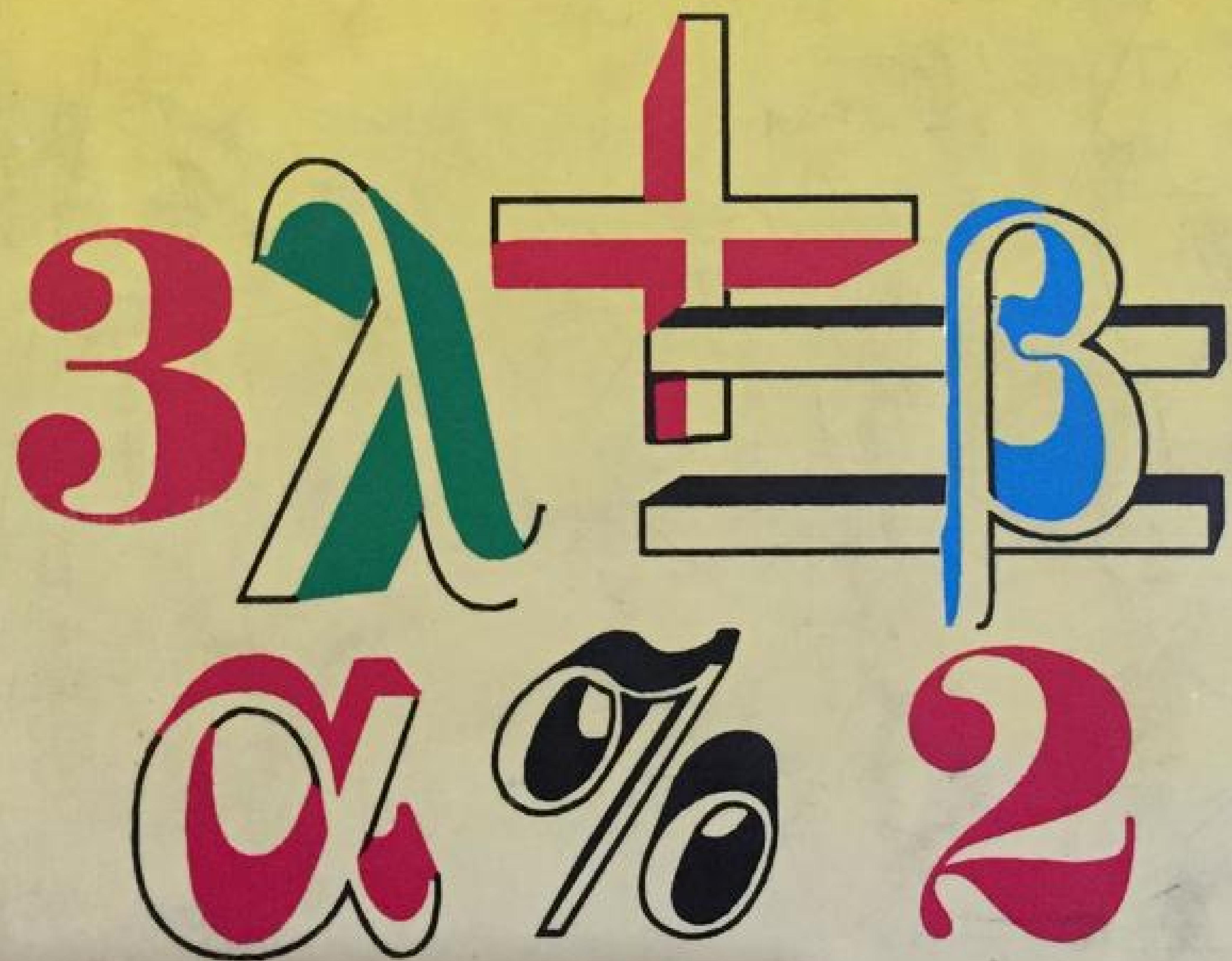


# ELEMENTS OF NUMERICAL ANALYSIS

$$(-) \times (-) = +$$

$$(+ \times (-) = -$$



*Dr. Faiz Ahmad  
Muhammad Afzal Rana*

## CONTENTS

### CHAPTER 1

#### ERRORS IN COMPUTATION

1.1	Introduction	1
1.2	Errors and their sources	2
1.3	Approximate numbers and significant figures	3
1.4	Rounding of numbers and around off errors	4
1.5	Truncation error	6
1.6	Absolute, relative, and percentage errors	7
1.7	The error of a sum	11
1.8	The error of a difference	13
1.9	The error of a product	14
1.10	The error of a quotient	16
1.11	The error of a power	17
1.12	The error of a root	18

### CHAPTER 2

#### LINEAR OPERATORS

2.1	Introduction	21
2.2	Functions of operators	22
2.3	Operator identities	23
2.4	Difference operators and the derivative operators	25

### CHAPTER 3

#### EIGENVALUES AND EIGENVECTORS OF MATRICES

3.1	Introduction	27
3.2	Eigenvalues and eigenvectors	27
3.3	Gershgorin's theorem	33
3.4	The Power method	36
3.5	Deflation	42

## CHAPTER 4

### INTERPOLATION WITH UNEQUALLY SPACED DATA

4.1	Introduction	8
4.2	Lagrange's formula	8
4.3	The error of the interpolation polynomial	8
4.4	Divided differences	8
4.4.1	Divided-difference table	8
4.4.2	Newton's interpolation polynomial	8

## CHAPTER 5

### INTERPOLATION WITH EQUALLY SPACED DATA

5.1	The difference table	62
5.2	Newton's forward and backward difference formulae	66
5.3	The Lozenge diagram	68
5.4	Gauss formulae	71
5.5	Stirling's interpolation formula	71
5.6	Bessel's interpolation formula	73
5.7	Everett's interpolation formula	74

## CHAPTER 6

### THE SOLUTION OF NONLINEAR EQUATIONS

6.1	Introduction	79
6.2	Bisection method	80
6.3	The method of false position or regula falsi method	83
6.4	Secant method	83
6.5	Newton-Raphson method	83
6.6	Fixed point iteration	84
6.7	Order of convergence	84

## CHAPTER 7

### DIFFERENCE EQUATIONS

7.1	Difference equations	102
7.2	Linear homogeneous difference equations	102
7.3	Linear inhomogeneous difference equations	102

## CHAPTER 4 INTERPOLATION WITH UNEQUALLY SPACED DATA

4.1	Introduction	8.1
4.2	Lagrange's formula	8.2
4.3	The error of the interpolation polynomial	8.3
4.4	Divided differences	8.3.1
4.4.1	Divided-difference table	8.3.2
4.4.2	Newton's interpolation polynomial	8.3.3

## CHAPTER 5

## INTERPOLATION WITH EQUALLY SPACED DATA

5.1	The difference table	62
5.2	Newton's forward and backward difference formulae	66
5.3	The Lozenge diagram	68
5.4	Gauss formulae	71
5.5	Stirling's interpolation formula	73
5.6	Bessel's interpolation formula	73
5.7	Everett's interpolation formula	74

## CHAPTER 6

## THE SOLUTION OF NONLINEAR EQUATIONS

6.1	Introduction	79
6.2	Bisection method	80
6.3	The method of false position or regula falsi method	83
6.4	Secant method	88
6.5	Newton-Raphson method	89
6.6	Fixed point iteration	92
6.7	Order of convergence	93

## CHAPTER 7

## DIFFERENCE EQUATIONS

7.1	Difference equations	10.2
7.2	Linear homogeneous difference equations	10.2
7.3	Linear inhomogeneous difference equations	10.3
		10.4
		10.5

## CHAPTER 8

### SOLUTION OF SYSTEMS OF LINEAR EQUATIONS (DIRECT METHODS)

8.1	Introduction	115
8.2	Gauss's elimination method	117
8.3	LU decomposition	122
8.3.1	Dolittle's method	123
8.3.2	Crout's method	126
8.3.3	Cholesky's method	128

## CHAPTER 9

### SOLUTION OF SYSTEMS OF LINEAR EQUATIONS (INDIRECT METHODS)

9.1	Iterative methods	133
9.1.1	Jacobi iterative method	134
9.1.2	Gauss Seidel iterative method	138
9.1.3	Successive over-relaxation	145
9.2	Residuals	147
9.3	Convergence of iterative methods	150
9.4	Ill-conditioned system	153

## CHAPTER 10

### NUMERICAL DIFFERENTIATION

10.1	Introduction	157
10.2	Numerical differentiation formulae based on equally spaced data	158
10.2.1	Numerical differentiation based on Newton's forward differences	158
10.2.2	Numerical differentiation based on Newton's backward differences	162
10.2.3	Numerical differentiation based on Stirling's formula	164
10.2.4	Numerical differentiation based on Bessel's formula	168
10.3	Numerical differentiation based on Lagrange's formula	170
10.4	Error analysis of differentiation formulae	172
10.5	Richardson extrapolation	174

## CHAPTER 11

### NUMERICAL INTEGRATION

11.1	Numerical integration	180
11.2	The trapezoidal rule with error term	184
11.3	Simpson's 1/3 rule with error term	187
11.4	Simpson's 3/8 rule with error term	189
11.5	Error estimation	193
11.6	Composite numerical integration	196
11.6.1	Composite trapezoidal rule	197
11.6.2	Composite Simpson's rule	199
11.7	Richardson's extrapolation	204
11.8	Newton-Cotes closed quadrature formulae	209
11.9	The method of undetermined co-efficients	215
11.10	Gaussian quadrature	217

## CHAPTER 12

### NUMERICAL SOLUTION OF ORDINARY DIFFERENTIAL EQUATIONS

12.1	Introduction	225
12.2	Euler's method	227
12.3	Improved Euler's method	231
12.4	Modified Euler method	233
12.5	Runge-Kutta methods	237
12.6	Predictor-Corrector (P-C) methods	244
12.6.1	Milne's method	246
12.6.2	Adams-Moulton method	246
12.7	Convergence and stability	250
12.7.1	Convergence and stability of single-step methods	251
12.7.2	Convergence and stability of multistep methods	255

## CHAPTER 13

### SYSTEMS OF DIFFERENTIAL EQUATIONS AND HIGHER ORDER EQUATIONS

13.1	Introduction	255
13.2	Systems of differential equations	259
13.3	High order equations - Initial value problems	262
13.4	Boundary value problems	264

**APPENDIX A**

- Some important theorems** 268

**APPENDIX B**

- Vector and matrix norms** 272

- References** 277

- Answers to selected exercises** 278

- Index** 290

## CHAPTER 1

### ERRORS IN COMPUTATION

#### 1.1 INTRODUCTION

Almost all calculations with numbers, whether performed by hand or by computer, are subject to numerical errors. This form of error is unavoidable in numerical computation and should not be confused with 'gross errors' arising from programming mistakes or from incorrect design or misuse of algorithms. For example, if we wish to add 1 to 7 we can obtain the exact answer '8', but if we wish to divide 1 by 7 the exact answer is a decimal fraction with an infinite number of figures. We have to be content with an approximate answer, say 0.143, which is in error by about 0.0001428.... We have incurred this error because, in practice, we can only cope with a finite number of significant figures: in this case three. All calculating devices, from pencil-and-paper to the most sophisticated computer, have this restriction. Typically, hand-held calculators are limited to eight or ten significant figures; large, high speed computers can handle a few more.

Whenever calculations are performed there are many possible sources of error, such as inaccurate data and inaccurate formulae, and errors will obviously affect the accuracy of the solution of a particular problem. It is therefore evident that the error in a computed result may be due to one or both of two sources: errors in the data and errors of calculation. Errors of the first type cannot be remedied, but those of the second type can usually be made as small as we please. Thus, when such a number as  $e$  is replaced by its approximate value in a computation, we can decrease the error due to the approximation by taking  $e$  to as many figures as desired.

Nearly all numerical calculations are in some way approximate, and the aim of the computer should be to obtain results consistent with the data with a minimum of labour and with better approximation. The object of the present chapter is to identify the sources of numerical errors, assess the effects on our computations and to set forth some basic ideas and methods relating to

approximate calculations and to give methods for estimating accuracy of the results obtained.

## 1.2 ERRORS AND THEIR SOURCES

The errors involved in mathematical problems may be divided into the following five groups.

1. Errors involved in the statement of the problem. More generally, there are almost always idealizations made in setting up a mathematical model before any computation can begin. The discrepancy between the model and the physical system is a source of error which can be termed as error of the problem. (Mathematical models and their relation to the real world are investigated in applied mathematics, physics, engineering, etc.)

It sometimes happens that it is either difficult or even impossible to solve a given problem when formulated precisely. If that is the case, it is replaced by an approximate problem yielding almost the same results. This is the source of an error regarded as the error of the method.

2. Errors occur due to the presence of infinite processes in mathematical analysis. Many numerical methods express the desired result as the limit of an infinite sequence or the sum of an infinite series. Since, generally speaking, an analysis infinite process cannot be completed in a finite number of steps, we are forced to stop at some term of the sequence and consider it to be an approximation to the required solution. Naturally, such a termination of the process gives rise to an error known as truncation error.

3. Errors associated with the system of numeration. In numerical computation there are fundamental limitations imposed by the finite nature of the computing medium at our disposal. Every possible means of computation is subject to the restrictions of finite space and finite time. Finite space gives rise to rounding error and finite time to truncation error. For example, if we wish to divide 2 by 3 the exact answer is a decimal fraction with an infinite number of figures,  $2/3=0.666\dots$ . We have to be content with an approximate number, say 0.67, which is in error by about -0.00333.... We have incurred this error because, in practice, we can only use a finite number of digits in our computations.

4. Errors due to operations involving approximate numbers (errors of operation). Experimental data involving real numbers (measurements of length, time etc.) are always subject to uncertainty. When performing computations with approximate numbers, we naturally carry (to some extent) the errors

the original data into the final results. In this way, errors of operation are inherent.

5. Errors due to numerical parameters (in formulae) whose values can only be determined approximately. Such, for instance, are all physical constants and mathematical constants such as  $\pi$  and  $e$ . Although such constants are defined precisely in mathematical terms, their representation requires an infinite sequence of digits, for example,

$$\pi = 3.1415926 \dots, e = 2.7182818 \dots$$

To use these in numerical computation we must approximate the values to a finite number of digits, for instance,

$$\pi \approx 3.142 \text{ (rounded to 3 decimal places or 4 significant figures),}$$

$$e \approx 2.71828 \text{ (rounded to 5 decimal places or 6 significant figures).}$$

The consequence is a data error which may be called the initial error.

In a specific problem, quite naturally, some errors are absent and others exert a negligible effect. But, generally, a complete analysis must deal with all types of errors.

### 1.3 APPROXIMATE NUMBERS AND SIGNIFICANT FIGURES

In the discussion of approximate computation, it is convenient to make a distinction between numbers which are exact and those which express approximate values.

**DEFINITION 1.1 (Approximate Numbers):** An approximate value of a number, called approximate number, is defined as a number which is used as an approximation to an exact number and differs only slightly from the exact number for which it stands.

The numbers such as 5,  $1/3$ , 200 etc. are exact numbers because there is no approximation or uncertainty associated with them. Although the numbers such as  $e$ ,  $\pi$ ,  $\sqrt{2}$ , etc. are exact numbers but they cannot be expressed exactly by a finite number of digits. When expressed in decimal form, they must be written as 2.7183, 3.1416, 1.4142, etc. Such numbers are therefore only approximations to the true values and are called approximate numbers.

The difference between the exact number  $x$  and its approximate value  $\bar{x}$  is called **error**, that is,

$$\text{Error} = \text{Exact number} - \text{Approximate number} = x - \bar{x}.$$

**DEFINITION 1.2 (Significant Figures):** Significant figures (digits) of an approximate number is any one of the digits 1, 2, ..., 9, in its decimal representation, or any zero except when it is used to fix the decimal point or to fill the places unknown or discarded digits.

