$	$\nabla^2 f_4$	$\nabla^3 f_4$	$\nabla^4 f_4$

Let us consider the following example:

**Example 1:** Evaluate the differences

$$(a) \nabla^3 [a_2x^2 + a_1x + a_0]$$

$$(b) \nabla^3 [a_3x^3 + a_2x^2 + a_1x + a_0]$$

**Solution:** (a)  $\nabla^3 [a_2x^2 + a_1x + a_0] = 0$

$$(b) \nabla^3 [a_3x^3 + a_2x^2 + a_1x + a_0]$$

$$= a_3 \nabla^3(x^3) + \nabla^3 [a_2x^2 + a_1x + a_0]$$

$$= a_3 \cdot 3! h^2$$

Note that the backward differences  $\nabla^k f_k$  lie on a straight line sloping upward to the right.

Also note that  $\Delta f_k = \nabla f_{k+1} = f_{k+1} - f_k$ .

Try to show that  $\Delta^4 f_0 = \nabla^4 f_4$ .

Let us now discuss about the central differences.

### 11.2.3 Central Differences

The first order central difference of  $f$  at  $x_k$ , denoted by  $\delta f_k$ , is defined as

$$\delta f = f(x + h/2) - f(x - h/2) = f_{k+1/2} - f_{k-1/2}$$

Operating with  $\delta$ , we obtain the higher order central differences as

$$\delta^s f = \delta^{s-1} (\delta f) = \delta^{s-1} f_{k+1/2} - \delta^{s-1} f_{k-1/2} \quad (13)$$

with  $\delta^s f_k = f_k$  when  $s = 0$ .

The second order central difference is given by

$$\delta^2 f_k = \delta[f_{k+1/2} - f_{k-1/2}] = \delta[f_{k+1/2}] - \delta[f_{k-1/2}]$$

$$= f_{k+1} - f_k - f_k + f_{k-1}$$

$$= f_{k+1} - 2f_k + f_{k-1}$$

Similarly,

$$\delta^3 f_k = f_{k+3/2} - 3f_{k+1/2} + 3f_{k-1/2} - f_{k-3/2}$$

and  $\delta^4 f_k = f_{k+2} - 4f_{k+1} + 6f_k - 4f_{k-1} + f_{k-2}$ .

Notice that the even order differences at a tabular value  $x_k$  are expressed in terms of tabular values of  $f$  and odd order differences at a tabular value  $x_k$  are expressed in terms of non-tabular value of  $f$ . Also note that the coefficients of  $\delta^s f_k$  are the same as those of the binomial expansion of  $(1-x)^s$ ,  $s = 1, 2, 3, \dots$ .

Since

$$\delta f_k = f_{k+1/2} - f_{k-1/2} = (E^{1/2} - E^{-1/2}) f_k$$

We have the operation relation

$$\delta = E^{1/2} - E^{-1/2} \quad (14)$$

The central differences at a non-tabular point  $x_{k+1/2}$  can be calculated in a similar way. For example,

$$\delta f_{k+1/2} = f_{k+1} - f_k$$

$$\delta^2 f_{k+1/2} = f_{k+3/2} - 2f_{k+1/2} + f_{k-1/2}$$

$$\delta^3 f_{k+1/2} = f_{k+2} - 3f_{k+1} + 3f_k - f_{k-1} \tag{15}$$

$$\delta^4 f_{k+1/2} = f_{k+3/2} - 4f_{k+3/2} + 6f_{k+1/2} - 4f_{k-1/2} + f_{k-3/2}$$

Relation (15) can be obtained easily by using the relation (14)

We have

$$\begin{aligned} \delta^s f_k &= [E^{1/2} - E^{-1/2}]^s f_k \\ &= \left[ \sum_{i=0}^s (-1)^i E^{-i/2} E^{(s-i)/2} \binom{s}{i} \right] f_k \\ &= \sum_{i=0}^s (-1)^i \binom{s}{i} f_{k+(s/2)-i} \end{aligned} \tag{16}$$

The following formulas can also be established:

$$f[x_0, \dots, x_{2m}] = \frac{1}{(2m)!h^{2m}} \delta^{2m} f_m \tag{17}$$

$$f[x_0, \dots, x_{2m+1}] = \frac{1}{(2m+1)!h^{2m+1}} \delta^{2m+1} f_{m+1/2} \tag{18}$$

$$f[x_{-m}, \dots, x_0, \dots, x_m] = \frac{1}{(2m)!h^{2m}} \delta^{2m} f_0 \tag{19}$$

$$f[x_{-m}, \dots, x_0, \dots, x_{m+1}] = \frac{1}{(2m+1)!h^{2m+1}} \delta^{2m+1} f_{1/2} \tag{20}$$

$$f[x_{-(m+1)}, \dots, x_0, \dots, x_m] = \frac{1}{(2m+1)!h^{2m+1}} \delta^{2m+1} f_{-1/2} \tag{21}$$

We now give below the central difference table with 5 nodes.

Table 3 : Central Difference Table

x	f	δf	δ²f	δ³f	δ⁴f
x <sub>-2</sub>	f <sub>-2</sub>				
		δf <sub>-3/2</sub>			
x <sub>-1</sub>	f <sub>-1</sub>		δ²f <sub>-1</sub>		
		δf <sub>-1/2</sub>		δ³f <sub>-1/2</sub>	
x <sub>0</sub>	f <sub>0</sub> .....		δ²f <sub>0</sub> .....		δ⁴f <sub>0</sub>
		δf <sub>1/2</sub>		δ³f <sub>1/2</sub>	
x <sub>1</sub>	f <sub>1</sub>		δ²f <sub>1</sub>		
		δf <sub>3/2</sub>			
x <sub>2</sub>	f <sub>2</sub>				

Note that the differences δ<sup>2m</sup>f<sub>0</sub> lie on a horizontal line shown by the dotted lines.



Table 4 : Central Difference Table

x	f	$\delta f$	$\delta^2 f$	$\delta^3 f$	$\delta^4 f$
$x_0$	$f_0$				
		$\delta f_{1/2}$			
$x_1$	$f_1$		$\delta^2 f_1$		
		$\delta f_{3/2}$		$\delta^3 f_{3/2}$	
$x_2$	$f_2$	.....	$\delta^2 f_2$	.....	$\delta^4 f_2$
		$\delta f_{5/2}$		$\delta^3 f_{5/2}$	
$x_3$	$f_3$		$\delta^2 f_3$		
		$\delta f_{7/2}$			
$x_4$	$f_4$				

Note that the differences  $\delta^{2m}f_2$  lie on a horizontal line.

We now define the mean operator  $\mu$  as follows

$$\begin{aligned} \mu f_k &= \frac{1}{2} [f_{k+1/2} + f_{k-1/2}] \\ &= \frac{1}{2} [E^{1/2} + E^{-1/2}] f_k. \end{aligned}$$

Hence

$$\mu = \frac{1}{2} [E^{1/2} + E^{-1/2}]$$

Relation Between the Operators  $\Delta$ ,  $\nabla$ ,  $\delta$  and  $\mu$

We have expressed  $\Delta$ ,  $\nabla$ ,  $\delta$  and  $\mu$  in terms of the operator  $E$  as follows

$$\Delta = E - 1$$

$$\nabla = 1 - E^{-1}$$

$$\delta = E^{1/2} - E^{-1/2}$$

$$\mu = \frac{1}{2} [E^{1/2} + E^{-1/2}]$$

$$\Delta = E(1 - E^{-1}) = E\nabla$$

$$= E^{1/2} (E^{1/2} - E^{-1/2}) = E^{1/2}\delta$$

Also  $E^{1/2} = \mu + \frac{\delta}{2}$

$$E^{-1/2} = \mu - \frac{\delta}{2}$$

Example 2: (a) Express  $\Delta^3 f_1$  as a backward difference.

(b) Express  $\Delta^3 f_1$  as a central difference.

(c) Express  $\delta^2 f_2$  as a forward difference.

Solution:

$$(a) \quad \Delta^3 f_1 = (E\nabla)^3 f_1 = E^3 \nabla^3 f_1 = \nabla^3 E^3 f_1 = \nabla^3 f_4 \quad (\Delta = E\nabla)$$

$$(b) \quad \Delta^3 f_1 = [E^{1/2}\delta]^3 f_1 = E^{3/2}\delta^3 f_1 = \delta^3 E^{3/2} f_1 = \delta^3 f_{5/2} \quad (\nabla = E^{1/2}\delta)$$

$$(c) \quad \delta^2 f_2 = [E^{-1/2}\delta]^2 f_2 = E^{-1}\Delta^2 f_2 = \Delta^2 E^{-1} f_2 = \Delta^2 f_1 \quad (\delta = E^{-1/2}\Delta)$$

Example 3: Prove that (a)  $\mu^2 = 1 + \frac{\delta^2}{4}$

$$(b) \quad \mu\delta = \frac{1}{2}(\Delta + \nabla)$$

$$(c) \quad \sqrt{1 + \mu^2\delta^2} = 1 + \frac{\delta^2}{2}$$

Solution: (a) We have  $\mu = \frac{1}{2}[E^{1/2} + E^{-1/2}]$

$$\begin{aligned} \mu^2 &= \frac{(E^{1/2} + E^{-1/2})^2}{4} = \frac{(E^{1/2} - E^{-1/2})^2 + 4}{4} \\ &= 1 + \frac{(E^{1/2} - E^{-1/2})^2}{4} \\ &= 1 + \frac{\delta^2}{4} \end{aligned}$$

(b) L.H.S.

$$\mu\delta = \frac{1}{2}(E^{1/2} + E^{-1/2})(E^{1/2} - E^{-1/2}) = \frac{1}{2}(E - E^{-1})$$

R.H.S.

$$\frac{1}{2}(\Delta + \nabla) = \frac{1}{2}[(E-1) + (1-E^{-1})] = \frac{1}{2}(E-E^{-1}).$$

Hence, the result.

(c) We have

$$\begin{aligned} \mu\delta &= \frac{1}{2}(E^{1/2} + E^{-1/2})(E^{1/2} - E^{-1/2}) = \frac{1}{2}(E - E^{-1}) \\ \therefore 1 + \mu^2\delta^2 &= 1 + \frac{(E - E^{-1})^2}{4} = \frac{(E - E^{-1})^2 + 4}{4} = \frac{(E + E^{-1})^2}{4} \\ \therefore \sqrt{1 + \mu^2\delta^2} &= \frac{E + E^{-1}}{2} = \frac{(E^{1/2} - E^{-1/2})^2 + 2}{2} \\ &= \frac{\delta^2 + 2}{2} = 1 + \frac{\delta^2}{2}. \end{aligned}$$

E1) Express  $\nabla^4 f_5$  in terms of function values.

E2) Show that  $(E + 1)\delta = 2(E - 1)\mu$ .

### 11.3 DIFFERENCE FORMULAS

We shall now derive different difference formulas using the results obtained in the preceding section (Section 11.2)

#### 11.3.1 Newton's Forward-Difference Formula

In Unit 10, we have derived Newton's form of interpolating polynomial (using divided

differences). We have also established in Sec. 11.2.1, the following relationship between divided differences and forward differences

$$f[x_k, \dots, x_{k+n}] = \frac{1}{n!h^n} \Delta^n f_k \quad (22)$$

Substituting the divided differences in terms of the forward differences in the Newton's form, and simplifying we get Newton's forward-difference form. The Newton's form of interpolating polynomial interpolating at  $x_k, x_{k+1}, \dots, x_{k+n}$  is

$$P_n(x) = \sum_{i=0}^n (x - x_k)(x - x_{k+1}) \dots (x - x_{k+i-1}) f[x_k, \dots, x_{k+i}]$$

Substituting (22), we obtain

$$P_n(x) = \sum_{i=0}^n (x - x_k)(x - x_{k+1}) \dots (x - x_{k+i-1}) \frac{1}{i!h^i} \Delta^i f_k \quad (23)$$

Setting  $k = 0$ , we have the form

$$\begin{aligned} P_n(x) &= \sum_{i=0}^n \frac{1}{i!h^i} (x - x_0)(x - x_1) \dots (x - x_{i-1}) \Delta^i f_0 \\ &= f_0 + \frac{(x - x_0)}{1!} \frac{\Delta f_0}{h} + \frac{(x - x_0)(x - x_1)}{h^2} \frac{\Delta^2 f_0}{h^2} + \dots + \frac{(x - x_0) \dots (x - x_{n-1})}{n!} \frac{\Delta^n f_0}{h^n} \quad (24) \end{aligned}$$

Using the transformation (2), we have

$$x - x_{k+j} = x_0 + sh - [x_0 + (k+j)h] = (s - k - j)h$$

Hence (23) can be rewritten as

$$\begin{aligned} P_n(x) &= P(x_0 + sh) = \sum_{i=0}^n \frac{1}{i!} (s - k)(s - k - 1) \dots (s - k - i + 1) \Delta^i f_k \\ &= \sum_{i=0}^n \Delta^i f_k \begin{bmatrix} s - k \\ i \end{bmatrix} \\ &= f_k + (s - k)\Delta f_k + \frac{(s - k)(s - k - 1)}{2!} \Delta^2 f_k + \dots + \frac{(s - k)(s - n - 1)}{n!} \Delta^n f_k \quad (25) \end{aligned}$$

of degree  $\leq n$ .

Setting  $k = 0$  in (25) we get the formula

$$P_n(x_0 + sh) = \sum_{i=0}^n \Delta^i f_0 \begin{bmatrix} s \\ i \end{bmatrix} \quad (26)$$

The form (23), (24), (25) or (26) is called the Newton's forward-difference formula.

The error term is now given by

$$E_n(x) = \begin{bmatrix} s \\ n+1 \end{bmatrix} h^{n+1} f^{(n+1)}(\xi)$$

**Example 4:** Find the Newton's forward-difference interpolating polynomial which agrees with the table of values given below. Hence obtain the value of  $f(x)$  at  $x = 1.5$ .

Interpolation

x	1	2	3	4	5	6
f(x)	10	19	40	79	142	235

Solution: We form a table of forward differences of f(x).

Table 5 : Forward Differences

x	f(x)	$\Delta f$	$\Delta^2 f$	$\Delta^3 f$
1	10			
		9		
2	19		12	
		21		6
3	40		18	
		39		6
4	79		24	
		63		6
5	142		30	
		93		
6	235			

Since the third order differences are constant, the higher order differences vanish and we can infer that f(x) is a polynomial of degree 3 and the Newton's forward-differences interpolation polynomial exactly represents f(x) and is not an approximation to f(x). The step length in the data is h = 1. Taking x<sub>0</sub> = 1 and the subsequent values of x as x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>5</sub>, the Newton's forward-differences interpolation polynomial.

$$f(x) = f_0 + (x-1)\Delta f_0 + \frac{(x-1)(x-2)}{2!} \Delta^2 f_0 + \frac{(x-1)(x-2)(x-3)}{3!} \Delta^3 f_0$$

becomes

$$f(x) = 10 + (x-1)(9) + \frac{(x-1)(x-2)}{2} (12) + \frac{(x-1)(x-2)(x-3)}{6} (6)$$

$$= 10 + (x-1) + 6(x-1)(x-2) + (x-1)(x-2)(x-3)$$

which on simplification gives

$$f(x) = x^3 + 2x + 7$$

$$\therefore f(1.5) = (1.5)^3 + 2(1.5) + 7$$

$$= 3.375 + 3 + 7 = 13.375$$

Note:

If we want only the value of f(1.5) and the interpolation polynomial is not needed, we can use the formula (26). In this case,

$$s = \frac{x - x_0}{h} = \frac{1.5 - 1}{1} = 0.5$$

and

$$f(1.5) = 10 + (0.5)(9) + \frac{(0.5)(-0.5)}{2} (12) + \frac{(0.5)(-0.5)(-1.5)}{6} (6)$$

$$= 10 + 4.5 - 1.5 + 0.375$$

$$= 13.375.$$

E3) The population of a town in the decennial census was given below. Estimate the population for the year 1915.

Year: x	1911	1921	1931	1941	1951
Population: y (in thousands)	46	66	81	93	101

Example 5: From the following table, find the number of students who obtained less than 45 marks.

Marks	30 - 40	40 - 50	50 - 60	60 - 70	70 - 80
No. of students	31	42	51	35	31

Solution: We form a table of the number of students  $f(x)$  whose marks are less than  $x$ . In other words, we form a cumulative frequency table.

Table 6 : Frequency Table

x	f(x)	$\Delta f$	$\Delta^2 f$	$\Delta^3 f$	$\Delta^4 f$
40	31				
		42			
50	73		9		
		51		-25	
60	124		-16		
		35		12	
70	159		-4		
		31			
80	190				

We have  $x_0 = 40$ ,  $x = 45$  and  $h = 10$

$$\therefore s = 0.5$$

$$\therefore f(45) = 31 + (0.5)(42) + \frac{(0.5)(-0.5)}{2}(9) + \frac{(0.5)(-0.5)(-1.5)}{6}(-25)$$

$$+ \frac{(0.5)(-0.5)(-1.5)(-2.5)}{24}(37)$$

$$= 31 + 21 - 1.125 - 1.5625 - 1.4453$$

$$= 47.8672 = 48$$

$\therefore$  The number of students who obtained less than 45 marks is approximately 48.

E4) From the following table, find the value of  $y(0.23)$ :

x	0.20	0.22	0.24	0.26	0.28	0.30
y	1.6596	1.6698	1.6804	1.6912	1.7024	1.7139

E5) Find the cubic polynomial which approximate  $y(x)$  given that

$$y(0) = 1, y(1) = 0, y(2) = 1 \text{ and } y(3) = 10.$$

E6) The following table gives the values of  $\tan x$  for  $0.1 \leq x \leq 0.3$ . Find the value of  $\tan(0.12)$ .

x	0.10	0.15	0.20	0.25	0.30
$\tan x$	0.1003	0.1511	0.2027	0.2553	0.3093

E7) The following table gives the population of a town in ten consecutive censuses. Calculate the population in the year 1915 and 1918. Hence obtain the increase in population during the period 1915 and 1918.

Year x	1911	1921	1931	1941	1951	1961
Population y (in thousands)	12	15	20	27	39	52

E8) Find the number of men getting wages between Rs. 10 and Rs. 15 from the following table.

Wages in Rs. x	0 - 10	10 - 20	20 - 30	30 - 40
No. of men y	9	30	35	42

E9) The following table shows the monthly premiums to be paid to a company at different ages. Find the premium to be paid at the age of 26 years.

Age	20	24	28	32	36
Premium in Rs.	14.27	15.81	17.72	19.96	22.48

E10) The area  $A$  of a circle of diameter  $d$  is given in the following table. Find the area of the circle when the diameter is 82 units.

d	80	85	90	95	100
A	5026	5674	6362	7088	7854

E11) In an examination, the number of candidates who secured marks in certain limits were as follows:

Marks	0 - 9	20 - 39	40 - 59	60 - 79	80 - 99
No. of candidates	41	62	65	50	17

Find the number of candidates whose marks are 25 or less.

E12) The following table gives the amount of a chemical dissolved in water at different temperatures.

Temperature	10°	15°	20°	25°	30°	35°
Amount dissolved	19.97	21.51	22.47	23.52	24.65	25.89

Compute the amount dissolved at 8°.

E13) Find a polynomial which fits the following data :

Interpolation at Equally Spaced Point

x	3	5	7	9	11
y	6	24	58	108	174

### 11.3.2 Newton's Backward-Difference Formula

Reordering the interpolating nodes as  $x_n, x_{n-1}, \dots, x_0$  and applying the Newton's divided difference form, we get

$$P_n(x) = f[x_n] + (x - x_n) f[x_{n-1}, x_n] + (x - x_n)(x - x_{n-1}) f[x_{n-2}, x_{n-1}, x_n] + \dots + (x - x) \dots (x - x_n) f[x_0, \dots, x_n] \quad (27)$$

We may also write

$$\begin{aligned} P_n(x) &= P_n \left[ x_n + \frac{x - x_n}{h} h \right] \\ &= P_n[x_n + sh] = \sum_{i=0}^n (x - x_n)(x - x_{n-1}) \dots (x - x_{n-i+1}) f[x_n, \dots, x_{n-i}] \\ &= \sum_{i=0}^n \frac{1}{i! h^i} (x - x_n)(x - x_{n-1}) \dots (x - x_{n-i+1}) \nabla^i f_n \end{aligned} \quad (28)$$

Set  $x = x_n + sh$ , then

$$x - x_i = x_n + sh - [x_n - (n-i)h] = (s + n - i)h$$

$$x - x_{n-j} = (s + n - n + j)h = (s + j)h$$

and

$$(x - x_n)(x - x_{n-1}) \dots (x - x_{n-i+1}) = s(s+1) \dots s(s+i-1)h^i$$

Equation (28) becomes

$$\begin{aligned} P_n(x) &= \sum_{i=0}^n \frac{1}{i!} s(s+1) \dots (s+i-1) \nabla^i f_n \\ &= f_n + s \nabla f_n + \frac{s(s+1)}{2!} \nabla^2 f_n + \frac{s(s+1) \dots (s+n-1)}{n!} \nabla^n f_n \end{aligned} \quad (29)$$

We have seen already that

$$\begin{bmatrix} -s \\ k \end{bmatrix} = (-1)^k \frac{s(s+1) \dots (s+k-1)}{k!}$$

Hence, equation (29) can be written as

$$\begin{aligned} P_n(x) &= f(x_n) + (-1) \begin{bmatrix} -s \\ 1 \end{bmatrix} \nabla f(x_n) + (-1)^2 \begin{bmatrix} -s \\ 2 \end{bmatrix} \nabla^2 f(x_n) \\ &+ \dots + (-1)^n \begin{bmatrix} -s \\ n \end{bmatrix} \nabla^n f(x_n) \end{aligned}$$

or

$$P_n(x) = \sum_{k=0}^n (-1)^k \begin{bmatrix} -s \\ k \end{bmatrix} \nabla^k f(x_n). \quad (30)$$

Equation (27), (28) or (29) is called the Newton's backward-difference form.

In this case error is given by

$$E_n(x) = (-1)^{n+1} \frac{s(s+1) \dots (s+n)}{(n+1)!} h^{n+1} f^{(n+1)}(\xi). \quad (31)$$

The backward-difference form is suitable for approximating the value of the function at  $x$  that lies towards the end of the table.

### Example 6

Find the Newton's backward differences interpolating polynomial for the data of Example 4.

**Solution:** We form the table of backward differences of  $f(x)$ .

Table 7 : Backward Difference Table

$x$	$f(x)$	$\nabla f$	$\nabla^2 f$	$\nabla^3 f$
1	10			
		9		
2	19		12	
		21		6
3	40		18	
		39		6
4	79		24	
		63		6
5	142		30	
		93		6
6	<u>235</u>			

Tables 5 and 7 are the same except that we consider the differences of Table 7 as backward differences. If we name the abscissas as  $x_0, x_1, \dots, x_5$ , then  $x_n = x_5 = 6$ ,  $f_n = f_5 = 235$ . With  $h = 1$ , the Newton's backward differences polynomial for the given data is given by

$$\begin{aligned} P(x) &= f_5 + (x-x_5) \nabla f_5 + \frac{(x-x_5)(x-x_4)}{2!} \nabla^2 f_5 + \frac{(x-x_5)(x-x_4)(x-x_3)}{3!} \nabla^3 f_5 \\ &= 235 + (x-6)(93) + \frac{(x-6)(x-5)}{2} (30) + \frac{(x-6)(x-5)(x-4)}{6} (6) \\ &= 235 + 93(x-6) + 15(x-6) + (x-4)(x-5)(x-6) \end{aligned}$$

which on simplification gives

$$P(x) = x^3 + 2x + 7,$$

which is the same as the Newton's forward differences interpolation polynomial in Example 4.



**Example 7:** Estimate the value of  $f(1.45)$  from the data given below:

Interpolation at Equally Spaced Point

x	1.1	1.2	1.3	1.4	1.5
f(x)	1.3357	1.5095	1.6984	1.9043	2.1293

**Solution:** We form the backward differences table for the data given.

Table 8 : Backward Differences Table

x	f(x)	$\nabla f$	$\nabla^2 f$	$\nabla^3 f$	$\nabla^4 f$
1.1	1.3357				
		0.1738			
1.2	1.5095		0.0151		
		0.1889		0.0019	
1.3	1.6984		0.0170		0.0002
		0.2059		0.0021	
1.4	1.9043		0.0191		
		0.2250			
1.5	2.1293				

Here  $x_n = 1.5$ ,  $x = 1.45$ ,  $h = 0.1$

$$\therefore s = \frac{x - x_n}{h} = \frac{1.45 - 1.5}{0.1} = -0.5$$

The Newton's backward differences interpolation formula gives

$$\begin{aligned} f(x) &= f_n + s\nabla f_n + \frac{s(s+1)}{2!} \nabla^2 f_n + \frac{s(s+1)(s+2)}{3!} \nabla^3 f_n + \frac{s(s+1)(s+2)(s+3)}{4!} \nabla^4 f_n \\ &= 2.1293 + (-0.5)(0.2250) + \frac{(-0.5)(0.5)}{2} (0.0191) + \\ &\quad \frac{(-0.5)(0.5)(1.5)}{6} (0.0021) + \frac{(-0.5)(0.5)(2.5)}{24} (0.0002) \\ &= 2.1293 - 0.1125 - 0.00239 - 0.00013 - 0.0000078 \\ &= 2.01427 \approx 2.0143 \end{aligned}$$

- E14) From the table of values of E4, find the value of  $y$  when  $x = 0.29$ .
- E15) Using the backward differences interpolation, find the polynomial which agrees with the values of  $y(x)$  where  
 $y(0) = 1$ ,  $y(1) = 0$ ,  $y(2) = 1$  and  $y(3) = 10$ .
- E16) From the table of values of E6, find the values of  $\tan(0.26)$  and  $\tan(0.40)$ .
- E17) From the data of E7, find the increase in population from 1954 to 1958 by applying the Newton's backward differences interpolation formula with 1961 as origin.
- E18) Find the area of the circle when the diameter is 98 from the data of E10.
- E19) In E11, find the number of candidates whose marks are less than or equal to (i) 70, (ii) 89.
- E20) Find the Newton's backward differences interpolating polynomial which fits the data of E13.

### 11.3.3 Stirling's Central Difference Form

A number of central difference formulas are available which can be used according to a situation to maximum advantage. But we shall consider only one such method known as Stirling's method. This formula is used whenever interpolation is required of  $x$  near the middle of the table of values.

For the central difference formulas, the origin  $x_0$ , is chosen near the point being approximated and points below  $x_0$  are labelled as  $x_1, x_2, \dots$  and those directly above as  $x_{-1}, x_{-2}, \dots$  (as in Table 3). Using this convention, Stirling's formula for interpolation is given by

$$\begin{aligned}
 P_n(x) = & f(x_0) + \frac{s}{2} [\delta f_{1/2} + \delta f_{-1/2}] + \frac{s^2}{2!} \delta^2 f_0 \\
 & + \frac{s(s^2 - 1^2)}{3!} \frac{1}{2} [\delta^3 f_{1/2} + \delta^3 f_{-1/2}] + \dots \\
 & + \frac{s(s^2 - 1^2) s(s^2 - 2^2) \dots [s^2 - (p-1)^2]}{(2p-1)!} \frac{1}{2} [\delta^{2p-1} f_{1/2} + \delta^{2p-1} f_{-1/2}] \\
 & + \frac{s^2(s^2 - 1^2) \dots [s^2 - (p-1)^2]}{(2p)!} \delta^{2p} f_0 \\
 & + \frac{s(s^2 - 1^2) \dots (s^2 - p^2)}{(2p+1)!} \frac{1}{2} [\delta^{2p+1} f_{1/2} + \delta^{2p+1} f_{-1/2}] \tag{32}
 \end{aligned}$$

where  $s = (x - x_0)/h$  and if  $n = 2p + 1$  is odd.

If  $n = 2p$  is even, then the same formula is used deleting the last term.

The Stirling's interpolation is used for calculation when  $x$  lies between  $x_0 - \frac{1}{4}h$  and  $x_0 + \frac{1}{4}h$ .

It may be noted from the Table 3, that the odd order differences at  $x_{-1/2}$  are those which lie along the horizontal line between  $x_0$  and  $x_{-1}$ . Similarly, the odd order differences at  $x_{1/2}$  are those which lie along the horizontal line between  $x_0$  and  $x_1$ . Even order differences at  $x_0$  are those which lie along the horizontal line through  $x_0$ .

**Example 8:** Using Stirling's formula, find the value of  $f(1.32)$  from the following table of values.

$x$	1.1	1.2	1.3	1.4	1.5
$f(x)$	1.3357	1.5095	1.6984	1.9043	2.1293

**Solution:**

Table 9: Central Difference

$x$	$f(x)$	$\delta f$	$\delta^2 f$	$\delta^3 f$	$\delta^4 f$
1.1	1.3357				
		0.1738			
1.2	1.5095		0.0151		
		0.1889		0.0019	
( $x_0$ ) 1.3	1.6984		0.0170		0.0002
		0.2059		0.0021	
1.4	1.9043		0.0191		
		0.2250			
1.5	2.1293				

Choose  $x_0 = 1.3$

$$\therefore s = \frac{(x-x_0)}{h} = \frac{1.32-1.3}{0.1} = 0.2.$$

From Eqn. (32), we have

$$f(x) = f_0 + \frac{s}{2} [\delta f_{-1/2} + \delta f_{1/2}] + \frac{s^2}{2!} \delta^2 f_0 + \frac{s(s^2-1^2)}{3!} \frac{1}{2} [\delta^3 f_{-1/2} + \delta^3 f_{1/2}] + \frac{s^2(s^2-1^2)}{4!} \delta^4 f_0.$$

Now,

$$\frac{1}{2} [\delta f_{-1/2} + \delta f_{1/2}] = \frac{1}{2}(0.1889 + 0.2059) = 0.1974$$

$$\frac{1}{2} [\delta^3 f_{-1/2} + \delta^3 f_{1/2}] = \frac{1}{2}(0.0019 + 0.0021) = 0.0020$$

Also  $\delta^2 f_0 = 0.0170$ ,  $\delta^4 f_0 = 0.0002$ .

Substituting in the above equation, we get

$$\begin{aligned} f(x) &= 1.6984 + (0.2)(0.1974) + \frac{0.04}{2}(0.0170) + \frac{(0.2)(-0.96)}{6}(0.0020) \\ &\quad + \frac{(0.04)(-0.96)}{24}(0.0002) \\ &= 1.6984 + 0.03948 + 0.00034 - 0.00006 - 0 \\ &= 1.73816 \approx 1.7382. \end{aligned}$$

In the following exercises, use the Stirling's interpolation formula.

E21) Find  $f(1.725)$  from the following table.

x	1.5	1.6	1.7	1.8	1.9
f(x)	4.4817	4.9530	5.4739	6.0496	6.6859

E22) Find the value of  $f(1.22)$  from the following table.

x	1.0	1.1	1.2	1.3	1.4
f(x)	0.8415	0.8912	0.9320	0.9636	0.9855

E23) Evaluate  $f(4.325)$  from the following.

x	4.1	4.2	4.3	4.4	4.5
f(x)	30.1784	33.3507	36.8567	40.7316	45.0141

E24) Find the value of  $y$  when  $x = 30$  from the table.

x	21	25	29	33	37
f(x)	18.4708	17.814	17.1070	16.3432	15.5134

E25) Find the approximate value of  $y(2.15)$  from the table.

x	0	1	2	3	4
y	6.9897	7.4036	7.7815	8.1281	8.4510

## 11.4 SUMMARY

In this unit, we have derived interpolation formulas for data with equally spaced values of the argument. We have seen how to find the value of  $f(x)$  for a given value of  $x$  by applying an appropriate interpolation formula derived in this section. The application of the formulas derived in this section is easier when compared to the application of the formulas derived in Units 9 and 10. However, the formulas derived in this unit can only be applied to data with equally spaced arguments whereas the formulas derived in Units 9 and 10 can be applied for data with equally spaced or unequally spaced arguments. Thus, the formulas derived in Units 9 and 10 are of a more general nature than those of Unit 11. The interpolation polynomial which fits a given data can be determined by using any of the formulas derived in this section which will be unique whatever be the interpolation formula that is used.

The interpolation formulas derived in this unit are listed below:

1. Newton's forward difference formula:

$$P_n(x) = P_n(x_0 + sh) = \sum_{i=0}^n \left[ \begin{matrix} s \\ i \end{matrix} \right] \Delta^i f_0$$

$$f_0 + s\Delta f_0 + \frac{s(s-1)}{2!} \Delta^2 f_0 + \dots + \frac{s(s-1)\dots s(s-n+1)}{n!} \Delta^n f_0$$

where  $s = (x - x_0)/h$ .

