An Introduction To Numerical Methods and Analysis Using MATLAB

DR. RIZWAN BUTT¹
Department of Mathematics, College of Science
King Saud University, P. O. Box 2455
Riyadh, 11451, Saudi Arabia.

May 5, 2015

¹E-mail: rizwanbu@ksu.edu.sa

ii Contents

Contents

1	Intr	roduction to Numerical Methods
	1.1	Introduction
	1.2	Error Analysis
	1.3	Sources of Errors
		1.3.1 Human Error
		1.3.2 Truncation Error
		1.3.3 Round-off Error
2	Solı	ution of Nonlinear Equations 7
	2.1	Introduction
	2.2	Method of Bisection
	2.3	Fixed-Point Method
	2.4	Newton's Method
	2.5	Secant Method
	2.6	Multiplicity of a Root
	2.7	Convergence of Iterative Methods
	2.8	Systems of Nonlinear Equations
	2.9	Exercises
3	Svs	tems of Linear Algebraic Equations 73
	3.1	Introduction
		3.1.1 Linear System in Matrix Notation
	3.2	Properties of Matrices and Determinant
		3.2.1 Introduction of Matrices
		3.2.2 Some Special Matrix Forms
		3.2.3 The Determinant of Matrix
		3.2.4 Matrix Inversion Method
	3.3	Solutions of Linear Systems of Equations
	3.4	Direct Numerical Methods for Linear Systems
		3.4.1 Gaussian Elimination Method
		3.4.2 Pivoting Strategies
		3.4.3 Gauss-Jordan Method
		3.4.4 LU Decomposition Method
	3.5	Norms of Vectors and Matrices
	3.6	Iterative Methods for Solving Linear Systems

		3.6.1 Jacobi Iterative Method
		3.6.2 Gauss-Seidel Iterative Method
		3.6.3 Matrix Forms of Iterative Methods for Linear System
		3.6.4 Convergence Criteria of Iterative Methods
	3.7	Errors in Solving Linear Systems
		3.7.1 Conditioning of Linear Systems
	3.8	Exercises
4	Pol	ynomial Interpolation and Approximation 181
	4.1	Introduction
	4.2	Polynomial Interpolation
		4.2.1 Lagrange Interpolating Polynomials
		4.2.2 Newton's General Interpolating Formula
	4.3	Interpolation with Spline Functions
		4.3.1 Piecewise Linear Interpolation
	4.4	Exercises
5	Nui	nerical Differentiation and Integration 239
	5.1	Introduction
	5.2	Numerical Differentiation
	5.3	Numerical Differentiation Formulas
		5.3.1 First Derivative Numerical Formulas
		5.3.2 Second Derivative Numerical Formula
	5.4	Formulas for Computing Derivatives
		5.4.1 Central Difference Formulas
		5.4.2 Forward and Backward Difference Formulas
	5.5	Numerical Integration
	5.6	Newton-Cotes Formulas
		5.6.1 Closed Newton-Cotes Formulas
	5.7	Exercises
6	Nui	nerical Solution of Ordinary Differential Equations 291
	6.1	Introduction
	6.2	Ordinary Differential Equations
		6.2.1 Classification of Differential Equations
	6.3	Numerical Methods for Solving IVP
	6.4	Single-Step Methods for IVP
		6.4.1 Euler's Method
		6.4.2 Analysis of the Euler's Method
		6.4.3 Higher-Order Taylor Methods
		6.4.4 Runge-Kutta Methods
	6.5	Exercises
In	dex	313

Chapter 1

Introduction to Numerical Methods

1.1 Introduction

I have written this book as an introductory course in numerical methods and numerical analysis for mathematicians, computer scientists, engineers, and other scientists. Numerical analysis is the branch of mathematics concerned with the theoretical foundations of numerical algorithms for the solution of problems arising in scientific applications. The subject addresses a variety of questions ranging from the approximation of functions and integrals to the approximate solution of algebraic, transcendental, differential and integral equations, with particular emphasis on the stability, accuracy, efficiency and reliability of numerical algorithms.