For example, in the number 0.00475 the significant figures are 4, 7, 5; the first three zeros are used merely to fix the position of the decimal point and indicate the place values of the other digits and are therefore not significant. In the number 830, however, all the digits, including zero, are significant figures. In the number 0.0006070 the first four zeros are not significant digits. All the remaining digits, including the other two zeros, are significant figures. If the last digit of 0.0006070 is not significant, then the number must be written as 0.000607. From this point of view, the numbers 0.0006070 and 0.000607 are the same, because the former has four significant figures and the latter only three.

When writing large numbers, the zeros on the right can serve both to indicate the significant figures and to fix the place values of the other digits. This can lead to misunderstanding when the numbers are written in the ordinary way. Consider, for example, the number 675,000. It is not clear how many significant figures there are, although we can say we have at least three. This ambiguity can be removed by using powers-of-ten notation.

(scientific notation) as  $6.75 \times 10^5$  if the number has three significant figures, or as  $6.750 \times 10^5$  if it has four significant figures, or as  $6.7500 \times 10^5$  if the number has five significant figures, etc. Speaking generally, this notation is convenient for numbers containing a large number of non-significant zeros, such as  $0.00000230 = 2.30 \times 10^{-6}$ , and the like.

#### 1.4 ROUNDING OF NUMBERS AND ROUND-OFF ERRORS

If we divide 13 by 7, we obtain

$$13/7 = 1.8571428571\dots$$

a quotient which never terminates. In order to use such a number in a practical computation, we must cut it down to a manageable form, such as 1.86, or 1.857, or 1.85714, or 1.857143, etc. The process of cutting off superfluous digits and retaining as desired is called rounding off.

**ROUNDING-OFF RULE**: significant figure in the nth place, or as place holder. In

- (a) if the first retained digit is even
- (b) if the first retained digit is odd
- (c) if the first non-zero digit is even
- (d) if the first non-zero digit is odd

**EXAMPLE 1.1:** Round off, and three significant figures:

**SOLUTION:** Rounding off, we obtain respectively 3.142, 3.14.

**EXAMPLE 1.2:** Round off, and three significant figures:

**SOLUTION:** We obtain

7.4727, 7.473, 7.473,

Let us now introduce

**DEFINITION 1.3 (Rounding-off Rule):** requires more significant figures than the resulting error must be rounded off. In other words, the result represented to its full accuracy.

If we attempt to divide a fraction with an infinite decimal expansion with an appropriate power of 10, the round-off error of the result will be

case, we do not know some calculating device that the result

**ROUNDING-OFF RULE.** To round off or simply round a number to  $n$  significant figures or digits, drop all digits to the right of the  $n$ th place, or replace them by zeros if the zeros are needed as place holder. In this rounding-off rule, note the following:

- (a) if the first of the discarded digits is less than 5, leave the remaining digits unchanged;
- (b) if the first discarded digit exceeds 5, add 1 to the last retained digit;
- (c) if the first discarded digit is exactly 5 and there are non-zero digits among those discarded, add 1 to the last retained digit;
- (d) if the first discarded digit, however, is exactly 5 and all other discarded digits are zeros, the last retained digit is left unchanged or is increased by 1 according as it is even or odd (the even-digit rule).

**EXAMPLE 1.1:** Round-off the number  $\pi = 3.14159265\dots$  to six, five, four, and three significant figures.

**SOLUTION:** Rounding the number to given significant figures we obtain respectively the approximate numbers 3.14159, 3.1416, 3.142, 3.14.

**EXAMPLE 1.2:** Round the following number to four significant figures:

7.4727, 76346, 15.235 and 15.245

**SOLUTION:** We obtain respectively the approximate numbers

7.473, 76350, 15.24, 15.24.

Let us now introduce the concept of the round-off error.

**DEFINITION 1.3 (Round-off error):** If the result of a calculation requires more significant figures than the number available, then the resulting error is called a round-off error, the exact value must be rounded off to a certain number of significant figures. In other words, the round-off error arises when a number is not represented to its full precision.

If we attempt to divide 1 by 7 the exact answer is a decimal fraction with an infinite number of figures. We have to be content with an approximate answer, say 0.143, which contains a round-off error of about 0.00014. But when, as is usually the case, we do not know the exact value of an answer obtained using some calculating device, we assume the worst. For example, suppose that the result of some calculation worked to three signifi-

cant figures is 0.527. In fact, the answer may be anything between 0.52650... and 0.52749.... In other words, 0.527 may be error due to round-off by  $\pm 0.0005$ , i.e. by five units in the position after the last significant figure which is retained.

**EXAMPLE 1.3:** Find the round-off errors introduced in Example 1.2.

**SOLUTION:** The round-off error introduced in each case of the example is as follows:

$$7.4727 - 7.473 = -0.0003$$

$$76346 - 76350 = -4$$

$$15.235 - 15.24 = -0.005$$

$$15.245 - 15.24 = 0.005$$

## 1.5 TRUNCATION ERROR

Many numerical methods express the desired results as the limit of an infinite sequence or the sum of an infinite series. In practice only a finite number of terms can be computed; the consequence of neglecting the remaining terms is the error known as the truncation error. A truncation error, therefore, leads to an approximate formula being used instead of an exact one. Consider a convergent infinite series

$$S = \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16} + \dots - \frac{1}{2^n} + \dots$$

Since we cannot include all of the infinite number of terms in the series, and since we see that successive terms are contributing progressively smaller and smaller amounts to  $S$ , we conclude that we can truncate the series after some finite number of terms ( $n$ ) and so obtain an approximation to  $S$ , which we might denote by  $S_n$ . The error  $S - S_n$  is called the truncation error in  $S_n$ .

It is important to note that it should be possible to reduce the truncation error by taking more terms in the series. Table 1.1 shows the values of  $S_n$  and the corresponding error for several values of  $n$ . The true value of  $S$  is 1. We can see that the error decreases as  $n$  increases, and that  $S_n$  is converging to a limit.

## 1.6 ABSOLUTE ERROR

**DEFINITION**  
number  $x$ ,

**DEFINITION**  
n decimal  
is the large

For  $\pi = 3.141592653589793$

and we see

**DEFINITION**  
error of an  
absolute error

Thus, if  $E$

number  $\bar{x}$  wh

From this it  
range

For brevity,

**Table 1.1 Truncation Error**

n	$S_n$	Error in $S_n$
3	0.875	0.125
5	0.96875	0.03125
10	0.99902	0.00098
15	0.99997	0.00003

## 1.6 ABSOLUTE, RELATIVE, AND PERCENTAGE ERRORS

**DEFINITION 1.4 (Absolute Error):** If  $\bar{x}$  is an approximation to a number  $x$ , then the absolute error  $E_a$  of  $\bar{x}$  is

$$E_a = |x - \bar{x}|. \quad (1.1)$$

**DEFINITION 1.5:** An approximation  $\bar{x}$  to  $x$  is said to be correct to  $n$  decimal places or positions (or agree to  $n$  decimal places) if  $n$  is the largest non-negative integer for which

$$|x - \bar{x}| < 0.5 \times 10^{-n} \quad (1.2)$$

For  $\pi = 3.14159\dots$  and  $\bar{\pi} = 3.142$ , the absolute error is

$$|\pi - \bar{\pi}| = 0.41 \times 10^{-3} < 0.5 \times 10^{-3},$$

and we see that  $\bar{\pi}$  is indeed correct to 3 decimal places.

**DEFINITION 1.6 (Limiting Absolute Error):** The limiting absolute error of an approximate number is any number not less than the absolute error of that number.

Thus, if  $E_{la}$  is the limiting absolute error of an approximate number  $\bar{x}$  which takes the place of the exact number  $x$ , then

$$E_a = |x - \bar{x}| \leq E_{la}. \quad (1.3)$$

From this it follows that the exact number  $x$  lies within the range

$$\bar{x} - E_{la} \leq x \leq \bar{x} + E_{la} \quad (1.4)$$

For brevity, we can then write  $x = \bar{x} \pm E_{la}$ .

**EXAMPLE 1.4:** Determine the limiting absolute error of the number  $\bar{x} = 3.14$  which is used instead of the number  $\pi$ .

**SOLUTION:** It follows from the inequality

$$3.14 < \pi < 3.15$$

that  $|x - \bar{x}| = |3.15 - 3.14| = 0.01$ ,

we can take  $E_{la} = 0.01$ .

We have a better estimate from the inequality

$$3.14 < \pi < 3.142$$

that

$$E_{la} = 0.002.$$

The absolute error (or the limiting absolute error) is not sufficient to describe the accuracy of a computation. Suppose that in measuring the lengths of two rods we get  $l_1 = 105.8\text{cm} \pm 0.1\text{cm}$  and  $l_2 = 200\text{cm} \pm 0.1\text{cm}$ . Despite the fact that the limiting absolute errors coincide, the second measurement is better than the first one. For this reason, an essential point in the accuracy of measurement is to define the relative error.

**DEFINITION 1.7 (Relative Error):** If  $\bar{x}$  is an approximation to  $x$ , then the relative error  $E_r$  of  $x$  is given by

$$E_r = \frac{|x - \bar{x}|}{|x|}, \quad (x \neq 0),$$

or

$$E_r = \frac{E}{|x|}. \quad (1.5)$$

**EXAMPLE 1.5:** a) If  $x = .3000 \times 10$  and  $\bar{x} = .3100 \times 10$ , the absolute error is .1 and the relative error is  $.3333 \times 10^{-1}$ .

b) If  $x = .3000 \times 10^{-3}$  and  $\bar{x} = .3100 \times 10^{-3}$  the absolute error is  $.1 \times 10^{-4}$  and the relative error is  $.3333 \times 10^{-1}$ .

c) If  $x = .3000 \times 10^4$  and  $\bar{x} = .3100 \times 10^4$  the absolute error is  $.1 \times 10^3$

and the relative error is  $.3333 \times 10^{-1}$ .

This example shows that the same relative error,  $.3333 \times 10^{-1}$ , occurs for widely varying absolute errors. Consequently, as a measure of accuracy, the absolute error may be misleading and the relative error more meaningful.

**DEFINITION 1.8:** An approximation  $\bar{x}$  to  $x$  is said to be correct to  $n$  significant figures (or digits) if  $n$  is the largest non-negative integer for which

$$|x - \bar{x}| \leq \frac{1}{2} 10^{m-n+1} \quad (1.6)$$

where  $m$ , an integer, is the highest power of ten when  $\bar{x}$  is represented in the form of decimal. For example, the number  $\bar{x} = 36.00$  is an approximation to the exact number  $x = 35.97$  correct to three figures, since

**PROOF:** Let  $x$  be the given approximate number. Consider  $|x - \bar{x}| = 0.03 < \frac{1}{2} \times 10^{1-3+1} = \frac{1}{2} \times 0.1$  ad of si If

**DEFINITION 1.9:** The limiting relative error  $E_{lr}$  of a given approximate number  $\bar{x}$  is any number not less than the relative error of that number. Thus, by definition we have

$$E_r \leq E_{lr} \quad (1.7)$$

or to the sum of  $\frac{E_a}{|x|} \leq E_{lr}$

and hence  $E_a \leq |x| E_{lr}$

Thus, for the limiting absolute error of a number  $\bar{x}$  we can take

$$E_{la} = |x| E_{lr} \quad (1.8)$$

Since, in practical situations,  $x \approx \bar{x}$ , in place of (1.8) one frequently uses

$$E_{la} = |\bar{x}| E_{lr} \quad (1.8')$$

**EXAMPLE 1.6:** The weight of  $1 \text{ dm}^3$  of water at  $0^\circ\text{C}$  is given as  $p = 999.847 \text{ gf} \pm 0.001 \text{ gf}$  ( $\text{gf} = \text{gram}(force)$ ). Determine the limiting relative error of the result of weighing the water.

**SOLUTION:** We have

$$E_{la} = 0.001 \text{ gf} \quad \text{and} \quad \bar{x} = p \leq 999.847 \text{ gf.}$$

Thus

$$E_{lr} = \frac{E_{la}}{\bar{x}} = \frac{0.001}{999.847} \approx 10^{-4}\%.$$

**DEFINITION 1.10:** The percentage error is 100 times the relative error. Thus

$$\text{percentage error of an approximate quantity} = 100E_r.$$

It is to be noted that relative and percentage errors are independent of the unit of measurement, whereas absolute errors are expressed in terms of the unit used.

### EXERCISE 1.1

- Round-off the following number to four significant figures:  
11.64498, 81.9772, 48.365, 67.495, 4.4995002, 39.63244.
- Round-off the number  $\pi = 3.1415926535\dots$  to three, four and five significant figures and also find the absolute error in each case.
- What will 7.52, 6.345, 7.4727, 76340, 15.235, 15.245, and 122377 become when rounded to one digit?
- What will 2.73235, 35.671 and 430050 become when rounded to two digits?
- What will 73.6547, 180.273 and 541500499 become when rounded to three digits?
- Find round-off errors in Ex. 3, 4 and 5.
- Find the relative error when  $c = 485165195.4\dots$  ( $c = e^{20}$ ) is rounded to 4 significant figures.

8. Round-off the following numbers:

- a) 8632574 to 4 significant figures
- b) 3.1415926 to 5 decimal places
- c) 8.5250 to 2 decimal places
- d) 1.6750 to 2 decimal places.

Find round-off error in each case.

9. A machinist measures a 0.5 inch bolt to the nearest thousandth of an inch, and a carpenter a 10 feet beam to the nearest eighth of an inch. Which measurement is more accurate?

### 1.7 THE ERROR OF A SUM

**THEOREM 1.1:** The absolute error of an algebraic sum of several approximate numbers does not exceed the sum of the absolute errors of the numbers.

**PROOF:** Let  $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$  be the given approximate numbers. Consider their algebraic sum

$$u = \pm \bar{x}_1 \pm \bar{x}_2 \pm \dots \pm \bar{x}_n$$

Then  $\Delta u = \pm \Delta \bar{x}_1 \pm \Delta \bar{x}_2 \pm \dots \pm \Delta \bar{x}_n$

We consider the  $\Delta \bar{x}_i$  to be the same for all terms. Then we have  
and, hence  $|\Delta u| = E_{\text{la}} \leq |\Delta \bar{x}_1| + |\Delta \bar{x}_2| + \dots + |\Delta \bar{x}_n|$ . (1.9)

**COROLLARY 1.1:** The limiting absolute error of an algebraic sum is equal to the sum of the limiting absolute errors of the terms. Thus

$$E_{\text{la}} = E_{1\text{la}} + E_{2\text{la}} + \dots + E_{n\text{la}} \quad (1.10)$$

Following is the practical rule for the addition of approximate numbers.

#### RULE FOR ADDITION:

- i) Find the approximate numbers with the least number of decimal places and leave them unchanged.
- ii) Round off the remaining approximate numbers, retaining one or two more decimal places than those with the small-

12

est number of decimals.

- iii) Add the numbers, taking into account all retained decimals.
- v) Round off the result, reducing it by one decimal.

**EXAMPLE 1.7:** Find the sum of the approximate numbers 0.348, 0.1834, 345.4, 235.2, 11.75, 9.27, 0.0849, 0.0214, 0.000354, each correct to the indicated significant figures.

**SOLUTION:** We find the least accurate numbers 345.4 and 235.2 whose absolute error may attain 0.1. Rounding the remaining numbers to two decimal places, we obtain

SUM A TO THE ERROR OF

345.4

235.2

11.75

9.27

0.08

0.02

0.35

0.18

0.00

602.25

602.25

Round the result to one decimal place by the even-digit rule. We get the approximate value of the sum; 602.2.