2. Newton's backward difference formula:

$$P_n(x) = P_n(x_n + sh) = \sum_{k=0}^n (-1)^k \left[ \begin{matrix} -s \\ k \end{matrix} \right] \nabla^k f_n \quad \text{where } s = (x - x_0)/h$$

3. Stirling's central difference formula:

$$P_n(x) = P_n(x_0 + sh) = f_0 + \frac{s}{2} [\delta f_{1/2} + \delta f_{-1/2}] + \frac{s^2}{2!} \delta^2 f_0 + \frac{s(s^2-1^2)}{3!} \frac{1}{2} [\delta^3 f_{1/2} + \delta^3 f_{-1/2}]$$

$$+ \dots + \frac{s(s^2-1^2)(s^2-2^2)\dots(s^2-(p-1)^2)}{(2p+1)!} \frac{1}{2} [\delta^{2p-1} f_{1/2} + \delta^{2p-1} f_{-1/2}]$$

$$+ \frac{s^2(s^2-1^2)\dots(s^2-(p-1)^2)}{(2p)!} s^2 f_0 + \frac{s^2(s^2-1^2)\dots(s^2-p^2)}{(2p+1)!} [\delta^{2p+1} f_{1/2} + \delta^{2p+1} f_{-1/2}]$$

if  $n = 2p + 1$  is odd. If  $n = 2p$  is even, the same formula is used deleting the last term.

## 11.5 SOLUTIONS/ANSWERS

E1) From Eqn. (12)  $\nabla^4 f_x = f_5 - 4f_4 + 6f_3 - 4f_2 + f_1$

E2) LHS =  $E^{1/2} (E^{1/2} + E^{-1/2}) = E^{1/2} 2\mu\delta = 2E^{1/2} \mu\delta$

RHS =  $2E^{1/2}(E^{1/2} - E^{-1/2})\mu = 2E^{1/2} \mu\delta$ .

- E3) The forward differences table is given below.

Table 10

x	y	$\Delta y$	$\Delta^2 y$	$\Delta^3 y$	$\Delta^4 y$
1911	46				
		20			
1921	66		-5		
		15		2	
1931	81		-3		-3
		12		1	
1941	93		-4		
		8			
1951	101				

Taking  $x_0 = 1911$ ,  $x = 1915$ ,  $h = 10$ , we get

$$s = \frac{1051 - 1911}{10} = 0.4$$

$$\begin{aligned} \therefore y(1915) &= 46 + (0.4)(20) + \frac{(0.4)(-0.6)}{2}(-5) + \frac{(0.4)(-0.6)(-1.6)}{6}(2) \\ &\quad + \frac{(0.4)(-0.6)(-1.6)(-2.6)}{24} \\ &= 46 + 8 + 0.6 + 0.128 + 0.1248 \\ &\approx 54.8528 \end{aligned}$$

or  $y(1915) \approx 54.85$  thousands.

- E4) 1.6751
- E5)  $y = x^3 - 2x^2 + 1$ ,  $y(4) = 33$
- E6) 0.1205
- E7) 12.54 thousands, 13.64 thousands, 1.1 thousands
- E8) 15
- E9) 16.25 Rs.
- E10) 5281
- E11) 58
- E12) 18.79 (Hint: Take all differences into consideration)
- E13)  $2x^2 - 7x + 9$
- E14) 1.7081
- E15)  $x^3 - 2x^2 + 1$
- E16) 0.2662, 0.4241
- E17) Population in 1954 is 43.33 thousands and the population in 1958 is 48.81. Hence the increase in population is approximately 5.48 thousands.
- E18) 7543
- E19) Hint: The number of candidates  $f(x)$  whose marks are less than or equal to  $x$  is as follows:

$x$	19	39	59	79	99
$f(x)$	41	103	168	218	235

- (i) Take 79 as origin and determine  $f(70)$

We get  $f(70) = 199$ .

- (ii) Take 99 as origin and obtain  $f(89) = 232$ .

E20)  $2x^2 - 7x + 9$

## Interpolation

E21)  $x = 1.725, x_0 = 1.7, h = 0.1$

$$\therefore s = \frac{1.725 - 1.7}{0.1} = 0.25$$

Table 11

x	f(x)	$\delta f$	$\delta^2 f$	$\delta^3 f$	$\delta^4 f$
1.5	4.4817				
		0.4713			
1.6	4.9530		0.0496		
		0.5209		0.0052	
1.7	5.4739		0.0548		0.0006
		0.5757		0.0058	
1.8	6.0496		0.0606		
		0.6363			
1.9	6.6849				

$$\mu\delta f_0 = \frac{1}{2} (0.5209 + 0.5757) = 0.5483$$

$$\mu\delta^3 f_0 = \frac{1}{2} (0.00552 + 0.0058) = 0.0055$$

$$\text{Also } \delta^2 f_0 = 0.0548, \delta^4 f_0 = 0.0006.$$

$$\begin{aligned} \therefore f(1.725) &= 5.4739 + (0.25)(0.5483) + \frac{(0.0625)}{2}(0.0548) \\ &\quad + \frac{(0.25)(-0.9375)}{6}(0.0055) + \frac{(0.0625)(-0.9375)}{24}(0.0006) \\ &= 5.4739 + 0.13708 + 0.00171 - 0.00021 - 0 \\ &= 5.61248 \approx 5.6125 \end{aligned}$$

E22) 0.9391

E23) 37.7894

E24) 16.9217

E25) 7.8352

NOTES

**NOTES**





UTTAR PRADESH  
RAJARSHI TANDON OPEN UNIVERSITY

# UGMM - 10

## Numerical Analysis

Block

# 4

### **NUMERICAL DIFFERENTIATION INTEGRATION AND SOLUTION OF DIFFERENTIAL EQUATIONS**

---

#### **UNIT 12**

**Numerical Differentiation** **5**

---

#### **UNIT 13**

**Numerical Integration** **27**

---

#### **UNIT 14**

**Numerical Solution of Ordinary Differential Equations** **45**

---

#### **UNIT 15**

**Numerical Solution of Differential Equations Using  
Runge-Kutta Methods** **65**

---

---

## Course Design Committee

---

Prof. S.R.K. Iyengar (*Editor*)  
Dept. of Mathematics  
I.I.T. Delhi

Prof. J.P. Agarwal  
Dept. of Mathematics  
I.I.T. Kharagpur

Dr. U. Anantha Krishnaiah  
Dept. of Mathematics  
KREC, Surathkal

Prof. R.K. Jain  
Dept. of Mathematics  
I.I.T. Delhi

Prof. C. Prabhakara Rao  
Dept. of Mathematics  
REC, Warangal

### Faculty Members School of Sciences, IGNOU

Prof. R.K. Bose

Dr. V.D Madan

Dr. Poornima Mital

Dr. Manik Patwardhan

Dr. Parvin Sinclair

Dr. Sujatha Varma

---

## Block Preparation Team

---

Prof. S.R.K. Iyengar (*Editor*)  
Dept. of Mathematics  
I.I.T. Delhi

Dr. U. Anantha Krishnaiah  
Dept. of Mathematics  
KREC, Surathkal

Dr. Poornima Mital  
School of Sciences  
IGNOU

**Course Coordinator : Dr. Poornima Mital**

---

## Production

---

Mr. Balakrishna Selvaraj  
Registrar (PPD)  
IGNOU

---

## Acknowledgements

Prof. R.K. Bose for comments on the manuscript. Kiran for typing the manuscript.

September, 1993

© Indira Gandhi National Open University, 1993

ISBN-81-7263-484-6

All rights reserved. No part of this work may be reproduced in any form, by mimeograph or any other means, without permission in writing from the Indira Gandhi National Open University.

Further information on the Indira Gandhi National Open University courses may be obtained from the University's office at Maidan Garhi, New Delhi-110 068.

Reproduced and reprinted with the permission of Indira Gandhi National Open University by Dr.A.K.Singh, Registrar, U.P.R.T.Open University, Allahabad (May, 2013)  
Reprinted by : Nitin Printers, 1 Old Katra, Manmohan Park, Allahabad.

---

## **BLOCK 4 NUMERICAL DIFFERENTIATION INTEGRATION AND SOLUTION OF DIFFERENTIAL EQUATIONS**

---

In Block 3, we developed some interpolation techniques for approximating a given function by a polynomial. In this block, we shall use these approximating polynomials for obtaining numerical methods to perform important mathematical operations, viz., differentiation and integration. In calculus you must have spent lot of time learning techniques to do these operations. However, if the function is not known explicitly, but only tabulated values are known or if the function is too complicated or the function is such that we cannot perform these operations by using calculus methods. In such cases, numerical methods permit us to perform these operations. These techniques are also useful in solving differential equation. This block consists of four units.

In Unit 12, which is the first unit of this block, we shall discuss, a few numerical differentiation methods, namely, the method based on undetermined coefficients, methods based on finite difference operators and methods based on interpolation.

In Unit 13, we shall derive, a few numerical integration methods, namely, the methods using Lagrange interpolation and methods using Newton's forward interpolation. We shall also discuss composite rules of integration to obtain more accurate results and improve the order of the results using Romberg integration in the unit.

In Unit 14, we begin with a recall of few basic concepts from the theory of differential equations. We shall then derive numerical methods for solving differential equations. We shall introduce here two such methods namely, Euler's method and Taylor series method to obtain numerical solution of ordinary differential equations (ODEs). We shall also introduce Richardson's extrapolation method to obtain higher order solutions to ODEs using lower order methods.

In Unit 15, which is the last unit of this block and also of this course, we shall discuss Runge-Kutta methods of second, third and fourth order for obtaining the solution of ODEs. We shall discuss Richardson's extrapolation method to extrapolate the solutions obtained by the Runge-Kutta methods also.



---

# UNIT 12 NUMERICAL DIFFERENTIATION

---

## Structure

- 12.1 Introduction
  - Objectives
- 12.2 Methods Based on Undetermined Coefficients
- 12.3 Methods Based on Finite Difference Operators
- 12.4 Methods Based on Interpolation
- 12.5 Richardson's Extrapolation
- 12.6 Optimum Choice of Step Length
- 12.7 Summary
- 12.8 Solutions/Answers

---

## 12.1 INTRODUCTION

---

Differentiation of a function  $f(x)$  is a fundamental and important concept in calculus. When the function is given explicitly its derivatives  $f'(x)$ ,  $f''(x)$ ,... etc. can be easily found using the methods of calculus. For example, if  $f(x) = x^2$ , we know that  $f'(x) = 2x$ ,  $f''(x) = 2$  and all the higher order derivatives are zero. However, if the function is not known explicitly but, we are given a table of values of  $f'(x)$  corresponding to a set of values of  $x$ , then we cannot find the derivatives by using calculus methods. For instance if  $f(x_k)$ , represents distance travelled by a car in time  $x_k$ ,  $k = 0, 1, 2, \dots$  seconds, and we require the velocity and acceleration of the car at any time  $x_k$ , then the derivatives  $f'(x)$  and  $f''(x)$  representing velocity and acceleration respectively, cannot be found analytically. Hence, the need arises to develop methods of differentiation to obtain the derivative of a given function  $f(x)$ , using the data given in the form of a table which might have been formed as a result of scientific experiments.

Numerical methods have the advantage that they are easily adaptable on calculators and computers. These methods make use of the interpolating polynomials, which we discussed in Block-3. We shall now discuss, in this unit, a few numerical differentiation methods, namely, the method based on undetermined coefficients, methods based on finite difference operators and methods based on interpolation.

### Objectives

After studying this unit you should be able to

- explain the importance of the numerical methods over the calculus methods;
- use the method of undetermined coefficients and methods based on finite difference operators to derive differentiation formulas and obtain the derivative of a function at step points;
- use the methods derived from the interpolation formulas to obtain the derivative of a function at off step points;
- use Richardson's extrapolation method for obtaining higher order solutions;
- obtain the optimal step length for the given formula.

---

## 12.2 METHODS BASED ON UNDETERMINED COEFFICIENTS

---

In Unit 1, we introduced you to the concepts of round-off and truncation errors. In the derivation of the methods of numerical differentiation, we shall be referring to these errors quite often. Let us first quickly recall these concepts here before going further.

**Definition :** The round-off error is the quantity  $R$  which must be added to the finite representation of a computed number in order to make it the true representation of that number. Thus

$$y(\text{machine representation}) + R = y(\text{true representation}).$$

**Definition :** The truncation error denoted by  $TE$  is the quantity which must be added to the finite representation of the computed quantity in order that the result be exactly equal to the quantity we are seeking to generate. Thus

$$y(\text{true representation}) + TE = y(\text{exact})$$

The total error  $E_n$  is then given by

$$\begin{aligned} |E_n| &= |y(\text{machine representation}) - y(\text{exact})| \\ &= |y(\text{machine representation}) - y(\text{true representation})| \\ &\quad + |y(\text{true representation}) - y(\text{exact})| \\ &\leq |y(\text{machine representation}) - y(\text{true representation})| \\ &\quad + |y(\text{true representation}) - y(\text{exact})| \\ &\leq |R| + |TE| \end{aligned}$$

**Defintion :** Let  $f(h)$  be the exact analytical value of a given problem obtained by using an analytical formula and  $f_h$  be the approximate value obtained by using a numerical method. If the error  $f(h) - f_h = C h^p$ , where  $C$  is a constant, then  $p$  is known as the order of the numerical method.

Let us consider a function  $f(x)$ , whose values are given at a set of tabular points. For developing numerical differentiation formulas for the derivatives  $f'(x)$ ,  $f''(x)$ , ... at a point  $x = x_k$ , we express the derivative  $f^q(x)$ ,  $q \geq 1$ , as a linear combination of the values of  $f(x)$  at an arbitrarily chosen set of tabular points. Here, we assume that the tabular points are equally spaced with the steplength  $h$  i.e. various step (nodal) points are  $x_m = x_0 \pm mh$ ,  $m = 0, 1, \dots$  etc. Then we write

$$h^q f^q(x_k) = \sum_{m=-s}^n \gamma_m f_{k+m} \quad (1)$$

where  $\gamma_i$ ,  $i = -s, -s+1, \dots, n$  are the unknowns to be determined and  $f_{k+m}$  denotes  $f(x_k + mh)$ . For example, when  $s = n = 1$  and  $q = 1$ , Eqn. (1) reduces to

$$h f'(x_k) = \gamma_{-1} f_{k-1} + \gamma_0 f_k + \gamma_1 f_{k+1}.$$

Similarly, when  $s = 1$ ,  $n = 2$  and  $q = 2$ , we have

$$h^2 f''(x_k) = \gamma_{-1} f_{k-1} + \gamma_0 f_k + \gamma_1 f_{k+1} + \gamma_2 f_{k+2}$$

Now suppose we wish to determine a numerical differentiation formula for  $f^q(x_k)$  of order  $p$  using the method of undetermined coefficients. In other words, we want our formula to give the exact derivative values when  $f(x)$  is a polynomial of degree  $\leq p$ , that is, for  $f(x) = 1, x, x^2, x^3, \dots, x^p$ . We then get  $p+1$  equations for the determination of the unknowns  $\gamma_i$ ,  $i = -s, -s+1, \dots, n$ . You know that if a method is of order  $p$ , then its TE is of the form  $C h^{p+1} f^{(p+1)}(\alpha)$ , for some constant  $C$ . This implies that if  $f(x) = x^m$ ,  $m = 0, 1, 2, \dots, p$  then the method gives exact results, since

$$\frac{d^{p+1}}{dx^{p+1}}(x^m) = 0, \text{ for } m = 0, 1, \dots, p.$$

Let us now illustrate this idea to find the numerical differentiation formula of  $O(h^4)$  for  $f''(x_k)$ .

### Derivation of $O(h^4)$ formula for $f''(x)$

Without loss of generality let us take  $x_k = 0$ . We shall take the points symmetrically, that is,  $x_m = mh$ ;  $m = 0, \pm 1, \pm 2$ .

Let  $f_{-2}, f_{-1}, f_0, f_1, f_2$  denote the values of  $f(x)$  at  $x = -2h, -h, 0, h, 2h$  respectively.

In this case the formula given by Eqn. (1) can be written as

$$h^2 f''(0) = \gamma_{-2} f_{-2} + \gamma_{-1} f_{-1} + \gamma_0 f_0 + \gamma_1 f_1 + \gamma_2 f_2 \quad (2)$$

Let us now make the formula exact for  $f(x) = 1, x, x^2, x^3, x^4$ . Then, we have

$$f(x) = 1, f''(0) = 0; f_{-2} = f_{-1} = f_0 = f_1 = f_2 = 1$$

$$f(x) = x, f''(0) = 0, f_{-2} = -2h; f_{-1} = -h; f_0 = 0; f_1 = h; f_2 = 2h;$$

$$f(x) = x^2, f''(0) = 2, f_{-2} = 4h^2 = f_2; f_{-1} = h^2 = f_1; f_0 = 0; \quad (3)$$

$$f(x) = x^3, f''(0) = 0, f_{-2} = -8h^3; f_{-1} = -h^3; f_0 = 0; f_1 = h^3, f_2 = 8h^3$$

$$f(x) = x^4, f''(0) = 0; f_{-2} = 16h^4 = f_2; f_{-1} = h^4 = f_1; f_0 = 0$$

Substituting these values in Eqn. (2), we obtain the following set of equations for determining  $\gamma_m, m = 0, 1, 2$ :

$$\begin{aligned} \gamma_{-2} + \gamma_{-1} + \gamma_0 + \gamma_1 + \gamma_2 &= 0 \\ -2\gamma_{-2} - \gamma_{-1} + \gamma_1 + 2\gamma_2 &= 0 \\ 4\gamma_{-2} + \gamma_{-1} + \gamma_1 + 4\gamma_2 &= 2 \\ -8\gamma_{-2} - \gamma_{-1} + \gamma_1 + 8\gamma_2 &= 0 \\ 16\gamma_{-2} + \gamma_{-1} + \gamma_1 + 16\gamma_2 &= 0 \end{aligned} \quad (4)$$

Thus we have a system of five equations for five unknowns. The solution of this system of Eqs. (4) is

$$\gamma_{-2} = \gamma_2 = -1/12; \gamma_{-1} = \gamma_1 = 16/12; \gamma_0 = 30/12;$$

Hence, the numerical differentiation formula of  $O(h^4)$  for  $f''(0)$  as given by Eqn. (2) is

$$f''(0) \approx f''_0 = \frac{1}{12h^2} [-f_{-2} + 16f_{-1} - 30f_0 + 16f_1 - f_2] \quad (5)$$

Now, we know that the TE of the formula (5) is given by the first non-zero term in the Taylor expression of

$$f''(x_0) - \frac{1}{12h^2} [-f(x_0 - 2h) + 16f(x_0 - h) - 30f(x_0) + 16f(x_0 + h) - f(x_0 + 2h)] \quad (6)$$

The Taylor series expansions give

$$f(x_0 - 2h) = f(x_0) - 2hf'(x_0) + 2h^2 f''(x_0) - \frac{4h^3}{3} f'''(x_0) + \frac{2h^4}{3} f^{IV}(x_0)$$

$$- \frac{4h^5}{15} f^V(x_0) + \frac{4h^6}{45} f^{VI}(x_0) - \dots$$

$$f(x_0 - h) = f(x_0) - hf'(x_0) + \frac{h^2}{2} f''(x_0) - \frac{h^3}{6} f'''(x_0) + \frac{h^4}{24} f^{IV}(x_0) - \frac{h^5}{120} f^V(x_0)$$

$$+ \frac{h^6}{720} f^{VI}(x_0) + \dots$$

$$f(x_0 + h) = f(x_0) + hf'(x_0) + \frac{h^2}{2} f''(x_0) + \frac{h^3}{6} f'''(x_0) + \frac{h^4}{24} f^{IV}(x_0) + \frac{h^5}{120} f^V(x_0)$$

$$+ \frac{h^6}{720} f^{VI}(x_0) + \dots$$

$$f(x_0 + 2h) = f(x_0) + 2hf'(x_0) + \frac{2h^2}{2} f''(x_0) + \frac{4h^3}{3} f'''(x_0) + \frac{2h^4}{3} f^{IV}(x_0) + \frac{4h^5}{15} f^V(x_0)$$

$$+ \frac{4h^6}{45} f^{VI}(x_0) + \dots$$

Substituting these expansions in Eqn. (6) and simplifying, we get the first non-zero term or the TE of the formula (5) as

$$TE = f''(x_0) - \frac{1}{12h^2} [-f(x_0 - 2h) + 16f(x_0 - h) - 30f(x_0) + 16f(x_0 + h) - f(x_0 + 2h)]$$

$$= - \frac{h^6}{90} f^{VI}(\alpha), 0 < \alpha < 1.$$

You may now try the following exercise.

E1) A differentiation rule of the form

$$f'_0 = \alpha_0 f_0 + \alpha_1 f_1 + \alpha_2 f_2$$

is given. Find  $\alpha_0, \alpha_1$  and  $\alpha_2$  so that the rule is exact for polynomials of degree 2.

You must have observed that in the numerical differentiation formula discussed above, we have to solve a linear system of equations. If the number of nodal points involved is large or if we have to determine a method of high order, then we have to solve a large system of linear equations, which becomes tedious. To avoid this, we can use finite difference operators to obtain the differentiation formulas, which we shall illustrate in the next section.

## 12.3 METHODS BASED ON FINITE DIFFERENCE OPERATORS

Recall that in Unit 10 of Block 3, we introduced the finite difference operators  $E$ ,  $\nabla$ ,  $\Delta$ ,  $\mu$  and  $\delta$ . There we also gave the relations among various operators.

In order to construct the numerical differentiation formulas using these operators, we shall first derive relations between the differential operator  $D$  where  $Df(x) = f'(x)$ , and the various difference operators.

By Taylor series, we have

$$\begin{aligned} f(x+h) &= f(x) + hf'(x) + \frac{h^2}{2}f''(x) + \dots \\ &= [1 + hD + \frac{h^2}{2}D^2 + \dots] f(x) \\ &= e^{hD} f(x) \end{aligned} \tag{7}$$

Since,  $Ef(x) = f(x+h)$

we obtain from Eqn. (7), the identity

$$E = e^{hD} \tag{8}$$

which gives the relations

$$hD = \log E = \log(1 + \Delta) = \Delta - \frac{\Delta^2}{2} + \frac{\Delta^3}{3} - \frac{\Delta^4}{4} + \dots \tag{9}$$

$$hD = \log E = -\log(1 - \nabla) = \nabla + \frac{\nabla^2}{2} + \frac{\nabla^3}{3} - \frac{\nabla^4}{4} + \dots \tag{10}$$

We can relate  $D$  with  $\delta$  as follows :

We know that  $\delta = E^{h/2} - E^{-h/2}$ . Using identity (8), we can write

$$\delta f(x) = [e^{hD/2} - e^{-hD/2}] f(x)$$

$$\text{Hence, } \delta = 2\sinh(hD/2)$$

$$\text{or } hD = 2\sinh^{-1}(\delta/2) \tag{11}$$

$$\text{Similarly } \mu = \cosh(hD/2) \tag{12}$$

$$\text{We also have } \mu\delta = \sinh(hD) \text{ or } hD = \sinh^{-1}(\mu\delta) \tag{13}$$

$$\text{and } \mu^2 = \cosh^2(hD/2) = 1 + \sinh^2(hD/2) = 1 + \frac{\delta^2}{4} \tag{14}$$

Using the Maclaurin's expansion of  $\sinh^{-1}x$ , in relation (11), we can express  $hD$  as an infinite series in  $\delta/2$ .

Thus, we have

$$\begin{aligned} hD &= 2\sinh^{-1}(\delta/2) \\ &= \delta - \frac{1^2\delta^3}{2^2 \cdot 3!} + \frac{1^2 \cdot 3^2 \delta^5}{2^4 \cdot 5!} + \frac{1^2 \cdot 3^2 \cdot 5^2 \delta^7}{2^6 \cdot 7!} + \dots \end{aligned} \tag{15}$$

**Notice** that this formula involves off-step points when operated on  $f(x)$ . The formula involving only the step points can be obtained by using the relation (13), i.e.,

$$\begin{aligned} hD &= \sinh^{-1}(\mu\delta) \\ &= \mu\delta - \frac{1^2\mu^3\delta^3}{3!} + \frac{1^2 \cdot 3^2 \cdot \mu^5 \delta^5}{5!} - \frac{1^2 \cdot 3^2 \cdot 5^2 \cdot \mu^7 \delta^7}{7!} + \dots \end{aligned} \tag{16}$$

Using the relation (14) in Eqn. (16), we obtain

$$hD = \mu \left[ \delta - \frac{\delta^3}{6} + \frac{\delta^5}{30} - \frac{\delta^7}{140} + \dots \right] \tag{17}$$

Thus, Eqns. (9), (10), (15) and (16) give us the relations between  $hD$  and various difference operators. Let us see how we can use these relations to derive numerical

$$E = \frac{\Delta + 1}{(1 - \nabla)^{-1}}$$

$$\mu = \frac{1}{2}(E^{h/2} + E^{-h/2})$$



differentiation formulas for  $f'_k, f''_k$  etc.

We first derive formulas for  $f'_k$ . From Eqn. (9), we get

$$hDf(x_k) = hf'_k = \left( \Delta - \frac{\Delta^2}{2} + \frac{\Delta^3}{3} - \frac{\Delta^4}{4} + \dots \right) f_k$$

Thus forward difference formulas of  $O(h), O(h^2), O(h^3)$  and  $O(h^4)$  can be obtained by retaining respectively 1, 2, 3, and 4 terms of the relation (9) as follows :

$0(h)$ method : $hf'_k = f_{k+1} - f_k$	(18)
$0(h^2)$ method : $hf'_k = \frac{1}{2}(-f_{k+2} + 4f_{k+1} - 3f_k)$	(19)
$0(h^3)$ method : $hf'_k = \frac{1}{6}(2f_{k+3} - 9f_{k+2} + 18f_{k+1} - 11f_k)$	(20)
$0(h^4)$ method : $hf'_k = \frac{1}{12}(-3f_{k+4} + 16f_{k+3} - 36f_{k+2} + 48f_{k+1} - 25f_k)$	(21)

TE of the formula (18) is

$$TE = f'(x_k) - \frac{1}{h}[f(x_{k+1}) - f(x_k)] = -\frac{h}{2}f''(\xi) \quad (22)$$

and that of formula (19) is

$$TE = f'(x_k) - \frac{1}{2h}[-f(x_{k+2}) + 4f(x_{k+1}) - 3f(x_k)]^2 = \frac{h^2}{3}f'''(\xi) \quad (23)$$

Similarly the TE of formulas (20) and (21) can be calculated. Backward difference formulas of  $O(h), O(h^2), O(h^3)$  and  $O(h^4)$  for  $f'_k$  can be obtained in the same way by using the equality (10) and retaining 1,2,3 or 4 terms. We are leaving it as an exercise for you to derive these formulas.

**E2) Derive backward difference formulas for  $f'_k$  of  $O(h), O(h^2), O(h^3)$  and  $O(h^4)$ .**

Central difference formulas for  $f'_k$  can be obtained by using the relation (17), i.e.,

$$hf'_k = \mu \left( \delta - \frac{\delta^3}{6} + \dots \right) f_k$$

Note that relation (17) will give us methods of  $O(h^2)$  and  $O(h^4)$ , on retaining 1 and 2 terms,

$0(h^2)$ method : $hf'_k = \frac{1}{2}(f_{k+1} - f_{k-1})$	(24)
$0(h^4)$ method : $hf'_k = \frac{1}{12}(-f_{k-2} + 8f_{k-1} + 8f_{k+1} + f_{k+2})$	(25)

We now illustrate these methods through an example.

**Example 1 :** Given the following table of values of  $f(x) = e^x$ , find  $f'(0.2)$  using formulas (18), (19), (24) and (25).

x :	0.0	0.1	0.2	0.3	0.4
f(x) :	1.000000	1.105171	1.221403	1.349859	1.491825

**Solution :** Here  $h = 0.1$  and exact value of  $e^x$  at  $x = 0.2$  is 1.221402758.

Using(18),  $f'(0.2) = \frac{f(0.3) - f(0.2)}{0.1}$

or  $f'(0.2) = \frac{1.349859 - 1.221403}{0.1} = 1.28456$

$TE = -\frac{h}{2}f''(0.2) = -\frac{1}{2}e^{0.2} = -0.061070$

$$\text{Actual error} = 1.221402758 - 1.28456 = -0.063157$$

$$\text{Using (19), } f'(0.2) = \frac{1}{0.2} [-f(0.4) + 4f(0.3) - 3f(0.2)] = 1.21701$$

$$\text{TE} = \frac{h^2}{3} f'''(0.2) = \frac{0.01}{3} e^{0.2} = 0.004071;$$

$$\text{Actual error} = 0.004393$$

$$\text{Using (24), } f'(0.2) = \frac{1}{0.2} [f(0.3) - f(0.1)] = 1.22344$$

$$\text{TE} = -\frac{h^2}{6} f'''(0.2) = -\frac{0.01}{6} e^{0.2} = -0.0020357;$$

$$\text{Actual error} = -0.002037$$

$$\text{Using (25), } f'(0.2) = \frac{1}{1.2} [-f(0.0) + 8f(0.1) - 8f(0.3) + f(0.4)] = 1.221399167$$

$$\text{TE} = \frac{h^4 f^{(4)}(0.2)}{30} = \frac{0.0001}{30} e^{0.2} = 0.4071 \times 10^{-5};$$

$$\text{Actual error} = 0.3591 \times 10^{-5}$$

Numerical differentiation formulas for  $f''_k$  can be obtained by considering

$$h^2 D^2 = \Delta^2 - \Delta^3 + \frac{11}{12} \Delta^4 - \frac{2}{3} \Delta^5 + \dots \quad (26)$$

$$= \Delta^2 + \Delta^3 + \frac{11}{12} \Delta^4 + \frac{2}{3} \Delta^5 + \dots \quad (27)$$

$$= \delta^2 - \frac{\delta^4}{12} + \frac{\delta^6}{90} - \dots \quad (28)$$

We can write the forward difference methods of  $O(h)$ ,  $O(h^2)$ ,  $O(h^3)$  and  $O(h^4)$  for  $f''_k$  by using Eqn. (26) and retaining 1, 2, 3 and 4 terms as follows :

$$O(h) \text{ method : } h^2 f''_k = f_{k+2} - 2f_{k+1} + f_k \quad (29)$$

$$O(h^2) \text{ method : } h^2 f''_k = -f_{k+3} + 4f_{k+2} - 5f_{k+1} + 2f_k \quad (30)$$

$$O(h^3) \text{ method : } h^2 f''_k = \frac{1}{12} (11f_{k+4} - 56f_{k+3} + 114f_{k+2} - 104f_{k+1} + 35f_k) \quad (31)$$

$$O(h^4) \text{ method : } h^2 f''_k = \frac{1}{12} (-8f_{k+5} + 51f_{k+4} - 136f_{k+3} + 194f_{k+2} - 144f_{k+1} + 43f_k) \quad (32)$$

Backward difference formulas can be written in the same way by using Eqn. (27). Central difference formulas of  $O(h^2)$  and  $O(h^4)$  for  $f''_k$  are obtained by using Eqn. (28) and retaining 1 or 2 terms in the form :

$$O(h^2) \text{ method : } h^2 f''_k = (f_{k-1} - 2f_k + f_{k+1}) \quad (33)$$

$$O(h^4) \text{ method : } h^2 f''_k = \frac{1}{12} (-f_{k-2} + 16f_{k-1} - 30f_k + 16f_{k+1} - f_{k+2}) \quad (34)$$

Let us consider an example.

**Example 2 :** For the table of values of  $f(x) = e^x$ , given in Example 1, find  $f''(0.2)$  using the formulas (33) and (34).

$$\text{Solution : Using (33), } f''(0.2) = \frac{1}{0.01} [f(0.1) - 2f(0.2) + f(0.3)] = 1.2224$$

$$\text{TE} = \frac{-h^2 f^{(4)}(0.2)}{2} = \frac{-(0.01) e^{0.2}}{12} = -0.0010178$$

$$\text{Actual error} = -0.0009972$$

Using Eqn.(34),  $f''(0.2) = \frac{[-f(0.0) + 16f(0.1) - 30f(0.2) + 16f(0.3) - f(0.4)]}{0.12} = 1.221375$

$TE = \frac{h^4 f^{(4)}(0.2)}{90} = 0.13571 \times 10^{-5}$

Actual error =  $0.27758 \times 10^{-4}$

And now the following exercises for you.