The intention of this book is to provide a gentle and sympathic introduction to many of the problems of scientific computing, and the wide variety of methods used for their solutions. The presentation of each numerical method is based on the successful teaching methodology of providing examples and geometric motivation for a method, and a concise statement of the steps to carry out the computation, before giving a mathematical derivation of the process or a discussion after more theoretical issues that are relevant to the use and understanding of the topic. Each topic illustrated by examples that range in complexity from very simple to moderate. Geometrical or graphical illustrations are included whenever they appropriate.

This book is concerned with the practical solution of problems on computers. In the process of problem solving, it is possible to distinguish several more or less distinct phases. The first phase is formulation. In formulating a mathematical model of a physical situation, the scientist should take into account beforehand the fact that he expects to solve his problem on a computer. He will therefore provide for specific objectives, proper input data, adequate checks, and for the type and amount of output. Once a problem has been formulated, numerical methods, together with a preliminary error analysis, must be devised for solving the problem. A numerical method which can be used to solve a problem will be called an algorithm. An algorithm is a complete and unambiguous set of procedures leading to the solution of a mathematical problem. The selection or construction of appropriate algorithms properly falls within the scope of numerical analysis. Having decided on a specific algorithm or set of algorithms for solving the problem, the numerical analyst should consider all the sources of error that may affect the results. He must consider how much accuracy is required, estimate the magnitude of the round-off and discretization error, determine an appropriate step size or the number of iterations required, provide for adequate checks on the accuracy, and make allowance for corrective action in cases of nonconvergence.

2 1.1 Introduction

The third phase of problem solving is programming. The programmer must transform the suggested algorithm into a set of unambiguous step-by-step instructions to the computer. the first step in this procedure is called *flow charting*. A flow chart is simply a set of procedures, usually in logical block form, which the computer will follow. It may be given in graphical or procedural statement form. The complexity of the problem and the amount of detail included. However, it should be possible for someone other than the programmer to follow the flow of information from the chart. The flow chart is an effective aid to the programmer, who must translate its major functions into machine code, and, at the same time, it is an effactive means of communication to others who wish to understand what the program does. In this book we sometimes use flow charts in graphical form, but more often in procedural statement form. Having produced a flow chart, the programmer must transform the indicated procedures into a set of machine instructions. This may be done directly in machine language or procedure-oriented language (sometimes called an algorithmic language). In this book MATLAB language is used exclusively. The practical justification of the methods is presented through computer examples through the use of MATLAB. In recent years, the number of MATLAB users has dramatically increased and now includes professionals who were trained in other high-level languages, Fortran, C, etc. but are now switching to MATLAB as well as students who are learning MATLAB as their first programming language. The surge of popularity in MATLAB is related to the increasing popularity of UNIX and computer graphics. To what extend numerical computations in the future will be programmed in MATLAB is uncertain. Nonetheless, there is no question that a need exists for comprehensive text especially geared to the requirements of those who want to learn, or use, numerical methods in MATLAB. This book has been written in response to this need.

The objectives of using MATLAB in this book include: (1) to be easily understood by undergraduate students with minimal knowledge of MATLAB, (2) to enable students to practice the methods in MATLAB, (3) to provide the short programs that can be easily used for scientific applications with or without modifications, and (4) to provide software that are easy to understand.

To provide maximum teaching flexibility, each chapter and each section begins with the basic, elementary material and gradually builds up to the more advanced material. The level of mathematical justification is determined largely by the desire to keep the mathematical prerequisites to a minimum. Thus, for example, no knowledge of linear algebra is assumed beyond the basic matrix algebra, and analytical results are based on a sound knowledge of the calculus.

In elementary calculus we learn how to differentiate and integrate to get exact answers to remarkably diverse range of realistic problems that could not be solved by purely algebraic methods. Unfortunately, from a practical point of view, the techniques of elementary (or even advanced) calculus alone are not adequate for solving calculus type problems such as solving polynomial equations of degree greater than four or even a simple equation such as

$$x = \cos x$$

also, to evaluate integrals of type

$$\int_a^b e^{x^2} dx$$
 and $\int_a^b \frac{\sin x}{x} dx$; etc.,

it is impossible to get the exact solutions of these problems. Even when an analytical solution can be found it may be of more theoretical than practical. Fortunately, one rarely needs exact answers.