The total absolute error  $E$  of the result is made up of three terms:

- i) the sum of limiting errors of the original data:

$$E_1 = 10^{-3} + 10^{-4} + 10^{-1} + 10^{-1} + 10^{-2} + 10^{-2} + 10^{-4} + 10^{-4} + 10^{-6}$$

$$= 0.221301 < 0.222.$$

- ii) the absolute value of the sum of the rounding errors of the terms:

$$E_2 = |0.0049 + 0.0014 - 0.002 + 0.0034 + 0.000354| = 0.008054 < 0.009.$$

- iii) the final rounding error of the result:

$$E_3 = |602.25 - 602.2| = 0.05.$$

Hence

and the

**EXAMPLE**

86939,

nificant

**SOLUTI**

hundred

give the

- 1000000

division

**EXAMPLE**

sum

all

**Round**

sum; 1545

1.8 THE E

rror

We cons

$\bar{x}_1, \bar{x}_2$ .

From e

difference

The limit

where  $x$  is  
erence betw

When one  
they must b  
tion.

$$\text{Hence } E_a = E_{1a} + E_{2a} + E_{3a} = 0.222 + 0.009 + 0.050 = 0.281 < 0.3$$

and thus the desired sum is  $602.2 \pm 0.3$ .

**EXAMPLE 1.8:** Find the sum of 36490, 994, 557.32, 295000, and 86939, assuming that the number 29500 is known to only three significant figures.

**SOLUTION:** Since the number 29500 is known only to the nearest hundred, we round off the others to the nearest ten, and add to give the sum as shown below:

$$\begin{array}{r} 29500 \\ 86940 \\ 36490 \\ 990 \\ 560 \\ \hline 154480 \end{array}$$

Round the sum to hundred, we get the approximate value of the sum; 154500 or  $1.545 \times 10^5$ .

### 1.8 THE ERROR OF A DIFFERENCE

We consider the difference  $u = \bar{x}_1 - \bar{x}_2$  of two approximate numbers  $\bar{x}_1, \bar{x}_2$ .

From equation (1.10), the limiting absolute error  $E_{1a}$  of the difference is

$$E_{1a} = E_{1la} + E_{2la}.$$

The limiting relative error of the difference is

$$E = \frac{E_{1a}}{x} = \frac{E_{1la} + E_{2la}}{x}$$

where  $x$  is the exact value of the absolute magnitude of the difference between the numbers  $\bar{x}_1$  and  $\bar{x}_2$ .

When one approximate number is to be subtracted from another, they must both be rounded off to the same place before subtraction.

**EXAMPLE 1.9:** Subtract 46.365 from 779.8, assuming that each number is approximate and correct only to its last figure.

**SOLUTION:** Since 779.8 is correct to one decimal place, therefore, we round off the number 46.365 to one decimal place to obtain 46.4. Thus

$$779.8 - 46.4 = 733.4.$$

### 1.9 THE ERROR OF A PRODUCT

**THEOREM 1.2:** The relative error of a product of several approximate nonzero numbers does not exceed the sum of the relative errors of the numbers.

**PROOF:** Let

$$u = \bar{x}_1 \cdot \bar{x}_2 \cdots \bar{x}_n$$

If all the approximate numbers  $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$  are positive, then

$$\ln u = \ln \bar{x}_1 + \ln \bar{x}_2 + \dots + \ln \bar{x}_n$$

Using the approximate formula  $\Delta \ln \bar{x} \approx d \ln \bar{x} = \frac{\Delta \bar{x}}{\bar{x}}$ , we obtain

$$\frac{\Delta u}{u} = \frac{\Delta \bar{x}_1}{\bar{x}_1} + \frac{\Delta \bar{x}_2}{\bar{x}_2} + \dots + \frac{\Delta \bar{x}_n}{\bar{x}_n}$$

Thus

$$\left| \frac{\Delta u}{u} \right| \leq \left| \frac{\Delta \bar{x}_1}{\bar{x}_1} \right| + \left| \frac{\Delta \bar{x}_2}{\bar{x}_2} \right| + \dots + \left| \frac{\Delta \bar{x}_n}{\bar{x}_n} \right|.$$

If  $x_i$  ( $i=1, 2, \dots, n$ ) are the exact values of the factors  $\bar{x}_i$ , we can write approximately

$$\left| \frac{\Delta \bar{x}_i}{\bar{x}_i} \right| \approx \left| \frac{\Delta x_i}{x_i} \right| = E_i \quad \text{and} \quad \left| \frac{\Delta u}{u} \right| = E_r$$

where  $E_i$  are the relative errors of the factors  $\bar{x}_i$  ( $i=1, 2, \dots, n$ ).

and  $E_r$  is the relative error of the product.

Consequently  $E_r \leq E_{1_r} + E_{2_r} + \dots + E_{n_r}$  (1.11)

It is clear that (1.11) also holds true if the factors  $\bar{x}_i$  ( $i = 1, 2, \dots, n$ ) have different signs.

**COROLLARY 1.2:** The limiting relative error of a product is equal to the sum of the limiting relative errors of the factors:

$$E_{lr} = E_{1_{lr}} + E_{2_{lr}} + \dots + E_{n_{lr}} \quad (1.12)$$

**EXAMPLE 1.10:** Find the product  $349.1 \times 863.4$  and state how many figures of the result are trustworthy if the factors are correct to all written figures.

**SOLUTION:** We have

$$E_{1_{la}} = \frac{1}{2} 10^{2-4+1} = 0.05 \text{ and } E_{2_{la}} = \frac{1}{2} \cdot 10^{2-4+1} = 0.05.$$

Hence

$$E_{lr} = \frac{0.05}{349.1} + \frac{0.05}{863.4} = .000143 + .0000579 = 0.00020.$$

Now,  $u = 349.1 \times 863.4 = 301412.94 \approx 301413$ .

The limiting absolute error of this product is

$$E_{la} = u E_{lr} = 301413 \times 0.00020 = 60.2826 \approx 60.$$

And so  $u$  is correct to only three figures and the result should be written as

$$u = 301413 \pm 60.$$

Thus the true result lies between 301473 and 301353. There is some uncertainty about the fourth figure.

It is interesting to note that when an approximate number is multiplied by an exact number  $k$ , the limiting relative error remains unchanged, while the limiting absolute error is increased by  $|k|$  times.

**REMARKS:** When it is desired to find the product of two or more approximate numbers of different accuracies, it is sufficient to:

1. round them off so that each contains one or two significant figures more than the least accurate factor;
2. in the final result, retain as many significant figures as there are correct figures in the least accurate factor (or keep one extra digit).

EXAMPLE 1.

56.3 / 45 , a  
gure but no**SOLUTION:**E<sub>1</sub>

Since

**EXAMPLE 1.II:** Find the product of the approximate numbers  $\bar{x}_1 = 2.5$  and  $\bar{x}_2 = 72.397$  correct to the number of digits written.

then

**SOLUTION:** Rounding  $\bar{x}_2$  to one decimal place, we have

$$\bar{x}_1 = 2.5, \quad \bar{x}_2 = 72.4. \text{ Thus}$$

$$\bar{x}_1 \bar{x}_2 = 2.5 \times 72.4 = 181 \approx 1.8 \times 10^2.$$

As this error

we take 25.2

**REMARKS:** If accurate than ded off so less accurate cant figures

## 1.10 THE ERROR OF A QUOTIENT

## 1.11 THE ERR

68.27

Here  $u = a$ 

and as to

as in the

Thus, we have

0.3862

correct 21 as

how many

Hence, the li

is m times the

true

EXAMPLE 1.

(0.3862)<sup>4</sup>, ass

its last figure

SOLUTION: Let

number a is

then

$$\ln u = \ln a - \ln b$$

and

$$\frac{\Delta u}{u} = \frac{\Delta a}{a} - \frac{\Delta b}{b}$$

Hence

$$\left| \frac{\Delta u}{u} \right| \leq \left| \frac{\Delta a}{a} \right| + \left| \frac{\Delta b}{b} \right|. \quad (1.13)$$

Thus, we have established the following theorem:

**THEOREM 1.3:** The relative error of a quotient does not exceed the sum of the relative errors of the dividend and the divisor.

**COROLLARY 1.3:** If  $u = \frac{a}{b}$ , then

$$E_{\frac{u}{lr}} = E_{\frac{a}{lr}} + E_{\frac{b}{lr}} = \text{limiting relative error of dividend} + \text{limiting relative error of divisor}$$

**EXAMPLE 1.12:** Find the number of correct figures in the quotient  $56.3/\sqrt{5}$ , assuming that the numerator is correct to its last figure but no further.

**SOLUTION:** We take  $\sqrt{5} = 2.236$ . Thus

$$E_{lr} = \frac{0.05}{56.3} + \frac{0.0005}{2.236} = .000888 + .000224 = 0.001$$

Since

$$u = \frac{56.3}{2.236} = 25.18 = 25.2,$$

then

$$E_{ua} = ux E_{lr} = 25.18 \times 0.001 = 0.025.$$

As this error does not affect the third figure of the quotient, we take 25.2 as the correct result.

**REMARKS:** If one of the numbers (divisor or dividend) is more accurate than the other, the more accurate number should be rounded off so as to contain one more significant figure than the less accurate one. The result should be given to as many significant figures as the less accurate number (or one digit more).

### 1.11 THE ERROR OF A POWER

Here  $u = a^m$  ( $m$  any natural number). Then  $\ln u = m \ln a$ , and

$$\left| \frac{\Delta u}{u} \right| = m \left| \frac{\Delta a}{a} \right|.$$

Thus, we have

$$E_{u_{lr}} = m E_{a_{lr}}. \quad (1.14)$$

Hence, the limiting relative error of the  $m$ th power of a number is  $m$  times the limiting relative error of the number.

**EXAMPLE 1.13:** Find the number of trustworthy figures in  $(0.3862)^4$ , assuming that the number in parentheses is correct to its last figure but no further.

**SOLUTION:** Let  $a = 0.3862$ . Then the limiting relative error of the number  $a$  is

$$E_{a_{lr}} = \frac{0.00005}{0.3862} \approx 0.00013.$$

18

The limiting relative error of the result is

$$E_{lr} = 4 \times E_{a_{lr}} = 4 \times 0.00013 = 0.00052.$$

$$\text{Also, } (0.3862)^4 = 0.022246.$$

$$\text{Therefore, } E_{la} = 0.022246 \times 0.00052 = 0.0000116.$$

Since this error affects the fourth significant figure of the result, the best we can do is to write.

$$(0.3862)^4 = 0.02225$$

### 1.12 THE ERROR OF A ROOT

Suppose  $u = \sqrt[m]{a}$ , then  $\ln u = \frac{1}{m} \ln a$  and, hence

$$\left| \frac{\Delta u}{u} \right| = \frac{1}{m} \left| \frac{\Delta a}{a} \right|$$

### 1.13 THE ERROR OF A QUOTIENT

Thus

$$E_{u_{lr}} = \frac{1}{m} E_{a_{lr}}. \quad (1.15)$$

That is, the limiting relative error of the  $m$ th root of an approximate number is only  $(1/m)$ th of the limiting relative error of the given number (i.e. radicand).

**EXAMPLE 1.14:** Find the number of trustworthy figures in  $(0.3862)^{1/4}$ , assuming that the number in parentheses is correct to its last figure but no farther.

**SOLUTION:** Here

$$E_{lr} = \frac{1}{4} (0.00013) = 0.000032, \quad (0.3862)^{1/4} = 0.78832.$$

Thus

$$E_{la} = 0.78832 \times 0.000032 = 0.000026.$$

Hence the fourth root 0.7883 is correct to four figures.

## EXERCISE 1.2

1. Find the sum of the approximate numbers 561.32, 86.954, 3.9462, and 491.6, each being correct to its last figure but no farther.
2. Compute the difference  $\sqrt{2.03} - \sqrt{2}$  correct to five significant figures.
3. Determine the product of the approximate numbers 12.2 and 73.56 and the number of correct digits in the product if the factors are correct to all written digits.
4. Evaluate the following expressions. Then, assuming maximum round-off error in each figure, estimate the limits between which your calculated value lies. Finally quote each result to as many significant figures as are reliable.

$$\text{a) } \frac{2.362 \times 1.76}{1.46 \times 0.785} \quad \text{b) } 7.62 + \frac{54}{15.3} \quad \text{c) } \frac{1.32 - 0.463}{16.1 \times 2.17}$$

- (1.15) 5. Find the sum of the approximate numbers 136.421, 28.3, 321, 68.243, 17.482 if the numbers are correct to their last digits only.
6. Determine the relative error and the number of correct figures in the product of the numbers 93.87 and 9.236.
7. The number 6852.4 and 48.392 are both approximate and true only to their last digits. Find their difference and state how many figures in the result are trustworthy.
8. Find the number of trustworthy figures in the quotient of  $876.3/494.2$ , assuming that both numbers are approximate and true only to the number of digits given.
9. Find the number of correct digits in the quotient  $u=25.7:3.6$  assuming the dividend and divisor are correct to the last digit given.
10. Evaluate the following expressions for the given values of  $x$ , first putting each in an appropriate form and estimate the maximum error in each result, assuming that all integer coefficients are exact but that all decimals are subject to the maximum round-off error.

## CHAPTER 2

### LINEAR OPERATORS

#### 2.1 INTRODUCTION

Let  $F$  be a linear real function space, that is,  $F$  is a set of functions  $f, g, \dots$  with the property that  $f \in F$  and  $g \in F$  implies  $\alpha f + \beta g \in F$ . Here  $\alpha$  and  $\beta$  are arbitrary constants. Let us define the following operators which map  $F$  into itself:

$$E f(x) = f(x + h) \quad (2.1)$$

$$\Delta f(x) = f(x + h) - f(x) \quad (2.2)$$

$$-\nabla f(x) = f(x) - f(x - h) \quad (2.3)$$

$$\delta f(x) = f(x + h/2) - f(x - h/2) \quad (2.4)$$

$$\mu f(x) = (1/2)[f(x + h/2) + f(x - h/2)] \quad (2.5)$$

The operators  $E, \Delta, \nabla, \delta$  and  $\mu$  are respectively called the **shifting operator**, the **forward difference operator**, the **backward difference operator**, the **central difference operator** and the **mean value operator**. All these operators are linear i.e. for any two functions  $f$  and  $g$ , arbitrary constants  $\alpha$  and  $\beta$  and  $P$  any of the above operators, we have

$$\{ P(\alpha f + \beta g) = \alpha Pf + \beta Pg \}$$

When  $P$  denotes any of the above operators, sums, differences and products of these operators are defined in the usual manner. An operator  $Q$  is said to be inverse of  $P$  if

$$QPf = f \quad \text{for all } f \in F.$$

The inverse of operator  $P$  is denoted by  $P^{-1}$ .