- E3) From the following table of values find  $f'(6.0)$  using an  $O(h)$  formula and  $f''(6.3)$  using an  $O(h^2)$  formula.

x	: 6.0	6.1	6.2	6.3	6.4
f(x)	: 0.1750	-0.1998	-0.2223	-0.2422	-0.2596

- E4) Calculate the first and second derivatives of  $\ln x$  at  $x = 500$  from the following table. Use  $O(h^2)$  forward difference method. Compute TE and actual errors.

x	: 500	510	520	530
f(x)	: 6.2146	6.2344	6.2538	6.2729

In Secs. 12.2 and 12.3, we have derived numerical differentiation formulas to obtain the derivative values at nodal points or step points, when the function values are given in the form of a table. However, these methods cannot be used to find the derivative values at off-step points. In the next section we shall derive methods which can be used for finding the derivative values at the off-step points as well as at step-points.

## 12.4 METHODS BASED ON INTERPOLATION

In these methods, given the values of  $f(x)$  at a set of points  $x_0, x_1, \dots, x_n$ , the general approach for deriving numerical differentiation formulas is to obtain the unique interpolating polynomial  $P_n(x)$  fitting the data. We then differentiate this polynomial  $q$  times ( $q \leq n$ ), to get  $(x)$ . The value  $P_n^{(q)}(x_k)$  then gives us the approximate value of  $f^{(q)}(x_k)$  where  $x_k$  may be a step point or an off-step point. We would like to point out here that even when the original data are known to be accurate i.e.  $P_n(x_k) = f(x_k)$ ,  $k = 0, 1, 2, \dots, n$ , yet the derivative values may differ considerably at these points. The approximations may further deteriorate while finding the values at off-step points or as the order of the derivative increases. However, these disadvantages are present in every numerical differentiation formula, as in general, one does not know whether the function representing a table of values has a derivative at every point or not.

We shall first derive differentiation formulas for the derivatives using non-uniform nodal points. That is, when the difference between any two consecutive points is not uniform.

### Non-uniform nodal points

Let the data  $(x_k, f_k)$ ,  $k = 0, 1, \dots, n$  be given at  $n + 1$  points where the step length  $x_i - x_{i-1}$  may not be uniform.

In Unit 9 you have seen that the Lagrange interpolating polynomial fitting the data  $(x_k, f_k)$ ,  $k = 0, 1, \dots, n$  is given by

$$P_n(x) = \sum_{k=0}^n L_k(x) f_k \tag{35}$$

where  $L_k(x)$  are the fundamental Lagrange polynomials given by

$$L_k(x) = \frac{\pi(x)}{(x - x_k) \pi'(x_k)} \tag{36}$$

and  $\pi(x) = (x - x_0)(x - x_1) \dots (x - x_n)$  (37)

$$\pi'(x_k) = (x_k - x_0)(x_k - x_1) \dots (x_k - x_{k-1})(x_k - x_{k+1}) \dots (x_k - x_n) \quad (38)$$

The error of interpolation is given by

$$E_n(x) = f(x) - P_n(x) = \frac{\pi(x)}{(n+1)!} f^{(n+1)}(\alpha), \quad x_0 < \alpha < x_n.$$

Differentiating  $P_n(x)$  w.r.t.  $x$ , we obtain

$$P'_n(x) = \sum_{k=0}^n L'_k(x) f_k \quad (39)$$

and the error is given by

$$E'_n(x) = \frac{1}{(n+1)!} \{ \pi'(x) f^{(n+1)}(\alpha) + \pi(x) (f^{(n+1)}(\alpha))' \} \quad (40)$$

Since in Eqn. (40), the function  $\alpha(x)$  is not known in the second term on the right hand side, we cannot evaluate  $E'_n(x)$  directly. However, since at a nodal point  $x_k$ ,  $\pi(x_k) = 0$ , we obtain

$$E'_n(x_k) = \frac{\pi'(x_k)}{(n+1)!} f^{(n+1)}(\alpha) \quad (41)$$

If we want to obtain the differentiation formulas for any higher order, say  $q$ th ( $1 \leq q \leq n$ ) order derivative, then we differentiate  $P_n(x)$ ,  $q$  times and get

$$f^{(q)}(x) \approx p_n^{(q)}(x) = \sum_{k=0}^n L_k^{(q)}(x) f_k \quad (42)$$

Similarly, the error term is obtained by differentiating  $E_n(x)$ ,  $q$  times.

Let us consider the following examples.

**Example 3 :** Find  $f'(x)$  and the error of approximation using Lagrange interpolation for the data  $(x_k, f_k)$ ,  $k = 0, 1$ .

**Solution :** We know that  $P_1(x) = L_0(x) f_0 + L_1(x) f_1$

$$\text{where } L_0(x) = \frac{x - x_1}{x_0 - x_1} \text{ and } L_1(x) = \frac{x - x_0}{x_1 - x_0}$$

Now,

$$P'_1(x) = L'_0(x) f_0 + L'_1(x) f_1$$

$$\text{and } L'_0(x) = \frac{1}{x_0 - x_1}, \quad L'_1(x) = \frac{1}{x_1 - x_0}$$

$$\text{Hence, } f'(x) \approx P'_1(x) = \frac{f_0}{x_0 - x_1} + \frac{f_1}{x_1 - x_0} = \frac{(f_1 - f_0)}{(x_1 - x_0)} \quad (43)$$

$$E'(x_0) = \frac{(x_0 - x_1)}{2} f''(\alpha) \text{ and } E'(x_1) = \frac{(x_1 - x_0)}{2} f''(\alpha), \quad x_0 < \alpha < x_1.$$

**Example 4 :** Find  $f'(x)$  and  $f''(x)$  given  $f_0, f_1, f_2$  at  $x_0, x_1, x_2$  respectively, using the Lagrange interpolation.

**Solution :** By Lagrange's interpolation formula

$$f(x) \approx P_2(x) = L_0(x) f_0 + L_1(x) f_1 + L_2(x) f_2$$

where,

$$L_0(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)}; \quad L'_0(x) = \frac{2x - x_1 - x_2}{(x_0 - x_1)(x_0 - x_2)}$$

$$L_1(x) = \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)}; \quad L'_1(x) = \frac{2x - x_0 - x_2}{(x_1 - x_0)(x_1 - x_2)}$$

$$L_2(x) = \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)}; \quad L'_2(x) = \frac{2x - x_0 - x_1}{(x_2 - x_0)(x_2 - x_1)}$$

Hence,  $f'(x) = P_2'(x) = L_0'(x) f_0 + L_1'(x) f_1 + L_2'(x) f_2$

and  $P_2''(x) = L_0''(x) f_0 + L_1''(x) f_1 + L_2''(x) f_2$   

$$= \frac{2f_0}{(x_0 - x_1)(x_0 - x_2)} + \frac{2f_1}{(x_1 - x_0)(x_1 - x_2)} + \frac{2f_2}{(x_2 - x_0)(x_2 - x_1)}$$

**Example 5 :** Given the following values of  $f(x) = \ln x$ , find the approximate value of  $f'(2.0)$  and  $f''(2.0)$ . Also find the errors of approximations.

x :	2.0	2.2	2.6
f(x) :	0.69315	0.78846	0.95551

**Solution :** Using the Lagrange's interpolation formula, we have

$$f'(x_0) = P_2'(x_0) = \frac{2x_0 - x_1 - x_2}{(x_0 - x_1)(x_0 - x_2)} f_0 + \frac{x_0 - x_2}{(x_1 - x_0)(x_1 - x_2)} f_1 + \frac{x_0 - x_1}{(x_2 - x_0)(x_2 - x_1)} f_2$$

∴ we get

$$f'(2.0) = \frac{4 - 2.2 - 2.6}{(2 - 2.2)(2 - 2.6)} (0.69315) + \frac{2 - 2.6}{(2.2 - 2)(2.2 - 2.6)} (0.78846) + \frac{2 - 2.2}{(2.6 - 2)(2.6 - 2.2)} (0.95551) = 0.49619$$

The exact value of  $f'(2.0) = 0.5$

Error is given by

$$E_2'(x_0) = \frac{1}{6} (x_0 - x_1)(x_0 - x_2) f'''(2.0) = \frac{1}{6} (2.0 - 2.2)(2.0 - 2.6) (-0.25) = -0.005$$

Similarly,

$$f''(x_0) = 2 \left[ \frac{f_0}{(x_0 - x_1)(x_0 - x_2)} + \frac{f_1}{(x_1 - x_0)(x_1 - x_2)} + \frac{f_2}{(x_2 - x_0)(x_2 - x_1)} \right]$$

$$\therefore f''(2.0) = 2 \left[ \frac{0.69315}{(2 - 2.2)(2 - 2.6)} + \frac{0.78846}{(2.2 - 2)(2.2 - 2.6)} + \frac{0.95551}{(2.6 - 2)(2.6 - 2.2)} \right] = -0.19642$$

The exact value of  $f''(2.0) = -0.25$ .

Error is given by

$$E_2''(x_0) = \frac{1}{3} (2x_0 - x_1 - x_2) f'''(2.0) + \frac{1}{24} (x_0 - x_1)(x_1 - x_2) [f^{IV}(2.0) + f^{IV}(2.0)] = -0.06917$$

You may now try the following exercise.

**E5)** Use Lagrange's interpolation to find  $f'(x)$  and  $f''(x)$  at  $x = 2.5, 5.0$  from the following table

x :	1	2	3	4
f(x)	1	16	81	256

Now let us consider the case of uniform nodal points.

### Uniform nodal points

When the difference between any two consecutive points is same, i.e., when we are given values of  $f(x)$  at equally spaced points, we can use Newton's forward or

backward interpolation formulas to find the unique interpolating polynomial  $P_n(x)$ . We can then differentiate this polynomial to find the derivative values either at the nodal points or at off-step points.

Let the data  $(x_k, f_k)$ ,  $k = 0, 1, \dots, n$  be given at  $(n + 1)$  points where the step points  $x_k$ ,  $k = 0, 1, \dots, n$  are equispaced with step length  $h$ . That is, we have

$$x_k = x_0 + kh, \quad k = 1, 2, \dots, n.$$

You know that by Newton's forward interpolation formula

$$f(x) = P_n(x) = f_0 + \frac{(x-x_0)}{h} \Delta f_0 + \frac{(x-x_0)(x-x_1)}{2!h^2} \Delta^2 f_0 + \dots + \frac{(x-x_0)(x-x_1)\dots(x-x_{n-1}) \Delta^n f_0}{n!h^n} \quad (44)$$

with error

$$E_n(x) = \frac{(x-x_0)(x-x_1)\dots(x-x_n)}{(n+1)!h^{n+1}} \Delta^{n+1} f(\alpha) \quad x_0 < \alpha < x_n. \quad (45)$$

If we put  $\frac{x-x_0}{h} = s$  or  $x = x_0 + sh$ , then Eqns. (44) and (45) reduce respectively to

$$f(x) = P_n(x) = f_0 + \frac{s \Delta f_0}{1!} + \frac{s(s-1) \Delta^2 f_0}{2!} + \frac{s(s-1)(s-2) \Delta^3 f_0}{3!} + \dots + \frac{s(s-1)\dots(s-n+1) \Delta^n f_0}{n!} \quad (46)$$

and

$$E_n(x) = \frac{s(s-1)\dots(s-n)}{(n+1)!} h^{(n+1)} f^{(n+1)}(\alpha) \quad (47)$$

Differentiation of  $P_n(x)$  w.r.t.  $x$  gives us

$$P'_n(x) = \frac{1}{h} \left[ \Delta f_0 + \frac{(2s-1)}{2} \Delta^2 f_0 + \frac{(3s^2-6s+2)}{6} \Delta^3 f_0 + \dots \right] \quad (48)$$

At  $x = x_0$ , we have  $s = 0$  and hence

$$f'(x_0) = \frac{1}{h} \left[ \Delta f_0 - \frac{\Delta^2 f_0}{2} + \frac{\Delta^3 f_0}{3} + \frac{\Delta^4 f_0}{4} + \dots \right]$$

which is same as formula (9) obtained in Sec. 12.3 by difference operator method. We can obtain the derivative at any step or off-step point by finding the value of  $s$  and substituting the same in Eqn. (48). The formula corresponding to Eqn. (47) in backward differences is

$$P'_n(x) = \frac{1}{h} \left[ \nabla f_n + \frac{(2s+1)}{2} \nabla^2 f_n + \frac{(2s^2+6s+2)}{6} \nabla^3 f_n + \dots \right] \quad (49)$$

where  $x = x_n + sh$ .

Formulas for higher order derivatives can be obtained by differentiating  $P'_n(x)$  further and the corresponding error can be obtained by differentiating  $E'_n(x)$ .

Let us illustrate the method through the following examples :

**Example 6 :** Find the first and second derivatives of  $f(x)$  at  $x = 1.1$  from the following tabulated values.

$x$ :	1.0	1.2	1.4	1.6	1.8	2.0
$f(x)$ :	0.0000	0.1280	0.5440	1.2960	2.4320	4.0000

**Solution :** Since we have to find the derivative at  $x = 1.1$ , we shall use the forward difference formula. The forward differences for the given data are given in Table 1.

Table 1

$x$	$f(x)$	$\Delta f(x)$	$\Delta^2 f(x)$	$\Delta^3 f(x)$	$\Delta^4 f(x)$	$\Delta^5 f(x)$
1.0	0.0					
		0.1280				
1.2	0.1280		0.2880			
		0.4160		0.0480		
1.4	0.5440		0.3360		0.0000	
		0.7520		0.0480		0.0000
1.6	1.2960		0.3840		0.0000	
		1.1360		0.0480		
1.8	2.4320		0.4320			
		1.5680				
2.0	4.0000					

Since,  $x = x_0 + s h$ ,  $x_0 = 1$ ,  $h = 0.2$  and  $x = 1.1$ , we have  $s = \frac{1.1-1}{0.2} = 0.5$

Substituting the value of  $s$  in formula (48), we get

$$f'(1.1) = \frac{1}{h} \left[ \Delta f_0 - \frac{0.25}{6} \Delta^3 f_0 \right] \quad (50)$$

Substituting the values of  $\Delta f_0$  and  $\Delta^3 f_0$  in Eqn. (50) from Table 1, we get

$$f'(1.1) = 0.63$$

To obtain the second derivative, we differentiate formula (48) and obtain

$$f''(x) = P''(x) = \frac{1}{h} \left[ \Delta^2 f_0 + (s-1) \Delta^3 f_0 \right]$$

Thus  $f''(1.1) = 6.6$

**Note :** If you construct a forward difference interpolating polynomial  $P(x)$ , fitting the data given in Table 1, you will find that  $f(x) = P(x) = x^3 - 3x + 2$ . Also,  $f'(1.1) = 6.3$ ,  $f''(1.1) = 6.6$ . The values obtained from this equation or directly as done above have to be same as the interpolating polynomial is unique.

**Example 7 :** Find  $f'(x)$  at  $x = 0.4$  from the following table of values.

$x$ :	0.1	0.2	0.3	0.4
$f(x)$ :	1.10517	1.22140	1.34986	1.49182

**Solution :** Since we are required to find the derivative at the end point, we will use the backward difference formula. The backward difference table for the given data is given by

Table 2

$x$ :	$f(x)$	$\Delta f(x)$	$\Delta^2 f(x)$	$\Delta^3 f(x)$
0.1	1.10517			
		0.11623		
0.2	1.22140		0.01223	
		0.12846		0.00127
0.3	1.34986		0.1350	
		0.14196		
0.4	1.49182			

Since  $x_n = 0.4$ ,  $h = 0.1$ ,  $x = 0.4$ , we get  $s = 0$

Substituting the value of  $s$  in formula (49), we get

$$\begin{aligned} f'(0.4) &= \frac{1}{h} \left[ \Delta f_3 + \frac{1}{2} \Delta^2 f_3 + \Delta^3 f_3 \right] \\ &= \frac{1}{0.1} \left[ 0.14196 + \frac{0.0135}{2} + \frac{0.00127}{3} \right] \\ &= 1.14913 \end{aligned}$$

How about trying a few exercises now ?

- E6) The position  $f(x)$  of a particle moving in a line at various times  $x_k$  is given in the following table. Estimate the velocity and acceleration of the particle at  $x = 1.5$  and  $3.5$

$x$ :	0	1	2	3	4
$f(x)$ :	-25	-9	0	7	15

- E7) Construct a difference table for the following data

$x$ :	1.3	1.5	1.7	1.9	2.1	2.3	2.5
$f(x)$ :	3.669	4.482	5.474	6.686	8.166	9.974	12.182

Taking  $h = 0.2$ , compute  $f'(1.5)$  and the error, if  $f(x) = e^x$ .

You must have by now observed that to obtain numerical differentiation methods of higher order, we require a large number of tabular points and thus a large number of function evaluations at these tabular points. Consequently, there is a possibility that the round-off errors may increase so much that the numerical results may become useless. However, it is possible to obtain higher order solutions by combining the computed values obtained by using the same method with two different step sizes. This technique is called extrapolation method or Richardson's extrapolation. We shall now discuss this method in the next section.

## 12.5 RICHARDSON'S EXTRAPOLATION

The underlying idea in this method is as follows :

Let  $f^{(q)}(h)$  denote the approximate value of  $f^{(q)}(x_k)$ , obtained by using a formula of order  $p$ , with steplength  $h$  and  $f^{(q)}(rh)$  denote the value of  $f^{(q)}(x_k)$  obtained by using the same method of order  $p$ , with steplength  $rh$ . Then,

$$f^{(q)}(h) = f^{(q)}(x_k) + Ch^p + O(h^{p+1}) \quad (51)$$

$$\text{and } f^{(q)}(rh) = f^{(q)}(x_k) + C(rh)^p + O(h^{p+1}) \quad (52)$$

Eliminating  $C$  between Eqns. (51) and (52), we get

$$f^{(q)}(x_k) = \frac{r^p f^{(q)}(h) - f^{(q)}(rh)}{r^p - 1} + O(h^{p+1}) \quad (53)$$

The new approximation to  $f^{(q)}(x_k)$  is therefore

$$f^{(q)}(x_k) = \frac{r^p f^{(q)}(h) - f^{(q)}(rh)}{r^p - 1} \quad (54)$$

The expression on the right hand side of Eqn. (54) for finding the value of the  $q$ th derivative by a certain method of order  $p$  has now become a method of order  $p + 1$ . This technique of combining two computed values obtained by using the same method



with two different step sizes, to obtain higher order solutions is called **Richardson's extrapolation method**.

We know that the truncation error of a numerical method of order  $p$  is given by

$$TE = C_1 h^p + O(h^{p+1}),$$

where  $C_1 \neq 0$ .

If, instead of denoting the higher order terms by  $O(h^{p+1})$ , we write down the actual terms, we have

$$TE = C_1 h^p + C_2 h^{p+1} + C_3 h^{p+2} + \dots$$

By repeated application of Richardson's extrapolation technique we can obtain solutions of higher orders, i.e.  $O(h^{p+1})$ ,  $O(h^{p+2})$ ,  $O(h^{p+3})$  etc. by eliminating  $C_1$ ,  $C_2$ ,  $C_3$  respectively. Let us see how this can be done.

Consider the central difference differentiation formula (24) of  $O(h^2)$  given by

$$f'_k = \frac{(f_{k+1} - f_{k-1}))}{2h}$$

Let  $g(x_k) = f'(x_k)$  be the exact value of the derivative, which is to be obtained and

$$g(h) = f'_k = \frac{(f_{k+1} - f_{k-1}))}{2h}$$

be the value given by the  $O(h^2)$  method. The truncation error of this method may be written as

$$g(h) = g(x_k) + C_1 h^2 + C_2 h^4 + C_3 h^6 + \dots \tag{55}$$

Let  $f'_k$  be evaluated with different step sizes  $\frac{h}{2^r}$ ,  $r = 0, 1, 2, \dots$

Then, we have

$$g\left(\frac{h}{2}\right) = g(x_k) + \frac{C_1 h^2}{4} + \frac{C_2 h^4}{16} + \frac{C_3 h^6}{64} + \dots \tag{56}$$

$$g\left(\frac{h}{4}\right) = g(x_k) + \frac{C_1 h^2}{16} + \frac{C_2 h^4}{256} + \frac{C_3 h^6}{4096} + \dots \tag{57}$$

Eliminating  $C_1$  from Eqns. (55) and (56), we get

$$g^{(1)}(h) = \frac{4g(h/2) - g(h)}{3} = g(x_k) - \frac{C_2 h^4}{4} - \frac{5C_3 h^6}{16} + \dots \tag{58}$$

Eliminating  $C_1$  from Eqns. (56) and (57), we obtain

$$g^{(1)}\left(\frac{h}{2}\right) = \frac{4g(h/4) - g(h/2)}{3} = g(x_k) - \frac{C_2 h^4}{64} - \frac{5C_3 h^6}{1024} + \dots \tag{59}$$

Notice that the methods  $g^{(1)}(h)$  and  $g^{(1)}(h/2)$  given by Eqns. (58) and (59) are  $O(h^4)$  approximations to  $g(x_k)$

Eliminating  $C_2$  from Eqns. (58) and (59), we get

$$g^{(2)}(h) = \frac{16g^{(1)}(h/2) - g^{(1)}(h)}{15} + \frac{C_3 h^6}{64} + \dots$$

which gives an  $O(h^6)$  approximation to  $g(x_k)$ . Generalising, we find that the successive higher order methods can be obtained from the formula.

$$g^{(m)}(h) = \frac{4^m g^{(m-1)}(h/2) - g^{(m-1)}(h)}{4^m - 1}$$

$m = 1, 2, \dots$ , with  $g^{(0)} = g(h)$  (60)

This procedure is known as the Richardson's repeated extrapolation to the limit.

These extrapolations can be stopped when

$$\left| g^{(k)}(h) - g^{(k-1)}\left(\frac{h}{2}\right) \right| < \epsilon \quad (61)$$

for a given error tolerance  $\epsilon$ .

Similarly, forward difference method of  $O(h^2)$  can be obtained by considering

$$g(h) = f'_k = \frac{(f_{k+1} - f_k)}{h}$$

and using Richardson's extrapolation technique in the form

$$g^{(1)}(h) = \frac{2g\left(\frac{h}{2}\right) - g(h)}{2 - 1} \quad (62)$$

This method is of  $O(h^2)$ .

You may note that in Richardson's extrapolation, each improvement made for forward (or backward) difference formula increases the order of solutions by one, whereas for central difference formula each improvement increases the order by two.

Let us now solve the following problems.

**Example 8 :** The following table of values of  $f(x) = x^4$ , is given :

x :	-1	1	2	3	4	5	6	7
f(x) :	1	1	16	81	256	625	1296	2401

Using the formula  $f'(x_1) = \left[ \frac{f(x_2) - f(x_0)}{2h} \right]$  and Richardson's extrapolation method find  $f'(3)$ .

**Solution :** Note that in this example  $x_1 = 3.0$ . The largest step  $h$  that can be taken is  $h = 4$ . Computations can also be done by using step lengths  $h_1 = \frac{h}{2} = 2$  and  $h_2 = h_1/2 = 1$ .

Using the formula

$$g(h) = f'(x_1) = \frac{f(x_2) - f(x_0)}{2h}$$

we get

$$g(h) = \frac{f(7) - f(-1)}{8} = 300 \quad O(h^2) \text{ method.}$$

$$g\left(\frac{h}{2}\right) = \frac{f(5) - f(1)}{4} = 156 \quad O(h^2) \text{ method.}$$

$$g\left(\frac{h}{4}\right) = \frac{f(4) - f(2)}{2} = 120 \quad O(h^2) \text{ method.}$$

Therefore, using the formula given by Eqn. (60), we have

$$g^{(1)}(h) = \frac{4g\left(\frac{h}{2}\right) - g(h)}{3} = \frac{624 - 300}{3} = 108 \quad O(h^4) \text{ method.}$$

$$g^{(1)}\left(\frac{h}{2}\right) = \frac{4g\left(\frac{h}{4}\right) - g\left(\frac{h}{2}\right)}{3} = \frac{480 - 156}{3} = 108 \quad O(h^4) \text{ method.}$$

$$g^{(2)}(h) = \frac{16g^{(1)}\left(\frac{h}{2}\right) - g^{(1)}(h)}{15} = 108$$

$O(h^6)$  method.

Writing in tabular form, we have

Step length	Second order method	Fourth order method	Sixth order method
h	$g(h) = 300$	$g^{(1)}(h) = 108$	
h/2	$g(h/2) = 156$		$g^{(2)}(h) = 108$
h/4	$g(h/4) = 120$	$g^{(1)}\left(\frac{h}{2}\right) = 108$	

Thus  $f'(3) = 108$ , must be the exact solution as we have

$$g^{(1)}(h) = g^{(1)}\left(\frac{h}{2}\right) = g^{(2)}(h).$$

**Example 9 :** Let  $f(x) = e^x$ . Using a central difference formula of  $O(h^2)$  find  $f''(1)$ . Improve this value using Richardson's extrapolation by taking  $h = 0.1$  and  $h = 0.05$ .

**Solution :** With  $h = 0.1$  and

$$f''_k = \frac{f_{k+1} - 2f_k + f_{k-1}}{h^2},$$

we get

$$f''(1) = \left( \frac{e^{1.1} - 2e^1 + e^{0.9}}{(0.1)^2} \right) = 2.720548$$

With  $h = 0.05$  we get  $f''(1) = \frac{e^{1.05} - 2e + e^{0.95}}{(0.05)^2} = 2.718848$

Both the solutions are  $O(h^2)$  approximations. Richardson's approximation using relation (54) with  $r = \frac{0.1}{0.05} = 2$  and  $p = 2$ , gives us

$$f''(1) = \left[ \frac{4(2.718848) - 2.720548}{3} \right] = 2.718281$$

The actual value is  $e = 2.718282$

Your may now try the following exercises :

- E8) Compute  $f''(0.6)$  from the following table using  $O(h^2)$  central difference formula. Improve it by Richardson's extrapolation method using step lengths  $h = 0.4, 0.2, 0.1$ .

x	: 0.2	0.4	0.5	0.6	0.7	0.8	1.0
f(x)	: 1.420072	1.881243	2.128147	2.386761	2.657971	2.942897	3.559753

- E9) Using central difference formula of  $O(h^2)$  find  $f''(0.3)$  from the given table and improve the accuracy using Richardson's extrapolation method using step lengths  $h = 0.1, 0.2$ .

x	: 0.1	0.2	0.3	0.4	0.5
f(x)	: 0.091	0.155	0.182	0.171	0.130

In the numerical differentiation methods, the truncation error is of the form  $Ch^p$  which tends to zero as  $h \rightarrow 0$ . However, the method which approximates  $f^{(q)}(x)$  contains  $h^q$  in

the denominator. As  $h$  is successively reduced to smaller values, the truncation error decreases but the round-off error in the method may increase as we are dividing by a small number. It may happen that after a certain critical value of  $h$ , the round-off error may become more dominant than the truncation error and the numerical results obtained may start worsening as  $h$  is further reduced. The problem of finding a steplength  $h$  small enough so that the truncation error is small, yet large enough so that round-off error does not dominate the actual error is referred to as the step size dilemma. Such a step length, if it can be determined is called the optimal steplength for that formula. We shall now discuss in the next section how to determine the optimal steplength.

## 12.6 OPTIMUM CHOICE OF STEPLENGTH

We begin by considering an example.

Consider the numerical differentiation formula

$$f'_k = \frac{1}{h} (f_{k+1} - f_k) \quad (63)$$

Let  $f(x) = e^x$  and we want to approximate  $f'(1)$  by taking  $h = \frac{2}{10^m}$ ,  $m = 1, 2, \dots, 7$ .

We have from the differentiation formula (63),

$$f'(1) = \frac{e^{1+h} - e}{h}$$

The exact solution is  $f'(1) = 2.718282$ . The actual error is  $e - f'(1)$  and the truncation error is  $-eh/2$ . With  $h = \frac{2}{10^m}$ ,  $m = 1, 2, \dots, 7$ , we have the results as given in Table 3.

Table 3

$h$	$f'(1)$	Actual error	Approximate Truncation error
$2 \times 10^{-1}$	3.009175	- 0.290893	- 0.271828
$2 \times 10^{-2}$	2.745650	- 0.027368	- 0.027183
$2 \times 10^{-3}$	2.721000	- 0.002718	- 0.002718
$2 \times 10^{-4}$	2.720000	- 0.017180	- $2.7 \times 10^{-4}$
$2 \times 10^{-5}$	2.700000	- 0.018280	- $2.7 \times 10^{-5}$
$2 \times 10^{-6}$	2.500000	- 0.218218	- $2.7 \times 10^{-6}$
$2 \times 10^{-7}$	0.000000	- 0.718282	- $2.7 \times 10^{-7}$

If you look at Table 3, you will observe that the improved accuracy of the formula, i.e.  $f'(1)$ , with decreasing  $h$  does not continue indefinitely. The truncation error agrees with the actual error till  $h = 2 \times 10^{-3} = 0.002$ . As  $h$  is further reduced, the truncation error ceases to approximate the actual error. This is because the actual error is dominated by round-off error rather than the truncation error. This effect gets worsened as  $h$  is reduced further. In such cases we determine the optimal steplength.

When  $f(x)$  is given in tabular form, these values may not be exact. These values contain round-off errors. In other words,  $f(x_k) = f_k + \epsilon_k$ , where  $f(x_k)$  is the exact value,  $f_k$  is the tabulated value and  $\epsilon_k$  is the round-off error. For the numerical differentiation formula (63), we have

$$f'(x_k) = \frac{(f_{k+1} - f_k)}{h} - \frac{h}{2} f''(\alpha), \quad x_k < \alpha < x_{k+1}$$

If the round-off errors in  $f_k$  and  $f_{k+1}$  are  $\epsilon_k$  and  $\epsilon_{k+1}$  then we have

$$f'(x_k) = \frac{1}{h} [(f_{k+1} + \epsilon_{k+1}) - (f_k + \epsilon_k)] - \frac{h}{2} f''(\alpha)$$

$$= \left( \frac{f_{k+1} - f_k}{h} \right) + \left( \frac{\epsilon_{k+1} - \epsilon_k}{h} \right) - \frac{h}{2} f''(\alpha)$$

$$= \left( \frac{f_{k+1} - f_k}{h} \right) + R + TE$$

where  $R$  is the round-off error and  $TE$  is the truncation error.

If we take  $\epsilon = \max(|\epsilon_k|, |\epsilon_{k+1}|)$  and  $M_2 = \max |f''(x)|$ , we find that

$$|R| \leq \frac{2\epsilon}{h} \text{ and } |TE| \leq \frac{h}{2} M_2$$

We define the optimum value of  $h$  as the one which satisfies either of the following conditions :

$$(i) \quad |R| = |TE| \quad (ii) \quad |R| + |TE| = \text{minimum} \quad (64)$$

By the first condition in (64), we have

$$\frac{2\epsilon}{h} = \frac{h}{2} M_2 \text{ or } h^2 = \frac{4\epsilon}{M_2} \text{ or } h = 2 \sqrt{\epsilon/M_2}$$

The value of the error is

$$|R| = |TE| = \sqrt{\epsilon M_2}$$

If we use the second condition  $|R| = |TE| = \min$ , we have

$$\frac{2\epsilon}{h} + \frac{h M_2}{2} = \min. \quad (65)$$

To find the minimum in Eqn. (65), we differentiate the left hand side of Eqn. (65) with respect to

$$-\frac{2\epsilon}{h^2} + \frac{M_2}{2} = 0 \text{ or } h^2 = \frac{4\epsilon}{M_2} \text{ or } h = 2 \sqrt{\epsilon/M_2}$$

$\therefore$  The minimum total error =  $|R| + |TE| = \sqrt{\epsilon M_2}$ .

Let us now consider an example.

**Example 10 :** For the method

$$f'_k = \frac{(-3f_k + 4f_{k+1} - f_{k+2})}{2h} + \frac{h^2}{3} f''(\alpha), \quad x_k < \alpha < x_{k+1}$$

determine the optimal value of  $h$  using the criteria

$$|R| = |TE| \text{ and (ii) } |R| + |TE| = \min.$$

Using this method and the first criterion, find the value of  $h$  and determine the value of  $f'(2.0)$ , from the following tabulated values of  $f(x) = \ln x$ . It is given that the maximum round-off error in the function evaluation is  $5 \times 10^{-6}$

$x$ :	2.0	2.01	2.02	2.06	2.12
$f(x)$ :	0.69315	0.69813	0.70310	0.72271	0.75142

**Solution :** If  $\epsilon_0, \epsilon_1$  and  $\epsilon_2$  are the round-off errors in the given function evaluations of  $f_0, f_1, f_2$  respectively, then we have

$$f'_0 = \frac{(-3f_0 + 4f_1 - f_2)}{2h} + \frac{(-3\epsilon_0 + 4\epsilon_1 - \epsilon_2)}{2h} + \frac{h^2}{3} f'''(\alpha)$$

$$\begin{aligned}
 &= \left( \frac{f_{k+1} - f_k}{h} \right) + \left( \frac{\epsilon_{k+1} - \epsilon_k}{h} \right) - \frac{h}{2} f''(\alpha) \\
 &= \left( \frac{f_{k+1} - f_k}{h} \right) + R + TE
 \end{aligned}$$

where R is the round-off error and TE is the truncation error.