Thus

$$E^{-1}f(x) = f(x - h)$$

We also define

$$E^{1/2}f(x) = f(x + h/2)$$

$$E^{-1/2}f(x) = f(x - h/2)$$

or we v

Two operators  $P_1$  and  $P_2$  are said to be equal if

$$P_1 f = P_2 f \quad \text{for all } f \in F$$

where e  
Similarly

**EXAMPLE 2.1:** Find  $\Delta e^x$ .

**SOLUTION:**

$$\Delta e^x = e^{x+h} - e^x$$

$$= (e^h - 1)e^x.$$

**EXAMPLE 2.2:** Find  $\Delta E^{-1}(\sin x)$ .

**SOLUTION:**

$$\Delta E^{-1} \sin x = \Delta \sin(x - h)$$

$$= \sin(x - h + h) - \sin(x - h)$$

$$= \sin x - \sin(x - h)$$

$$= 2 \sin \frac{(x - x + h)}{2} \cos \frac{(x + x - h)}{2}$$

$$= 2 \sin \frac{h}{2} \cos \frac{(2x - h)}{2}$$

etc. An c  
ingful on  
tion  $f \in F$ If  $F$  is th  
ingful. How  
since

**EXAMPLE 2.3:** Find  $\mu^2 f(x)$ .

$$\begin{aligned} \text{SOLUTION: } \mu^2 f(x) &= \mu(\mu f(x)) = \mu(1/2)[f(x + h/2) + f(x - h/2)] \\ &= (1/2)[\mu[f(x + h/2) + f(x - h/2)]] \\ &= (1/2)[(1/2)(f(x+h)+f(x))+(1/2)(f(x)+f(x-h))] \\ &= (1/4)[f(x + h) + 2 f(x) + f(x - h)]. \end{aligned}$$

and the seri

## 2.3 OPERATOR

Any of the  
other operato

## 2.2 FUNCTIONS OF OPERATORS

Suppose every  $f \in F$  can be expanded in Taylor series:

$$f(x+h) = f(x) + hf'(x) + (h^2/2)f''(x) + \dots \quad (2)$$

Denoting the derivative operator by  $D$ , we can write (2.6) in the form

$$E f(x) = (1 + hD + (h^2/2)D^2 + \dots) f(x)$$

we have

Thus we have the operator identity

$$E = 1 + hD + (h^2/2)D^2 + \dots$$

or we write

$$E = e^{hD}$$

where  $e^{hD}$  is defined by means of the power series in (2.7). Similarly one can define, for any operator  $P$ ,

$$\sin P = P - \frac{P^3}{3!} + \frac{P^5}{5!} - \dots$$

$$\log(1+P) = P - \frac{P^2}{2} + \frac{P^3}{3} - \dots$$

etc. An operator defined by means of an infinite series is meaningful only if the series obtained when it operates on any function  $f \in F$  is convergent. For example, let  $Q$  be defined by

$$Q = 1 + D + \frac{D^2}{2} + \frac{D^3}{3} + \dots$$

If  $F$  is the space of polynomials, the above definition is meaningful. However if  $e^x \in F$  then  $Qe^x$  does not have any meaning, since

$$Qe^x = e^x(1 + 1 + 1/2 + 1/3 + \dots),$$

and the series on the right side does not converge.

### 2.3 OPERATOR IDENTITIES

Any of the above operators can be expressed in terms of any other operator. For example the relations

$$\underline{\Delta} = E - 1 \quad (2.8)$$

$$\underline{\nabla} = 1 - E^{-1} \quad (2.9)$$

$$\delta = E^{1/2} - E^{-1/2} \quad (2.10)$$

$$\mu = (1/2)[E^{1/2} + E^{-1/2}] \quad (2.11)$$

follow from the definition of these operators. From

$$\begin{aligned} \delta &= E^{1/2} - E^{-1/2}, \\ \text{we have } \delta^2 &= (E^{1/2} - E^{-1/2})(E^{1/2} - E^{-1/2}) \\ &= E + E^{-1} - 2 \end{aligned} \quad (2.12)$$

and

$$\mu = (1/2)[E^{1/2} + E^{-1/2}],$$

gives

$$\mu^2 = (1/4)[E + E^{-1} + 2]. \quad (2.13)$$

Thus

$$\delta^2 + 2 = 4\mu^2 - 2.$$

This gives the useful relation

$$\mu^2 = 1 + \frac{\delta^2}{4} \quad (2.14)$$

**EXAMPLE 2.4:** Prove that  $\delta^2 = \Delta^2 (1 + \Delta)^{-1}$ .**SOLUTION:**  $\delta = E^{1/2} - E^{-1/2}$ 

Therefore

$$\delta^2 = E + E^{-1} - 2$$

But

$$E = \Delta + 1,$$

hence

$$\delta^2 = \Delta + 1 + \frac{1}{\Delta+1} - 2$$

$$= \frac{(\Delta + 1)^2 + 1 - 2(\Delta + 1)}{\Delta + 1}$$

$$= \Delta^2 / (\Delta + 1).$$

**EXAMPLE 2.5:** Show that  $\mu^2 = (1 + \Delta/2)^2 (1 + \Delta)^{-1}$ .**SOLUTION:** Since  $\mu = (1/2)[E^{1/2} + E^{-1/2}]$ ,

therefore

$$\mu^2 = (1/4)[E + E^{-1} + 2]$$

But

$$E = 1 + \Delta,$$

therefore

$$\mu^2 = (1/4)[1 + \Delta + \frac{1}{1+\Delta} + 2]$$

$$= (1/4)[\frac{(1 + \Delta)^2 + 1 + 2(1 + \Delta)}{1 + \Delta}]$$

$$= (1/4)[2 + \Delta]^2 / (1 + \Delta)$$

$$= (1 + \Delta/2)^2 (1 + \Delta)^{-1}.$$

**EXAMPLE 2.6:** Prove that  $E = 1 + \frac{\delta^2}{2} + \delta \left[ 1 + \frac{\delta^2}{4} \right]^{1/2}$ .**SOLUTION:** We know  $\delta = E^{1/2} - E^{-1/2}$ .

and

$$\mu = (1/2)[E^{1/2} + E^{-1/2}]$$

which give

$$2E^{1/2} = \delta + 2\mu = \delta + 2 \left[ 1 + \frac{\delta^2}{4} \right]^{1/2}.$$

On squaring, we get the desired result.

## 2.4 DIFFERENCE OPERATORS AND THE DERIVATIVE OPERATORS

(2.13)

From the relation

$$E = e^{hD}$$

we have

(2.14)

$$\Delta = E - 1 = e^{hD} - 1,$$

$$\nabla = 1 - E^{-1} = 1 - e^{-hD},$$

$$\delta = E^{1/2} - E^{-1/2} = e^{hD/2} - e^{-hD/2}, \\ = 2 \sinh(hD),$$

$$\mu = (1/2)[E^{1/2} + E^{-1/2}] = \cosh(hD).$$

Later on, in chapter 10, we shall derive the following formulae:

$$hD = \Delta - \frac{\Delta^2}{2} + \frac{\Delta^3}{3} - \frac{\Delta^4}{4} + \dots \quad (2.15)$$

$$h^2 D^2 = \Delta^2 - \Delta^3 + \frac{11}{12} \Delta^4 - \dots \quad (2.16)$$

$$hD = \nabla + \frac{\nabla^2}{2} + \frac{\nabla^3}{3} + \frac{\nabla^4}{4} + \dots \quad (2.17)$$

$$h^2 D^2 = \nabla^2 + \nabla^3 + \frac{11}{12} \nabla^4 + \dots \quad (2.18)$$

## EXERCISE 2

1. Evaluate  $\underline{\Delta(2^x)}$ ,  $\nabla(\cos x)$ ,  $\mu\delta(x)$ .

2. Prove the following identities:

$$(a) \quad \delta^2 = \nabla^2(1 - \nabla)^{-1} = \Delta^2(1 + \Delta)^{-1}$$

$$(b) \quad \mu\delta = (1/2)[E - E^{-1}] = (1/2)[\Delta + \nabla]$$

$$(c) \quad \mu^2 = [1 - \nabla/2]^2(1 - \nabla)^{-1}$$

$$(d) \quad \Delta = (1/2)\delta^2 + \delta \left[ 1 + \frac{\delta^2}{4} \right]^{1/2}$$

$$(e) \quad \nabla = - (1/2)\delta^2 + \delta \left[ 1 + \frac{\delta^2}{4} \right]^{1/2}$$

$$(f) \quad \mu\delta = \Delta(1 + \Delta/2)(1 + \Delta)^{-1}$$

$$(g) \quad 1 + \mu^2 \delta^2 = (1 + \delta^2/2)^2$$

$$(h) \quad (1 + \Delta/2)(1 + \Delta)^{-1/2} = (1 - \nabla/2)(1 - \nabla)^{-1/2}$$

## CHAPTER 3

### EIGENVALUES AND EIGENVECTORS OF MATRICES

#### 3.1 INTRODUCTION

Eigenvalues play a very important role in obtaining deep understanding of many physical system. For example, the stability of aircraft in flight and the modes of vibration of a bridge are governed by eigenvalues. They are also important in the study of the growth of certain types of population. Suppose  $X$  denotes the number of females of a certain age in a given population. If the population is in the state of age stability, the vector  $X$  satisfies an equation of the form

$$AX = \lambda X$$

where  $A$  is a square matrix and  $\lambda$  is a positive real number which indicates whether the population increases or not. If  $\lambda < 1$  the population decreases, If  $\lambda > 1$  it increases, and if  $\lambda = 1$  it remains the same.

In fact, for a given square matrix  $A$  there are only a limited number of values of  $\lambda$  which lead to a non-trivial solution. These are known as the **eigenvalues** of the matrix and the corresponding solution vectors are known as the **eigenvectors** of  $A$ .

#### 3.2 EIGENVALUES AND EIGENVECTORS

Let  $A$  be a square matrix of order  $n$ . A number  $\lambda$  is called an eigenvalue or characteristic value of  $A$  if an  $n \times 1$  non zero column vector  $X$  can be found such that

$$AX = \lambda X. \quad (3.1)$$

The above equation may also be written as

$$(A - \lambda I)X = 0$$

If  $X^T = (x_1, x_2, \dots, x_n)$ , then we have a system of  $n$  equations in  $n$  unknowns:

$$(a_{11} - \lambda)x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n = 0,$$

$$a_{21}x_1 + (a_{22} - \lambda)x_2 + \dots + a_{2n}x_n = 0,$$

.....

.....

$$a_{n1}x_1 + a_{n2}x_2 + a_{n3}x_3 + \dots + (a_{nn} - \lambda)x_n = 0.$$

It is well-known that a non-trivial solution exists if and only if the determinant of the system vanishes i.e.,

$$\begin{vmatrix} a_{11} - \lambda & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} - \lambda & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} - \lambda \end{vmatrix} = 0 \quad (3.2)$$

This gives an equation of degree  $n$  in  $\lambda$ :

$$\lambda^n - (a_{11} + a_{22} + \dots + a_{nn})\lambda^{n-1} + \dots + (-1)^n k = 0, \quad (3.3)$$

where  $k$  is yet to be found. It is clear that if we put  $\lambda = 0$  in the left hand side of (3.2), we are left with  $\det A$ . Hence  $k$  in (3.3) must be equal to  $\det A$ . The equation

$$\lambda^n - (a_{11} + a_{22} + \dots + a_{nn})\lambda^{n-1} + \dots + (-1)^n \det A = 0, \quad (3.4)$$

is called the **characteristic equation** of the matrix  $A$  and any root of this equation will be an eigenvalue of  $A$ . Thus a matrix of order  $n$  has  $n$  eigenvalues (although all of them may not be distinct).

**EXAMPLE 3.1:** Find eigenvalues and corresponding eigenvectors of the matrix

$$A = \begin{bmatrix} 2 & 3 \\ 4 & 3 \end{bmatrix}$$

$$\text{SOLUTION: } |A - \lambda I| = \begin{vmatrix} 2-\lambda & 3 \\ 4 & 3-\lambda \end{vmatrix} = (2 - \lambda)(3 - \lambda) - 12 \\ = \lambda^2 - 5\lambda - 6.$$

Thus  $\lambda^2 - 5\lambda - 6 = 0$  is the characteristic equation of  $A$  and  $\lambda = 6, -1$  are the two eigenvalues. To find eigenvectors we have to solve

$$(2 - \lambda)x_1 + 3x_2 = 0, \\ 4x_1 + (3 - \lambda)x_2 = 0.$$

First if  $\lambda = 6$ , then the system becomes

$$4x_1 - 3x_2 = 0, \\ 4x_1 - 3x_2 = 0.$$

So  $x_2$  may be chosen arbitrarily. Let  $x_2 = 1$ , then  $x_1 = 3/4$ . So

$$x_1 = \begin{bmatrix} 3/4 \\ 1 \end{bmatrix}$$

is an eigenvector corresponding to  $\lambda_1 = 6$ . Similarly we find that

$$x_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

is an eigenvector corresponding to  $\lambda_2 = -1$ .

**EXAMPLE 3.2:** Find the eigenvalues and the corresponding eigenvectors of

$$A = \begin{bmatrix} -2 & 6 & -24 \\ 0 & -3 & 10 \\ 1 & -4 & 13 \end{bmatrix}$$

$$\text{SOLUTION: } |A - \lambda I| = \begin{vmatrix} -2-\lambda & 6 & -24 \\ 0 & -3-\lambda & 10 \\ 1 & -4 & 13-\lambda \end{vmatrix} = -\lambda^3 + 8\lambda^3 - 5\lambda - 14.$$

The characteristic equation is

$$\lambda^3 - 8\lambda^2 + 5\lambda + 14 = 0,$$

with roots  $\lambda = -1, 2, 7$ . Hence A has eigenvalues -1, 2 and 7. To find the eigenvectors we proceed as follows. Corresponding to -1 we have to solve

$$AX = -X$$

Let  $X = (x_1, x_2, x_3)$ . Thus we have to solve the system of equations

$$-2x_1 + 6x_2 - 24x_3 = -x_1,$$

$$-3x_2 + 10x_3 = -x_2,$$

$$x_1 - 4x_2 + 13x_3 = -x_3,$$

or

$$x_1 - 6x_2 + 24x_3 = 0,$$

$$2x_2 - 10x_3 = 0,$$

$$x_1 - 4x_2 + 14x_3 = 0.$$

Since the equations are not linearly independent, we can choose

$x_3$  arbitrarily. Let  $x_3 = 1$ , then  $x_2 = 5$ ,  $x_1 = 6$ . Thus  $\begin{bmatrix} 6 \\ 5 \\ 1 \end{bmatrix}$  is an eigenvector corresponding to the eigenvalue  $\lambda_1 = -1$ . Similarly we find that  $\begin{bmatrix} -3 \\ 2 \\ 1 \end{bmatrix}$  and  $\begin{bmatrix} -2 \\ 1 \\ 1 \end{bmatrix}$  respectively are the eigenvectors corresponding to the eigenvalues  $\lambda_2 = 2$  and  $\lambda_3 = 7$ .

We now give some simple results about eigenvalues of a matrix. Suppose  $\lambda_1, \lambda_2, \dots, \lambda_n$  are the (not necessarily distinct) eigenvalues of a square matrix A.

**THEOREM 3.1:**  $\lambda_1 + \lambda_2 + \dots + \lambda_n = a_{11} + a_{22} + \dots + a_{nn}$

$$\lambda_1 \lambda_2 \dots \lambda_n = \det A$$

This result follows from equation (3.4). The number  $\lambda_1 + \lambda_2 + \dots + \lambda_n$  is called the **trace** of A and is denoted by  $\text{tr}(A)$ .

**THEOREM 3.2:** If X is an eigenvector of A corresponding to an eigenvalue  $\lambda$  then so is  $cX$  where c is any non zero constant.

**PROOF:**

$$\text{or } AX = \lambda X \Rightarrow cAX = c\lambda X$$

$$A(cX) = \lambda(cX).$$

**DEFINITION 3.1:** Two column vectors  $X_1$  and  $X_2$  are said to be orthogonal if  $\bar{X}_1 X_2 = 0$ .

Here  $\bar{X}_1 = (X_1^T)^*$ , for example if  $X_1 = \begin{bmatrix} i \\ 2-i \\ 3 \end{bmatrix}$ , then  $\bar{X}_1 = (-i, 2+i, 3)$ .

**DEFINITION 3.2:** A matrix is said to be hermitian if  $\bar{A} = A$ .

**EXAMPLE 3.3:** If  $A = \begin{bmatrix} 1 & 1-i \\ 1+i & 3 \end{bmatrix}$ , then  $A^T = \begin{bmatrix} 1 & 1+i \\ 1-i & 3 \end{bmatrix}$

and  $(A^T)^* = \bar{A} = \begin{bmatrix} 1 & 1-i \\ 1+i & 3 \end{bmatrix} = A$ .

Thus  $A$  is a hermitian matrix.

**THEOREM 3.3:** A hermitian matrix has only real eigenvalues.

**PROOF:** Let  $A$  be hermitian and  $\lambda$  be an eigenvalue and  $X$  the corresponding eigenvector i.e.

$$AX = \lambda X.$$

By definition  $X$  is not the zero vector. Taking the transpose of both sides we have

$$(AX)^T = \lambda X^T,$$

$$X^T A^T = \lambda X^T.$$

Take complex conjugates

$$\bar{X} \bar{A} = \bar{\lambda} \bar{X},$$

$$\bar{X} A^T = \bar{\lambda} \bar{X},$$

or

because  $A$  is hermitian. Postmultiply with  $X$

$$\bar{X} A X = \bar{\lambda} \bar{X} X,$$

$$\bar{X} \lambda X = \bar{\lambda} \bar{X} X,$$

or

$$(\lambda - \bar{\lambda}) \bar{X} X = 0.$$

But  $\bar{X} X \neq 0$ , since  $X$  is assumed to be a non zero vector. Hence  $\lambda = \bar{\lambda}$  which shows that  $\lambda$  is real.

**THEOREM 3.4:** Eigenvectors corresponding to distinct eigenvalues of a hermitian matrix are orthogonal.

**PROOF:** Let  $AX_1 = \lambda_1 X_1$ ,  
and  $AX_2 = \lambda_2 X_2$ ,

where  $\lambda_1 \neq \lambda_2$ , and A is a hermitian matrix. Proceeding as above we have

$$\bar{X}_1 A X_2 = \lambda_1 \bar{X}_1 X_2$$

where we used the fact that  $\lambda_1$  is real. The above equation implies

$$(\lambda_2 - \lambda_1) \bar{X}_1 X_2 = 0.$$

Since  $\lambda_1 \neq \lambda_2$ , we have

$$\bar{X}_1 X_2 = 0,$$

which proves the theorem.

An important corollary of the above theorem is as follows.

**COROLLARY:** A symmetric real matrix has only real eigenvalues and the eigenvectors corresponding to distinct eigenvalues are orthogonal.

It is easily shown [8], (page 84) that if  $X_1, X_2, \dots, X_n$  are eigenvectors of A corresponding to distinct eigenvalues then the set  $\{X_1, \dots, X_n\}$  is linearly independent. Hence every  $n \times 1$  column vector X can be expressed as a linear combination of these vectors. However if A is hermitian, it can be shown that [1] even if the eigenvalues are not distinct an orthogonal set of linearly independent eigenvectors  $\{X_r\}_{r=1}^n$  can always be found.

**EXAMPLE 3.4:** The unit matrix of order 3 has an eigenvalue repeated three times. However

$$X_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, X_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, X_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

are three linearly independent orthonormal vectors such that

$$IX_i = X_i, \quad i = 1, 2, 3.$$

**EXAMPLE 3.5:** The matrix

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

has the eigenvalue unity repeated three times. However there is only one linearly independent eigenvector

$$X = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

### 3.3 GERSHGORIN'S THEOREM

Now we consider some estimates for the absolute value  $|\lambda|$  where  $\lambda$  is any eigenvalue of a square matrix  $A = (a_{ij})$ . If

$X = (x_1, x_2, \dots, x_n)^T$  is an eigenvector corresponding to  $\lambda$ , then we have

$$\lambda x_i = \sum_{r=1}^n a_{ir} x_r, \quad i = 1, 2, \dots, n. \quad (3.5)$$

If  $|x_N| = \max_r |x_r|$ , then we have

$$\begin{aligned} |\lambda x_N| &= \left| \sum_{r=1}^n a_{Nr} x_r \right| \leq \sum_{r=1}^n |a_{Nr}| |x_r| \\ &\leq \sum_{r=1}^n |a_{Nr}| |x_N| \\ &= |x_N| \sum_{r=1}^n |a_{Nr}| \end{aligned}$$

Hence  $|\lambda| \leq \sum_{r=1}^n |a_{Nr}|$ .

Since  $N$  is, in general, not known we replace the above estimate by a weaker one

$$|\lambda| \leq \max_i \sum_{r=1}^n |a_{ir}| \quad (3.6)$$

Also the transposed matrix has the same set of eigenvalues, hence

$$|\lambda| \leq \max_i \sum_{r=1}^n |a_{ri}|. \quad (3.7)$$

From (3.5) we have

$$(\lambda - a_{NN})x_N = \sum_{r \neq N}^n a_{Nr} x_r.$$

Thus  $|\lambda - a_{NN}| |x_N| \leq \sum_{r \neq N}^n |a_{Nr}| |x_r|$

$$\text{or } |\lambda - a_{NN}| \leq |x_N| \sum_{r \neq N}^n |a_{Nr}|.$$

or  $|\lambda - a_{NN}| \leq \sum_{r \neq N}^n |a_{Nr}|.$

As before we replace the above result by a weaker one: " $\lambda$ " lies in one of the discs with centre at  $a_{ii}$  and radius

$$\sum_{\substack{r=1 \\ r \neq i}}^n |a_{ir}|, \quad i = 1, 2, \dots, n.$$

$$|\lambda - a_{ii}| \leq \sum_{r \neq i}^n |a_{ir}|, \quad i = 1, 2, \dots, n. \quad (3.8)$$

and from the transposed matrix, we have, " $\lambda$ " lies in one of the discs,

$$|\lambda - a_{ii}| \leq \sum_{r \neq i}^n |a_{ri}|, \quad i = 1, 2, \dots, n. \quad (3.9)$$

All the above estimates hold simultaneously, therefore the best of them is preferred. The above results are known as Gershgorin's Theorem: Here we

**THEOREM 3.5:** If  $\lambda$  is an eigenvalue of  $A = \{a_{ij}\}$  then it satisfies the estimates (3.6) - (3.9).

**EXAMPLE 3.6:** Consider

$$A = \begin{bmatrix} -2 & 6 & -24 \\ 0 & -3 & 10 \\ 1 & -4 & 13 \end{bmatrix}$$

Here

$$\sum_{r=1}^3 |a_{1r}| = 2 + 6 + 24 = 32$$

$$\sum_{r=1}^3 |a_{2r}| = 0 + 3 + 10 = 13$$

$$\sum_{r=1}^3 |a_{3r}| = 1 + 4 + 13 = 18$$

Therefore (3.6) gives

$$|\lambda| \leq 32. \quad (3.10)$$

Similarly (3.7) gives

$$|\lambda| \leq 47 \quad (3.11)$$

From (3.8) we get

$$|\lambda + 2| \leq 30 \quad \text{or} \quad |\lambda + 3| \leq 10 \quad \text{or} \quad |\lambda - 13| \leq 5. \quad (3.12)$$

Finally (3.9) gives

$$|\lambda + 2| \leq 1 \quad \text{or} \quad |\lambda + 3| \leq 10 \quad \text{or} \quad |\lambda - 13| \leq 34 \quad (3.13)$$

In (3.12) the disc  $|\lambda + 2| \leq 30$  contains the other two and in (3.13)  $|\lambda - 13| \leq 34$  does the same. Thus

$|\lambda + 2| \leq 30$  is the best estimate.

**EXAMPLE 3.7:** Consider

$$A = \begin{bmatrix} 1 & 2 & 2 \\ 2 & 3 & -2 \\ -5 & 3 & 8 \end{bmatrix}$$

Here we have  $|\lambda| \leq 16$  and  $|\lambda| \leq 12$ .

Also  $|\lambda - 1| \leq 4$  or  $|\lambda - 3| \leq 4$  or  $|\lambda - 8| \leq 8$ .  
 and  $|\lambda - 1| \leq 7$  or  $|\lambda - 3| \leq 5$  or  $|\lambda - 8| \leq 4$ .

Thus all eigenvalues lie in the union of  $|\lambda - 1| \leq 7$  and  $|\lambda - 8| \leq 4$ .  
 The matrix has eigenvalues 3, 4 and 5.

### 3.4 THE POWER METHOD

The power method is an iterative technique to find the dominant eigenvalue of a given matrix A. Assume that a matrix has eigenvalues which are so arranged that  $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$  and  $u_1, u_2, \dots, u_n$  are n linearly independent eigenvectors. If  $u$  is any arbitrary  $n \times 1$  vector then it can be expressed as a linear combination of the eigenvectors and we write

$$u = a_1 u_1 + a_2 u_2 + \dots + a_n u_n.$$

Then  $Au = a_1 A u_1 + a_2 A u_2 + \dots + a_n A u_n,$   
 $= a_1 \lambda_1 u_1 + a_2 \lambda_2 u_2 + \dots + a_n \lambda_n u_n.$

It is clear

$$A^k u = a_1 (\lambda_1^k) u_1 + a_2 (\lambda_2^k) u_2 + \dots + a_n (\lambda_n^k) u_n$$

$$= \lambda_1^k \left[ a_1 u_1 + a_2 \left( \frac{\lambda_2}{\lambda_1} \right)^k u_2 + \dots + a_n \left( \frac{\lambda_n}{\lambda_1} \right)^k u_n \right] \quad (3.14)$$

Since  $\left| \frac{\lambda_r}{\lambda_1} \right| < 1, \quad r = 2, 3, \dots, n,$

the vector in brackets will tend to  $a_1 u_1$  as  $k$  is taken large. If  $(V)_i$  denotes the i-th component of the column vector V, we have

$$\frac{(A^{k+1} u)_i}{(A^k u)_i} = \lambda_1 \left[ \frac{[a_1 u_1 + a_2 \left( \frac{\lambda_2}{\lambda_1} \right)^{k+1} u_2 + \dots + a_n \left( \frac{\lambda_n}{\lambda_1} \right)^{k+1} u_n]_i}{[a_1 u_1 + a_2 \left( \frac{\lambda_2}{\lambda_1} \right)^k u_2 + \dots + a_n \left( \frac{\lambda_n}{\lambda_1} \right)^k u_n]_i} \right]$$

Hence  $\lim_{k \rightarrow \infty} \frac{(A^{k+1}u)_1}{(A^k u)_1} = \lambda_1 \frac{a_1(u_1)_1}{a_1(u_1)_1} = \lambda_1$  (3.15)

From (3.15) we have for large  $k$

$$A^{k+1}u \approx \lambda_1 A^k u$$

Thus

$$A(A^k u) \approx \lambda_1 (A^k u)$$

Therefore  $A^k u$  is an approximate eigenvector corresponding to the eigenvalue  $\lambda_1$ . The method based on the above analysis, for finding the dominant (in absolute value) eigenvalue of a matrix, is called the power method. We proceed as follows:

- Choose a suitable arbitrary column vector  $u$ . Usually  $u$  is chosen as  $(1, 1, \dots, 1)^T$ .
- Form  $v = Au$ .
- Multiply  $v$  by a number  $k_i$  so as to make the largest element equal to 1. If one is working with an automatic machine, it is easier to multiply, at every step, with

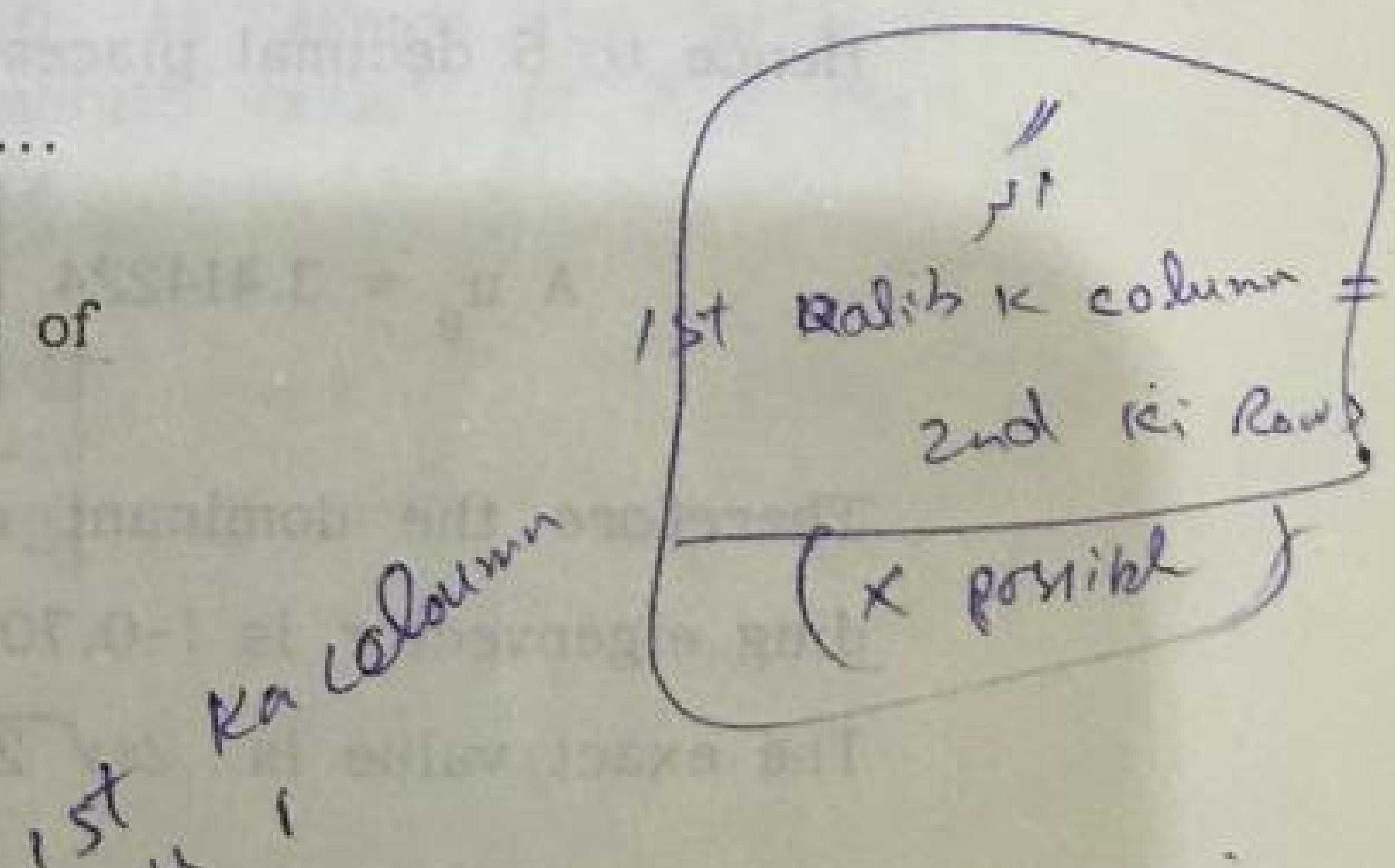
$$k_i = \left( |v_1|^2 + |v_2|^2 + \dots + |v_n|^2 \right)^{-1/2}$$

We do this because if  $v$  is an eigenvector then so is  $kv$ . This is advisable to stop the numbers from getting too large.

- Iterate  $v_{i+1} = Ak_i v_i$   $i = 1, 2, \dots$

**EXAMPLE 3.8:** Find the dominant eigenvalue of

$$A = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$



**SOLUTION:** Let

$$u = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \text{ then } v = Au = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \quad u_1 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

adjoint S. J. S. S. S. S.