If we take  $\epsilon = \max(|\epsilon_k|, |\epsilon_{k+1}|)$  and  $M_2 = \max |f''(x)|$ , we find that

$$|R| \leq \frac{2\epsilon}{h} \text{ and } |TE| \leq \frac{h}{2} M_2$$

We define the optimum value of h as the one which satisfies either of the following conditions :

$$(i) \quad |R| = |TE| \qquad (ii) \quad |R| + |TE| = \text{minimum} \qquad (64)$$

By the first condition in (64), we have

$$\frac{2\epsilon}{h} = \frac{h}{2} M_2 \text{ or } h^2 = \frac{4\epsilon}{M_2} \text{ or } h = 2 \sqrt{\epsilon/M_2}$$

The value of the error is

$$|R| = |TE| = \sqrt{\epsilon M_2}$$

If we use the second condition  $|R| = |TE| = \text{min}$ , we have

$$\frac{2\epsilon}{h} + \frac{h M_2}{2} = \text{min.} \qquad (65)$$

To find the minimum in Eqn. (65), we differentiate the left hand side of Eqn. (65) with respect to

$$-\frac{2\epsilon}{h^2} + \frac{M_2}{2} = 0 \text{ or } h^2 = \frac{4\epsilon}{M_2} \text{ or } h = 2\sqrt{\epsilon/M_2}$$

$\therefore$  The minimum total error =  $|R| + |TE| = \sqrt{\epsilon M_2}$

Let us now consider an example.

**Example 10 :** For the method

$$f'_k = \frac{(-3f_k + 4f_{k+1} - f_{k+2})}{2h} + \frac{h^2}{3} f''(\alpha), \quad x_k < \alpha < x_{k+1}$$

determine the optimal value of h using the criteria

$$|R| = |TE| \text{ and (ii) } |R| + |TE| = \text{min.}$$

Using this method and the first criterion, find the value of h and determine the value of  $f'(2.0)$ , from the following tabulated values of  $f(x) = \ln x$ . It is given that the maximum round-off error in the function evaluation is  $5 \times 10^{-6}$

x :	2.0	2.01	2.02	2.06	2.12
f(x) :	0.69315	0.69813	0.70310	0.72271	0.75142

**Solution :** If  $\epsilon_0, \epsilon_1$  and  $\epsilon_3$  are the round-off errors in the given function evaluations of  $f_0, f_1, f_2$  respectively, then we have

$$f'_0 = \frac{(-3f_0 + 4f_1 - f_2)}{2h} + \frac{(-3\epsilon_0 + 4\epsilon_1 - \epsilon_2)}{2h} + \frac{h^2}{3} f'''(\alpha)$$

determine  $h_{\text{opt}}$  using the criteria.

- (i)  $|R| = |TE|$  and
- (ii)  $|R| + |TE| = \text{minimum}$ .

Using this method and the second criterion, find  $h_{\text{opt}}$  for  $f(x) = \ln x$  and determine the value of  $f'(2.03)$  from the following table of values of  $f(x)$ , if it is given that the maximum round-off error in the function evaluation is  $5 \times 10^{-6}$

x	:	0.2	2.01	2.02	2.03	2.04	2.06
f(x)	:	0.69315	0.69813	0.70310	0.70804	0.71295	0.72271

We now end this unit by giving a summary of what we have covered in it.

## 12.7 SUMMARY

In this unit we have covered the following :

- 1) If a function  $f(x)$  is not known explicitly but a table of values of  $f(x)$  corresponding to a set of values of  $x$  is given then its derivatives can be obtained by numerical differentiation methods.
- 2) Numerical differentiation formulas using
  - (i) the method of undetermined coefficients and
  - (ii) methods based on finite difference operators can be obtained for the derivatives of a function at nodal or step points when the function is given in the form of table.
- 3) When it is required to find the derivative of a function at off-step points then the methods mentioned in (2) above cannot be used. In such cases, the methods derived from the interpolation formulas are useful.
- 4) Higher order solutions can be obtained by Richardson's extrapolation method which uses the lower order solutions. These results are more accurate than the results obtained directly from higher order differentiation formulas.
- 5) Round-off errors play a very important role in numerical differentiation. Sometimes, if the step size is too small, the round-off errors gets magnified unmanageably. In such cases the optimal step length for the given formula could be used, provided that it can be determined.

## 12.8 SOLUTIONS/ANSWERS

E1) Let  $f'(x) = \alpha_0 f_0 + \alpha_1 f_1 + \alpha_2 f_2$ . Setting  $f(x) = 1, x, x^2$ , we obtain

$$\alpha_0 + \alpha_1 + \alpha_2 = 0$$

$$(\alpha_1 + 2\alpha_2)h = 1$$

$$(\alpha_1 + 4\alpha_2)h^2 = 1$$

$$\text{Solving we obtain } \alpha_0 = -\frac{3}{2h}, \alpha_1 = \frac{2}{h}, \alpha_2 = -\frac{1}{2h}$$

$$\text{Hence, } f'_0 = \left( \frac{-3f_0 + 4f_1 - f_2}{2h} \right)$$

E2)  $0(h)$  method :  $hf'_k = (f_k - f_{k-1})$

$$0(h^2) \text{ method : } hf'_k = \left( \frac{3f_k - 2f_{k-1} + f_{k-2}}{2} \right)$$

$$O(h^3) \text{ method : } hf'_k = \left( \frac{11f_k - 22f_{k-1} + 9f_{k-2} - 6f_{k-3}}{6} \right)$$

$$O(h^4) \text{ method : } hf'_k = \left( \frac{25f_k - 56f_{k-1} + 36f_{k-2} - 24f_{k-3} + 3f_{k-4}}{12} \right)$$

E3) Using formula (18), we have

$$f'(6.0) = \left[ \frac{f(6.1) - f(6.0)}{0.1} \right] = -3.7480$$

Using formula (33),

$$f'(6.3) = \left[ \frac{f(6.4) - 2f(6.3) + f(6.2)}{(0.1)^2} \right] = 0.25$$

E4) Using formula (20), we have

$$f'(500) = \left[ \frac{-3f(500) + 4f(510) - f(520)}{2h} \right] = 0.002$$

Using (32), we have

$$f''(500) = \left[ \frac{2f(500) - 5f(510) + 4f(520) - f(530)}{h^2} \right] = -0.5 \times 10^{-5}$$

$$\text{Exact value } f'(x) = 1/x = 0.002; f''(x) = -1/x^2 = -0.4 \times 10^{-5}$$

Actual error in  $f'(500)$  is 0, whereas in  $f''(500)$  it is  $0.1 \times 10^{-5}$ . Truncation error in  $f'(x)$  is  $\frac{-h^2 f'''}{3} = -5.33 \times 10^{-7}$  and in  $f''(x)$  it is  $\frac{11h^2 f^{IV}}{12} = 8.8 \times 10^{-9}$

E5) In the given problem  $x_0 = 1, x_1 = 2, x_2 = 3, x_3 = 4$  and  $f_0 = 1, f_1 = 16, f_2 = 81$  and  $f_3 = 256$ .

Constructing the Lagrange fundamental polynomials, we get

$$L_0(x) = \left( \frac{x^3 - 9x^2 + 26x - 24}{6} \right); L_1(x) = \left( \frac{x^3 - 8x^2 + 19x - 12}{2} \right)$$

$$L_2(x) = \left( \frac{x^3 - 7x^2 + 14x - 8}{2} \right); L_3(x) = \left( \frac{x^3 - 6x^2 + 11x - 6}{6} \right)$$

$$P_3(x) = L_0'(x) f_0 + L_1'(x) f_1 + L_2'(x) f_2 + L_3'(x) f_3$$

$$P_3'(x) = L_0'(x) f_0 + L_1'(x) f_1 + L_2'(x) f_2 + L_3'(x) f_3$$

$$P_3''(x) = L_0''(x) f_0 + L_1''(x) f_1 + L_2''(x) f_2 + L_3''(x) f_3$$

We obtain after substitution,

$$P_3'(2.5) = 62.4167; P_3''(2.5) = 79; P_3'(5) = 453.667; P_3''(5) = 234.$$

The exact values of  $f'(x)$  and  $f''(x)$  are (from  $f(x) = x^4$ )

$$f'(2.5) = 62.5, f'(5) = 500; f''(2.5) = 75; f''(5) = 300.$$

E6) We are required to find  $f'(x)$  and  $f''(x)$  at  $x = 1.5$  and  $3.5$  which are off-step points. Using the Newton's forward difference formula with  $x_0 = 0, x = 1.5, s = 1.5$ , we get  $f'(1.5) = 8.7915$  and  $f''(1.5) = -4.0834$ .

Using the backward difference formula with  $x_n = 4, x = 3.5, s = -0.5$ , we get  $f'(3.5) = 7.393$  and  $f''(3.5) = 1.917$ .



E7) The difference table for the given problem is :

x	f(x)	$\Delta f$	$\Delta^2 f$	$\Delta^3 f$	$\Delta^4 f$
1.3	3.669				
		0.813			
1.5	4.482		0.179		
		0.992		0.41	
1.7	5.574		0.220		0.007
		1.212		0.48	
1.9	6.686		0.268		0.012
		1.480		0.060	
2.1	8.166		0.328		0.012
		1.808		0.072	
2.3	9.974		0.400		
		2.208			
2.5	12.182				

Taking  $x_0 = 1.5$  we see that  $s = 0$  and we obtain from the interpolation formula

$$f'(1.5) = \frac{1}{h} \left[ \Delta f_0 - \frac{\Delta^2 f_0}{2} + \frac{\Delta^3 f_0}{3} - \frac{\Delta^4 f_0}{4} + \dots \right]$$

$$= \left[ 0.992 - \frac{0.220}{2} + \frac{0.048}{3} - \frac{0.012}{4} \right] = 4.475$$

Exact value is  $e^{1.5} = 4.4817$  and error is  $= 0.0067$

E8) Use the  $O(h^2)$  formula (33). With  $h = 0.1$ ,  $f''(0.6) = 1.2596$ ,  $h = 0.2$ ,  $f''(0.6) = 1.26545$ ,  $h = 0.4$ ,  $f''(0.6) = 1.289394$ .

Using Richardson's extrapolation formula,

$$g^{(1)}(0.2) = \left[ \frac{4g(0.1) - g(0.2)}{3} \right] = 1.25765$$

$$g^{(1)}(0.4) = \left[ \frac{4g(0.2) - g(0.4)}{3} \right] = 1.257468$$

These two results are of  $O(h^4)$ . To get  $O(h^6)$  result we repeat the extrapolation technique and obtain

$$g^{(2)}(0.4) = \left[ \frac{16g^{(1)}(0.2) - g^{(1)}(0.4)}{15} \right] = 1.257662$$

E9) Using (24) with  $h = 0.1, 0.2$ , we have

$$g(0.1) = f''(0.3) = -3.8; g(0.2) = -3.575$$

$$g^{(1)}(0.2) = \left[ \frac{4g(0.1) - g(0.2)}{3} \right] = 3.875$$

E10) If  $\epsilon_{-1}, \epsilon_0, \epsilon_1$  are the round-off errors in the given function evaluations  $f_{-1}, f_0, f_1$  respectively, and if  $\epsilon = \max(|\epsilon_{-1}|, |\epsilon_0|, |\epsilon_1|)$  and  $M_3 = \max |f'''(x)|$  then

$$|R| \leq \frac{\epsilon}{h} \text{ and } |TE| \leq \frac{h^2}{3} M_3.$$

If we use  $|R| = |TE|$ , we get

$$h_{\text{opt}} \left( \frac{3\epsilon}{M_3} \right)^{1/3}$$

and error is given by

$$|R| = |TE| = \epsilon^{2/3} \left( \frac{M_3}{3} \right)^{1/3}$$

If we use  $|R| = |TE| = \min$ , then

$$h_{\text{opt}} \left( \frac{3\epsilon}{2M_3} \right)^{1/3}$$

$$\text{and error is } = \frac{M^{1/3}}{3} \left( \frac{3\epsilon}{2} \right)^{2/3}$$

For  $f(x) = \ln x$  and using the second criterion, we get

$$h_{\text{opt}} = \left( 30 \times 10^{-6} \right)^{1/3} = 0.03.$$

For  $h = 0.03$ , we get

$$f'(2.03) = \frac{0.72271 - 0.69315}{0.06} = 0.492667.$$

If we take  $h = 0.01$ , we get

$$f'(2.03) = 0.4925.$$

The exact value of  $f'(2.03) = 0.492611$ .

The result deteriorate for  $h < h_{\text{opt}}$ .

---

# UNIT 13 NUMERICAL INTEGRATION

---

## Structure

- 13.1 Introduction
  - Objectives
- 13.2 Methods Based on Interpolation
  - Methods Using Lagrange Interpolation
  - Methods Using Newton's Forward Interpolation
- 13.3 Composite Integration
- 13.4 Romberg Integration
- 13.5 Summary
- 13.6 Solutions/Answers

---

## 13.1 INTRODUCTION

---

In Unit 12, we developed methods of differentiation to obtain the derivative of a function  $f(x)$ , when its values are not known explicitly, but are given in the form of a table. In this unit, we shall now derive numerical methods for evaluating the definite integrals of such functions  $f(x)$ . You may recall that in calculus, the definite integral of  $f(x)$  over the interval  $[a, b]$  is defined as

$$\int_a^b f(x) dx = \lim_{h \rightarrow 0} R[h]$$

where  $R[h]$  is the left-end Riemann sum for  $n$  subintervals of length  $h = \frac{(b-a)}{n}$  and is given by

$$R[h] = \sum_{k=0}^{n-1} h f(x_k)$$

The need for deriving accurate numerical methods for evaluating the definite integral arises mainly, when the integral is either

- i) a complicated function such as  $f(x) = e^{-x^2}$ ,  $f(x) = \frac{\sin(x)}{x}$  etc. which have no anti-derivatives expressible in terms of elementary functions, or
- ii) when the integrand is given in the form of tables.

Many scientific experiments lead to a table of values and we may not only require an approximation to the function  $f(x)$  but also may require approximate representation of the integral of the function. Moreover, analytical evaluation of the integral may lead to transcendental, logarithmic or circular functions. The evaluation of these functions for a given value of  $x$  may not be an accurate process. This motivates us to study numerical integration methods which can be easily implemented on calculators.

In this unit we shall develop numerical integration methods wherein the integral is approximated by a linear combination of the values of the integrand i.e.,

$$\int_a^b f(x) dx = \beta_0 f(x_0) + \beta_1 f(x_1) + \dots + \beta_n f(x_n) \quad (1)$$

where  $x_0, x_1, \dots, x_n$  are the points which divide the interval  $[a, b]$  into  $n$  sub-intervals and  $\beta_0, \beta_1, \dots, \beta_n$  are the weights to be determined. We shall discuss in this unit, a few techniques to determine the unknowns in Eqn. (1).

## Objectives

After studying this unit you should be able to

- use trapezoidal and Simpson's rules of integration to integrate functions given in the form of tables and find the errors in these rules;
- improve the order of the results using Romberg integration or its accuracy, by composite rules of integration.

## 13.2 METHODS BASED ON INTERPOLATION

In Block 3, you have studied several interpolation formulas, which fits the given data  $(x_k, f_k)$ ,  $k = 0, 1, 2, \dots, n$ . We shall now see how these interpolation formulas can be used to develop numerical integration methods for evaluating the definite integral of a function which is given in a tabular form. The problem of numerical integration is to approximate the definite integral as a linear combination of the values of  $f(x)$  in the form

$$\int_a^b f(x) dx \approx \sum_{k=0}^n \beta_k f_k \quad (2)$$

where the  $n + 1$  distinct points  $x_k$ ,  $k = 0, 1, 2, \dots, n$  are called the **nodes** or **abscissas** which divide the interval  $[a, b]$  into  $n$  sub-intervals  $(x_0 < x_1 < x_2 < \dots < x_n)$  and  $\beta_k$ ,  $k = 0, 1, \dots, n$  are called the **weights** of the **integration rule** or **quadrature formula**. We shall denote the exact value of the definite integral by  $I$  and denote the rule of integration by

$$I_h[f] = \sum_{k=0}^n \beta_k f_k \quad (3)$$

The error of approximating the integral  $I$  by  $I_h[f]$  is given by

$$E_h[f] = \int_a^b f(x) dx - \sum_{k=0}^n \beta_k f_k \quad (4)$$

The order of the integration method (3) is defined as follows :

**Definition :** An integration method of the form (3) is said to be of order  $p$  if it produces exact results for all polynomials of degree less than or equal to  $p$ .

In Eqn. (3) we have  $2n + 2$  unknowns viz.,  $n + 1$  nodes  $x_k$  and the  $n + 1$  weights  $\beta_k$  and the method can be made exact for polynomials of degree  $\leq 2n + 1$ . Thus, the method of the form (3) can be of maximum order  $2n + 1$ . But, if some of the nodes are prescribed in advance, then the order will be reduced. If all the  $n + 1$  nodes are prescribed, then we have to determine only  $n + 1$  weights and the corresponding method will be of maximum order  $n$ .

We first derive the numerical method based on Lagrange interpolation.

### 13.2.1 Methods Using Lagrange Interpolation

Suppose we are given the  $n + 1$  abscissas  $x_k$ 's and the corresponding values  $f_k$ 's. We know that the unique Lagrange interpolating polynomial  $P_n(x)$  of degree  $\leq n$ , satisfying the interpolatory conditions  $P_n(x_k) = f(x_k)$ ,  $k = 0, 1, 2, \dots, n$ , is given by

$$f(x) \approx P_n(x) = \sum_{k=0}^n L_k(x) f_k \quad (5)$$

with the error of interpolation

$$E_{n+1}[P_n(x)] = \frac{\pi(x)}{(n+1)!} f^{(n+1)}(\alpha) \quad x_0 < \alpha < x_n \quad (6)$$

$$\text{where } L_k(x) = \frac{(x-x_0)(x-x_1)\dots(x-x_{k-1})(x-x_{k+1})\dots(x-x_n)}{(x_k-x_0)(x_k-x_1)\dots(x_k-x_{k-1})(x_k-x_{k+1})\dots(x_k-x_n)}$$

and  $\pi(x) = (x-x_0)(x-x_1)\dots(x-x_n)$ .

We replace the function  $f(x)$  in the definite integral (2) by the Lagrange interpolating polynomial  $P_n(x)$  given by Eqn. (5) and obtain

$$\begin{aligned} I_h[f] &= \int_a^b P_n(x) dx = \sum_{k=0}^n \int_a^b L_k(x) f_k dx \\ &= \sum_{k=0}^n \beta_k f_k \end{aligned} \quad (7)$$

where

$$\beta_k = \int_a^b L_k(x) dx. \quad (8)$$

The error in the integration rule is

$$E_n[f] = \int_a^b E_{n+1}[P_n(x)] dx = \int_a^b \frac{\pi(x)}{(n+1)!} f^{(n+1)}(\alpha) dx \quad (9)$$

We have

$$|E_n[f]| \leq \frac{M_{n+1}}{(n+1)!} \int_a^b |\pi(x)| dx \quad (10)$$

where  $M_{n+1} = \max_{x_0 < x < x_n} |f^{(n+1)}(x)|$

Let us consider now the case when the nodes  $x_k$ 's are equispaced with  $x_0 = a$ ,  $x_n = b$ , and the length of each subinterval is  $h = \frac{b-a}{n}$ . The numerical integration methods given by (7) are then known as **Newton-Cotes formulas** and the weights  $\beta_k$ 's given by (8) are known as **Cotes numbers**. Any point  $x \in [a, b]$  can be written as  $x = x_0 + sh$ .

With this substitution, we have

$$\pi(x) = h^{n+1} s(s-1)(s-2)\dots(s-n)$$

$$L_k(x) = \frac{(-1)^{n-k} s(s-1)\dots(s-k+1)(s-k-1)\dots(s-n)}{k!(n-k)!} \quad (11)$$

Using  $x = x_0 + sh$  and changing the variable of integration from  $x$  to  $s$ , we obtain

$$\beta_k = \frac{(-1)^{n-k}}{k!(n-k)!} h \int_0^n s(s-1)(s-2)\dots(s-k+1)(s-k-1)\dots(s-n) ds \quad (12)$$

$$\text{and } |E_n[f]| \leq \frac{h^{n+2} M_{n+1}}{(n+1)!} \int_a^b s(s-1)(s-2)\dots(s-n) ds \quad (13)$$

We now derive some of the Newton Cotes formulas viz. trapezoidal rule and Simpson's rule by using first and second degree Lagrange polynomials with equally spaced nodes. You might have studied these rules in your calculus course.

### Trapezoidal Rule

When  $n = 1$ , we have  $x_0 = a$ ,  $x_n = b$  and  $h = b-a$ . Using Eqn. (12) the Cotes numbers can be found as

$$\beta_0 = -h \int_0^1 (s-1) ds = \frac{h}{2};$$

$$\text{and } \beta_1 = h \int_0^1 s ds = \frac{h}{2}.$$

Substituting the values of  $\beta_0$  and  $\beta_1$  in Eqn. (7), we get

$$I_T[f] = \frac{h}{2} [f_0 + f_1] \quad (14)$$

The error of integration is

$$|E_T[f]| \leq \frac{h^3}{2} M_2 \int_0^1 s(s-1) ds = -\frac{h^3}{12} M_2 \quad (15)$$

where  $M_2 = \max_{x_0 < x < x_1} |f''(x)|$  (16)

Thus, by trapezoidal rule,  $\int_a^b f(x) dx$  is given by

$$I[f] = \frac{h}{2} (f_0 + f_1) - \frac{h^3}{12} M_2$$

The reason for calling this formula the trapezoidal rule is that geometrically when  $f(x)$  is a function with positive value then  $\frac{h}{2} (f_0 + f_1)$  is the area of the trapezium with height  $h = b - a$  and parallel sides as  $f_0$  and  $f_1$ . This is an approximation to the actual area under the curve  $y = f(x)$  above the  $x$ -axis bounded by the ordinates  $x = x_0, x = x_1$  (see Fig. 1). Since the error given by Eqn. (15) contains the second derivative, trapezoidal rule integrates exactly polynomials of degree  $\leq 1$ .

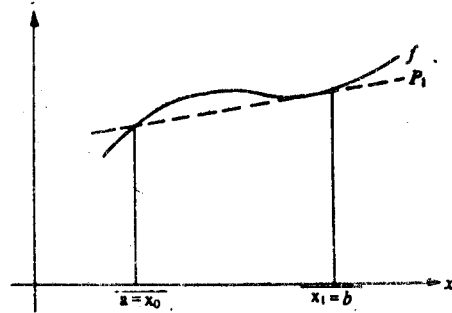


Fig. 1

Let us now consider an example.

**Example 1 :** Find the approximate value of

$$I = \int_0^1 \frac{dx}{1+x}$$

using trapezoidal rule and obtain a bound for the error. The exact value of  $I = \ln 2 = 0.693147$  correct to six decimal places.

**Solution :** Here  $x_0 = 0, x_1 = 1$  and  $h = 1 - 0 = 1$ . Using Eqn. (14), we get

$$I_T[f] = \frac{1}{2} \left( 1 + \frac{1}{2} \right) = 0.75$$

Actual error =  $0.75 - 0.693147 = 0.056853$ .

The error in the trapezoidal rule is given by

$$|E_T[f]| \leq \frac{1}{12} = \max_{0 \leq x \leq 1} \left| \frac{2}{(1+x)^3} \right| = \frac{1}{6} = 0.166667.$$

Thus, the error bound obtained is much greater than the actual error.

We now derive the Simpson's rule.

**Simpson's Rule**

For  $n = 2$ , we have  $h = \frac{b-a}{2}, x_0 = a, x_1 = \frac{a+b}{2}$  and  $x_2 = b$ .

From (12), we find the Cotes numbers as

$$\beta_0 = \frac{h}{2} \int_0^2 (s-1)(s-2) ds = \frac{h}{3}$$

$$\beta_1 = h \int_0^2 s(s-2) ds = \frac{4h}{3}, \beta_2 = \frac{h}{2} \int_0^2 s(s-1) ds = \frac{h}{3}.$$

Eqn. (7) in this case reduces to

$$I_s[f] = \frac{h}{3} [f_0 + 4f_1 + f_2] \quad (17)$$

Eqn. (17) is the Simpson's rule for approximating  $I = \int_a^b f(x) dx$ .

The magnitude of the error of integration is

$$\begin{aligned} |E_s[f]| &\leq \frac{h^4 M_3}{3!} \int_0^2 |s(s-1)(s-2)| ds \\ &= \frac{h^4 M_3}{3!} \left[ \int_0^1 s(s-1)(s-2) ds + \int_1^2 s(s-1)(s-2) ds \right] \\ &= \frac{h^4 M_3}{3!} \left[ \left( \frac{s^4}{4} - s^3 + s^2 \right)_0^1 + \left( \frac{s^4}{4} - s^3 + s^2 \right)_1^2 \right] \\ &= \frac{h^4 M_3}{3!} \left[ \frac{1}{4} - \frac{1}{4} \right] = 0 \end{aligned}$$

This indicates that Simpson's rule integrates polynomials of degree 3 also exactly. Hence, we have to write the error expression (13) with  $n = 3$ . We find

$$\begin{aligned} |E_s[f]| &\leq \frac{h^5 M_4}{24} \int_0^2 s(s-1)(s-2)(s-3) ds \\ &= \frac{h^5 M_4}{24} \left[ \int_0^1 s(s-1)(s-2)(s-3) ds + \int_1^2 s(s-1)(s-2)(s-3) ds \right] \\ &= \frac{-h^5 M_4}{90} \end{aligned} \quad (18)$$

where  $M_4 = \max_{x_0 < x < x_3} |f^{IV}(x)|$

Since the error in Simpson's rule contains the fourth derivative, Simpson's rule integrates exactly all polynomials of degree  $\leq 3$ .

Thus, by Simpson's rule,  $\int_a^b f(x) dx$  is given by

$$I[f] = \frac{h}{3} [f_0 + 4f_1 + f_2] - \frac{h^5}{90} M_4$$

Geometrically,  $\frac{h}{3} [f_0 + 4f_1 + f_2]$  represents the area bounded by the quadratic curve passing through  $(x_0, f_0)$ ,  $(x_1, f_1)$  and  $(x_2, f_2)$  above the  $x$ -axis and lying between the ordinates  $x = x_0$ ,  $x = x_2$  (see Fig. 2).

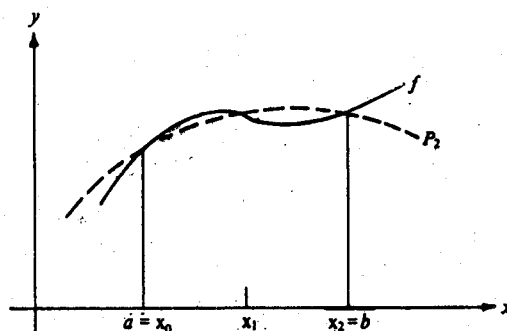


Fig. 2

In case we are given only one tabulated value in the interval  $[a, b]$ , then  $h = b - a$ , and the interpolating polynomial of degree zero is  $P_0(x) = f_k$ . In this case, we obtain the rectangular integration rule given by

$$I_R[f] = \int_a^b f_k dx \approx hf_k \quad (19)$$

The error in the integration rule is obtained from Eqn. (13) as

$$E_R[f] \leq \frac{h^2 M_1}{2} \quad (20)$$

where  $M_1 = \max_{a < x < b} |f'(x)|$

If the given tabulated value in the interval  $[a, b]$  is the value at the mid-point, then we have  $x_k = \frac{(a+b)}{2}$ , and  $f_k = f_{k+\frac{1}{2}}$ . In this case  $h = b - a$  and we obtain the integration rule as

$$I_M[f] = \int_a^b f_{k+\frac{1}{2}} dx \approx hf_{k+\frac{1}{2}} \quad (21)$$

Rule (21) is called the mid-point rule. The error in the rule calculated from (13) is

$$E_M[f] = \frac{h^2}{2} \int_{-\frac{1}{2}}^{\frac{1}{2}} s ds = 0.$$

This shows that the mid-point rule integrates polynomials of degree one exactly. Hence the error for the mid-point rule is given by

$$E_M[f] \leq \frac{h^3 M_2}{2} \int_{-\frac{1}{2}}^{\frac{1}{2}} s(s-1) ds = \frac{h^3 M_2}{24} \quad (22)$$

where  $M_2 = \max_{a < x < b} |f''(x)|$  and  $h = b - a$

We now illustrate these methods through an example.

**Example 2 :** Evaluate  $\int_0^1 e^{-x^2} dx$  using

a) rectangular rule b) mid-point rule c) trapezoidal rule and d) Simpson's rule.

If the exact value of the integral is 0.74682 correct to 5 decimal places, find the error in these rules.

**Solution :** The values of the function  $f(x) = e^{-x^2}$  at  $x = 0, 0.5$  and  $1$  are

$$f(0) = 1, f(0.5) = 0.7788, f(1) = 0.36788$$

Taking  $h = 1$  and using

a)  $I_R[f] = hf_0$ , we get  $I_R[f] = 1$ .

b)  $I_M[f] = hf_{1/2}$ , we get  $I_M[f] = 0.7788$ .

c)  $I_T[f] = \frac{h}{2} [f_0 + f_1]$ , we get  $I_T[f] = \frac{1}{2} (1 + 0.36788) = 0.68394$ , Taking  $h = 0.5$  and

using Simpson's rule, we get

d)  $I_S[f] = \frac{h}{3} [f_0 + 4f_1 + f_2]$

$$= \frac{h}{3} [f(0) + 4f(0.5) + f(1)]$$

$$= 0.74718.$$

Exact value of the integral is 0.74682.



The errors in these rules are given by

$$E_R[f] = -0.25318, E_M[f] = -0.03198$$

$$E_T[f] = 0.06288, E_{11}[f] = -0.00036.$$

You may now try the following exercise :

E1) Use the trapezoidal and Simpson's rule to approximate the following integrals. Compare the approximations to the actual value and find a bound for the error in each case.

a)  $\int_1^2 \ln x \, dx$

b)  $\int_0^{0.1} x^{1/3} \, dx$

c)  $\int_0^{\pi/4} \tan x \, dx$

We now derive integration methods using Newton's forward interpolation formula.

### 13.2.2 Methods Using Newton's Forward Interpolation

Let the data be given at equi-spaced nodal points  $x_k = x_0 + sh, s = 0, 1, 2, \dots, n$ , where  $x_0 = a$  and  $x_n = x_0 + nh = b$ .

The step length is given by  $h = \frac{b-a}{n}$ .

The Newton's forward finite difference interpolation formula interpolating this data is given by

$$f(x) \approx P_n(x) = f_0 + s\Delta f_0 + s(s-1)\frac{\Delta^2 f_0}{2} + \dots + \frac{s(s-1)(s-2)\dots(s-n+1)\Delta^n f_0}{n!} \quad (23)$$

with the error of interpolation

$$E_{n+1}[f] = \frac{h^{n+1} s(s-1)(s-2)\dots(s-n)}{(n+1)!} f^{(n+1)}(\alpha)$$

Integrating both sides of Eqn. (23) w.r.t.  $x$  between the limits  $a$  and  $b$ , we can approximate the definite integral  $I$  by the numerical integration rule

$$I_h[f] = \int_a^b P_n(x) \, dx = h \int_0^1 \left[ f_0 + s\Delta f_0 + \frac{s(s-1)}{2}\Delta^2 f_0 + \dots \right] ds \quad (24)$$

The error of interpolation of (24) is given by

$$|E_h(f)| \leq \frac{h^{n+2} M_{n+1}}{(n+1)!} \int_0^1 s(s-1)(s-2)\dots(s-n) \, ds$$

We can obtain the trapezoidal rule (14) from (24) by using linear interpolation i.e.,  $f(x) \approx P_1(x) = f_0 + s\Delta f_0$ . We then have

$$\begin{aligned} I_T[f] &= h \int_0^1 [f_0 + s\Delta f_0] \, ds \\ &= h \left[ s f_0 + \frac{s^2}{2} \Delta f_0 \right]_0^1 \\ &= h \left[ f_0 + \frac{\Delta f_0}{2} \right] = \frac{h}{2} [f_0 + f_1] \end{aligned}$$

with the error of integration given by (15).