Denote by  $u_i$  the normalized column vector  $k v_i$  i.e. the vector whose largest element in absolute value has been made equal to 1.

largest e  
Let  $\lambda$  be  
eigenvecto

$$v_2 = A u_1 = \begin{bmatrix} 2 \\ -2 \\ 2 \end{bmatrix} \xrightarrow{\text{① value } \leq 1} u_2 = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \xrightarrow{\text{② } \frac{1}{2}, \frac{-1}{2}, \frac{1}{2}}$$

$$v_3 = A u_2 = \begin{bmatrix} 3 \\ -4 \\ 3 \end{bmatrix} \xrightarrow{\frac{3}{4}, \frac{-4}{4}} u_3 = \begin{bmatrix} -0.75 \\ 1 \\ -0.75 \end{bmatrix}$$

$$v_4 = A u_3 = \begin{bmatrix} -2.5 \\ 3.5 \\ -2.5 \end{bmatrix} \xrightarrow{\frac{-2.5}{5}, \frac{3.5}{5}} u_4 = \begin{bmatrix} -0.75 \\ 1 \\ -0.75 \end{bmatrix}$$

$$v_5 = A u_4 = \begin{bmatrix} -2.428 \\ 3.428 \\ -2.428 \end{bmatrix} \xrightarrow{\frac{-2.428}{5}, \frac{3.428}{5}} u_5 = \begin{bmatrix} -0.7083 \\ 1 \\ -0.7083 \end{bmatrix}$$

$$v_6 = A u_5 = \begin{bmatrix} -2.4166 \\ 3.4166 \\ -2.4166 \end{bmatrix} \xrightarrow{\frac{-2.4166}{5}, \frac{3.4166}{5}} u_6 = \begin{bmatrix} -0.7073 \\ 1 \\ -0.7073 \end{bmatrix}$$

$$v_7 = A u_6 = \begin{bmatrix} -2.4146 \\ 3.4146 \\ -2.4146 \end{bmatrix} \xrightarrow{\frac{-2.4146}{5}, \frac{3.4146}{5}} u_7 = \begin{bmatrix} -0.70714 \\ 1 \\ -0.70714 \end{bmatrix}$$

$$v_8 = A u_7 = \begin{bmatrix} -2.41428 \\ 3.41428 \\ -2.41428 \end{bmatrix} \xrightarrow{\frac{-2.41428}{5}, \frac{3.41428}{5}} u_8 = \begin{bmatrix} -0.707112 \\ 1 \\ -0.707112 \end{bmatrix}$$

$$v_9 = A u_8 = \begin{bmatrix} -2.414224 \\ 3.414224 \\ -2.414224 \end{bmatrix} \xrightarrow{\frac{-2.414224}{5}, \frac{3.414224}{5}} u_9 = \begin{bmatrix} -0.70711 \\ 1 \\ -0.70711 \end{bmatrix}$$

Subtracting  
Thus  $\lambda - \alpha$   
ated vecto  
element o  
subtracting  
unchanged.  
largest or  
value.

Suppose  
eigenvalue  
found, add  
lue of A.  
that found

## EXAMPLE 3

SOLUTION:  
Thus they  
eigenvalue is  
smallest eigen

Hence to 5 decimal places

$$A u_8 = 3.414224 \begin{bmatrix} -0.70711 \\ 1 \\ -0.70711 \end{bmatrix} = 3.414224 u_9$$

Therefore the dominant eigenvalue is 3.414224 and the corresponding eigenvector is  $(-0.70711, 1, -0.70711)^T$ .

The exact value is  $2 + \sqrt{2} = 3.414214$ .

Apply the pow

## The smallest and the largest eigenvalues:

We can modify the power method to find the smallest and

Let u

largest eigenvalues.

Let  $\lambda$  be an eigenvalue of a matrix  $A$  and  $X$  be the corresponding eigenvector. Then

$$AX = \lambda X$$

Subtracting  $\alpha X$  from both sides of the above equation we get

$$\begin{aligned} AX - \alpha X &= \lambda X - \alpha X \\ (A - \alpha I)X &= (\lambda - \alpha)X \end{aligned}$$

Thus  $\lambda - \alpha$  is an eigenvalue of the matrix  $A - \alpha I$  and  $X$  is the associated vector. In otherwords, the effect of subtracting  $\alpha$  from each element of the leading diagonal of a matrix has the effect of subtracting  $\alpha$  from each eigenvalue but leaves the eigenvectors unchanged. Therefore, we can make, by a proper choice of  $\alpha$ , the largest or the smallest eigenvalue the dominant one in absolute value.

Suppose that  $\alpha$  is taken equal to or greater than the dominant eigenvalue of  $A$ . Then, if the dominant eigenvalue of  $A - \alpha I$  is found, adding  $\alpha$  to it will give the numerically smallest eigenvalue of  $A$ . But the correspondig eigenvector will be the same as that found for the dominant eigenvalue of  $A - \alpha I$ .

**EXAMPLE 3.9:** Find the smallest eigenvalue of

$$\begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$

**SOLUTION:** It is clear that all eigenvalues satisfy  $|\lambda - 2| \leq 2$ . Thus they must lie in the interval  $[0, 4]$ . Hence the largest eigenvalue is also the dominant (in absolute value). To find the smallest eigenvalue we consider the matrix

$$B = A - 4I = \begin{bmatrix} -2 & -1 & 0 \\ -1 & -2 & -1 \\ 0 & -1 & -2 \end{bmatrix}$$

Apply the power method to  $B$ .

$$\text{Let } u = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \text{ then } v = Au = \begin{bmatrix} -3 \\ -4 \\ -3 \end{bmatrix} = -4 \begin{bmatrix} 0.75 \\ 1 \\ 0.75 \end{bmatrix} = -4u_1$$

$$v_2 = Au_1 = \begin{bmatrix} -2.5 \\ -3.5 \\ -2.5 \end{bmatrix} = -3.5 \begin{bmatrix} 0.7143 \\ 1 \\ 0.7143 \end{bmatrix} = -3.5 u_2$$

$$v_3 = Au_2 = \begin{bmatrix} -2.4286 \\ -3.4286 \\ -2.4286 \end{bmatrix} = -3.4286 \begin{bmatrix} 0.70833 \\ 1 \\ 0.70833 \end{bmatrix} = -3.4286 u_3$$

$$v_4 = Au_3 = \begin{bmatrix} -2.41667 \\ -3.41667 \\ -2.41667 \end{bmatrix} = -3.41667 \begin{bmatrix} 0.70732 \\ 1 \\ 0.70732 \end{bmatrix} = -3.41667 u_4$$

$$v_5 = Au_4 = \begin{bmatrix} -2.41464 \\ -3.41464 \\ -2.41464 \end{bmatrix} = -3.41464 \begin{bmatrix} 0.7071434 \\ 1 \\ 0.7071434 \end{bmatrix} = -3.41464 u_5$$

$$v_6 = Au_5 = \begin{bmatrix} -2.4142868 \\ -3.4142868 \\ -2.4142868 \end{bmatrix} = -3.4142868 \begin{bmatrix} 0.7071131 \\ 1 \\ 0.7071131 \end{bmatrix} = -3.4142868 u_6$$

$$v_7 = Au_6 = \begin{bmatrix} -2.414226127 \\ -3.414226127 \\ -2.414226127 \end{bmatrix} = -3.414226 \begin{bmatrix} 0.707108 \\ 1 \\ 0.707108 \end{bmatrix} = -3.414226 u_7$$

**EXAMPLE 3.10:**

By power method, we find that the approximate eigenvalue for  $A$  is 0.25. The corresponding eigenvector is  $(1, -0.5)^T$ .

For  $B$ , the approximate eigenvalue is 4, 9. Thus both eigenvalues are real.

**Rayleigh Quotient:**

$u_1, u_2, \dots, u_n$  are the orthonormal eigenvectors of  $A$ .

Then  $\bar{u}_n = u_1 + u_2 + \dots + u_n$  is also an eigenvector of  $A$ .

Thus an approximate eigenvalue for  $B$  is -3.414226. The corresponding eigenvalue for  $A$  is  $-3.414226 + 4 = 0.585774$ .

The exact value is  $2\sqrt{2} = 0.58578644$ . The approximate corresponding eigenvector for this eigenvalue is

Hence using the

$$\begin{bmatrix} 0.707108 \\ 1 \\ 0.707108 \end{bmatrix}$$

The question arises that whether the smallest eigenvalue can be found without finding the largest one. Suppose

$$AX = \lambda X$$

then

$$X = A^{-1}\lambda X = \lambda A^{-1}X \quad (\text{As } \lambda \text{ is scalar})$$

or

$$A^{-1}X = \lambda^{-1}X.$$

This equation shows that  $1/\lambda$  is an eigenvalue of  $A^{-1}$ . If power method is used to find the largest eigenvalue  $1/\lambda$  of  $A^{-1}$ , then the smallest eigenvalue of  $A$  is  $\lambda$ . Moreover, if  $X$  is an eigenvector of  $A^{-1}$  corresponding to eigenvalue  $1/\lambda$ , then  $X$  is also the eigenvector of  $A$  corresponding to eigenvalue  $\lambda$ .

$$\gamma_p = \frac{\bar{v}_{p+1} v_p}{\bar{v}_p v_p}$$

Thus we see that

much faster as compared to the power method. This is called the sequence inverse iteration method.

$$\gamma_1 = -3$$

EXAMPLE 3.10: Let  $A = \begin{bmatrix} 5 & 2 \\ 2 & 8 \end{bmatrix}$ , then

$$A^{-1} = \frac{1}{36} \begin{bmatrix} 8 & -2 \\ -2 & 5 \end{bmatrix} = \begin{bmatrix} 0.2222 & -0.0556 \\ -0.0556 & 0.1389 \end{bmatrix}$$

By power method we can find the largest eigenvalue of  $A^{-1}$  which is 0.25. The corresponding eigenvector is  $(1, -0.5)^T$ . The eigenvalue of  $A$  is  $1/0.25$  or 4 and the corresponding eigenvector is  $(1, -0.5)^T$ . The eigenvalues of  $A$  calculated analytically are 4, 9. Thus both the results agree.

**Rayleigh Quotient:** If a matrix is Hermitian its eigenvectors  $u_1, u_2, \dots, u_n$  are orthogonal. We suppose that we have also made them orthonormal. Let  $u$  be an arbitrary column vector such that

$$\bar{u}_1 u_1 \neq 0. \text{ Now}$$

$$u = c_1 u_1 + c_2 u_2 + \dots + c_n u_n$$

$$v_p = A^p u = c_1 (\lambda_1)^p u_1 + c_2 (\lambda_2)^p u_2 + \dots + c_n (\lambda_n)^p u_n$$

$$\text{and } v_{p+1} = c_1 (\lambda_1)^{p+1} u_1 + \dots + c_n (\lambda_n)^{p+1} u_n.$$

Hence using the orthonormality property of the eigenvectors i.e.

$$\bar{u}_i u_j = \delta_{ij}, \text{ we find}$$

$$\gamma_p = \frac{\bar{v}_{p+1} v_p}{\bar{v}_p v_p} = \lambda_1 \frac{|c_1|^2 + |c_2|^2 \left(\frac{\lambda_2}{\lambda_1}\right)^{2p+1} + \dots + |c_n|^2 \left(\frac{\lambda_n}{\lambda_1}\right)^{2p+1}}{|c_1|^2 + |c_2|^2 \left(\frac{\lambda_2}{\lambda_1}\right)^{2p} + \dots + |c_n|^2 \left(\frac{\lambda_n}{\lambda_1}\right)^{2p}}$$

Thus we see that the convergence of the sequence  $\{\gamma_p\}_{p=1}^{\infty}$  to  $\lambda_1$  is much faster as compared with the sequence given by (3.15). This is called the sequence of Rayleigh quotients. For example in the example 3.9,

42

$$\gamma_2 = -3.4117647, \dots, \gamma_6 = -3.414213473$$

and we see that  $\gamma_6$  is much closer to the exact value  $-3.41421356$  as compared with the value  $-3.414226$  obtained by the earlier method.

Now we for

### 3.5 DEFLECTION

Suppose we have found an eigenvalue  $\lambda_1$  and the corresponding eigenvector  $X_1$  of an  $n \times n$  matrix  $A$ . Deflation means finding a matrix of order  $(n-1) \times (n-1)$  whose eigenvalues coincide with the rest of the eigenvalues of  $A$ . Suppose  $X_1 = (x_1, x_2, \dots, x_n)^T$ . Assume that  $x_1 \neq 0$  and multiply  $X_1$  by  $1/x_1$  to get an eigenvector  $X_1^*$  whose top element is unity i.e.

$$X_1^* = (1, x_2^*, \dots, x_n^*)^T.$$

Note that  $X_1^*$  has only zero entries except the first one which is unity.

**THEOREM 3:** If  $A$  has  $n$  eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_n$

**PROOF:** Suppose

where  $X_1^*$  has only zero entries except the first one which is unity. Then  $X_2^*, \dots, X_n^*$  has its top elements unity.

If the  $i$ -th row of  $A$  is denoted by  $R_i$ , we see that

$$AX_1^* = \begin{bmatrix} R_1 X_1^* \\ R_2 X_1^* \\ R_3 X_1^* \\ \vdots \\ R_n X_1^* \end{bmatrix}.$$

Hence  $O$  is an eigenvector of  $A$  with eigenvalue  $\lambda_1$  since its top element is unity.

Since  $AX_1^* = \lambda_1 X_1^* = (\lambda_1, \lambda_1 x_2^*, \dots, \lambda_1 x_n^*)^T$ , it is clear that

$$R_1 X_1^* = \lambda_1. \quad (3.1)$$

If rest of the eigenvectors have also been normalized in the same manner i.e. their top element is 1, then

$$R_i X_1^* = \lambda_i.$$

In the above, we learn from (3.17) that

If, however  $X_1$  is such that its top element is zero, then

$$R_1 X_i^* = 0. \quad (3.16c)$$

Now we form the matrix B defined as

$$B = A - X_1^* R_1.$$

Note that  $X_1^* R_1$  is an  $n \times n$  matrix whose first row is  $R_1$ , hence B has only zeros in its first row. We show that with the exception of  $\lambda_1$ , matrices A and B have all eigenvalues common.

**THEOREM 3.6:** If A has eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_n$  then B has eigenvalues 0,  $\lambda_2, \dots, \lambda_n$ .

**PROOF:** Suppose  $A X_i^* = \lambda_i X_i^*$  ( $i = 1, 2, \dots, n$ )

where  $X_1^*$  has its top element equal to unity and each of  $X_2^*, X_3^*, \dots, X_n^*$  has its top element either 1 or 0. Now

$$\begin{aligned} BX_1^* &= (A - X_1^* R_1)X_1^* \\ &= AX_1^* - X_1^*(R_1 X_1^*), \\ &= \lambda_1 X_1^* - X_1^* \lambda_1, \\ &= 0. \end{aligned}$$

Hence 0 is an eigenvalue of B with eigenvector  $X_1^*$ . If  $X_i^*$  has its top element unity then

$$\begin{aligned} B(X_i^* - X_1^*) &= (A - X_1^* R_1)(X_i^* - X_1^*) \\ &= AX_i^* - AX_1^* - X_1^*(R_1 X_i^*) + X_1^*(R_1 X_1^*) \\ &= \lambda_i X_i^* - \lambda_1 X_i^* - \lambda_1 X_1^* + \lambda_1 X_1^*, \\ &= \lambda_i (X_i^* - X_1^*). \end{aligned} \quad (3.17)$$

In the above, we have made use of (3.16a) and (3.16b). Now it is clear from (3.17) that  $\lambda_i$  is an eigenvalue of B and  $X_i^* - X_1^*$  is the

corresponding eigenvector. Finally if  $X_1^*$  has its top element equal to zero, then

$$\begin{aligned} BX_j^* &= AX_j^* - X_1^* R_1 X_j^* \\ &= AX_j^*. \end{aligned}$$

Since  $R_1 X_j^* = 0$  by (3.16c). Thus again we have

$$BX_j^* = \lambda_j X_j^* \quad (3.18)$$

This completes the proof.