Similarly Simpson's rule (16) can be obtained from (24) by using quadratic interpolation i.e.,  $f(x) \approx P_2(x)$ .

Taking  $x_0 = a$ ,  $x_1 = x_0 + h$ ,  $x_2 = x_0 + 2h = b$ , we have

$$\begin{aligned} I_s[f] &= \int_a^b f(x) dx \approx h \int_0^2 \left[ f_0 + s \Delta f_0 + \frac{s(s-1)}{2} \Delta^2 f_0 \right] ds \\ &= h \left[ 2f_0 + 2\Delta f_0 + \frac{\Delta^2 f_0}{3} \right] \\ &= \frac{h}{3} [f_0 + 4f_1 + f_2]. \end{aligned}$$

The error of interpolation is given by Eqn. (18).

**Example 3 :** Find the approximate value of  $I = \int_0^1 \frac{dx}{1+x}$  using

Simpson's rule. Obtain the error bound and compare it with the actual error. Also compare the result obtained here with the one obtained in Example 1.

**Solution :** Here  $x_0 = 0$ ,  $x_1 = 0.5$ ,  $x_2 = 1$  and  $h = \frac{1}{2}$ .

Using Simpson's rule, we have

$$I_s[f] = \frac{h}{3} [f(0) + 4f(0.5) + f(1)] = \frac{1}{6} \left[ 1 + \frac{8}{3} + 0.5 \right] = 0.694445$$

Exact value of  $I = \ln 2 = 0.693147$ .

Actual error = 0.001297. The bound for the error is given by

$$|E_s[f]| \leq \frac{h^5}{90} M_4 = 0.00833, \text{ where } M_4 = \max \left| \frac{24}{(1+x)^5} \right| = 24.$$

Here too the actual error is less than the given bound.

Also actual error obtained here is much less than that obtained in Example 1.

You may now try the following exercise.

E2) Find an approximation to  $\int_{1.1}^{1.5} e^x dx$ , using

- the trapezoidal rule with  $h = 0.4$
- Simpson's rule with  $h = 0.2$

The Newton-Cotes formulas as derived above are generally unsuitable for use over large integration intervals. Consider for instance, an approximation to

$\int_0^4 e^x dx$ , using Simpson's rule with  $h = 2$ . Here

$$\int_0^4 e^x dx \approx \frac{2}{3} (e^0 + 4e^2 + e^4) = 56.76958.$$

Since the exact value in this case is  $e^4 - e^0 = 53.59815$ , the error is  $-3.17143$ . This error is much larger than what we would generally regard as acceptable. However, large error is to be expected as the step length  $h = 2.0$  is too large to make the error expression meaningful. In such cases, we would be required to use higher order formulas. An alternate approach to obtain more accurate results while using lower order methods is the use of composite integration methods, which we shall discuss in the next section.

### 13.3 COMPOSITE INTEGRATION

In composite integration we divide the given interval  $[a, b]$  into a number of subintervals and evaluate the integral in each of the subintervals using one of the

integration rules. We shall construct composite rules of integration for trapezoidal and Simpson's methods and find the corresponding errors of integration when these composite rules are used.

### Composite Trapezoidal Rule

We divide the interval  $[a, b]$  into  $N$  subintervals of length  $h = \frac{(b-a)}{N}$ . We denote the subintervals as

$(x_{k-1}, x_k)$ ,  $k = 1, 2, \dots, N$  where  $x_0 = a$ ,  $x_N = b$ . Then

$$I = \int_a^b f(x) dx = \sum_{k=1}^N \int_{x_{k-1}}^{x_k} f(x) dx \quad (25)$$

Evaluating each of the integrals on the right hand side by trapezoidal rule, we have

$$\begin{aligned} I_T[f] &= \sum_{k=1}^N \frac{h}{2} [f_{k-1} + f_k] \\ &= \frac{h}{2} [f_0 + f_N + 2(f_1 + f_2 + \dots + f_{N-1})] \end{aligned} \quad (26)$$

The method (26) is known as **composite trapezoidal rule**. The error is given by

$$E_T[f] = -\frac{h^3}{12} \left[ \sum_{i=1}^N f''(\alpha_i) \right], \quad x_{k-1} < \alpha_i < x_k, \quad k = 1, \dots, N.$$

Now since  $f$  is a continuous function on the interval  $[a, b]$ , we have as a consequence of Intermediate-value theorem

$$\sum_{i=1}^N f''(\alpha_i) = f''(\xi) \sum_{i=1}^N 1, \quad \text{where } a < \xi < b,$$

$$\begin{aligned} \therefore E_T[f] &= -\frac{h^3}{12} f''(\xi) N, \quad a < \xi < b, \\ &= -\frac{Nh}{12} h^2 f''(\xi) \\ &= -\frac{(b-a)h^2}{12} f''(\xi). \end{aligned}$$

If  $M_2 = \max_{a < \xi < b} |f''(\xi)|$ . Then

$$|E_T[f]| \leq \frac{(b-a)h^2}{2} M_2 \quad (27)$$

The error is of order  $h^2$  and it decreases as  $h$  decreases.

Composite trapezoidal rule integrates exactly polynomials of degree  $\leq 1$ . We can try to remember the formula (26) as

$$I_T[f] = \left(\frac{h}{2}\right) [\text{first ordinate} + \text{last ordinate} + 2(\text{sum of the remaining ordinates})].$$

### Composite Simpson's Rule

In using Simpson's rule of integration (17), we need three abscissas. Hence, we divide the interval  $[a, b]$  into an even number of subintervals of equal length giving an odd

number of abscissas in the form  $a = x_0 < x_1 < x_2 < \dots < x_{2N} = b$  with  $h = \frac{b-a}{2N}$  and

$x_k = x_0 + kh, k = 0, 1, 2, \dots, 2N$ . We then write

$$I = \int_a^b f(x) dx = \sum_{k=1}^N \int_{x_{2k-2}}^{x_{2k}} f(x) dx \quad (28)$$

Evaluating each of the integrals on the right hand side of Eqn. (28) by the Simpson's rule, we have

$$\begin{aligned} I_s[f] &= \sum_{k=1}^N \frac{h}{3} [f_{2k-2} + 4f_{2k-1} + f_{2k}] \\ &= \frac{h}{3} [f_0 + f_{2N} + 4(f_1 + f_3 + \dots + f_{2N-1}) + 2(f_2 + f_4 + \dots + f_{2N})] \end{aligned} \quad (29)$$

The formula (29) is known as the composite Simpson's rule of numerical integration. The error in (29) is obtained from (18) by adding up the errors. Thus we get

$$\begin{aligned} E_s[f] &= -\frac{h^5}{90} \left[ \sum_{k=1}^N f^{IV}(\alpha_k) \right], x_{2k-2} < \alpha_k < x_{2k} \\ &= -\frac{h^5}{90} f^{IV}(\xi) \sum_{i=1}^N 1, a < \xi < b \\ &= -\frac{Nh^5}{90} f^{IV}(\xi) \\ &= -\frac{(b-a)h^4}{180} f^{IV}(\xi). \end{aligned}$$

If  $M_4 = \max_{a \leq \xi \leq b} |f^{IV}(\xi)|$ , we can write using  $h = \frac{(b-a)}{2N}$

$$|E_s[f]| \leq \frac{(b-a)}{180} h^4 M_4 = \frac{(b-a)^5 M_4}{2880N^4} \quad (30)$$

The error is of order  $h^4$  and it approaches zero very fast as  $h \rightarrow 0$ . The rule integrates exactly polynomials of degree  $\leq 3$ . We can remember the composite Simpson's rule as

$$I_s[f] = \left(\frac{h}{3}\right) [\text{first ordinate} + \text{last ordinate} + 2(\text{sum of even ordinates}) + 4(\text{sum of the remaining odd ordinates})]$$

We now illustrate composite trapezoidal and Simpson's rule through examples.

**Example 4 :** Evaluate  $\int_0^1 \frac{dx}{1+x}$  using

(a) composite trapezoidal rule and (b) composite Simpson's rule with 2, 4 and 8 subintervals.

**Solution :** We give in Table 1 the values of  $f(x)$  with  $h = \frac{1}{8}$  from  $x = 0$  to  $x = 1$ .

**Table 1**

$x$ :	0	1/8	2/8	3/8	4/8	5/8	6/8	7/8	1
$f(x)$ :	1	8/9	8/10	8/11	8/12	8/13	8/14	8/15	8/16
	$f_0$	$f_1$	$f_2$	$f_3$	$f_4$	$f_5$	$f_6$	$f_7$	$f_8$

If  $N = 2$  then  $h = 0.5$  and the ordinates  $f_0, f_4$  and  $f_8$  are to be used.

We get

$$I_T[f] = \frac{1}{4} [f_0 + 2f_4 + f_8] = \frac{17}{24} = 0.708333$$

$$I_s[f] = \frac{1}{6} [f_0 + 4f_1 + f_8] = \frac{23}{36} = 0.694444$$

If  $N = 4$  then  $h = 0.25$  and the ordinates  $f_0, f_2, f_4, f_6, f_8$  are to be used.

We have

$$I_T[f] = \frac{1}{8} [f_0 + f_8 + 2(f_2 + f_4 + f_6)] = 0.697024$$

$$I_s[f] = \frac{1}{12} [f_0 + f_8 + 4(f_2 + f_6) + 2f_4] = 0.693254$$

If  $N = 8$  then  $h = 1/8$  and all the ordinates in Table 1 are to be used.

We obtain

$$I_T[f] = \frac{1}{16} [f_0 + f_8 + 2(f_1 + f_2 + f_3 + f_4 + f_5 + f_6 + f_7)] = 0.694122$$

$$I_s[f] = \frac{1}{24} [f_0 + f_8 + 4(f_1 + f_3 + f_5 + f_7) + 2(f_2 + f_4 + f_6)] = 0.693147$$

The exact value of the given integral correct to six decimal places is  $ln2 = 0.693147$ .

We now give the actual errors in Table 2 below.

**Table 2**

N	$E_T [f]$	$E_s [f]$
2	0.015186	0.001297
4	0.003877	0.000107
8	0.000975	0.000008

Note that as  $h$  decreases the errors in both trapezoidal and Simpson's rule also decreases.

Let us consider another example.

**Example 5 :** Find the minimum number of intervals required to evaluate  $\int_0^1 \frac{dx}{1+x}$  with an accuracy  $10^{-6}$ , by using the Simpson rule.

**Solution :** In Example 4 you may observe from Table 2 that  $N \approx 8$  gives  $10^{-6}$  (1.E - 06) accuracy. We shall now determine  $N$  from the theoretical error bound for Simpson's rule which gives 1.E - 06 accuracy. Now

$$|E_s[f]| \leq \frac{(b-a)^5 M_4}{2880N^4}$$

where

$$M_4 = \max_{0 < x < 1} |f^{IV}(x)|$$

$$= \max_{0 < x < 1} \left| \frac{24}{(1+x)^5} \right| = 24$$

To obtain the required accuracy we should therefore have

$$\frac{24}{2880N^4} \leq 10^{-6}, \text{ or } N^4 \geq \frac{24 \times 10^6}{2880} = 8333.3333$$

$$\therefore N \geq 9.5$$

We find that we cannot take  $N = 9$  since to make use of Simpson's rule we should have even number of intervals. We therefore conclude that  $N = 10$  should be the minimum number of subintervals to obtain the accuracy 1.E - 0.6 (i.e.,  $10^{-6}$ )

You may now try the following exercises :

E3) Evaluate  $\int_0^1 \frac{dx}{1+x^2}$  by subdividing the interval (0, 1) into 6 equal parts and using

(a) Trapezoidal rule (b) Simpson's rule. Hence find the value of  $\pi$  and actual errors.

E4) A function  $f(x)$  is given by the table

$x$	1.0	1.5	2.0	2.5	3.0
$f(x)$	1.000	2.875	7.000	14.125	25.000

Find the integral of  $f(x)$  using (a) trapezoidal rule (b) Simpson's rule.

E5) The speedometer reading of a car moving on a straight road is given. Estimate the distance travelled by the car in 12 minutes using (a) Trapezoidal rule (b) Simpson's rule.

Time : (minutes)	0	2	4	6	8	10	12
Speedometer Reading	0	15	25	40	45	20	0

E6) Evaluate  $\int_{0.2}^{0.4} (\sin x - \ln x + e^x) dx$  using (a) Trapezoidal rule (b) Simpson's rule taking  $h = 0.1$ . Find the actual errors.

E7) Determine  $N$  so that the composite trapezoidal rule gives the value of  $\int_0^1 e^{-x^2} dx$  correct upto 3 digits after the decimal point, assuming that  $e^{-x^2}$  can be calculated accurately.

You must have realised that though the trapezoidal rule is the easiest Newton-Cotes formula to apply but it lacks the degree of accuracy generally required. There is a way to improve the accuracy of the results obtained by the trapezoidal and Simpson rules. This method is known as **Romberg integration**, or as **extrapolation to the limit**. Richardson's extrapolation technique (ref. Sec. 12.5 of Unit 12) applied to the integration methods is called Romberg integration. We shall now discuss this technique in the next section.

### 13.4 ROMBERG INTEGRATION

In Romberg integration, first we find the power series expansion of the error term in the integration method. Then by eliminating the leading terms in the error expression, we obtain new values which are of higher order than the previously computed values.

If  $F_0(h)$  denotes the approximate value obtained by using the composite trapezoidal rule, then

$$I = F_0(h) + C_1 h^2 + C_2 h^4 + C_3 h^6 + \dots$$

where  $I$  is the exact value of the integral.

Let the integral be evaluated with the step lengths  $h$ ,  $\frac{h}{2}$  and  $\frac{h}{4}$ .

Then

$$I = F_0(h) + C_1 h^2 + C_2 h^4 + \dots \quad (31)$$

$$I = F_0\left(\frac{h}{2}\right) + \frac{C_1}{4} h^2 + \frac{C_2}{16} h^4 + \dots \quad (32)$$

Eliminating  $C_1$  from Eqns. (31) and (32), we get

$$I \approx \frac{1}{3} \left[ 4F_0\left(\frac{h}{2}\right) - F_0(h) \right] = F_1(h) \quad (33)$$

Note that this value is of  $O(h^4)$ . Similarly,

$$I \approx \frac{1}{3} \left[ 4F_0\left(\frac{h}{4}\right) - F_0\left(\frac{h}{2}\right) \right] = F_1\left(\frac{h}{2}\right) \quad (34)$$

etc.

Applying this method repeatedly by eliminating  $C_2$ , then  $C_3$  etc. we get the Romberg integration formula

$$F_m(h) = \frac{4^m F_{m-1}\left(\frac{h}{2}\right) - F_{m-1}(h)}{4^m - 1}, m = 1, 2, \dots \tag{35}$$

In the same way if  $G_0(h)$  denote the value of the integral obtained by using the Simpson's rule, then

$$I = G_0(h) + d_1 h^4 + d_2 h^6 + d_3 h^8 + \dots$$

where  $I$  is the exact value of the integral.

Let the integral be evaluated with step lengths  $h$ ,  $h/2$  and  $h/4$ .

Then, we have

$$I = G_0(h) + d_1 h^4 + d_2 h^6 + \dots \tag{36}$$

$$I = G_0\left(\frac{h}{2}\right) + \frac{d_1}{16} h^4 + \frac{d_2}{64} h^6 + \dots \tag{37}$$

Eliminating  $d_1$  from Eqns. (36) and (37), we get

$$I \approx \frac{4^2 G_0\left(\frac{h}{2}\right) - G_0\left(\frac{h}{2}\right)}{4^2 - 1} = G_1(h) \tag{38}$$

Similarly,

$$I \approx \frac{4^2 G_0\left(\frac{h}{4}\right) - G_0\left(\frac{h}{2}\right)}{4^2 - 1} = G_1\left(\frac{h}{2}\right) \tag{39}$$

etc.

Note that these values are of order  $h^6$ .

Applying extrapolation technique repeatedly, we get

$$G_m(h) = \frac{4^{m+1} G_{m-1}\left(\frac{h}{2}\right) - G_{m-1}(h)}{4^{m+1} - 1}, m = 1, 2, \dots \tag{40}$$

We now illustrate this technique through an example.

**Example 6 :** Find the value of the integral  $I = \int_0^1 f(x) dx$  where  $f(x) = \frac{1}{1+x}$  using

(a) composite trapezoidal and (b) composite Simpson's rules, with 3, 5 and 9 nodes. Use extrapolation technique to improve the results.

**Solution :** We take the computed values from Example 4. The Romberg integration values are given in Tables 3 and 4 for composite trapezoidal and composite Simpson's rules respectively.

**Table 3**

N	h	$I_T [f] = F_0(h)$	$F_1(h)$	$F_2(h)$
2	$\frac{1}{2}$	0.708333	0.693155	0.693148
4	$\frac{1}{4}$	0.697024	0.693254	
8	$\frac{1}{8}$	0.694122		

Note that

$$F_1\left(\frac{1}{2}\right) = \frac{4F_0\left(\frac{1}{4}\right) - F_0\left(\frac{1}{2}\right)}{3}$$

$$F_1\left(\frac{1}{4}\right) = \frac{4F_0\left(\frac{1}{8}\right) - F_0\left(\frac{1}{4}\right)}{3}$$

and

$$F_2\left(\frac{1}{2}\right) = \frac{16F_1\left(\frac{1}{4}\right) - F_1\left(\frac{1}{2}\right)}{15}$$

Table 4

N	h	$I_s [f] = G_0 (h)$	$G_1 (h)$	$G_2 (h)$
2	$\frac{1}{2}$	0.694444	0.693175	0.693148
4	$\frac{1}{4}$	0.693254	0.693148	
8	$\frac{1}{8}$	0.693155		

Note that

$$G_1\left(\frac{1}{2}\right) = \frac{16G_0\left(\frac{1}{4}\right) - G_0\left(\frac{1}{2}\right)}{15}$$

$$G_1\left(\frac{1}{4}\right) = \frac{16G_0\left(\frac{1}{8}\right) - G_0\left(\frac{1}{4}\right)}{15}$$

and

$$G_2\left(\frac{1}{2}\right) = \frac{64G_1\left(\frac{1}{4}\right) - G_1\left(\frac{1}{2}\right)}{63}$$

Suppose that we wish to evaluate the integral in the above example directly by the trapezoidal and Simpson's rules to an accuracy  $1.0E - 06$ . What should be the maximum value of step length to be chosen to achieve this accuracy ?

To answer this question let us calculate the error bound for trapezoidal rule.

$$|E_T[f]| \leq \frac{h^2}{12} \max_{0 < x < 1} |f''(x)| = \frac{h^2}{12} \max_{0 < x < 1} \left| \frac{2}{(1-x)^3} \right| > \frac{h^2}{48}$$

Hence

$$\frac{h^2}{48} < 1.E-06 \text{ or } h \approx 0.007$$

or  $N \approx 145$ .

Thus to obtain  $1.E - 06$  accuracy by trapezoidal rule we need to use 145 subintervals, i.e., 146 function evaluations. But by extrapolation we have used only 9 evaluations and improved these values.

Let us consider another example.

**Example 7 :** Use composite trapezoidal rule to find  $\int_1^{2.2} \ln x \, dx$  with  $N = 3, 6, 12$  and improve the accuracy by Romberg integration.

**Solution :** We give the result in the form of the following table.

Table 5

N	h	$I_T [f] = F_0 (h)$	$F_1 (h)$	$F_2 (h)$
3	0.4	0.527395	0.534605	0.534606
6	0.2	0.534591	0.534591	
12	0.1	0.534152		



You may now try the following exercise.

- E8) The following table gives the values of  $\ln x$  for  $x = 1, 2, \dots, 11$ . Evaluate the integral of the tabulated function using Trapezoidal rule with  $h = 1, 2$ . Use Richardson's extrapolation technique to improve the accuracy and obtain the actual error. Compare the results obtained by using Simpson's rule with  $h = 1$ .

x	:	1	2	3	4	5	6
ln x	:	0.0000	0.6931	1.0986	1.3863	1.6094	1.7918
x	:	7	8	9	10	11	
ln x	:	1.9459	2.0974	2.1972	2.3026	2.3979	

We now end this unit by giving a summary of what we have covered in it.

### 13.5 SUMMARY

In this unit, we have learnt the following :

- 1) If a function  $f(x)$  is not known explicitly but a table of values of  $x$  is given or when it has no anti-derivative expressible in terms of elementary functions then its integral cannot be obtained by calculus methods. In such cases numerical integration methods are used to find the definite integral of  $f(x)$  using the given data.
- 2) The basic idea of numerical integration methods is to approximate the definite integral as a linear combination of the values of  $f(x)$  in the form

$$\int_a^b f(x) dx \approx \sum_{k=0}^n \beta_k f(x_k) \quad (\text{see Eqn. (2)})$$

where the  $(n + 1)$  distinct nodes  $x_k, k = 0, 1, \dots, n, x_0 < x_1 < x_2 < \dots < x_n$  divide the integral  $[a, b]$  into  $n$  subintervals and  $\beta_k, k = 0, 1, \dots, n$  are the weights of the integration rule. The error of the integration methods is then given by

$$|E_n[f]| = \left| \int_a^b f(x) dx - \sum_{k=0}^n \beta_k f(x_k) \right| \quad (\text{see Eqn. (4)})$$

- 3) For equispaced nodes, the integration formulas derived by using Lagrange interpolating polynomials  $P_n(x)$  of degree  $\leq n$ , satisfying the interpolatory conditions  $P_n(x_k) = f(x_k), k = 0, 1, \dots, n$  are known as Newton-Cotes formulas. Corresponding to  $n = 1$  and  $n = 2$ , Newton-Cotes formulas viz., trapezoidal rule and Simpson's rule are obtained.
- 4) For large integration intervals, the Newton-Cotes formulas are generally unsuitable for they give large errors. Composite integration methods can be used in such cases by dividing the interval into a large number of subintervals and evaluating the integral in each of the subintervals using one of the integration rules.
- 5) For improving the accuracy of the trapezoidal or Simpson's rules and to obtain higher order solutions, Romberg's integration can be used.

### 13.6 SOLUTIONS/ANSWERS

$$\begin{aligned} \text{E1) a) } I_T[f] &= \frac{h}{2} [f_0 + f_1] = 0.346574 \\ I_S[f] &= \frac{h}{3} [f_0 + 4f_1 + f_2] \\ &= \frac{0.5}{3} [4 \ln 1.5 + \ln 2] = 0.385835 \end{aligned}$$

Exact value of  $I = 0.386294$

Actual error in  $I_T[f] = 0.03972$

Actual error in  $I_S[f] = 0.000459$

Also

$$|E_T[f]| \leq -\frac{h^3}{12} \max_{1 < x < 2} \left| \frac{1}{x^2} \right| = -\frac{1}{12} = -0.083334$$

$$|E_S[f]| \leq -\frac{h^5}{90} \max_{1 < x < 2} \left| \frac{6}{x^4} \right| = -0.002083$$

b)  $I_T[f] = 0.023208$ ,  $|E_T[f]| = \text{none}$ .

$I_S[f] = 0.032296$ ,  $|E_S[f]| = \text{none}$ .

Exact value = 0.034812.

c)  $I_T[f] = 0.39270$ ,  $|E_T[f]| = 0.161$

$I_S[f] = 0.34778$ ,  $|E_S[f]| = 0.00831$

Exact value = 0.34657.

E2)  $I_T[f] = 1.49718$

$I_S[f] = 1.47754$ .

E3) With  $h = 1/6$ , the values of  $f(x) = \frac{1}{1+x^2}$

from  $x = 0$  to  $1$  are

x	:	0	1/6	2/6	3/6	4/6	5/6	1
f(x)	:	1	0.972973	0.9	0.8	0.692308	0.590164	0.5

Now

$$I_T[f] = \frac{h}{2} [f_0 + f_6 + f_6 + 2(f_1 + f_2 + f_3 + f_4 + f_5)]$$

$$= 0.784241$$

$$I_S[f] = \frac{h}{3} [f_0 + f_6 + 4(f_1 + f_3 + f_5) + 2(f_2 + f_4)]$$

$$= 0.785398$$

$$\int_0^1 \frac{dx}{1+x^2} = [\tan^{-1}x]_0^1 = \frac{\pi}{4}. \quad \text{Exact } \pi = 3.141593$$

Value of  $\pi$  from  $I_T[f] = 4 \times 0.784241 = 3.136963$

Error in calculating  $\pi$  by  $I_T[f]$  is  $E_T[f] = 0.004629$

Value of  $\pi$  from  $I_S[f] = 4 \times 0.785398 = 3.141592$

Error in  $\pi$  by  $I_S[f]$  is  $E_S[f] = 1.0 \times 10^{-6}$ .

E4)  $I_T[f] = \left(\frac{h}{2}\right) [f_0 + f_4 + 2(f_1 + f_2 + f_3)]$

$$= (1/4) [1 + 25 + 2(2.875 + 7 + 14.125)] = 18.5$$

$$I_S[f] = \left(\frac{h}{3}\right) [f_0 + f_4 + 2f_2 + 4(f_1 + f_3)]$$

$$= (1/6) [1 + 25 + 2 \times 7 + 4(2.875 + 14.125)] = 18$$

E5) Let  $v_0 = 0, v_1 = 15, v_2 = 25, v_3 = 40, v_4 = 45, v_5 = 20, v_6 = 0$ . Then

$$I = \int_0^{12} v \, dt, I_T[v] = \left(\frac{h}{2}\right) [v_0 + v_6 + 2(v_1 + v_2 + v_3 + v_4 + v_5)] = 290$$

$$I_s[v] = \frac{880}{30} = 293.33.$$

E6) The values of  $f(x) = \sin x = \ln x + e^x$  are

$$f(0.2) = 3.02951, f(0.3) = 2.849352, f(0.4) = 2.797534$$

$$I_T[f] = \left(\frac{0.1}{2}\right) [f(0.2) + 2f(0.3) + f(0.4)] = 0.57629$$

$$I_s[f] = \left(\frac{0.1}{3}\right) [f(0.2) + 4f(0.3) + f(0.4)] = 0.574148$$

$$\text{Exact value} = 0.574056$$

$$E_T = 2.234 \times 10^{-03}$$

$$E_S = 9.2 \times 10^{-05}$$

E7) Error in composite trapezoidal rule

$$E_T[f] = -\frac{(b-a)^3}{12N^2} M_2, M_2 = \max_{0 < x < 1} |f''(x)|$$

Thus

$$|E_T[f]| \leq \frac{1}{12N^2} \max_{0 < x < 1} |f''(x)|$$

$$f(x) = e^{-x^2}, f''(x) = e^{-x^2} (4x^2 - 2)$$

$$f'''(x) = e^{-x^2} 4x(3-2x^2) = 0 \text{ when } x = 0, x = \sqrt{1.5}$$

$$\max \{|f''(0)|, |f''(1)|\} = \max [2, 2e^{-1}] = 2$$

For getting the correct value upto 3 digits, we must have

$$\frac{2}{12N^2} < 10^{-03} \text{ or } N^2 > \frac{10^3}{6} = \frac{10^4}{60}$$

or

$$N > \frac{100}{\sqrt{60}} = 12.9.$$

The interger value is  $N = 13$ .

E8) With  $h = 1$ , using trapezoidal rule

$$I_T[f] = F(h) = \left(\frac{h}{2}\right) \left[ f_0 + f_{10} + 2 \sum_{k=1}^9 f_k \right]$$

$$= 16.32125$$

With  $h = 2$ , -

$$I_T[f] = F(2h) = \left(\frac{h}{3}\right) \left[ f_0 + f_{10} + 2 \sum_{k=1}^5 f_{2k} \right]$$

$$= 16.1001$$

By extrapolation

$$F_1(h) = \frac{1}{3} [4F(h) - F(2h)] = 16.39496667$$

By Simpson's rule

$$I_S[f] = \left(\frac{h}{3}\right) \left[ f_0 + f_{10} + 4 \sum_{k=1}^5 f_{2k-1} + 2 \sum_{k=1}^4 f_{2k} \right]$$

$$= 16.39496667 \text{ (which is same as the value obtained by extrapolation)}$$

Exact value of the integral = 16.376848

Actual error = 0.01811867.

---

# UNIT 14 NUMERICAL SOLUTION OF ORDINARY DIFFERENTIAL EQUATIONS

---

## Structure

- 14.1 Introduction
  - Objectives
- 14.2 Basic Concepts
- 14.3 Taylor Series Method
- 14.4 Euler's Method
- 14.5 Richardson's Extrapolation
- 14.6 Summary
- 14.7 Solutions/Answers

---

## 14.1 INTRODUCTION

---

In the previous two units, you have seen how a complicated or tabulated function can be replaced by an approximating polynomial so that the fundamental operations of calculus viz., differentiation and integration can be performed more easily. In this unit we shall solve a differential equation, that is, we shall find the unknown function which satisfies a combination of the independent variable, dependent variable and its derivatives. In physics, engineering, chemistry and many other disciplines it has become necessary to build mathematical models to represent complicated processes. Differential equations are one of the most important mathematical tools used in modelling problems in the engineering and physical sciences. As it is not always possible to obtain the analytical solution of differential equations recourse must necessarily be made to numerical methods for solving differential equations. In this unit, we shall introduce two such methods namely, Euler's method and Taylor series method to obtain numerical solution of ordinary differential equations (ODEs). We shall also introduce Richardson's extrapolation method to obtain higher order solutions to ODEs using lower order methods. To begin with, we shall recall few basic concepts from the theory of differential equations which we shall be referring quite often.

### Objectives

After studying this unit you should be able to :

- identify the initial value problem for the first order ordinary differential equations;
- obtain the solution of the initial value problems by using Taylor series method and Euler's method;
- use Richardson's extrapolation technique for improving the accuracy of the result obtained by Euler's method.

---

## 14.2 BASIC CONCEPTS

---

In this section we shall state a few definitions from the theory of differential equations and define some concepts involved in the numerical solution of differential equations.

**Definition :** An equation involving one or more unknown functions (dependent variables) and its derivatives with respect to one or more known functions (independent variables) is called a **differential equation**.

For example,

$$x \frac{dy}{dx} = 2y \quad (1)$$

$$x \frac{\partial z}{\partial x} + y \frac{\partial z}{\partial y} - z = 0 \tag{2}$$

are differential equations.

Differential equations of the form (1), involving derivatives w.r.t. a single independent variable are called **ordinary differential equations (ODEs)** whereas, those involving derivatives w.r.t. two or more independent variables are **partial differential equations (PDEs)**. Eqn. (2) is an example of PDE.

**Definition :** The **order** of a differential equation is the order of the highest order derivative appearing in the equation and its **degree** is the highest exponent of the highest order derivative after the equation has been rationalised i.e., after it has been expressed in the form free from radicals and any fractional power of the derivatives or negative power. For example equation

$$\left(\frac{d^3y}{dx^3}\right)^2 + 2 \frac{d^2y}{dx^2} - \frac{dy}{dx} + x^2 \left(\frac{dx}{dx}\right)^3 = 0 \tag{3}$$

is of **third order** and **second degree**. Equation

$$y = x \frac{dy}{dx} + \frac{a}{dy/dx}$$

is of **first order** and **second degree** as it can be written in the form

$$y \frac{dy}{dx} = x \left(\frac{dx}{dx}\right)^2 + a \tag{4}$$

**Definition :** When the dependent variable and its derivatives occur in the first degree only and not as higher powers or products, the equation is said to be **linear**; otherwise it is **nonlinear**.

Equation  $\frac{d^2y}{dx^2} + y = x^2$  is a linear ODE, whereas,  $(x+y)^2 \frac{dy}{dx} = 1$  is a nonlinear ODE.

Similarly,  $\frac{\partial^2 z}{\partial x^2} + \frac{\partial^2 z}{\partial y^2} - \left(\frac{\partial^2 z}{\partial x \partial y}\right) = 0$ , is a nonlinear PDE.

In this unit we shall be concerned only with the ODEs.

The general form of a linear ODE of order n can be expressed in the form

$$L[y] = a_0(t) y^{(n)}(t) + a_1(t) y^{(n-1)}(t) + \dots + a_{n-1}(t) y'(t) + a_n(t) y(t) = r(t) \tag{5}$$

where  $r(t)$ ,  $a_i(t)$ ,  $i = 1, 2, \dots, n$  are known functions of t and

$$L = a_0(t) \frac{d^n}{dt^n} + a_1(t) \frac{d^{n-1}}{dt^{n-1}} + \dots + a_{n-1}(t) \frac{d}{dt} + a_n(t),$$

is the linear differential operator. The general nonlinear ODE of order n can be written as

$$F(t, y, y', y'', \dots, y^{(n)}) = 0 \tag{6}$$

$$\text{or, } y^{(n)} = f(t, y, y', y'', \dots, y^{(n-1)}) \tag{7}$$

Eqn. (7) is called a canonical representation of Eqn. (6). In such a form, the highest order derivative is expressed in terms of lower order derivatives and the independent variable.

The **general solution** of an nth order ODE contains n arbitrary constants. In order to determine these arbitrary constants, we require n conditions. If these conditions are given at one point, then these conditions are known as **initial conditions** and the differential equation together with the initial conditions is called an **initial value problem (IVP)**. The nth order IVP can be written as

$$y^{(n)}(t) = f(t, y, y', y'', \dots, y^{(n-1)})$$

$$y^{(p)}(t_0) = y_0^{(p)}, p = 0, 1, 2, \dots, n-1. \tag{8}$$

$$y_{(0)}^{(k)} = \frac{d^k y}{dt^k}$$

$k = 1, 2, \dots, n$

If the  $n$  conditions are prescribed at more than one point then these conditions are known as **boundary conditions**. The differential equation together with the boundary conditions is then known as a **boundary value problem (BVP)**.

The  $n$ th order IVP (8) is equivalent to the following system of  $n$  first order equations :

Set  $y = y_1$ . Then

$$\begin{array}{ll} y' = y_1' = y_2 & y_1(t_0) = y_0 \\ y_2' = y_3 & y_2(t_0) = y_0' \\ \dots\dots\dots & \dots\dots\dots \\ y_{n-1}' = y_n & y_{n-1}(t_0) = y_0^{(n-2)} \\ y_n' = f(t, y_1, y_2, \dots, y_n) & y_n(t_0) = y_0^{(n-1)}; \end{array}$$

In vector notation, this system can be written as a single equation as

$$\frac{dy}{dt} = f(t, y), \quad y(t_0) = \alpha \tag{9}$$

where  $y = (y_1, y_2, \dots, y_n)^T$ ,  $f(t, y) = (y_2, y_3, \dots, f(t, y_1, \dots, y_n))^T$

$$\alpha = (y_0, y_0', \dots, y_0^{(n-1)})^T$$

Hence, it is sufficient to study numerical methods for the solution of the first order IVP.

$$y' = f(t, y), \quad y(t_0) = y_0 \tag{10}$$

The vector form of these methods can then be used to solve Eqn. (9). Before attempting to obtain numerical solutions to Eqn. (10), we must make sure that the problem has a unique solution. The following theorem ensures the existence and uniqueness of the solution to IVP (10).

**Theorem 1 :** If  $f(t, y)$  satisfies the conditions

- i)  $f(t, y)$  is a real function
- ii)  $f(t, y)$  is defined and continuous for  $t \in [t_0, b]$ ,

$$y \in ]-\infty, \infty [$$

- iii) there exists a constant  $L$  such that for any  $t \in [t_0, b]$  and for any two numbers  $y_1$  and  $y_2$

$$|f(t, y_1) - f(t, y_2)| \leq L |y_1 - y_2|$$

then for any  $y_0$ , the IVP (10) has a unique solution. This condition is called the **Lipschitz condition** and  $L$  is called the **Lipschitz constant**.

We assume the existence and uniqueness of the solution and also that  $f(t, y)$  has continuous partial derivatives w.r.t.  $t$  and  $y$  of as high order as we desire.

Let us assume that  $[t_0, b]$  be an interval over which the solution of the IVP (10) is required. If we subdivide the interval  $[t_0, b]$  into  $n$  subintervals using a stepsize

$h = \left[ \frac{t_n - t_0}{n} \right]$ , where  $t_n = b$ , we obtain the **mesh points** or **grid points**  $t_0, t_1, t_2, \dots, t_n$  as shown in Fig. 1.

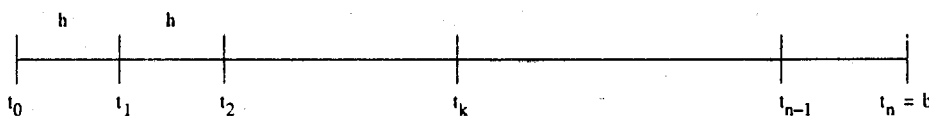


Fig. 1

We can then write  $t_k = t_0 + kh$ ,  $k = 0, 1, \dots, n$ . A numerical method for the solution of the IVP (10), will produce approximate values  $y_k$  at the grid points  $t_k$ .

Remember that the approximate values  $y_k$  may contain the truncation and round-off errors. We shall now discuss the construction of numerical methods and related basic concepts with reference to a simple ODE.

$$\frac{dy}{dt} = \lambda y, t \in [t_0, b]$$

$$y(t_0) = y_0 \tag{11}$$

Let the grid points be defined by

$$t_j = t_0 + jh, j = 0, 1, \dots, N$$

where  $t_0 = a$  and  $t_0 + Nh = b$ .

Separating the variables and integrating, we find that the exact solution of Eqn. (11) is

$$y(t) = y(t_0) e^{\lambda(t-t_0)} \tag{12}$$

In order to obtain a relation connecting two successive solution values, we set  $t = t_n$  and  $t_{n+1}$  in Eqn. (12). Thus we get

$$y(t_n) = y(t_0) e^{\lambda(t_n-t_0)}$$

and

$$y(t_{n+1}) = y(t_0) e^{\lambda(t_{n+1}-t_0)}$$

Dividing, we get

$$\frac{y(t_{n+1})}{y(t_n)} = \frac{e^{\lambda t_{n+1}}}{e^{\lambda t_n}} = e^{\lambda(t_{n+1}-t_n)}$$

Hence we have

$$y(t_{n+1}) = e^{\lambda h} y(t_n), n = 0, 1, \dots, N-1 \tag{13}$$

Eqn. (13) gives the required relation between  $y(t_n)$  and  $y(t_{n+1})$ .

Setting  $n = 0, 1, 2, \dots, N-1$ , successively, we can find  $y(t_1), y(t_2), \dots, y(t_N)$  from the given value  $y(t_0)$ .

An approximate method or a numerical method can be obtained by approximating  $e^{\lambda h}$  in Eqn. (13). For example, we may use the following polynomial approximations.

$$e^{\lambda h} = 1 + \lambda h + 0 (|\lambda h|^2) \tag{14}$$

$$e^{\lambda h} = 1 + \lambda h + \frac{\lambda^2 h^2}{2} + 0 (|\lambda h|^3) \tag{15}$$

$$e^{\lambda h} = 1 + \lambda h + \frac{\lambda^2 h^2}{2} + \frac{\lambda^3 h^3}{6} + 0 (|\lambda h|^4) \tag{16}$$

and so on.

Let us retain  $(p+1)$  terms in the expansion of  $e^{\lambda y}$  and denote the approximation to  $e^{\lambda y}$  by  $E(\lambda h)$ . The numerical method for obtaining the approximate values  $y_n$  of  $y(t_n)$  can then be written as

$$y_{n+1} = E(\lambda h) y_n, n = 0, 1, \dots, N-1 \tag{17}$$

The truncation error (TE) of the method is defined by

$$TE = y(t_{n+1}) - y_{n+1}$$

Since  $(p+1)$  terms are retained in the expansion of  $e^{\lambda h}$ , we have

$$TE = \left( 1 + \lambda h + \dots + \frac{(\lambda h)^p}{p!} + \frac{(\lambda h)^{p+1}}{(p+1)!} e^{\theta \lambda h} \right) - \left( 1 + \lambda h + \dots + \frac{(\lambda h)^p}{p!} \right)$$



$$= \frac{(\lambda h)^{p+1}}{(p+1)!} e^{\theta \lambda h}, \quad 0 < \theta < 1.$$

The TE is of order  $p+1$ . The integer  $p$  is then called the order of the method.

We say that a numerical method is **stable** if the error at any stage, i.e.  $y_n - y(t_n) = \epsilon_n$  remains bounded as  $n \rightarrow \infty$ . Let us examine the stability of the numerical method (17). Putting  $y_{n+1} = y(t_{n+1}) + \epsilon_{n+1}$  and  $y_n = y(t_n) + \epsilon_n$  in Eqn. (17), we have

$$\begin{aligned} y(t_{n+1}) + \epsilon_{n+1} &= E(\lambda h) [y(t_n) + \epsilon_n] \\ \epsilon_{n+1} &= E(\lambda h) [y(t_n) + \epsilon_n] - y(t_{n+1}) \\ &= E(\lambda h) [y(t_n) + \epsilon_n] - e^{\lambda h} y(t_n) \text{ (using Eqn. (13))} \\ \therefore \epsilon_{n+1} &= [E(\lambda h) - e^{\lambda h}] y(t_n) + E(\lambda h) \epsilon_n \end{aligned} \quad (18)$$

We note from Eqn. (18) that the error at  $t_{n+1}$  consists of two parts. The first part  $E[\lambda h] - e^{\lambda h}$  is the **local truncation error** and can be made as small as we like by suitably determining  $E[\lambda h]$ . The second part  $|E(\lambda h)| \epsilon_n$  is the **propagation error** from the previous step  $t_n$  to  $t_{n+1}$  and will not grow if  $|E(\lambda h)| < 1$ . If  $|E(\lambda h)| < 1$ , then as  $n \rightarrow \infty$  the propagation error tends to zero and method is said to be absolutely stable. Formally we give the following definition.

**Definition :** A numerical method (17) is called **absolutely stable** if  $|E(\lambda h)| \leq 1$ .

You may also observe here that the exact value  $y(t_n)$  given by Eqn. (13) increases if  $\lambda > 0$  and decreases if  $\lambda < 0$ , with the growth factor  $e^{\lambda h}$ . The approximate value  $y_n$  given by Eqn. (17) grows or decreases with the factor  $|E(\lambda h)|$ . Thus, in order to have meaningful numerical results, it is necessary that the growth factor of the numerical method should not increase faster than the growth factor of exact solution when  $\lambda > 0$  and should decay at least as fast as the growth factor of the exact solution when  $\lambda < 0$ . Accordingly, we give here the following definition.

**Definition :** A numerical method is said to be **relatively stable** if  $|E(\lambda h)| \leq e^{\lambda h}$ ,  $\lambda > 0$ .

The polynomial approximations (14), (15) and (16) always give relatively stable methods. Let us now find when the methods  $y_{n+1} = E(\lambda h) y_n$  are absolutely stable where  $E(\lambda h)$  is given by (14), (15) or (16).

This methods are given by

**First order :**  $y_{n+1} = (1 + \lambda h) y_n$

**Second order :**  $y_{n+1} = \left( 1 + \lambda h + \frac{\lambda^2 h^2}{2} \right) y_n$

**Third order :**  $y_{n+1} = \left( 1 + \lambda h + \frac{\lambda^2 h^2}{2} + \frac{\lambda^3 h^3}{6} \right) y_n$

These methods are absolutely stable when

**First order :**  $|1 + \lambda h| \leq 1$

or  $-1 \leq \lambda h \leq 1$

or  $-2 \leq \lambda h \leq 2$

**Second order :**  $\left| 1 + \lambda h + \frac{\lambda^2 h^2}{2} \right| \leq 1$

or  $-1 \leq 1 + \lambda h + \frac{\lambda^2 h^2}{2} \leq 1$

The right inequality gives

$$\lambda h \left( 1 + \frac{\lambda h}{2} \right) \leq 0$$

i.e.,  $\lambda h \leq 0$  and  $1 + \frac{\lambda h}{2} \geq 0$ .

The second condition gives  $-2 \leq \lambda h$ . Hence the right inequality gives  $-2 \leq \lambda h \leq 0$ . The left inequality gives

$$2 + \lambda h + \frac{\lambda^2 h^2}{2} \geq 0.$$

For  $-2 \leq \lambda h \leq 0$ , this equation is always satisfied. Hence the stability condition is

$$-2 \leq \lambda h \leq 0$$

Third order :  $\left| 1 + \lambda h + \frac{\lambda^2 h^2}{2} + \frac{\lambda^3 h^3}{6} \right| \leq 1$

Using the right and left inequalities, we get

$$-2.5 \leq \lambda h \leq 0.$$

These intervals for  $\lambda h$  are known as stability intervals.

Numerical methods for finding the solution of IVP given by Eqn. (10) may be broadly classified as

- i) Singlestep methods
- ii) Multistep methods

**Singlestep methods** enable us to find  $y_{n+1}$ , an approximation to  $y(t_{n+1})$ , if  $y_n, y_n'$  and  $h$  are known.

**Multistep methods** enable us to find  $y_{n+1}$ , an approximation to  $y(t_{n+1})$ , if  $y_i, y_i', i = n, n-1, \dots, n-m+1$  and  $h$  are known. Such methods are called  $m$ -step multistep methods.

In this course we shall be discussing about the singlestep methods only.

A singlestep method for the solution of the IVP

$$y' = f(t, y), y(t_0) = y_0, t \in (t_0, b)$$

is a recurrence relation of the form

$$y_{n+1} = y_n + h \phi(t_n, y_n, h) \tag{19}$$

where  $\phi(t_n, y_n, h)$  is known as the increment function

If  $y_{n+1}$  can be determined from Eqn. (19) by evaluating the right hand side, then the singlestep method is known as an **explicit method**, otherwise it is known as an **implicit method**. The local truncation error of the method (19) is defined by

$$TE = y(t_{n+1}) - y(t_n) - h \phi(t_n, y_n, h) \tag{20}$$

The largest integer  $p$  such that

$$|h^{-1} TE| = O(h^p) \tag{21}$$

is called the order of the singlestep method.

Let us now take up an example to understand how the singlestep method works.

**Example 1 :** find the solution of the IVP  $y' = \lambda y, y(0) = 1$  in  $0 < t \leq 0.5$ , using the first order method

$$y_{n+1} = (1 + \lambda h) y_n \text{ with } h = 0.1 \text{ and } \lambda = \pm 1.$$

**Solution :** Here the number of intervals are  $N = \frac{0.5}{h} = \frac{0.5}{0.1} = 5$

We have  $y_0 = 1$

$$y_1 = (1 + \lambda h) y_0 = (1 + \lambda h) = (1 + 0.1\lambda)$$

$$y_2 = (1 + \lambda h) y_1 = (1 + \lambda h)^2 = (1 + 0.1\lambda)^2$$

---


$$y_5 = (1 + \lambda h)^5 = (1 + 0.1\lambda)^5$$

The exact solution is  $y(t) = e^{\lambda t}$ .

We now give in Table 1 the values of  $y_n$  for  $\lambda = \pm 1$  together with exact values.

**Table 1**

Solution of $y' = \lambda y, y(0) = 1, 0 \leq t \leq 0.5$ with $h = 0.1$				
$\lambda = 1$			$\lambda = -1$	
t	First Order method	Exact Solution	First Order method	Exact Solution
0	1	1	1	1
0.1	1.1	1.10517	0.9	0.90484
0.2	1.21000	1.22140	0.81	0.81873
0.3	1.33100	1.34986	0.729	0.74082
0.4	1.46410	1.49182	0.6561	0.67032
0.5	1.61051	1.64872	0.59049	0.60653

In the same way you can obtain the solution using the second order method and compare the results obtained in the two cases.

E1) Find the solution of the IVP

$$y' = \lambda y, y(0) = 1$$

in  $0 \leq t \leq 0.5$  using the second order method

$$y_{n+1} = \left( 1 + \lambda h + \frac{\lambda^2 h^2}{2} \right) y_n \text{ with } h = 0.1 \text{ and } \lambda = 1.$$

We are now prepared to consider numerical methods for integrating differential equations. The first method we discuss is the Taylor series method. It is not strictly a numerical method, but it is the most fundamental method to which every numerical method must compare.

### 14.3 TAYLOR SERIES METHOD

Let us consider the IVP given by Eqn. (10), i.e.,

$$y' = f(t, y), y(t_0) = y_0, \quad t \in [t_0, b]$$

The function  $f$  may be linear or nonlinear, but we assume that  $f$  is sufficiently differentiable w.r.t. both  $t$  and  $y$ .

The Taylor series expansion of  $y(t)$  about any point  $t_k$  is given by

$$y(t) = y(t_k) + (t-t_k) y'(t_k) + \frac{(t-t_k)^2}{2!} y''(t_k) + \dots + \frac{(t-t_k)^p}{p!} y^{(p)}(t_k) + \dots \quad (22)$$

Substituting  $t = t_{k+1}$  in Eqn. (22), we have

$$y(t_{k+1}) = y(t_k) + h y'(t_k) + \frac{h^2 y''(t_k)}{2!} + \dots + \frac{h^p y^{(p)}(t_k)}{p!} + \dots \quad (23)$$

where  $t_{k+1} = t_k + h$ . Neglecting the terms of order  $h^{p+1}$  and higher order terms, we have the approximation

$$y_{k+1} = y_k + hy'_k + \frac{h^2}{2!} y''_k + \dots + \frac{h^p}{p!} y_k^{(p)}$$

$$= y_k + h \phi(t_k, y_k, h) \tag{24}$$

where  $\phi(t_k, y_k, h) = y'_k + \frac{h}{2!} y''_k + \dots + \frac{h^{p-1}}{p!} y_k^{(p)}$

This is called the Taylor Series method of order  $p$ . The truncation error of the method is given by

$$TE = y(t_{k+1}) - y(t_k) - h\phi(t_k, y(t_k), h)$$

$$= \frac{h^{p+1}}{(p+1)!} y^{(p+1)}(t_k + \theta h), 0 < \theta < 1 \tag{25}$$

When  $p = 1$ , we get from Eqn. (24)

$$y_{k+1} = y_k + hy'_k \tag{26}$$

which is the Taylor series method of order one.

To apply (24), we must know  $y(t_k), y'(t_k), y''(t_k), \dots, y^{(p)}(t_k)$ .

However,  $y(t_k)$  is known to us and if  $f$  is sufficiently differentiable, then higher order derivatives can be obtained by calculating the total derivative of the given differential equation w.r.t.  $t$ , keeping in mind that  $y$  is itself a function of  $t$ . Thus we obtain for the first few derivatives as :

$$y' = f(t, y)$$

$$y'' = f_t + f_y y'$$

$$y''' = f_{tt} + 2f_{ty} y' + f_{yy} (y')^2 + f_y (f_t + f_y y')$$
 etc.

where  $f_t = \partial f / \partial t, f_{tt} = \partial^2 f / \partial t^2$  etc.

The number of terms to be included in the method depends on the accuracy requirements.

Let  $p = 2$ . Then the Taylor Series method of  $O(h^2)$  is

$$y_{k+1} = y_k + hy'_k + \frac{h^2}{2} y''_k \tag{27}$$

with the TE =  $\frac{h^3}{6} y'''(\alpha), t_n < \alpha < t_{n+1}$

The Taylor series method of  $O(h^3), (p=3)$  is

$$y_{k+1} = y_k + hy'_k + \frac{h^2}{2} y''_k + \frac{h^3}{6} y'''_k \tag{28}$$

with the TE =  $\frac{h^4}{24} y^{(IV)}(\alpha), t_n < \alpha < t_{n+1}$ .

Let us consider the following examples.

**Example 2 :** Using the third order Taylor series method find the solution of the differential equation

$$xy' = x - y, y(2) = 2 \text{ at } x = 2.1 \text{ taking } h=0.1$$

**Solution :** We have the derivatives and their values at  $x=2, y=2$  as follows :

$$y' = 1 - \frac{y}{x} \qquad y'(2) = 0$$

$$y'' = -\frac{y'}{x} + \frac{y}{x^2} \qquad y''(2) = 1/2$$

$$y''' = \frac{-y''}{x} + \frac{2y'}{x^2} - \frac{2y}{x^3} \qquad y'''(2) = -3/4$$

Using Taylor series method of  $O(h^3)$  given by Eqn. (28), we obtain

$$y(2.1) = 2 + 0.0025 - 0.000125 = 2.002375.$$

**Example 3 :** Solve the equation  $x^2y' = 1 - xy - x^2y^2$ ,  $y(1) = -1$  from  $x=1$  to  $x=2$  by using Taylor series method of  $O(h^2)$  with  $h = 1/3$  and  $1/4$  and find the actual error at  $x=2$  if the exact solution is  $y = -1/x$ .

**Solution :** From the given equation, we have  $y' = \frac{1}{x^2} - \frac{y}{x} - y^2$

Differentiating it w.r.t.  $x$ , we get

$$y'' = \frac{-2}{x^3} - \frac{y'}{x} + \frac{y}{x^2} - 2yy'$$

Using the second order method (27),

$$y_{k+1} = y_k + hy'_k + \frac{h^2}{2} y''_k$$

we have the following results

$$y(1) = -1, \quad y'(1) = 1, \quad y''(1) = -2$$

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

$$x_1 = \frac{4}{3}, \quad y(x_1) = -0.7778, \quad y'(x_1) = 0.5409, \quad y''(x_1) = -0.8455$$

$$x_2 = \frac{5}{3}, \quad y(x_2) = -0.6445, \quad y'(x_2) = 0.3313, \quad y''(x_2) = -0.4358$$

$$x_3 = 2, \quad y(x_3) = -0.5583 = y(2)$$

$$h = \frac{1}{4}$$

$$x_1 = \frac{5}{4}, \quad y(x_1) = -0.8125, \quad y'(x_1) = 0.6298, \quad y''(x_1) = -1.0244$$

$$x_2 = \frac{3}{2}, \quad y(x_2) = -0.6871, \quad y'(x_2) = 0.4304, \quad y''(x_2) = -0.5934$$

$$x_3 = \frac{7}{4}, \quad y(x_3) = -0.5980, \quad y'(x_3) = 0.3106, \quad y''(x_3) = -0.3745$$

$$x_4 = 2, \quad y(x_4) = -0.5321 = y(2)$$

Since the exact value is  $y(2) = -0.5$ , we have the actual errors as

$$e_1 = 0.0583 \text{ with } h = \frac{1}{3}$$

$$e_2 = 0.0321 \text{ with } h = \frac{1}{4}$$

Note that error is small when the step size  $h$  is small.

You may now try the following exercises:

Write the Taylor series method of order four and solve the IVPs E2) and E3).

E2)  $y' = x - y^2$ ,  $y(0) = 1$ . Find  $y(0.1)$  taking  $h = 0.1$ .

E3)  $y' = x^2 + y^2, y(0) = 0.5$ . Find  $y(0.4)$  taking  $h = 0.2$ .

E4) Using second order Taylor series method solve the IVP

$y' = 3x + \frac{y}{2}, y(0) = 1$ . Find  $y(0.6)$  taking  $h = 0.2$  and  $h = 0.1$ .

Find the actual error at  $x = 0.6$  if the exact solution is  $y = -6x - 12$ .

Notice that though the Taylor series method of order  $p$  gives us results of desired accuracy in a few number of steps, it requires evaluation of the higher order derivatives and becomes tedious to apply if the various derivatives are complicated. Also, it is difficult to determine the error in such cases. We now consider a method, the Euler's method which can be regarded as Taylor series method of order one and avoids these difficulties.

## 14.4 EULER'S METHOD

Let the given IVP be

$$y' = f(t, y), y(t_0) = y_0$$

Let  $[t_0, b]$  be the interval over which the solution of the given IVP is to be

determined. Let  $h$  be the steplength. Then the nodal points are defined by  $t_k = t_0 + kh$ ,  $k = 0, 1, 2, \dots, N$  with  $t_N = t_0 + Nh = b$ .

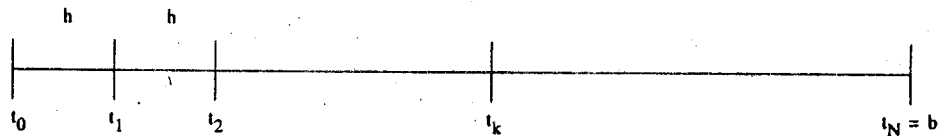


Fig. 1

The exact solution  $y(t)$  at  $t = t_{k+1}$  can be written by Taylor series as

$$y(t_k + h) = y(t_k) + hy'(t_k) + \left(\frac{h^2}{2}\right)y''(t_k) + \dots \quad (29)$$

Neglecting the term of  $O(h^2)$  and higher order terms, we get

$$y_{k+1} = y_k + hy'_k \quad (30)$$

with  $TE = \left(\frac{h^2}{2}\right)y''(\alpha), t_k < \alpha < t_{k+1} \quad (31)$

From the given IVP,  $y'(t_k) = f(t_k, y_k) = f_k$

We can rewrite Eqn. (30) as

$$y_{k+1} = y_k + h f_k \quad (32)$$

for  $k = 0, 1, \dots, N$

Eqn. (32) is known as the Euler's method and it calculates recursively the solution at the nodal points  $t_k, k = 0, 1, \dots, N$ .

Since the truncation error (31) is of order  $h^2$ , Euler's method is of first order. It is also called an  $O(h)$  method.

Let us now see the geometrical representation of the Euler's method.

### Geometrical Interpretation

Let  $y(t)$  be the solution of the given IVP, Integrating  $\frac{dy}{dt} = f(t, y)$  from  $t_k$  to  $t_{k+1}$ , we get

$$\int_{t_k}^{t_{k+1}} \frac{dy}{dt} dt = \int_{t_k}^{t_{k+1}} f(t, y) dt = y(t_{k+1}) - y(t_k) \quad (33)$$

We know that geometrically  $f(t, y)$  represents the slope of the curve  $y(t)$ . Let us approximate the slope of the curve between  $t_k$  and  $t_{k+1}$  by the slope at  $t_k$  only. If we approximate  $y(t_{k+1})$  and  $y(t_k)$  by  $y_{k+1}$  and  $y_k$  respectively, then we have

$$\begin{aligned} y_{k+1} - y_k &= f(t_k, y_k) \int_{t_k}^{t_{k+1}} dt \\ &= (t_{k+1} - t_k) f(t_k, y_k) \\ &= hf(t_k, y_k) \end{aligned} \quad (34)$$

$$\therefore y_{k+1} = y_k + hf(t_k, y_k), \quad k = 0, 1, 2, \dots, N.$$

Thus in Euler's method the actual curve is approximated by a sequence of line segments and the area under the curve is approximated by the area of the quadrilateral. (see Fig.3)

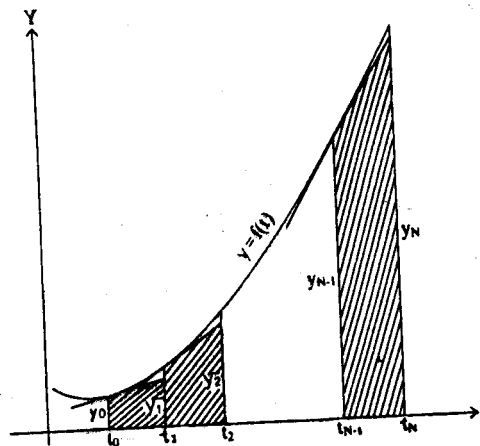


Fig. 2 : Geometrical representation of Euler's method.

Let us now consider the following examples.

**Example 4 :** Use Euler method to find the solution of  $y' = t + |y|$ , given  $y(0) = 1$ . Find the solution on  $[0, 0.8]$  with  $h = 0.2$ .

**Solution :** We have

$$\begin{aligned} y_{n+1} &= y_n + hf_n \\ y(0.2) \approx y_1 &= y_0 + (0.2) f_0 \\ &= 1 + (0.2) [0+1] = 1.2 \\ y(0.4) \approx y_2 &= y_1 + (0.2) f_1 \\ &= 1.2 + (0.2) [0.2 + 1.2] \\ &= 1.48 \\ y(0.6) \approx y_3 &= y_2 + (0.2) f_2 \\ &= 1.48 + (0.2) [0.4 + 1.48] \\ &= 1.856 \end{aligned}$$

$$\begin{aligned}
 y(0.8) \approx y_4 &= y_3 + (0.2) f_3 \\
 &= 1.856 + (0.2) [0.6 + 1.856] \\
 &= 2.3472
 \end{aligned}$$

**Example 5 :** Solve the differential equation  $y' = t+y$ ,  $y(0) = 1$ ,  $t \in [0,1]$  by Euler's method using  $h = 0.1$ . If the exact value is  $y(1) = 3.436564$ , find the exact error.

**Solution :** Euler's method is

$$y_{n+1} = y_n + hy'_n$$

For the given problem, we have

$$y_{n+1} = y_n + h [t_n + y_n]$$

$$= (1 + h)y_n + ht_n$$

$$h = 0.1, y(0) = 1,$$

$$y_1 = y_0 = (1 + 0.1) + (0.1)(0) = 1.1$$

$$y_2 = (1.1)(1.1) + (0.1)(0.1) = 1.22, y_3 = 1.362$$

$$y_4 = 1.5282, y_5 = 1.72102, y_6 = 1.943122,$$

$$y_7 = 2.197434, y_8 = 2.487178, y_9 = 2.815895$$

$$y_{10} = 3.187485 \approx y(1)$$

$$\text{actual error} = y(1) - y_{10} = 3.436564 - 3.187485 = 0.2491.$$

**Remark :** Since Euler's method is of  $O(h)$ , it requires  $h$  to be very small to attain the desired accuracy. Hence, very often, the number of steps to be carried out becomes very large. In such cases, we need higher order methods to obtain the required accuracy in a limited number of steps.

Euler's method construct  $y_k \approx y(t_k)$  for each  $k = 1, 2, \dots, N$ ,

where

$$y_{n+1} = y_k + hf(t_k, y_k)$$

This equation is called the **difference equation** associated with Euler's method. A difference equation of order  $N$  is a relation involving  $y_n, y_{n+1}, \dots, y_{n+N}$  some simple difference equations are

$$\left. \begin{aligned}
 y_{n+1} - y_n &= 1 \\
 y_{n+1} - y_n &= n \\
 y_{n+1} - (n+1)y_n &= 0
 \end{aligned} \right\} \quad (35)$$

where  $n$  is an integer.

A difference equation is said to be **linear** if the unknown functions  $y_{n+k}$  ( $k = 0, 1, \dots, N$ ) appear linearly in the difference equation. The general form of a linear nonhomogeneous difference equation of order  $N$  is

$$y_{n+N} + a_{N-1}y_{n+N-1} + \dots + a_0y_n = b \quad (36)$$

where the coefficients  $a_{N-1}, a_{N-2}, \dots, a_0$  and  $b$  may be functions of  $n$  but not of  $y$ . All the Eqns. (35) are linear. It is easy to solve the difference Eqn. (36), when the coefficients are constant or a function of  $n$  say linear or a quadratic function of  $n$ .

The general solution of Eqn. (36) can be written in the form

$$y_n = y_n(c) + y_n(p),$$

where  $y_n(c)$  is the complementary solution of the homogeneous equation associated with Eqn. (36) and  $y_n(p)$  is a particular solution of Eqn. (36). To obtain the complementary solution of the homogeneous equations, we start with a solution in the form  $y_n = \beta^n$  and substitute it in the given equation. This gives us a polynomial of degree  $N$ . We assume that its roots  $\beta_1, \beta_2, \dots, \beta_N$  are all real and distinct.



Then by linearity it follows that

$$y_n = C_1 \beta_1^n + C_2 \beta_2^n + \dots + C_N \beta_N^n$$

for arbitrary constants  $C_i$ , is a solution of the homogeneous equation associated with Eqn. (36). A particular solution of Eqn. (36) when  $b$  is a constant can be obtained by setting  $y_n(p) = A$  (a constant) in Eqn. (36) and determining the value of  $A$ . For detail, you can refer to elementary numerical analysis by Conte- deBoor. We illustrate this method by considering a few examples.

**Example 6 :** Find the solution of the initial-value difference equations

$$y_{n+2} - 4y_{n+1} + 3y_n = 2^n, y_0 = 0, y_1 = 1$$

**Solution :** The homogeneous equation of the given problem is

$$y_{n+2} - 4y_{n+1} + 3y_n = 0 \tag{37}$$

Let  $y_n = \beta^n$ . Then Eqn. (37) reduces to

$$\beta^{n+2} - 4\beta^{n+1} + 3\beta^n = 0.$$

Dividing by  $\beta^n$ , we obtain the characteristic equation

$$\beta^2 - 4\beta + 3 = 0$$

i.e.,  $\beta = 1, 3$

$$\begin{aligned} \therefore y_n(C) &= C_1 (1)^n + C_2 (3)^n \\ &= C_1 + 3^n C_2 \end{aligned} \tag{38}$$

For obtaining the particular solution we try  $y_n(p) = A2^n$ .