Now define an  $(n - 1) \times (n - 1)$  matrix  $A_1$  obtained by deleting the first row and the first column of  $B$ . Let

$$B = \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ b_{21} & b_{22} & b_{23} & \dots & b_{2n} \\ b_{31} & b_{32} & b_{33} & \dots & b_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & b_{n3} & \dots & b_{nn} \end{bmatrix}.$$

Then

$$A_1 = \begin{bmatrix} b_{22} & b_{23} & b_{24} & \dots & b_{2n} \\ b_{32} & b_{33} & b_{34} & \dots & b_{3n} \\ b_{42} & b_{43} & b_{44} & \dots & b_{4n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ b_{n2} & b_{n3} & b_{n4} & \dots & b_{nn} \end{bmatrix}.$$

Since

$$|B - \lambda I_n| = -\lambda |A_1 - \lambda I_{n-1}|, \quad (3.19)$$

it is obvious that  $B$  and  $A_1$  have all eigenvalues common with exception of  $\lambda = 0$ . In (3.19)  $I_n$  denotes the unit matrix of order  $n$ .  $A_1$  is the so called deflated matrix we were seeking.

### EXAMPLE 3.11:

$A$  has an eigenvalue

$$X_1^* = (1, 1, 1/2)^T$$

Now

$$\text{and } B = A - X_1^* R_1 X_1^*$$

$$\text{Since } B = A - X_1^* R_1 X_1^*$$

The matrix  $A_1$  obtained by deleting the first row and the first column of  $B$  is found to be

$$A_1 =$$

Now

$$|A_1 - \lambda I_{n-1}|$$

Thus  $A_1$  has eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_{n-1}$  which are the eigenvalues of  $A$ .

We can recover the deflated matrix  $A_1$  in (3.19).

$$(y_1, y_2, \dots, y_{n-1})$$

$$(0, y_1, y_2, \dots, y_{n-1})$$

denote this vector by

is  $\lambda_1$ , i.e.

**EXAMPLE 3.11:** Let  $A = \begin{bmatrix} 1 & 2 & 2 \\ 2 & 3 & -2 \\ -5 & 3 & 8 \end{bmatrix}$ .

$A$  has an eigenvalue 4 with  $(2, 2, 1)^T$  as an eigenvector. Here  $\lambda_1 = 4$ ,  $X_1^* = (1, 1, 1/2)^T$  and  $R_1 = (1, 2, 2)$ .

Now

$$X_1^* R_1 = \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 2 \\ 1/2 & 1 & 1 \end{bmatrix},$$

and  $B = A - X_1^* R_1$  comes out as

$$B = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 1 & -4 \\ -\frac{11}{2} & 2 & 7 \end{bmatrix}$$

The matrix  $A_1$  obtained by deleting the first row and the first column of  $B$  is found to be

$$A_1 = \begin{bmatrix} 1 & -4 \\ 2 & 7 \end{bmatrix}$$

Now  $|A_1 - \lambda I| = \lambda^2 - 8\lambda + 15$   
 $= (\lambda - 5)(\lambda - 3)$ .

Thus  $A_1$  has eigenvalues 5 and 3. These are also the remaining eigenvalues of  $A$ .

We can recover eigenvectors of the matrix  $A$  from those of deflated matrix  $A_1$  in a simple manner. First we note that if

$(y_1, y_2, \dots, y_{n-1})^T$  is an eigenvector of matrix  $A_1$  then  $(0, y_1, y_2, \dots, y_{n-1})^T$  is an eigenvector of the matrix  $B$ . Let us denote this vector by  $Y_1$  and suppose the corresponding eigenvalue is  $\lambda_1$ , i.e.

(3.20)

$$By_1 = \lambda_1 Y_1$$

However from (3.17) and (3.18) we have

$$B(X_1^* - X_i^*) = \lambda_i(X_1^* - X_i^*) \quad (3.21)$$

if

$$R_1 X_1^* = \lambda_i \neq 0$$

and

$$BX_i^* = \lambda_i X_i^*$$

if

$$R_1 X_i^* = 0.$$

If (3.21) holds then  $Y_i$  must be proportional to  $X_1^* - X_i^*$  since these are eigenvectors of  $B$  corresponding to the same eigenvalue. Thus for some  $\alpha$

$$X_1^* - X_i^* = \alpha Y_i. \quad (3.23)$$

Multiply (3.23) with  $R_1$ , we get

$$R_1 X_1^* - R_1 X_i^* = \alpha R_1 Y_i,$$

or

$$\lambda_1 - \lambda_i = \alpha R_1 Y_i.$$

This gives  $\alpha = (\lambda_1 - \lambda_i)/R_1 Y_i$  and (3.23) gives

$$X_1^* = X_1^* - (\lambda_1 - \lambda_i)Y_i / R_1 Y_i. \quad (3.24)$$

If  $R_1 X_i^* = 0$ , i.e. (3.22) holds, then

$$Y_i = \beta X_i^*$$

for some  $\beta \neq 0$ , since  $X_i^*$  and  $Y_i$  happen to be eigenvectors of  $B$  corresponding to the same eigenvalue. In this case  $Y_i$  itself is an eigenvector of  $A$ . Actually we do not know  $X_i^*$ , however  $R_1 Y_i^*$  implies  $R_1 X_1^* \neq 0$  and  $R_1 Y_i^* = 0$  implies  $R_1 X_i^* = 0$ .

### EXAMPLE 3.1

We know that matrix  $A_1$  can

$A_1$  has an eigenvalue 3 with a

Since  $R_1 Y_2 = 0$   
let  $\lambda_3 = 3$  and

Here  $R_1 Y_3 = 2$  and

Thus we have found  
and 3 with respect to  
 $(1,0,1)^T$ .

In the beginning element of  $X_1$  is determined  
make the second element a suitable constant and modifications.

In practice we can find the corresponding eigenvector

**EXAMPLE 3.12:** Let  $A = \begin{bmatrix} 1 & 2 & 2 \\ 2 & 3 & -2 \\ -5 & 3 & 8 \end{bmatrix}$ .

We know that  $\lambda_1 = 4$ ,  $X_1^* = (1, 1, 1/2)^T$ . We found that the deflated matrix  $A_1$  came out to be

$$A_1 = \begin{bmatrix} 1 & -4 \\ 2 & 7 \end{bmatrix}$$

$A_1$  has an eigenvalue 5 with an eigenvector  $(1, -1)^T$  and an eigenvalue 3 with an eigenvector  $(2, -1)^T$ . Let us take  $\lambda_2 = 5$  and

$$Y_2 = \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}$$

Since  $R_1 Y_2 = 0$ , we conclude that  $Y_2$  is an eigenvector of  $A$ . Now let  $\lambda_3 = 3$  and

$$Y_3 = \begin{bmatrix} 0 \\ 2 \\ -1 \end{bmatrix}$$

Here  $R_1 Y_3 = 2$  and (3.24) gives

$$\begin{aligned} X_3^* &= \begin{bmatrix} 1 \\ 1 \\ 1/2 \end{bmatrix} - \frac{(4-3)}{2} \begin{bmatrix} 0 \\ 2 \\ -1 \end{bmatrix}, \\ &= \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}. \end{aligned}$$

Thus we have finally found that the matrix  $A$  has eigenvalues 4, 5 and 3 with respective eigenvectors  $(1, 1, 1/2)^T$ ,  $(0, 2, -1)^T$  and  $(1, 0, 1)^T$ .

In the beginning of this section we assumed that the top element of  $X_1$  is different from zero. If it happens to be zero we make the second element from the top unity by multiplying with a suitable constant and carry through with the analysis with minor modifications.

In practice we can find the dominant eigenvalue and the corresponding eigenvector of a given matrix  $A$  by the power method, then

find the matrix  $A_1$  by deflation, apply the power method to  $A_1$  to find the next dominant eigenvalue and so on.

In the case of  $3 \times 3$  matrix we can find the dominant eigenvalue and the corresponding eigenvector of the matrix by the power method, then the smallest eigenvalue and the corresponding vector. The third eigenvalue follows immediately as the sum of the eigenvalues is equal to the trace of the given matrix..

The power method fails if the dominant eigenvalues are a pair of complex conjugate numbers. Also the convergence is poor if ratio of the two dominant eigenvalues is close to unity. In such a case Jacobi's method may be used. For details see Fröberg [4].

### EXERCISE 3

- ✓ 1. Find the eigenvalues and the corresponding eigenvectors of the following matrices analytically.

$$A = \begin{bmatrix} 1 & 2 & -8 \\ -2 & 2 & -2 \\ 1 & -4 & 10 \end{bmatrix}, B = \begin{bmatrix} 0 & 2 & -10 \\ -3 & 1 & -3 \\ 1 & -5 & 11 \end{bmatrix}, C = \begin{bmatrix} -2 & 2 & 2 \\ 3 & -1 & 3 \\ 1 & 1 & -3 \end{bmatrix}.$$

- ✓ 2. Find the dominant eigenvalue and the corresponding eigenvector for each of the following matrices. Give your answers correct to three decimal places.

$$(a) \begin{bmatrix} 1 & 0 & 2 \\ 0 & 1 & 3 \\ 1 & 1 & 1 \end{bmatrix} \quad (b) \begin{bmatrix} -2 & 1 & 0 \\ 1 & -2 & 1 \\ 0 & 1 & -2 \end{bmatrix}.$$

- ✓ 3. Find all the eigenvalues of the matrix

$$\begin{bmatrix} 3 & 0 & 1 \\ 0 & 3 & 2 \\ 1 & 2 & 3 \end{bmatrix}$$

by calculating the roots of its characteristic equation. Hence determine the associated eigenvectors and normalize them. Give your answers correct to four decimal places.

- ✓ 4. Taking  $u = (0,0,1)^T$ , use the power method to evaluate the dominant eigenvalue and corresponding eigenvector of the matrix  $A$  given in Exercise 3. Give your answer correct to four decimal places and compare them with those obtained previously.

ELEMENTS  
Using  
deflate  
values  
those  
Using  
find t  
1. Consid  
✓ having  
corre  
calcu  
vecto

5. Using the estimate of eigenvector obtained in Exercise 4, deflate the matrix A and hence calculate the remaining eigenvalues and associated eigenvectors. Compare your results with those obtained in Exercise 3.
6. Using the estimate of eigenvector obtained in Exercise 4, find the smallest eigenvalue and associated eigenvector of A.
7. Consider the matrix

~~X~~

$$A = \begin{bmatrix} 0 & 2 & -8 \\ -2 & 1 & -2 \\ 1 & -4 & 9 \end{bmatrix}$$

having  $\lambda_1 = -1$  as an eigenvalue and  $x_1 = (1, 3/2, 1/2)^T$  the corresponding eigenvector. Deflate the matrix and hence calculate the remaining eigenvalues and associated eigenvectors.

## CHAPTER 4

### INTERPOLATION WITH UNEQUALLY SPACED DATA

#### 4.1 INTRODUCTION

Suppose we are given the values of a function  $f(x)$  at certain points  $x_0, x_1, \dots, x_n$  and we want to estimate  $f(\alpha)$  where  $\alpha$  is any point in the interval  $[x_0, x_n]$ . We assume that the function can be approximated by a polynomial. This is called polynomial interpolation.

There are a number of interpolation formulae, most of which have certain advantages over the others in certain situations, but no one of which is preferable to all others in all respects. In the general case, when the abscissas are not equally spaced, the use of divided differences is convenient.

#### 4.2 LAGRANGE'S FORMULA

We assume that the polynomial

$$P(x) = a_0 + a_1 x + \dots + a_n x^n$$

is such that  $P(x_r) = y_r$ ,  $r = 0, 1, 2, \dots, n$ . Thus the curves representing graphs of the two functions  $y = f(x)$  and  $y = P(x)$  pass through the same  $n+1$  points  $(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)$ . If the function  $f(x)$  has continuous derivatives of all orders, it is reasonable to suppose that  $f(x)$  and  $P(x)$  will be fairly close at other points of the interval  $[x_0, x_n]$ . To determine  $a_0, a_1, \dots, a_n$  we have to solve the system of equations:

$$\left. \begin{array}{l} a_0 + a_1 x_0 + a_2 x_0^2 + \dots + a_n x_0^n = y_0, \\ a_0 + a_1 x_1 + a_2 x_1^2 + \dots + a_n x_1^n = y_1, \\ \dots \\ \dots \\ a_0 + a_1 x_n + a_2 x_n^2 + \dots + a_n x_n^n = y_n. \end{array} \right\} \quad (4.1)$$

The determinant of the system is

$1$	$x_0$	$x_0^2$	$\dots$	$x_0^n$
$1$	$x_1$	$x_1^2$	$\dots$	$x_1^n$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$1$	$x_n$	$x_n^2$	$\dots$	$x_n^n$

Since  $x_0, x_1, \dots, x_n$  are all distinct, this determinant is not zero and the system (4.1) has a unique solution. However one or more, but not all, of  $a_n$  may turn out to be zero. Thus we have shown that a polynomial of degree  $\leq n$  exists such that the curve representing  $y = P(x)$  passes through the given  $n+1$  points. Also this polynomial is unique. Let us define

$$L_k(x) = \frac{(x-x_0)(x-x_1)\dots(x-x_{k-1})(x-x_{k+1})\dots(x-x_n)}{(x_k-x_0)(x_k-x_1)\dots(x_k-x_{k-1})(x_k-x_{k+1})\dots(x_k-x_n)}.$$

We see that

$$L_k(x_r) = \begin{cases} 0, & \text{if } r \neq k, \\ 1, & \text{if } r = k. \end{cases} \quad (4.2)$$

If we write

$$P(x) = L_0(x)y_0 + L_1(x)y_1 + \dots + L_n(x)y_n, \quad (4.3)$$

We see that  $P(x)$  is a polynomial of degree  $\leq n$  (since  $L_0(x), L_1(x), \dots, L_n(x)$  etc. are polynomials of degree  $n$  and coefficients of  $x^n, x^{n-1}, \dots$

may add up to zero), also

$$P(x_0) = y_0, P(x_1) = y_1, \dots, P(x_n) = y_n.$$

Thus  $P(x)$  is the required polynomial. Equation (4.3) is called the **Lagrange's interpolation formula**. Note that, although we proved the existence of such a polynomial by a different method, we seldom use the determinant method to actually evaluate the constants  $a_0, a_1, \dots, a_n$ . Later we shall discuss other methods for finding the interpolation polynomial. Since such a polynomial is unique, it does not matter how we find it. The result is always the same.

**EXAMPLE 4.1:** Let values of  $y = f(x)$  be as given in the following table. Find a polynomial which interpolates the given data

x	y
-1	-8
1	-2
2	4
3	28

**SOLUTION:** Here  $x_0 = -1, x_1 = 1, x_2 = 2, x_3 = 3,$

$$y_0 = -8, y_1 = -2, y_2 = 4, y_3 = 28.$$

Thus

$$L_0(x) = \frac{(x-1)(x-2)(x-3)}{(-1-1)(-1-2)(-1-3)},$$

$$L_1(x) = \frac{(x-(-1))(x-2)(x-3)}{(1-(-1))(1-2)(1-3)} \text{ etc.}$$

Here

$$P(x) = \frac{(x-1)(x-2)(x-3)}{(-2)(-3)(-4)}(-8) + \frac{(x+1)(x-2)(x-3)}{2(-1)(-2)}(-2)$$

$$+ \frac{(x-(-1))(x-1)(x-3)}{(2-(-1))(2-1)(2-3)}(4) + \frac{(x-(-1))(x-1)(x-2)}{(3-(-1))(3-1)(3-2)}(28)$$

$$= (x^3 - 6x^2 + 11x - 6)(1/3) + (x^3 - 4x^2 + x + 6)(-1/2) \\ + (x^3 - 3x^2 - x + 3)(-4/3) + (x^3 - 2x^2 - x + 2)(7/2)$$

$$= 2x^3 - 3x^2 + x - 2.$$

## 4.3 THE ERROR OF THE INTERPOLATION POLYNOMIAL

Suppose we want to consider the error  $R(x) = f(x) - P(x)$  at a point  $\alpha$  in the interval  $[x_0, x_n]$ . Of course  $R=0$  if  $\alpha = x_0, x_1, \dots, x_n$ . So let  $\alpha$  be a point distinct from them. Let us define

$$F(x) = (x - x_0)(x - x_1) \dots (x - x_n)$$

and write

$$g(x) = f(x) - p(x) - KF(x). \quad (4.4)$$

In (4.4)  $K$  is a constant to be chosen in a manner to make  $g(\alpha)=0$ . Note that  $g(x_0) = g(x_1) = \dots = g(x_n) = 0$ .