This gives

$$2^{n+2}A - 4 \times 2^{n+1}A + 3 \times 2^n A = 2^n$$

or,  $A = -1$

Therefore, the general solution of the given problem is

$$y_n = C_1 + 3^n C_2 - 1$$

Using conditions for  $n = 0, 1$ , we obtain

$$C_1 + C_2 = 1$$

$$C_1 + 3C_2 = 2$$

$$\therefore C_2 = 1/2, C_1 = 1/2 \text{ and}$$

$$y_n = \frac{3^n - 1}{2} \tag{39}$$

which is the required solution.

**Note :** In the above method we can obtain  $y_n$  for all  $n$  from one formula given by Eqn. (39). Whereas, in the Euler's method for obtaining the value at each iteration, we require the previous iterated value. We illustrate it by considering another example.

**Example 7 :** Using difference method find the solution of  $y_{k+1} = y_k + h(5+3y_k)$ , given  $y(0) = 1$ . Find the solution  $y(0.6)$  with  $h = 0.1$ .

**Solution :** We have

$$y_{k+1} - (1 + 3h) y_k = 5h \tag{40}$$

Solution of the homogeneous equation is

$$y_k(c) = C(1 + 3h)^k.$$

For obtaining the particular solution we try  $y_k(p) = Ah$ .

This give

$$A = -\frac{5}{3}.$$

Therefore, the general solution of the given problem is

$$y_k = C(1 + 3h)^k - \frac{5}{3}.$$

Using the condition  $y(0) = 1$ , we obtain  $C = 8/3$ .

Thus

$$y_k = \frac{8}{3}(1 + 3h)^k - \frac{5}{3}. \quad (41)$$

Eqn. (41) gives the formula for obtaining  $y_k$   $\forall$   $k$ .

$$\begin{aligned} y_6 = y(0.6) &= \frac{8}{3}(1 + 3 \times 0.1)^6 - \frac{5}{3} \\ &= 11.204824. \end{aligned}$$

Now Euler's method is

$$y_{k+1} = (1 + 3h)y_k + 5h$$

and we get for  $h = 0.1$

$$y_1 = 1.8, y_2 = 2.84, y_3 = 4.192, y_4 = 5.9496, y_5 = 8.23448, y_6 = 11.204824.$$

You may now try the following exercises

Solve the following IVPs using Euler's method

E5)  $y' = 1 - 2xy, y(0.2) = 0.1948$ . Find  $y(0.4)$  with  $h = 0.2$

E6)  $y' = \frac{1}{x^2 - 4y}, y(4) = 4$ . Find  $y(4.1)$  taking  $h = 0.1$

E7)  $y' = \frac{y-x}{y+x}, y(0) = 1$ . Find  $y(0.1)$  with  $h = 0.1$

E8)  $y' = 1 + y^2, y(0) = 1$ . Find  $y(0.6)$  taking  $h = 0.2$  and  $h = 0.1$ .

You may recall that in Unit 12 we studied Richardson's extrapolation technique to increase the order of a numerical differentiation formula without increasing the function evaluations. In Unit 13, we introduced Romberg integration which is the Richardson's extrapolation technique applied to the integration rules. In both the cases, the order of the numerical value was improved by the application of the Richardson's extrapolation. In the next section we shall use this technique to obtain higher order solutions to differential equations using lower order methods.

## 14.5 RICHARDSON'S EXTRAPOLATION

Consider the Euler's method

$$y_{k+1} = y_k + hf'_k$$

which is an  $O(h)$  method. Let  $F(h)$  and  $F(h/2)$  be the solutions obtained by using step lengths  $h$  and  $h/2$  respectively.

Recall that the Richardson's extrapolation method of combining two computed values with two different step sizes, to obtain a higher order method is given by (ref. Formula (54) Unit 12)

$$F^{(1)}(h) = \frac{F(2h) - r^p F(h)}{1 - r^p} \quad (42)$$

where  $p$  is the order of the method.

Thus, in the case of Euler's method which is of first order, once we know the values  $F(h)$  and  $F(h/2)$  at two different step sizes  $h$  and  $h/2$ , Formula (42) for  $r = 1/2$ ,  $p=1$ , reduces to

$$F^{(1)}(h) = \frac{2F(h/2) - F(h)}{2 - 1} \quad (43)$$

Before illustrating this technique, we give you a method of determining numerically, the order of a method.

Let  $y_1(t_k)$  and  $y_2(t_k)$  be the two values obtained by a numerical method of order  $p$  with step sizes  $h_1$  and  $h_2$ . If  $e_1$  and  $e_2$  are the corresponding errors, then

$$\frac{e_1}{e_2} = \frac{Ch_1^p}{Ch_2^p} = \left(\frac{h_1}{h_2}\right)^p$$

By taking logarithms, we get

$$p \ln \left(\frac{h_1}{h_2}\right) = \ln \left(\frac{e_1}{e_2}\right)$$

Hence the order  $p$  of the method is

$$p = \frac{\ln(e_1/e_2)}{\ln(h_1/h_2)}$$

Let us now consider the following examples.

**Example 6 :** Using the Euler's method tabulate the solution of the IVP

$$y' = -2t y^2, y(0) = 1$$

in the interval  $[0, 1]$  taking  $h = 0.2, 0.1$ . Using Richardson's extrapolation technique obtain the improved value at  $t = 1$ .

**Solution :** Euler's method gives

$$\begin{aligned} y_{k+1} &= y_k + h f_k \text{ where } f_k = -2t_k y_k^2 \\ &= y_k - 2h t_k y_k^2. \end{aligned}$$

Starting with  $t_0 = 0, y_0 = 1$ , we obtain the following table of values for  $h = 0.2$ .

**Table 2 :  $h = 0.2$**

t	y(t)
0.2	1
0.4	0.92
0.6	0.78458
0.8	0.63684
1.0	0.50706

Thus,  $y(1.0) = 0.50706$  with  $h = 0.2$

Similarly, starting with  $t_0 = 0, y_0 = 1$ , we obtain the following table of values for  $h = 0.1$ .

**Table 3 :  $h = 0.1$**

t	y(t)	t	y(t)
0.1	1.0	0.6	0.75715
0.2	0.98	0.7	0.68835
0.3	0.94158	0.8	0.62202
0.4	0.88839	0.9	0.56011
0.5	0.82525	1.0	0.50364

$y(1.0) = 0.50364$  with  $h = 0.1$

Using formula (43), the extrapolated value at  $y(1)$  is given by

$$\begin{aligned} F^{(1)}(0.1) &= \frac{2F(0.1) - F(0.2)}{1} \\ &= 2(0.50364) - (0.50706) \\ &= 0.50022 \end{aligned}$$

Let us consider another example

**Example 7 :** Use Euler's method to solve numerically the initial value problem  $y' = t + y$ ,  $y(0) = 1$  with  $h = 0.2, 0.1$  and  $0.05$  in the interval  $[0, 0.6]$ . Apply Richardson's extrapolation technique to compute  $y(0.6)$ .

**Solution :** Euler's method gives

$$\begin{aligned} y_{k+1} &= y_k + h f_k \\ &= y_k + h (t_k + y_k) \\ &= (1+h) y_k + h t_k \end{aligned}$$

Starting with  $t_0 = 0, y_0 = 1$ , we obtain the following table of values.

**Table 4 :  $h = 0.2$**

t	y(t)
0.2	1.2
0.4	1.48
0.6	1.856

$\therefore y(0.6) = 1.856$  with  $h = 0.2$

**Table 5 :  $h = 0.1$**

t	y(t)
0.1	1.1
0.2	1.22
0.3	1.362
0.4	1.5282
0.5	1.72102
0.6	1.943122

$\therefore y(0.6) = 1.943122$  with  $h = 0.1$

**Table 6 :  $h = 0.05$**

t	y(t)	t	y(t)
0.05	1.05	0.35	1.46420
0.1	1.105	0.4	1.55491
0.15	1.16525	0.45	1.65266
0.2	1.23101	0.5	1.75779
0.25	1.30256	0.55	1.87068
0.3	1.38019	0.6	1.99171

$\therefore y(0.6) = 1.99171$  with  $h = 0.05$

By Richardson's extrapolation method (43), we have

$$F^{(1)}(0.05) = 2F(0.05) - F(0.1) \\ = 2.040298$$

$$F^{(1)}(0.1) = 2F(0.1) - F(0.2) \\ = 2.030244$$

Repeating Richardson's technique and using formula (42) with  $p=2$ , we obtain

$$F^{(2)}(0.05) = \frac{(2)^2 F^{(1)}(0.05) - F^{(1)}(0.1)}{(2)^2 - 1} \\ = \frac{4(2.040298) - 2.030244}{3} \\ = 2.043649$$

The exact solution is  $y = -(1+t) + 2e^t$

Hence  $y(0.6) = 2.044238$ .

The actual error of the extrapolated value is

$$\text{error} = y(0.6) - F^{(2)}(0.05) \\ = 2.044238 - 2.043649 \\ = 0.000589$$

And now a few exercises for you

E9) The IVP

$$y' = 3t + \frac{y}{2}, y(0) = 1.$$

is given. Find  $y(0.6)$  with  $h = 0.2$  and  $h = 0.1$ , using Euler's method and extrapolate the value  $y(0.6)$ . Compare with the exact solution.

E10) Extrapolate the value  $y(0.6)$  obtained in E8).

We now end this unit by giving a summary of what we have covered in it.

## 14.6 SUMMARY

In this unit, we have covered the following

1) Taylor series method of order  $p$  for the solution of the IVP

$$y' = f(t, y), y(t_0) = y_0, t \in [t_0, b] \text{ (see Eqn. (10))}$$

is given by

$$y_{k+1} = y_k + h \phi [t_k, y_k, h]$$

where  $\phi [t_k, y_k, h] = y'_k + \frac{h}{2!} y''_k + \dots + \frac{h^{p-1}}{p!} y^{(p)}_k$  and  $t_k = t_0 + kh, k = 0, 1,$

$2, \dots, N, t_N = b$ . The error of approximation is given by

$$TE = \frac{h^{p+1}}{(p+1)!} y^{(p+1)}(t_k + \theta h), 0 < \theta < 1.$$

2) Euler's method is the Taylor series method of order one. The steps involved in solving the IVP given by (10) by Euler's method are as follows :

**Step 1 :** Evaluate  $f(t_0, y_0)$

**Step 2 :** Find  $y_1 = y_0 + h f(t_0, y_0)$

**Step 3 :** If  $t_0 < b$ , change  $t_0$  to  $t_0 + h$  and  $y_0$  to  $y_1$  and repeat steps 1 and 2

Step 4: If  $t_0 = b$ , write the value of  $y_1$ .

3) Richardson's extrapolation method given by Eqn. (42) can be used to improve the values of the function evaluated by the Euler's method.

## 14.7 SOLUTIONS/ANSWERS

E1) We have  $y_0 = 1, \lambda = 1, h = 0.1$

$$y_1 = \left( 1 + 0.1 + \frac{(0.1)^2}{2} \right)$$

$$y_2 = (1.105)^5$$

$$y_5 = (1.105)^5$$

Table giving the values of  $y_n$  together with exact values is

Table 7

t	Second order method	Exact solution
0	1	1
0.1	1.105	1.10517
0.2	1.22103	1.22140
0.3	1.34923	1.34986
0.4	1.49090	1.49182
0.5	1.64745	1.64872

E2) Taylor series method of  $O(h^4)$  to solve  $y' = x - y^2, y(0) = 1$  is

$$y_{n+1} = y_n + hy'_n + \frac{h^2}{2} y''_n + \frac{h^3}{6} y'''_n + \frac{h^4}{24} y^{iv}_n$$

$$y' = x - y^2$$

$$y'(0) = -1$$

$$y'' = 1 - 2yy'$$

$$y''(0) = 3$$

$$y''' = -2yy'' - 2(y')^2$$

$$y'''(0) = -8$$

$$y^{iv} = -2yy''' - 6y'y''$$

$$y^{iv}(0) = 34$$

Substituting

$$y(0.1) = 1 - (0.1)(-1) + \frac{(0.1)^2}{2}(3) + \frac{(0.1)^3}{6}(-8) + \frac{(0.1)^4}{24}(34) = 0.9138083.$$

E3) Taylor series method :

$$y' = x^2 + y^2,$$

$$y(0) = 0.5,$$

$$y'(0) = 0.25,$$

$$y'(0.2) = 0.35175$$

$$y'' = 2x + 2yy'$$

$$y''(0) = 0.25,$$

$$y''(0.2) = 0.79280$$

$$y''' = 2 + 2yy'' + 2(y')^2$$

$$y'''(0) = 2.375,$$

$$y'''(0.2) = 3.13278$$

$$y^{iv} = 2yy''' + 6y'y''$$

$$y^{iv}(0) = 2.75,$$

$$y^{iv}(0.2) = 5.17158$$

$$y(0.2) = 0.55835,$$

$$y(0.4) = 0.64908$$

E4) Second order Taylor's method is

$$y_{n+1} = y_n + hy'_n + \frac{h^2}{2} y''_n.$$

$h = 0.2$

$$y(0) = 1, \quad y'(0) = 0.5, \quad y''(0) = 3.25$$

$$y(0.2) = 1.165, \quad y'(0.2) = 1.1825, \quad y''(0.2) = 3.59125$$

$$y(0.4) = 1.47333, \quad y'(0.4) = 1.93667, \quad y''(0.4) = 3.96833$$

$$y(0.6) = 1.94003$$

$$h = 0.1$$

$$y(0.1) = 1.06625, \quad y'(0.1) = 0.83313, \quad y''(0.1) = 3.41656$$

$$y(0.2) = 1.16665, \quad y'(0.2) = 1.18332, \quad y''(0.2) = 3.59167$$

$$y(0.3) = 1.46457, \quad y'(0.3) = 1.63228, \quad y''(0.3) = 3.81614$$

$$y(0.4) = 1.64688, \quad y'(0.4) = 2.02344, \quad y''(0.4) = 4.01172$$

$$y(0.5) = 1.86928, \quad y'(0.5) = 2.43464, \quad y''(0.5) = 4.21732$$

$$y(0.6) = 2.13383$$

E5) Euler's method is  $y_{k+1} = y_k + hf_k = y_k + h(1 - 2x_k y_k)$

$$y(0.4) = 0.1948 + (0.2)(1 - 2 \times 0.2 \times 0.1948) \\ = 0.379216.$$

E6)  $y' = \frac{1}{x^2 + y}, \quad y(4) = 4, \quad y'(4) = 0.05$

$$y(4.1) = y(4) + hy'(4) \\ = 4 + (0.1)(0.05) = 4.005.$$

E7) Euler's method  $y' = (y-x)/(y+x), y(0) = 1, y'(0) = 1$

$$y(0.1) = 1 + (0.1)(1) = 1.1$$

E8) Euler's method is

$$y_{k+1} = h + y_k + hy_k^2$$

Starting with  $t_0 = 0$  and  $y_0 = 1$ , we have the following tables of values

**Table 8 : h = 0.2**

t	y(t)
0.2	1.4
0.4	1.992
0.6	2.9856

$$\therefore y(0.6) = 2.9856$$

**Table 9 : h = 0.1**

t	y(t)
0.1	1.2
0.2	1.444
0.3	1.7525
0.4	2.1596
0.5	2.7260
0.6	3.5691

$$\therefore y(0.6) = 3.5691$$

E9)  $y_{k+1} = \left(1 + \frac{h}{2}\right)y_k + 3ht_k$

Starting with  $t_0 = 0, y_0 = 1$ , we have the following table of values

**Table 10 : h = 0.2**

t	y(t)
0.2	1.1
0.4	1.33
0.6	1.703

$$\therefore y(0.6) = 1.703$$

Table 11 :  $h = 0.1$

t	y(t)
0.1	1.05
0.2	1.1325
0.3	1.2491
0.4	1.40156
0.5	1.59164
0.6	1.82122

$$\therefore y(0.6) = 1.82122$$

Using formula (42), we have

$$\begin{aligned} F^{(1)}(0.1) &= 2F(0.1) - F(0.2) \\ &= 1.93944 \end{aligned}$$

Exact solution is

$$y = -6(t+2) + 13e^{\frac{1}{2}t}$$

Hence

$$y(0.6) = 1.948164$$

The actual error of the extrapolated value is

$$\begin{aligned} \text{error} &= y(0.6) - F^{(1)}(0.1) \\ &= 1.948164 - 1.93944 \\ &= 0.008724 \end{aligned}$$

E10) From E8), we have  $F(0.1) = 3.5691$  and  $F(0.2) = 2.9856$

$$\therefore F^{(1)}(0.1) = 4.1526$$



---

# UNIT 15 SOLUTION OF ORDINARY DIFFERENTIAL EQUATIONS USING RUNGE-KUTTA METHODS

---

## Structure

- 15.1 Introduction
  - Objectives
- 15.2 Runge-Kutta Methods
  - Runge-Kutta Methods of Second Order
  - Runge-Kutta Methods of Third Order
  - Runge-Kutta Methods of Fourth Order
- 15.3 Richardson's Extrapolation
- 15.4 Summary
- 15.5 Solutions/Answers

---

## 15.1 INTRODUCTION

---

In Unit 14, we considered the IVPs

$$y' = f(t, y), \quad y'(t_0) = y_0 \quad (1)$$

and developed Taylor series method and Euler's method for its solution. As mentioned earlier, Euler's method being a first order method, requires a very small step size for reasonable accuracy and therefore may require lot of computations. Higher order Taylor series methods require evaluation of higher order derivatives either manually or computationally. For complicated functions, finding second, third and higher order total derivatives is very tedious. Hence Taylor series methods of higher order are not of much practical use in finding the solution of IVPs of the form given by Eqn. (1).

In order to avoid this difficulty, at the end of nineteenth century, the German mathematician, Runge observed that the expression for the increment function  $\phi(t, y, h)$  in the singlestep methods [see Eqn. (24) of Sec. 14.3, Unit 14]

$$y_{n+1} = y_n + h \phi(t_n, y_n, h) \quad (2)$$

can be modified to avoid evaluation of higher order derivatives. This idea was further developed by Runge and Kutta (another German mathematician) and the methods given by them are known as Runge-Kutta methods. Using their ideas, we can construct higher order methods using only the function  $f(t, y)$  at selected points on each subinterval. We shall, in the next section, derive some of these methods.

### Objectives

After studying this unit, you should be able to :

- obtain the solution of IVPs using Runge-Kutta methods of second, third and fourth order;
- compare the solutions obtained by using Runge-Kutta and Taylor series methods;
- extrapolate the approximate value of the solutions obtained by the Runge-Kutta methods of second, third and fourth order.

---

## 15.2 RUNGE-KUTTA METHODS

---

We shall first try to discuss the basic idea of how the Runge- Kutta methods are developed.

Consider the  $O(h^2)$  singlestep method

$$y_{n+1} = y_n + hy'_n + \frac{h^2}{2} y''_n \quad (3)$$

If we write Eqn. (3) in the form of Eqn. (2) i.e., in terms of  $\phi [t_n, y_n, h]$  involving partial derivatives of  $f(t, y)$ , we obtain

$$\phi(t, y, h) = f(t_n, y_n) + \frac{h}{2} [f_t(t_n, y_n) + hf(t_n, y_n) f_y(t_n, y_n)] \quad (4)$$

Runge observed that the r.h.s. of Eqn. (4) can also be obtained using the Taylor series expansion of  $f(t_n + ph, y_n + qhf_n)$  as

$$f(t_n + ph, y_n + qhf_n) \approx f_n + ph f_t(t_n, y_n) + qhf_n f_y(t_n, y_n) \quad (5)$$

Comparing Eqns. (4) and (5) we find that  $p = q = 1/2$  and the Taylor series method of  $O(h^2)$  given by Eqn. (3) can also be written as

$$y_{n+1} = y_n + hf\left(t_n + \frac{h}{2}, y_n + \frac{h}{2} f_n\right) \quad (6)$$

Since (5) is of  $O(h^2)$ , the value of  $y_{n+1}$  in (6) has the TE of  $O(h^3)$ . Hence the method (6) is of  $O(h^2)$  which is same as that of (3).

The advantage of using (6) over Taylor series method (3) is that we need to evaluate the function  $f(t, y)$  only at two points  $(t_n, y_n)$  and  $(t_n + \frac{h}{2}, y_n + \frac{h}{2} f_n)$ . We observe that  $f(t_n, y_n)$  denotes the slope of the solution curve to the IVP (1) at  $(t_n, y_n)$ . Further,  $f\left[t_n + \frac{h}{2}, y_n + \left(\frac{h}{2} f_n\right)\right]$  denotes an approximation to the slope of the solution curve at the point  $\left[t_n + \frac{h}{2}, y\left(t_n + \frac{h}{2}\right)\right]$ . Eqn. (6) denotes geometrically, that the slope of the solution curve in the interval  $[t_n, t_{n+1}]$  is being approximated by an approximation to the slope at the middle points  $t_n + \frac{h}{2}$ . This idea can be generalised and the slope of the solution curve in  $[t_n, t_{n+1}]$  can be replaced by a weighted sum of slopes at a number of points in  $[t_n, t_{n+1}]$  (called off- step points). This idea is the basis of the Runge-Kutta methods.

Let us consider for example, the weighted sum of the slopes at the two points  $[t_n, y_n]$  and  $[t_n + ph, y_n + qhf_n]$ ,  $0 < p, q < 1$  as

$$\phi(t_n, y_n, h) = W_1 f(t_n, y_n) + W_2 f[t_n + ph, y_n + qhf_n] \quad (7)$$

We call  $W_1$  and  $W_2$  as weights and  $p$  and  $q$  as scale factors. We have to determine the four unknowns  $W_1, W_2, p$  and  $q$  such that  $\phi(t_n, y_n, h)$  is of  $O(h^2)$ . Substituting Eqn. (5) in (7), we have

$$\phi(t_n, y_n, h) = W_1 f_n + W_2 [f_n + phf_t(t_n, y_n) + qhf_n f_y(t_n, y_n)] \quad (8)$$

and the method (2) reduces to

$$\begin{aligned} y_{n+1} &= y_n + h [W_1 f_n + W_2 \{f_n + phf_t(t_n, y_n) + qhf_n f_y(t_n, y_n)\}] \\ &= y_n + h(W_1 + W_2) f_n + h^2 W_2 (pf_t + qf_n f_y)_n \end{aligned} \quad (9)$$

where  $( )_n$  denotes that the quantities inside the brackets are evaluated at  $(t_n, y_n)$ .

Comparing the r.h.s. of Eqn. (9) with Eqn. (3), we find that

$$\left. \begin{aligned} W_1 + W_2 &= 1 \\ W_2 p &= W_2 q = \frac{1}{2} \end{aligned} \right\} \quad (10)$$

In the system of Eqns. (10), since the number of unknowns is more than the number of equations, the solution is not unique and we have infinite number of solutions. The solution of Eqn. (10) can be written as

$$\begin{aligned} W_1 &= 1 - W_2 \\ p &= q = 1/(2W_2) \end{aligned} \quad (11)$$

By choosing  $W_2$  arbitrarily we may obtain infinite number of second order Runge-Kutta methods. If  $W_2 = 1$ ,  $p = q = \frac{1}{2}$  and  $W_1 = 0$ , then we get the method (6). Another choice is  $W_2 = \frac{1}{2}$  which gives  $p = q = 1$  and  $W_1 = \frac{1}{2}$ . With this choice we obtain from (7), the method

$$y_{n+1} = y_n + \frac{h}{2} [f(t_n, y_n) + f(t_n + h, y_n + hf_n)] \quad (12)$$

which is known as **Heun's method**.

**Note** that when  $f$  is a function of  $t$  only, the method (12) is equivalent to the trapezoidal rule of integration, whereas the method (6) is equivalent to the midpoint rule of integration. Both the methods (6) and (12) are of  $O(h^2)$ . The methods (6) and (12) can easily be implemented to solve the IVP (1). Method (6) is usually known as **improved tangent method** or **modified Euler method**. Method (12) is also known as **Euler-Cauchy method**.

We shall now discuss the Runge-Kutta methods of  $O(h^2)$ ,  $O(h^3)$  and  $O(h^4)$ .

### 15.2.1 Runge-Kutta Methods of Second Order

The general idea of the Runge-Kutta (R-K) methods is to write the required methods as

$$\begin{aligned} y_{n+1} &= y_n + h \text{ (weighted sum of the slopes).} \\ &= y_n + \sum_{i=1}^m W_i K_i \end{aligned} \quad (13)$$

where  $m$  slopes are being used. These slopes are defined by

$$\begin{aligned} K_1 &= hf(t_n, y_n), \\ K_2 &= hf(t_n + C_2h, y_n + a_{21}K_1), \\ K_3 &= hf(t_n + C_3h, y_n + a_{31}K_1 + a_{32}K_2), \\ K_4 &= hf(t_n + C_4h, y_n + a_{41}K_1 + a_{42}K_2 + a_{43}K_3), \end{aligned}$$

etc. In general, we can write

$$K_i = f \left[ t_n + C_i h, \sum_{j=1}^{i-1} a_{ij} K_j \right], \quad i = 1, 2, \dots, m \text{ with } C_1 = 0 \quad (14)$$

The parameters  $C_i, a_{ij}, W_j$  are unknowns and are to be determined to obtain the Runge-Kutta methods.

We shall now derive the second order Runge-Kutta methods.

Consider the method as

$$y_{n+1} = y_n + W_1 K_1 + W_2 K_2 \quad (15)$$

where

$$\begin{aligned} K_1 &= hf(t_n, y_n) \\ K_2 &= hf(t_n + C_2h, y_n + a_{21}K_1) \end{aligned} \quad (16)$$

where the parameters  $C_2, a_{21}, W_1$  and  $W_2$  are chosen to make  $y_{n+1}$  closer to  $y(t_{n+1})$ .

The exact solution satisfies the Taylor series

$$y(t_{n+1}) = y(t_n) + (t_n) h y' + \frac{h^2}{2} y''(t_n) + \frac{h^3}{6} y'''(t_n) + \dots \quad (17)$$

where

$$y' = f(t, y)$$

$$y'' = f_t + f f_y$$

$$y''' = f_{tt} + 2f f_{ty} + f_{yy} f^2 + f_y(f_t + f f_y)$$

We expand  $K_1$  and  $K_2$  about the point  $(t_n, y_n)$

$$K_1 = hf(t_n, y_n) = hf_n$$

$$K_2 = hf(t_n + C_2 h, y_n + a_{21} h f_n)$$

$$= h \left\{ f(t_n, y_n) + (C_2 h f_t + a_{21} h f_n f_y) + \frac{1}{2!} (C_2^2 h^2 f_{tt} + 2C_2 a_{21} h^2 f_n f_{ty} + a_{21}^2 h^2 f_n^2 f_{yy}) + \dots \right\}$$

Substituting these values of  $K_1$  and  $K_2$  in Eqn. (15), we have

$$y_{n+1} = y_n + (W_1 + W_2) hf_n + h^2 [W_2 C_2 f_t + W_2 a_{21} f_n f_y] + \frac{h^3}{2} W_2 (C_2^2 f_{tt} + 2C_2 a_{21} f_n f_{ty} + a_{21}^2 f_n^2 f_{yy}) + \dots \quad (18)$$

Comparing Eqn. (18) with (17), we have

$$W_1 + W_2 = 1$$

$$C_2 W_2 = \frac{1}{2}$$

$$a_{21} W_2 = \frac{1}{2}$$

From these equations we find that if  $C_2$  is chosen arbitrarily we have

$$a_{21} = C_2, W_2 = 1/(2C_2), \quad W_1 = 1 - 1/(2C_2) \quad (19)$$

The R-K method is given by

$$y_{n+1} = y_n + h [W_1 f(t_n, y_n) + W_2 f(t_n + C_2 h, y_n + C_2 h f_n)]$$

and Eqn. (18) becomes

$$y_{n+1} = y_n + hf_n + \frac{h^2}{2} (f_t + f_n f_y) + \frac{C_2 h^3}{4} (f_{tt} + 2f_n f_{ty} + f_n^2 f_{yy}) + \dots \quad (20)$$

Subtracting Eqn. (20) from the Taylor series (17), we get the truncation error as

$$\begin{aligned} TE &= y(t_{n+1}) - y_{n+1} \\ &= h^2 \left[ \left( \frac{1}{6} - \frac{C_2}{4} \right) (f_{tt} + 2f_n f_{ty} + f_n^2 f_{yy}) + \frac{1}{6} f_y (f_t + f_n f_y) \right] + \dots \\ &= \frac{h^3}{12} [(2 - 3C_2) y'' + 3C_2 f_y y''] + \dots \end{aligned} \quad (21)$$

Since the TE is of  $O(h^3)$ , all the above R-K methods are of second order. Observe that no choice of  $C_2$  will make the leading term of TE zero for all  $f(t, y)$ . The local TE depends not only on derivatives of the solution  $y(t)$  but also on the function  $f(t, y)$ . This is typical of all the Runge-Kutta methods. Generally,  $C_2$  is chosen between 0 and 1 so that we are evaluating  $f(t, y)$  at an off-step point in  $[t_n, t_{n+1}]$ . From the definition, every Runge-Kutta formula must reduce to a quadrature formula of the same order or greater if  $f(t, y)$  is independent of  $y$ , where  $W_i$  and  $C_i$  will be weights and abscissas of the corresponding numerical integration formula.

Best way of obtaining the value of the arbitrary parameter  $C_2$  in our formula is to

- i) choose some of  $W_i$ 's zero so as to minimize the computations.
- ii) choose the parameter to obtain least TE,
- iii) choose the parameter to have longer stability interval.

Methods satisfying either of the condition (ii) or (iii) are called optimal Runge-Kutta methods.

We made the following choices :

i)  $C_2 = \frac{1}{2}, \therefore a_{21} = \frac{1}{2}, W_1 = 0, W_2 = 1$ , then

$$y_{n+1} = y_n + K_2,$$

$$K_1 = hf(t_n, y_n),$$

$$K_2 = hf\left(t_n + \frac{h}{2}, y_n + \frac{K_1}{2}\right) \tag{22}$$

which is the same as improved tangent or modified Euler's method.

ii)  $C_2 = 1, \therefore a_{21} = 1, W_1 = W_2 = \frac{1}{2}$ , then

$$y_{n+1} = y_n + \frac{1}{2}(K_1 + K_2),$$

$$K_1 = hf(t_n, y_n),$$

$$K_2 = hf(t_n + h, y_n + K_1) \tag{23}$$

which is same as the Euler-Cauchy method or Heun's method.

iii)  $C_2 = \frac{2}{3}, \therefore a_{21} = \frac{2}{3}, W_1 = \frac{1}{4}, W_2 = \frac{3}{4}$ , then

$$y_{n+1} = y_n + \frac{1}{4}(K_1 + 3K_2),$$

$$K_1 = hf(t_n, y_n),$$

$$K_2 = hf\left(t_n + \frac{2h}{3}, y_n + \frac{2K_1}{3}\right) \tag{24}$$

which is the optimal R-K method.

Method (24) is optimal in the sense that it has minimum TE. In other words, with the above choice of unknowns, the leading term in the TE given by Eqn. (21) is minimum. Though several other choices are possible, we shall limit our discussion with the above three methods only.

In order to remember the weights  $W_i$  and scale factors  $C_i$  and  $a_{ij}$  we draw the following tables :

$C_2$	$a_{21}$	
	$W_1$	$W_2$

General form

$1/2$	$1/2$	
	0	1

Improved tangent method

1	1	
	$1/2$	$1/2$

Heun's method

$2/3$	$2/3$	
	$1/4$	$1/4$

Optimal method

We now illustrate these methods through an example.

**Example 1 :** Solve the IVP  $y' = -t y^2$ ,  $y(2) = 1$  and find  $y(2.1)$  and  $y(2.2)$  with  $h = 0.1$  using the following R-K methods of  $O(h^2)$

- Improved tangent method [modified Euler method (22)]
- Heun's method [Euler-Cauchy method (23)]
- Optimal R-K method [method (24)]
- Taylor series method of  $O(h^2)$ .