Consider  $g(x)$  on the interval  $[x_0, x_1]$ . Now  $g(x_0)=g(x_1)=0$ , also  $g$  is continuous on  $[x_0, x_1]$  and differentiable on  $(x_0, x_1)$  hence, by Rolle's theorem, there is a point  $\beta_1$  in  $(x_0, x_1)$  such that  $g'(\beta_1) = 0$ . Applying the same argument to  $g(x)$  on intervals  $[x_1, x_2], \dots, [x_{n-1}, x_n]$  we find that there exist points  $\beta_2, \beta_3, \dots, \beta_{n+1}$  such that the derivative of  $g(x)$  vanishes at each of these points. Now we apply Rolle's theorem to  $g'(x)$  and find that there exist points  $\gamma_1, \gamma_2, \dots, \gamma_n$  such that  $g''(x)$  vanishes at each of these points. Proceeding in this manner, we find that there exists a point  $\xi$  in  $(x_0, x_n)$  such that  $g^{(n+1)}(\xi)=0$ . Now from

(4.4)

$$g^{n+1}(\xi) = f^{(n+1)}(\xi) - P^{(n+1)}(\xi) - K(n+1)!$$

But  $P^{n+1}(\xi)=0$ , because  $P(x)$  is a polynomial of degree at most  $n$ . Hence

$$K = \frac{f^{(n+1)}(\xi)}{(n+1)!}$$

From (4.4) we get, on putting  $x = \alpha$  (remember  $g(\alpha) = 0$ )

$$f(\alpha) = P(\alpha) + \frac{f^{(n+1)}(\xi)}{(n+1)!} F(\alpha). \quad (4.5)$$

In the beginning we assumed that  $\alpha \neq x_0, x_1, \dots, x_n$ . However (4.5) is obviously satisfied even if we remove this restriction. Since  $\alpha$  is an arbitrary point we can write

$$f(x) = P(x) + \frac{f^{(n+1)}(\xi)}{(n+1)!} F(x), \quad x \in [x_0, x_n].$$

Note that  $\xi$  also depends on  $x$ . We have shown that the error  $f(x) - P(x)$  committed by approximating the function  $f(x)$  by the interpolation polynomial is

$$R = \frac{f^{(n+1)}(\xi) F(x)}{(n+1)!}, \quad x_0 < \xi < x_n. \quad (4.6)$$

#### 4.4 DIVIDED DIFFERENCES

Let the values of a function  $f$  be given at  $n+1$  points  $x_0, x_1, \dots, x_n$ . We define the first divided difference of  $f(x)$  between  $x_0$  and  $x_1$  as

$$f(x_0, x_1) = \frac{f(x_1) - f(x_0)}{x_1 - x_0}. \quad (4.7)$$

Note that  $f(x_0, x_1) = f(x_1, x_0)$ . We define the second divided difference,  $f(x_0, x_1, x_2)$ , as

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

Similarly the  $n$ -th divided difference is defined as

$$f(x_0, x_1, \dots, x_n) = \frac{f(x_1, x_2, \dots, x_n) - f(x_0, x_1, \dots, x_{n-1})}{x_n - x_0}. \quad (4.8)$$

Now we shall prove the following theorem:

**THEOREM 4.1:**

$$f(x_0, x_1, \dots, x_n) = \sum_{r=0}^n \frac{f(x_r)}{(x_r - x_0)(x_r - x_1) \dots (x_r - x_{r-1})(x_r - x_{r+1}) \dots (x_r - x_n)}$$

**PROOF:** Obviously the result is true for  $n = 1$ , (see 4.7). Let the result be true for  $n = k$ . Thus we have

$$\begin{aligned} f(x_0, x_1, \dots, x_k) &= \sum_{r=0}^k \frac{f(x_r)}{(x_r - x_0)(x_r - x_1) \dots (x_r - x_{r-1})(x_r - x_{r+1}) \dots (x_r - x_k)} \quad (4.9) \end{aligned}$$

(4.6)

Now by definition

$$f(x_0, x_1, \dots, x_{k+1}) = \frac{f(x_1, \dots, x_{k+1}) - f(x_0, \dots, x_k)}{x_{k+1} - x_0}$$

$$= \frac{1}{(x_{k+1} - x_0)} \left[ \sum_{r=1}^{k+1} \frac{f(x_r)}{(x_r - x_1) \dots (x_r - x_{k+1})} - \sum_{r=0}^k \frac{f(x_r)}{(x_r - x_0) \dots (x_r - x_k)} \right]$$

where we have used Eq. (4.9),

(4.7)

$$\begin{aligned} &= \frac{1}{(x_{k+1} - x_0)} \left[ \frac{f(x_{k+1})}{(x_{k+1} - x_1) \dots (x_{k+1} - x_k)} - \frac{f(x_0)}{(x_0 - x_1)(x_0 - x_2) \dots (x_0 - x_k)} \right. \\ &\quad \left. + \sum_{r=1}^k \frac{[x_r - x_0 - (x_r - x_{k+1})]f(x_r)}{(x_r - x_0) \dots (x_r - x_{k+1})} \right] \end{aligned}$$

$$= \frac{f(x_0)}{(x_0 - x_1) \dots (x_0 - x_{k+1})} + \sum_{r=1}^k \frac{f(x_r)}{(x_r - x_0) \dots (x_r - x_{k+1})}$$

$$+ \frac{f(x_{k+1})}{(x_{k+1} - x_0) \dots (x_{k+1} - x_k)}$$

$$f(x_0, x_1, \dots, x_{k+1}) = \sum_{r=0}^{k+1} \frac{f(x_r)}{(x_r - x_0) \dots (x_r - x_{k+1})}$$

(4.8)

Thus we see that, if the result is true for  $n = k$ , then it is true for  $n = k+1$ . Since it is true for  $n = 1$ , it is true for all integral  $n$ .

We see that the right hand side of Eq. (4.8) remains the same even if we interchange any two of the arguments  $x_0, x_1, \dots, x_n$ . Thus the divided differences are symmetric functions of their arguments.

If any two of  $x_0, x_1, \dots, x_n$  happen to be equal, we define the divided differences as follows:

$$f(x_0, x_0) = \lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0} = f'(x_0).$$

Similarly,

$$f(x_0, x_0, x_0) = f''(x_0).$$

#### 4.4.1 DIVIDED-DIFFERENCE TABLE

The first, second and higher divided differences are usually arranged in the form of a table like the one given below:

$x$	$f(\ )$	$f(,)$	$f(,,)$	$f(,,,)$
$x_0$	$f(x_0)$			
		$f(x_0, x_1)$		
$x_1$	$f(x_1)$		$f(x_0, x_1, x_2)$	
		$f(x_1, x_2)$		$f(x_0, x_1, x_2, x_3)$
$x_2$	$f(x_2)$		$f(x_1, x_2, x_3)$	
		$f(x_2, x_3)$		$f(x_1, x_2, x_3, x_4)$
$x_3$	$f(x_3)$		$f(x_2, x_3, x_4)$	
		$f(x_3, x_4)$		
$x_4$	$f(x_4)$			

Thus if we draw a line starting at  $f(x_1)$  and sloping towards the right, it passes through successive divided differences  $f(x_1, x_2)$ ,  $f(x_1, x_2, x_3)$  and so on.

EXAMPLE 4.2: Form the divided difference table for the data

x	1	2	4	7	8
f(x)	-9	-41	-189	9	523

SOLUTION: The divided-difference table will be as follows

x	f(x)	f(,)	f(,,)	f(,,,)	f(,,,,)
1	-9				
2	-41	-32			
4	-189		-14		
7	9			7	
8	523				1
		66		14	
			112		
			514		

From the above table, we can easily write down any difference we want. For example to write down  $f(4,7,8)$  we draw a diagonal from  $f(4) = -189$ . Thus  $f(4,7) = 66$  and  $f(4,7,8) = 112$ .

#### 4.4.2 NEWTON'S INTERPOLATION POLYNOMIAL

From the definition and symmetry of  $f(x_0, x_1, x_2)$  etc., we can write

$$f(x) = f(x_0) + (x - x_0) f(x, x_0)$$

$$f(x, x_0) = f(x_0, x_1) + (x - x_1) f(x, x_0, x_1)$$

$$f(x, x_0, x_1) = f(x_0, x_1, x_2) + (x - x_2) f(x, x_0, x_1, x_2)$$

... ... ... ... ...

... ... ... ... ...

... ... ... ... ...

$$f(x, x_0, x_1, \dots, x_{n-1}) = f(x_0, x_1, \dots, x_n) + (x - x_n) f(x, x_0, x_1, \dots, x_n)$$

Multiply the second equation by  $(x - x_0)$ , the third by  $(x - x_0)(x - x_1)$ , the fourth by  $(x - x_0)(x - x_1)(x - x_2)$  and so on, and add. We get

$$f(x) = f(x_0) + (x-x_0)f(x_0, x_1) + (x-x_0)(x-x_1)f(x_0, x_1, x_2) + \dots + (x-x_0)(x-x_1)\dots(x-x_{n-1})f(x, x_0, x_1, \dots, x_n) + R_n \quad (4.10)$$

$$\text{where } R_n = (x - x_0)(x - x_1)\dots(x - x_n)f(x, x_0, x_1, \dots, x_n) \quad (4.11)$$

Let us write (4.10) as

$$f(x) = P_n(x) + R_n(x)$$

$$P_n = f(x_0) + (x-x_0)f(x_0, x_1) + \dots + (x-x_0)\dots(x-x_{n-1})f(x, x_0, x_1, \dots, x_n) \quad (4.12)$$

It is clear that  $P_n(x)$  is a polynomial of degree at most  $n$ , also  $R_n(x_r) = 0$  for  $r=0,1,2,\dots,n$ . Hence we see that  $P_n(x_r) = f(x_r)$ ,  $r=0,1,2,\dots,n$ . Therefore  $P_n(x)$  is a polynomial of degree at most  $n$  that passes through the  $n+1$  given points  $(x_r, f(x_r))$ ,  $r = 0, 1, \dots, n$ . We have already seen that such a polynomial is unique. Thus  $R_n(x)$  must be identical with the error term found in (4.6)

$$R_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} F(x). \quad (4.13)$$

where, as before,

$$F(x) = (x - x_0)(x - x_1)\dots(x - x_n).$$

From (4.11) and (4.13) we also see that

$$f(x, x_0, x_1, \dots, x_n) = \frac{f^{(n+1)}(\xi)}{(n+1)!}$$

The polynomial  $P_n(x)$  is called Newton's interpolation polynomial with divided differences. We emphasize the point that the Lagrangian and Newton's polynomials are identical. However calculations in the case of Newton's polynomial are simpler. Also  $x_r$  ( $r=1,2,\dots,n$ ) does not appear in any of the differences which involve  $x_0, x_1, \dots$ , upto  $x_{r-1}$ . Suppose we have found the Newton's

(4.10)

polynomial  $P_n(x)$  which passes through  $n + 1$  given points and now wish to add a new point  $(x_{n+1}, f(x_{n+1}))$  to the data. All we have to do is to calculate

$$(x - x_0)(x - x_1) \dots (x - x_n) f(x_0, x_1, \dots, x_{n+1})$$

and add this to  $P_n(x)$ :

$$P_{n+1}(x) = P_n(x) + (x - x_0)(x - x_1) \dots (x - x_n) f(x_0, x_1, \dots, x_{n+1}).$$

However, for the Lagrangian polynomial we have to calculate every term afresh and the labour is increased manyfold.

**EXAMPLE 4.3:** Find a 4th degree polynomial which passes through the 5 points given below.

x	f(x)
1.0	-9
2.0	-41
4.0	-189
7.0	9
8.0	523

**SOLUTION:** The table of divided differences is as follows:

x	f(x)	f(., )	f(., , )	f(., , , )	f(., . . . , )
1.0	-9				
2.0	-41	-32			
4.0	-189		-74	7	1
7.0	9			112	
8.0	523		514		

$$P_4(x) = f(x_0) + (x - x_0)f(x_0, x_1) + (x - x_0)(x - x_1)f(x_0, x_1, x_2) + \dots$$

$$= -9 + (x-1)(-32) + (x-1)(x-2)(-14) + (x-1)(x-2)(x-4) \quad (7)$$

$$+ (x-1)(x-2)(x-4)(x-7) \quad (1)$$

$$= -9 - 32(x-1) - 14(x^2 - 3x + 2) + 7(x^3 - 7x^2 + 14x - 8)$$

$$+ (x^4 - 14x^3 + 63x^2 - 106x + 56)$$

$$= -5 + 2x + 0x^2 - 7x^3 + x^4$$

So,  $P_4(x) = x^4 - 7x^3 + 2x - 5.$

Now if we add a point  $(5, -173)$  to the data and ask for the polynomial which involves the 6 given points, all we have to do is to calculate one line in the divided difference table and add one term to the already calculated  $P_4(x)$ . The last line in the new table is

$$5 \quad -173 \quad 232 \quad 141 \quad 29 \quad 5 \quad 1$$

Hence the new polynomial will be

$$\begin{aligned} P_5(x) &= P_4(x) + (x - 1)(x - 2)(x - 4)(x - 7)(x - 8)(1) \\ &= x^4 - 7x^3 + 2x - 5 + (x^5 - 22x^4 + 175x^3 - 610x^2 + 904x - 448) \\ &= x^5 - 21x^4 + 168x^3 - 610x^2 + 906x - 453. \end{aligned}$$

#### EXERCISE 4

1. A function  $y = f(x)$  is given by the following table.

x	0	2.5069	5.0154	7.52270
f(x)	0.3989423	0.3988169	0.3984408	0.3978138

- a) Construct the divided differences table for the above data.  
 b) Find  $f(3.7608)$  using Newton's interpolation polynomial.
2. The following table gives the function  $y = f(x)$ . Calculate  $f(323.5)$  by Newton's interpolation polynomial.

x	321.0	322.8	324.2	325.0
f(x)	2.50651	2.50893	2.51081	2.51188

Compare your result with the exact value if  $f(x) = \log_{10} x$ .

3. The relation between steam pressure  $P$  and temperature  $T$  is given by the following table. Evaluate the pressure at temperature  $372.1^\circ$ .

T	$361^\circ$	$367^\circ$	$378^\circ$	$387^\circ$	$399^\circ$
P	154.9	167.0	191.0	212.5	244.2

4. Compute the value of  $\log_{10} 323.5$  by Lagrange's formula taking the data of problem 2. Compare your result with those obtained previously.
5. Construct the Lagrange interpolation polynomial for the function  $y = \sin \pi x$ , choosing the points
- $$x_0 = 0, x_1 = 1/6, x_2 = 1/2.$$
6. Find the second degree Lagrange interpolating polynomial passing through the three points in the following table.

i	$x_i$	$f(x_i)$
0	0	-5
1	1	1
2	3	25

7. Use appropriate Newton's interpolating polynomials of degree one, two, and three to approximate  $f(2.5)$  if

$$f(2.0) = .5103757 \quad f(2.6) = .4813306$$

$$f(2.2) = .5207843 \quad f(2.8) = .4359160$$

$$f(2.4) = .5104147$$

8. Use the Newton's interpolating polynomial of degree 3 or less to approximate  $\cos(.750)$  using the values given below.

$$\begin{aligned} \cos(.698) &= .7661 \\ \cos(.733) &= .7432 \end{aligned}$$

$$\begin{aligned} \cos(.768) &= .7193 \\ \cos(.803) &= .6946. \end{aligned}$$