Compare the results with the exact solution

$$y(t) = \frac{2}{t^2 - 2}$$

**Solution :** We have the exact values

$$y(2.1) = 0.82988 \text{ and } y(2.2) = 0.70422$$

a) Improved tangent method is

$$y_{n+1} = y_n + K_2$$

$$K_1 = hf(t_n, y_n)$$

$$K_2 = hf\left(t_n + \frac{h}{2}, y_n + \frac{K_1}{2}\right)$$

For this problem  $f(t, y) = -t y^2$  and

$$K_1 = (0.1) [(-2)(1)] = -0.2$$

$$K_2 = (0.1) [(-2.05)(1 - 0.1)^2] = -0.16605$$

$$y(2.1) = 1 - 0.16605 = 0.83395$$

Taking  $t_1 = 2.1$  and  $y_1 = 0.83395$ , we have

$$K_1 = hf(t_1, y_1) = (0.1) [(-2.1)(0.83395)^2] = -0.146049$$

$$K_2 = hf\left(t_1 + \frac{h}{2}, y_1 + \frac{K_1}{2}\right)$$

$$= (0.1) [-(2.15)(0.83395 - 0.0730245)^2] = -0.124487$$

$$y(2.2) = y_1 + K_2 = 0.83395 - 0.124487 = 0.70946$$

b) Heun's method is :

$$y_{n+1} = y_n + \frac{1}{2}(K_1 + K_2)$$

$$K_1 = hf(t_n, y_n) = -0.2$$

$$K_2 = hf(t_n + h, y_n + K_1) = -0.1344$$

$$y(2.1) = 0.8328$$

Taking  $t_1 = 2.1$  and  $y_1 = 0.8328$ , we have

$$K_1 = -0.14564, \quad K_2 = -0.10388$$

$$y(2.2) = 0.70804$$

c) Optimal method is :

$$y_{n+1} = y_n + \frac{1}{4}(K_1 + 3K_2)$$

$$K_1 = hf(t_n, y_n) = -0.2$$

$$K_2 = hf\left(t_n + \frac{2h}{3}, y_n + \frac{2K_1}{3}\right) = 0.15523$$

$$y(2.1) = 0.83358$$

Taking  $t_1 = 2.1$  and  $y_1 = 0.83358$ , we have  
 $K_1 = -0.1459197$ ,  $K_2 = -0.117463$

$$y(2.2) = 0.7090$$

d) Taylor series method of  $O(h^2)$  :

$$y_{n+1} = y_n + hy'_n + \frac{h^2}{2} y''_n$$

$$y' = -ty^2, y'' = -y^2 - 2tyy'$$

$$y(2) = 1, y'(2) = -2, y''(2) = 7$$

$$y(2.1) = 0.8350$$

With  $t_1 = 2.1$ ,  $y_1 = 0.835$ , we get

$$y'(2.1) = -1.4641725, y''(2.1) = 4.437627958$$

$$y(2.2) = 0.71077$$

We now summarise the results obtained and give them in Table 1.

**Table 1**

Solutions and errors in solution of  $y' = -ty^2$ ,  $y(2) = 1$ ,  $h = 0.1$ . Numbers inside brackets denote the errors.

t	Method (22)	Method (23)	Method (24)	Method Taylor $O(h^2)$	Exact Solution
2.1	0.83395 (0.00405)	0.8328 (0.0029)	0.83358 (0.00368)	0.8350 (0.0051)	0.8299
2.2	0.70746 (0.0033)	0.70804 (0.00384)	0.7090 (0.0048)	0.71077 (0.00657)	0.7042

You may observe here that all the above numerical solutions have almost the same error.

You may now try the following exercises :

Solve the following IVPs using Heun's method of  $O(h^2)$  and the optimal R-K method of  $O(h^2)$ .

E1)  $10y' = t^2 + y^2$ ,  $y(0) = 1$ . Find  $y(0.2)$  taking  $h = 0.1$ .

E2)  $y' = 1 + y^2$ ,  $y(0) = 0$ . Find  $y(0.4)$  taking  $h = 0.2$ . Given that the exact solution is  $y(t) = \tan t$ , find the errors.

Also compare the errors at  $t = 0.4$ , obtained here with the one obtained by Taylor series method of  $O(h^2)$ .

E3)  $y' = 3t + \frac{1}{2}y$ ,  $y(0) = 1$ . Find  $y(0.2)$  taking  $h = 0.1$ . Given  $y(t) = 13e^{t/2} - 6t - 12$ , find the errors.

Let us now discuss the R-K methods of third order.

### 15.2.2 Runge-Kutta Methods of Third Order

Here we consider the method as

$$y_{n+1} = y_n + W_1 K_1 + W_2 K_2 + W_3 K_3 \tag{25}$$

where

$$K_1 = h f(t_n, y_n)$$

$$K_2 = h f(t_n + C_2 h, y_n + a_{21} K_1)$$

$$K_3 = h f(t_n + C_3 h, y_n + a_{31} K_1 + a_{32} K_2)$$

Expanding  $K_2$ ,  $K_3$  and  $y_{n+1}$  into Taylor series, substituting their values in Eqn. (25) and comparing the coefficients of powers of  $h$ ,  $h^2$  and  $h^3$ , we obtain

$$\begin{aligned} a_{21} &= C_2 & C_2 W_2 + C_3 W_3 &= \frac{1}{2} \\ a_{31} + a_{32} &= C_3 & C_2^2 W_2 + C_3^2 W_3 &= \frac{1}{3} \\ W_1 + W_2 + W_3 &= 1 & C_2 a_{32} W_3 &= \frac{1}{6} \end{aligned} \quad (26)$$

We have 6 equations to determine the 8 unknowns. Hence the system has two arbitrary parameters. Eqns. (26) are typical of all the R-K methods. Looking at Eqn. (26), you may note that the sum of  $a_{ij}$ 's in any row equals the corresponding  $C_i$ 's and the sum of the  $W_i$ 's is equal to 1. Further, the equations are linear in  $W_2$  and  $W_3$  and have a solution for  $W_2$  and  $W_3$  if and only if

$$\begin{vmatrix} C_2 & C_3 & -1/2 \\ C_2^2 & C_3^2 & -1/3 \\ 0 & C_2 a_{32} & -1/6 \end{vmatrix} = 0$$

(Ref. Sec. 8.4.2, Unit 8, Block-2, MTE-02).

Expanding the determinant and simplifying we obtain

$$C_2(2 - 3C_2) a_{32} - C_3(C_3 - C_2) = 0, C_2 \neq 0 \quad (27)$$

Thus we choose  $C_2$ ,  $C_3$  and  $a_{32}$  satisfying Eqns. (27).

Since two parameters of this system are arbitrary, we can choose  $C_2$ ,  $C_3$  and determine  $a_{32}$  from Eqn. (27) as

$$a_{32} = \frac{C_3(C_3 - C_2)}{C_2(2 - 3C_2)}$$

If  $C_3 = 0$ , or  $C_2 = C_3$  then  $C_2 = \frac{2}{3}$  and we can choose  $a_{32} \neq 0$ , arbitrarily. All  $C_i$ 's should be chosen such that  $0 < C_i < 1$ . Once  $C_2$  and  $C_3$  are prescribed,  $W_i$ 's and  $a_{ij}$ 's can be determined from Eqns. (26).

We shall list a few methods in the following notation.

$C_2$	$a_{21}$		
$C_3$	$a_{31}$	$a_{32}$	
	$W_1$	$W_2$	$W_3$

i) Classical third order R-K method

$1/2$	$1/2$		
$1$	$-1$	$2$	
	$1/6$	$4/6$	$1/6$

$$y_{n+1} = y_n + \frac{1}{6} (K_1 + 4K_2 + K_3) \quad (28)$$

$$K_1 = hf(t_n, y_n)$$

$$K_2 = hf\left(t_n + \frac{h}{2}, y_n + \frac{K_1}{2}\right)$$



$$K_3 = hf \left( t_n + \frac{h}{2}, y_n - K_1 + 2K_2 \right)$$

ii) Heun's Method

1/3	1/3		
2/3	0	2/3	
	1/4	0	3/4

$$y_{n+1} = y_n + \frac{1}{4} (K_1 + 3K_3) \quad (29)$$

$$K_1 = h f(t_n, y_n)$$

$$K_2 = h f \left( t_n + \frac{h}{3}, y_n + \frac{K_1}{3} \right)$$

$$K_3 = h f \left( t_n + \frac{2h}{3}, y_n + \frac{2K_2}{3} \right)$$

iii) Optimal method

1/2	1/2		
3/4	0	3/4	
	2/9	3/9	4/9

$$y_{n+1} = y_n + \frac{1}{9} (2K_1 + 3K_2 + 4K_3) \quad (30)$$

$$K_1 = h f(t_n, y_n),$$

$$K_2 = h f \left( t_n + \frac{h}{2}, y_n + \frac{K_1}{2} \right),$$

$$K_3 = h f \left( t_n + \frac{3h}{4}, y_n + \frac{3K_2}{4} \right).$$

We now illustrate the third order R-K methods by solving the problem considered in Example 1, using (a) Heun's method (b) optimal method

a) Heun's method

$$y_{n+1} = y_n + \frac{1}{4} (K_1 + 3K_3)$$

$$K_1 = h f(t_n, y_n) \\ = -0.2$$

$$K_2 = h f \left( t_n + \frac{h}{3}, y_n + \frac{K_1}{3} \right) \\ = -0.17697$$

$$K_3 = h f \left( t_n + \frac{2h}{3}, y_n + \frac{2K_2}{3} \right) \\ = -0.16080$$

$$y(2.1) = 0.8294$$

Taking  $t_1 = 2.1$  and  $y_1 = 0.8294$ , we have

$$K_1 = -0.14446$$

$$K_2 = -0.13017$$

$$K_3 = -0.11950$$

$$y(2.2) = 0.70366$$

b) Optimal method

$$y_{n+1} = y_n + \frac{1}{9} (2K_1 + 3K_2 + 4K_3)$$

$$K_1 = -0.2$$

$$K_2 = -0.16605$$

$$K_3 = -0.15905$$

$$y(2.1) = 0.8297$$

Taking  $t_1 = 2.1$  and  $y_1 = 0.8297$ , we have

$$K_1 = -0.14456$$

$$K_2 = -0.12335$$

$$K_3 = -0.11820$$

$$y(2.2) = 0.70405$$

You can now easily find the errors in these solutions and compare the results with those obtained in Example 1.

And now here is an exercise for you.

E4) Solve the IVP

$$y' = y - t, \quad y(0) = 2$$

using third order Heun's or optimal R-K methods. Find  $y(0.2)$  taking  $h = 0.1$ . Given the exact solution to be  $y(t) = 1 + t + e^t$ , find the errors at  $t = 0.2$ .

We now discuss the fourth order R-K methods.

### 15.2.3 Runge-Kutta Methods of Fourth Order

Consider the method as

$$y_{n+1} = y_n + W_1 K_1 + W_2 K_2 + W_3 K_3 + W_4 K_4 \quad (31)$$

$$K_1 = h f(t_n, y_n),$$

$$K_2 = h f(t_n + C_2 h, y_n + a_{21} K_1),$$

$$K_3 = h f(t_n + C_3 h, y_n + a_{31} K_1 + a_{32} K_2),$$

$$K_4 = h f(t_n + C_4 h, y_n + a_{41} K_1 + a_{42} K_2 + a_{43} K_3).$$

Since the expansions of  $K_2, K_3, K_4$  and  $y_{n+1}$  in Taylor series are complicated, we shall not write down the resulting system of equations for the determination of the unknowns. It may be noted that the system of equations has 3 arbitrary parameters. We shall state directly a few R-K methods of  $O(h^4)$ . The R-K methods (31) can be denoted by

$C_2$	$a_{21}$			
$C_3$	$a_{31}$	$a_{32}$		
$C_4$	$a_{41}$	$a_{42}$	$a_{43}$	
	$W_1$	$W_2$	$W_3$	$W_4$

For different choices of these unknowns we have the following methods :

i) Classical R-K method

$\frac{1}{2}$	$\frac{1}{2}$			
$\frac{1}{2}$	0	$\frac{1}{2}$		
1	0	0	1	
	$\frac{1}{6}$	$\frac{2}{6}$	$\frac{2}{6}$	$\frac{1}{6}$

$$y_{n+1} = y_n + \frac{1}{6} (K_1 + 2K_2 + 2K_3 + K_4) \quad (32)$$

$$K_1 = h f(t_n, y_n),$$

$$K_2 = h f\left(t_n + \frac{h}{2}, y_n + \frac{K_1}{2}\right),$$

$$K_3 = h f\left(t_n + \frac{h}{2}, y_n + \frac{K_2}{2}\right),$$

$$K_4 = h f(t_n + h, y_n + K_3).$$

This is the widely used method due to its simplicity and moderate order. We shall also be working out problems mostly by the classical R-K method unless specified otherwise.

ii) Runge-Kutta-Gill method

$\frac{1}{2}$	$\frac{1}{2}$			
$\frac{1}{2}$	$\frac{(\sqrt{2}-1)}{2}$	$\frac{(2-\sqrt{2})}{2}$		
1	0	$-\frac{\sqrt{2}}{2}$	$1 + \frac{\sqrt{2}}{2}$	
	$\frac{1}{6}$	$\frac{(2-\sqrt{2})}{6}$	$\frac{(2+\sqrt{2})}{6}$	$\frac{1}{6}$

$$y_{n+1} = y_n + \frac{1}{6} (K_1 + (2-\sqrt{2}) K_2 + (2+\sqrt{2}) K_3 + K_4) \quad (33)$$

$$K_1 = h f(t_n, y_n),$$

$$K_2 = h f\left(t_n + \frac{h}{2}, y_n + \frac{K_1}{2}\right),$$

$$K_3 = h f\left(t_n + \frac{h}{2}, y_n + \left(\frac{\sqrt{2}-1}{2}\right) K_1 + \left(\frac{2-\sqrt{2}}{2}\right) K_2\right),$$

$$K_4 = h f\left(t_n + h, y_n - \frac{\sqrt{2}}{2} K_2 + \left(1 + \frac{\sqrt{2}}{2}\right) K_3\right).$$

The Runge-Kutta-Gill method is also used widely. But, in this unit, we shall mostly work out problems with the classical R-K method of  $O(h^4)$ . Hence, whenever we refer to R-K method of  $O(h^4)$  we mean only the classical R-K method of  $O(h^4)$  given by (32). We shall now illustrate this method through examples.

**Example 2 :** Solve the IVP  $y' = t + y, y(0) = 1$  by Runge-Kutta method of  $O(h^4)$  for  $t \in [0, 0.5]$  with  $h = 0.1$ . Also find the error at  $t = 0.5$ , if the exact solution is  $y(t) = 2e^t - t - 1$ .

**Solution :** We use the R-K method of  $O(h^4)$  given by (32).

Initially,  $t_0 = 0, y_0 = 1$ .

We have

$$K_1 = hf(t_0, y_0) = (0.1) [0+1] = 0.1$$

$$K_2 = hf\left(t_0 + \frac{h}{2}, y_0 + \frac{K_1}{2}\right) = (0.1) [0.05 + 1 + 0.05] = 0.11$$

$$K_3 = hf\left(t_0 + \frac{h}{2}, y_0 + \frac{K_2}{2}\right) = (0.1) [0.05 + 1 + 0.055] = 0.1105$$

$$K_4 = hf(t_0 + h, y_0 + K_3) = (0.1) [0.1 + 1 + 0.1105] = 0.12105$$

$$y_1 = y_0 + \frac{1}{6}(K_1 + 2K_2 + 2K_3 + K_4)$$

$$= 1 + \frac{1}{6}[1 + 0.22 + 0.2218 + 0.2205] = 1.11034167$$

Taking  $t_1 = 0.1$  and  $y_1 = 1.11034167$ , we repeat the process.

$$K_1 = hf(t_1, y_1) = (0.1)[0.1 + 1.11034167] = 0.121034167$$

$$K_2 = hf\left(t_1 + \frac{h}{2}, y_1 + \frac{K_1}{2}\right) = (0.1)\left[0.1 + 0.05 + 1.11034167 + \frac{(0.121034167)}{2}\right]$$

$$= 0.132085875$$

$$K_3 = hf\left(t_1 + \frac{h}{2}, y_1 + \frac{K_2}{2}\right) = (0.1)\left[0.1 + 0.05 + 1.11034167 + \frac{(0.132085875)}{2}\right]$$

$$= 0.132638461$$

$$K_4 = hf\left(t_1 + h, y_1 + K_3\right) = (0.1)[0.1 + 0.05 + 1.11034167 + 0.132638461]$$

$$= 0.144303013$$

$$y_2 = y_1 + \frac{1}{6}(K_1 + 2K_2 + 2K_3 + K_4)$$

$$= 1.11034167 + \frac{1}{6}[(0.121034167 + 2(0.132085875) + 2(0.132638461)$$

$$+ 0.144303013] = 1.24280514$$

Rest of the values  $y_3, y_4, y_5$  we give in Table 2.

Table 2

$t_n$	$y_n$
0.0	1
0.1	1.11034167
0.2	1.24280514
0.3	1.39971699
0.4	1.58364848
0.5	1.79744128

Now the exact solution is

$$y(t) = 2e^t - t - 1$$

Error at  $t = 0.5$  is

$$y(0.5) - y_5 = (2e^{0.5} - 0.5 - 1) - 1.79744128$$

$$= 1.79744254 - 1.79744128$$

$$= 0.000001261$$

$$= 0.13 \times 10^{-05}$$

Let us consider another example

**Example 3 :** Solve the IVP

$$y' = 2y + 3e^t, y(0) = 0 \text{ using}$$

a) classical R-K method of  $O(h^4)$

b) R-K Gill method of  $O(h^4)$ ,

Find  $y(0.1), y(0.2), y(0.3)$  taking  $h = 0.1$ . Also find the errors at  $t = 0.3$ , if the exact solution is  $y(t) = 3(e^{2t} - e^t)$ .

Solution : a) Classical R-K method is

$$y_{n+1} = y_n + \frac{1}{6} (K_1 + 2K_2 + 2K_3 + K_4)$$

Here  $t_0 = 0$ ,  $y_0 = 0$ ,  $h = 0.1$

$$K_1 = h f(t_0, y_0) = 0.3$$

$$K_2 = h f\left(t_0 + \frac{h}{2}, y_0 + \frac{K_1}{2}\right) = 0.3453813289$$

$$K_3 = h f\left(t_0 + \frac{h}{2}, y_0 + \frac{K_2}{2}\right) = 0.3499194618$$

$$K_4 = h f(t_0 + h, y_0 + K_3) = 0.4015351678$$

$$y_1 = 0.3486894582$$

Taking  $t_1 = 0.1$ ,  $y_1 = 0.3486894582$ , we repeat the process and obtain

$$K_1 = 0.4012891671, \quad K_2 = 0.4584170812$$

$$K_3 = 0.4641298726, \quad K_4 = 0.6887058455$$

$$y(0.2) = 0.8112570941$$

Taking  $t_2 = 0.2$ ,  $y_2 = 0.837870944$  and repeating the process we get

$$K_1 = 0.53399502, \quad K_2 = 0.579481565$$

$$K_3 = 0.61072997, \quad K_4 = 0.694677825$$

$$\therefore y(0.3) = 1.416807999$$

b) R-K-Gill method is

$$y_{n+1} = y_n + \frac{1}{6} (K_1 + (2-\sqrt{2}) K_2 + (2+\sqrt{2}) K_3 + K_4)$$

Taking  $t_0 = 0$ ,  $y_0 = 1$  and  $h = 0.1$ , we obtain

$$K_1 = 0.3, \quad K_2 = 0.3453813289$$

$$K_3 = 0.3480397056, \quad K_4 = 0.4015351678$$

$$y(0.1) = 0.3486894582$$

Taking  $t_1 = 0.1$ ,  $y_1 = 0.3486894582$ , we obtain

$$K_1 = 0.4012891671, \quad K_2 = 0.4584170812$$

$$K_3 = 0.4617635569, \quad K_4 = 0.5289846936$$

$$y(0.2) = 0.8112507529$$

Taking  $t_2 = 0.2$ ,  $y_2 = 0.8112507529$ , we obtain

$$K_1 = 0.528670978, \quad K_2 = 0.6003248734$$

$$K_3 = 0.6045222614, \quad K_4 = 0.6887058455$$

$$y(0.3) = 1.416751936$$

From the exact solution we get

$$y(0.3) = 1.416779978$$

Error in classical R-K method (at  $t = 0.3$ ) =  $0.2802 \times 10^{-04}$

Error in R-K-Gill method (at  $t = 0.3$ ) =  $0.2804 \times 10^{-04}$ .

You may now try the following exercises.

Solve the following IVPs using R-K method of  $O(h^4)$

E5)  $y' = \frac{y-t}{y+t}$ ,  $y(0) = 1$ . Find  $y(0.5)$  taking  $h = 0.5$ .

E6)  $y' = 1 - 2ty$ ,  $y(0.2) = 0.1948$ . Find  $y(0.4)$  taking  $h = 0.2$ .

- E7)  $10ty' + y^2 = 0, y(4) = 1$ . Find  $y(4.2)$  taking  $h = 0.2$ . Find the error given the exact solution is  $y(t) = \frac{1}{c + 0.1 \ln t}$ , where  $c = 0.86137$
- E8)  $y' = \frac{1}{t^2} - \frac{y}{t} - y^2, y(1) = -1$ . Find  $y(1.3)$  taking  $h = 0.1$ . Given the exact solution to be  $y(t) = \frac{1}{t}$ , find the error at  $t = 1.3$ .

In the next section, we shall study the application of Richardson's extrapolation to the solutions of ordinary differential equations.

### 15.3 RICHARDSON'S EXTRAPOLATION

You know that Richardson's extrapolation technique improves the approximate value of  $y(t_h)$  and the order of this improved value of  $y(t_h)$  exceeds the order of the method by one.

Here we shall first calculate the solutions  $F(h_1)$  and  $F(h_2)$  of the given IVP with steplengths  $h_1$  and  $h_2$  where  $h_2 = h_1/2$  at a given point using a Runge-Kutta method. Then by Richardson's extrapolation technique we have for the second order method

$$F^{(1)}(h) = \frac{4F(h/2) - F(h)}{4 - 1} = \left[ \frac{4F(h/2) - F(h)}{3} \right] \quad (34)$$

and for the fourth order method

$$F^{(1)}(h) = \frac{16F(h/2) - F(h)}{16 - 1} = \frac{1}{15} \left[ 16F\left(\frac{h}{2}\right) - F(h) \right] \quad (35)$$

as the improved solution at that point, which will be of higher order than the original method. We shall now illustrate the technique through an example.

**Example 4 :** Using Runge-Kutta method of  $O(h^2)$  find the solution of the IVP  $y' = t + y, y(0) = 1$  using  $h = 0.1$  and  $0.2$  at  $t = 0.4$ . Use extrapolation technique to improve the accuracy. Also find the errors if the exact solution is  $y(t) = 2e^t - t - 1$ .

**Solution :** We shall use Heun's second order method (23) to find the solution at  $t = 0.4$  with  $h = 0.1$  and  $0.2$ . The following Table 3 gives values of  $y(t)$  at  $t = 0.2$  and  $t = 0.4$  with  $h = 0.1$  and  $0.2$ .

Table 3

$t_h$	$F_1 = F(0.1)$	$F_2 = F(0.2)$	Extrapolated value = $\frac{1}{3}(4F_1 - F_2)$	Errors
0.2	1.24205	1.24	1.242733	$0.725 \times 10^{-4}$
0.4	1.58180	1.5768	1.583472	$0.177 \times 10^{-3}$

You may now try the following exercises :

- E9) Solve E2) taking  $h = 0.1$  and  $0.2$  using  $O(h^2)$  Heun's method. Extrapolate the value at  $t = 0.4$ . Also find the error at  $t = 0.4$ .
- E10) Solve E6), taking  $h = 0.1$  and  $0.2$  using  $O(h^2)$  Heun's method. Extrapolate the value at  $t = 0.4$ . Compare this solution with the solution obtained by the classical  $O(h^4)$  R-K method.

We now end this unit by giving a summary of what we have covered in it.

## 15.4 SUMMARY

In this unit we have learnt the following :

- 1) Runge-Kutta methods being singlestep methods are self- starting methods.
- 2) Unlike Taylor series methods, R-K methods do not need calculation of higher order derivatives of  $f(t, y)$  but need only the evaluation of  $f(t, y)$  at the off-step points.
- 3) For a given IVP of the form

$$y' = f(t, y), \quad y'(t_0) = y_0, \quad t \in [t_0, b]$$

where the mesh points are  $t_j = t_0 + jh, j = 0, 1, \dots, n$ .

$t_n = b = t_0 + nh$ , R-K methods are obtained by writing

$$y_{n+1} = y_n + h \text{ (weighted sum of the slopes)}$$

$$= y_n + \sum_{i=1}^m W_i K_i$$

where  $m$  slopes are used. These slopes are defined by

$$K_i = f \left[ t_n + C_i h, \sum_{j=1}^{i-1} a_{ij} k_j \right], \quad i = 1, 2, \dots, m, \quad C_1 = 0.$$

The unknowns  $C_i, a_{ij}$  and  $W_j$  are then obtained by expanding  $K_i$ 's and  $y_{n+1}$  in Taylor series about the point  $(t_n, y_n)$  and comparing the coefficients of different powers of  $h$ .

- 4) Richardson's extrapolation technique can be used to improve the approximate value of  $y(t_n)$  obtained by  $O(h^2), O(h^3)$  and  $O(h^4)$  methods and obtain the method of order one higher than the method.

## 15.5 SOLUTIONS/ANSWERS

E1) Heun's method :  $y_{n+1} = y_n + \frac{1}{2}(K_1 + K_2)$

Starting with  $t_0 = 0, y_0 = 1, h = 0.1$

$$\therefore K_1 = 0.01$$

$$K_2 = 0.010301$$

$$y(0.1) = 1.0101505$$

Taking  $t_1 = 0.1, y_1 = 1.0101505$

$$K_1 = 0.0103040403$$

$$K_2 = 0.0181327468$$

$$y(0.2) = 1.020709158$$

Optimal R-K, method :  $y_{n+1} = y_n + \frac{1}{4}(K_1 + 3K_2)$

$$t_0 = 0, y_0 = 1, h = 0.1$$

$$K_1 = 0.01, \quad K_2 = 0.01017823$$

$$y(0.1) = 1.010133673$$

$$t_1 = 0.1, y_1 = 1.010133673$$

$$K_1 = 0.0103037, \quad K_2 = 0.010620$$

$$y(0.2) = 1.020675142$$

E2) Heun's method :

$$K_1 = 0.2, \quad K_2 = 0.208$$

$$y(0.2) = 0.204$$

$$K_1 = 0.2083232, \quad K_2 = 0.2340020843$$

$$y(0.4) = 0.4251626422$$

Optimal R-K, method :

$$\begin{aligned} K_1 &= 0.2, \quad K_2 = 0.2035556 \\ y(0.2) &= 0.2026667 \\ K_1 &= 0.2082148, \quad K_2 = 0.223321245 \\ y(0.4) &= 0.422211334 \end{aligned}$$

Taylor series method

$$\begin{aligned} y' &= 1 + y^2, \quad y'' = 2yy' \\ y(0) &= 0, \quad y'(0) = 1, \quad y''(0) = 0 \\ y(0.2) &= 0.2 \\ y'(0.2) &= 1.04, \quad y''(0.2) = 0.416 \\ y(0.4) &= 0.41632 \end{aligned}$$

Now the exact solution is  $y(t) = \tan t$

$$\text{Exact } y(0.4) = 0.422793219$$

$$\text{Error in Heun's method} = 0.236 \times 10^{-2}$$

$$\text{Error in Optimal R-K method} = 0.582 \times 10^{-3}$$

$$\text{Error in Taylor series method} = 0.647 \times 10^{-2}$$

E3) Heun's method :

$$\begin{aligned} K_1 &= 0.05, \quad K_2 = 0.0825 \\ y(0.1) &= 1.06625 \\ K_1 &= 0.0833125, \quad K_2 = 0.117478125 \\ y(0.2) &= 1.166645313 \end{aligned}$$

Optimal R-K, method :

$$\begin{aligned} K_1 &= 0.05, \quad K_2 = 0.071666667 \\ y(0.1) &= 1.06625 \\ K_1 &= 0.0833125; \quad K_2 = 0.106089583 \\ y(0.2) &= 1.166645313 \\ \text{Exact } y(0.2) &= 1.167221935 \end{aligned}$$

$$\text{Error in both the methods is same and} = 0.577 \times 10^{-3}$$

E4) Heun's method :  $y_{n+1} = y_n + \frac{1}{4}(K_1 + 3K_2)$

Starting with  $t_0 = 0, y_0 = 2, h = 0.1$ , we have

$$\begin{aligned} K_1 &= 0.2, \quad K_2 = 0.203334, \quad K_3 = 0.206889 \\ y(0.1) &= 2.205167 \\ t_1 &= 0.1, \quad y_1 = 2.205167 \text{ we have} \\ K_1 &= 0.210517, \quad K_2 = 0.214201, \quad K_3 = 0.218130 \\ y(0.2) &= 2.421393717 \end{aligned}$$

Optimal R-K method :  $y_{n+1} = y_n + \frac{1}{9}(2K_1 + 3K_2 + 4K_3)$

$$\begin{aligned} K_1 &= 0.2, \quad K_2 = 0.205, \quad K_3 = 0.207875 \\ y(0.1) &= 2.205167 \\ t_1 &= 0.1, \quad y_1 = 2.205167 \\ K_1 &= 0.2105167, \quad K_2 = 0.2160425, \quad K_3 = 0.219220 \\ y(0.2) &= 2.421393717 \\ \text{exact } y(0.2) &= 2.421402758 \end{aligned}$$

Since  $y(0.2)$  is same by both the methods

$$\text{Error} = 0.9041 \times 10^{-5} \text{ in both the methods at } t = 0.2.$$

E5)  $K_1 = 0.5, \quad K_2 = 0.333333$

$$\begin{aligned} K_3 &= 0.3235204118, \quad K_4 = 0.2258064816 \\ y(0.5) &= 1.33992199. \end{aligned}$$



- E6)  $K_1 = 0.184416$ ,  $K_2 = 0.16555904$   
 $K_3 = 0.1666904576$ ,  $K_4 = 0.1421615268$   
 $y(0.4) = 0.3599794203$ .
- E7)  $K_1 = -0.005$ ,  $K_2 = -0.004853689024$   
 $K_3 = -0.0048544$ ,  $K_4 = -0.004715784587$   
 $y(4.2) = 0.9951446726$ .  
 Exact  $y(4.2) = 0.995145231$ , Error =  $0.559 \times 10^{-6}$

- E8)  $K_1 = 0.1$ ,  $K_2 = 0.09092913832$   
 $K_3 = 0.09049729525$ ,  $K_4 = 0.08260717517$   
 $y(1.1) = -0.909089993$   
 $K_1 = 0.08264471138$ ,  $K_2 = 0.07577035491$   
 $K_3 = 0.07547152415$ ,  $K_4 = 0.06942067502$   
 $y(1.2) = -0.8333318022$   
 $K_1 = 0.06944457204$ ,  $K_2 = 0.06411104536$   
 $K_3 = 0.06389773475$ ,  $K_4 = 0.0591559551$   
 $y(1.3) = -0.7692287876$   
 Exact  $y(1.3) = -0.7692307692$   
 Error =  $0.19816 \times 10^{-5}$

- E9) Heun's method :  
 with  $h = 0.1$

$$K_1 = 0.1, \quad K_2 = 0.101$$

$$y(0.1) = 0.1005$$

$$K_1 = 0.101010, \quad K_2 = 0.104061$$

$$y(0.2) = 0.203035$$

$$K_1 = 0.1041223, \quad K_2 = 0.1094346$$

$$y(0.3) = 0.309813$$

$$K_1 = 0.1095984, \quad K_2 = 0.1048047$$

$$F\left(\frac{h}{2}\right) = y(0.4) = 0.417014563$$

with  $h = 0.2$

$$F(h) = y(0.4) = 0.4251626422 \text{ [see E2]}$$

Now

$$F^{(1)}(0.4) = \frac{4F(h/2) - F(h)}{3}$$

$$= 0.414298537$$

$$\text{Exact } y(0.4) = 0.422793219$$

$$\text{Error} = 0.8495 \times 10^{-2}$$

- E10) Heun's method : with  $h = 0.1$

$$y(0.2) = 0.1948$$

$$K_1 = 0.092208, \quad K_2 = 0.08277952$$

$$y(0.3) = 0.28229375$$

$$K_1 = 0.083062374, \quad K_2 = 0.07077151$$

$$y(0.4) = 0.359210692$$

Heun's method with  $h = 0.2$

$$K_1 = 0.184416, \quad K_2 = 0.13932541$$

$$y(0.4) = 0.35667072$$

$$F^{(1)}(0.4) = 0.360057349$$

Result obtained by classical R-K method of  $O(h^4)$  is  
 $y(0.4) = 0.3599794203$  (see E6)

**NOTES**