Indeed, in the real world the problems themselves are usually inexact because they are generally possessed in terms of parameters that are measured, hence only approximate. What we are likely to require in a realistic situation is not an exact answer but rather one having a prescribed accuracy. The basic approach used to solve problems in numerical analysis is the algorithm which is used to describe a step-by-step procedure and requires a finite number of steps. So a numerical method is an algorithm which consists of a sequence of arithmetic and logical operations and which produces an approximate solution to within any prescribed accuracy. There are different numerical methods for the solution of one problem but the particular method chosen depends on the context from which the problem is taken.

Types of Numerical Methods

There are two basic types of numerical methods, direct numerical and indirect (iterative) numerical methods.

Direct methods compute the solution to a problem in a finite number of steps. These methods would give the precise answer if they were performed in infinite precision arithmetic. Examples include Gaussian elimination, the LU factorization method for solving systems of linear equations. In practice, finite precision is used and the result is an approximation of the true solution (assuming stability). In the absence of rounding errors, direct methods would deliver an exact solution.

In contrast to direct methods, iterative methods are not expected to terminate in a finite number of steps. Starting from an initial guess, iterative methods form successive approximations that converge to the exact solution only in the limit. In computational mathematics, an iterative method is a mathematical procedure that generates a sequence of improving approximate solutions for a class of problems. A specific implementation of an iterative method, including the termination criteria, is an algorithm of the iterative method. An iterative method is called convergent if the corresponding sequence converges for given initial approximations. A mathematically rigorous convergence analysis of an iterative method is usually performed. A convergence test, often involving the residual, is specified in order to decide when a sufficiently accurate solution has (hopefully) been found. Even using infinite precision arithmetic these methods would not reach the solution within a finite number of steps (in general). Iterative methods are often the only choice for nonlinear equations. However, iterative methods are often useful even for linear problems involving a large number of variables (sometimes of the order of millions), where direct methods would be prohibitively expensive (and in some cases impossible) even with the best available computing power. Examples include Newton's method, bisection method, and Jacobi iteration. In computational matrix algebra, iterative methods are generally needed for large problems.

An iterative method for the given problem converges means:- approximate values should come in side the given interval I- difference between two successive approximations should be small. Otherwise diverges. An iterative process may converge or diverge. If the divergence occurs, the procedure should be terminated because there may be no solution. We can restart the procedure by changing the initial approximation if necessary. But in the case of convergence we have to apply some stopping procedures to end the computations. In the following there are some more stopping criterion that can be used, each of them can be apply to any iterative technique considered in this chapter. By selecting a tolerance $\epsilon > 0$ and generate approximate solutions x_1, x_2, \ldots, x_n until one of the

4 1.2 Error Analysis

following conditions is satisfied:

$$|x_n - x_{n-1}| < \epsilon$$
 or $\frac{|x_n - x_{n-1}|}{|x_n|} < \epsilon$, $x_n \neq 0$.

Sometimes difficulties can arise using any of these stopping criteria. For example, there exist sequence $\{x_n\}_0^{\infty}$ with the property that the differences $(x_n - x_{n-1})$ converge to zero while the sequence itself diverges. It is also possible for $f(x_n)$ to be close to zero while x_n differs significantly from α . Without additional knowledge about f(x) or α , the above second inequality is the best stooping criterion to apply because it tests relative error. Also, one of the other stopping criteria is to use a fixed number of iterations, and then the final approximation x_n may be considered as the value of the required root. This type of stopping criteria is helpful when the convergence is very slow. It is important to note that in considering whether an iteration converges or not, it may be necessary to ignore the first few iterations since the procedure may appear diverge initially, even though it ultimately converges.

Iterative methods are more common than direct methods in numerical analysis. Some methods are direct in principle but are usually used as though they were not. For these methods the number of steps needed to obtain the exact solution is so large that an approximation is accepted in the same manner as for an iterative method. The numerical methods deal with numbers. We exam the sources of various types of computational errors.

1.2 Error Analysis

In generally numerical methods give an approximate solution (in number) of the given problem. How good is the approximate answer, we have to check by using the error analysis theory. There are basically two ways to know about the resulting approximation:- by using actual (exact) error and the error bound (upper bound) formulas of the using numerical methods.

What is Error

An approximate number p is a number that differs but slightly from an exact number α . We write

$$p \approx \alpha$$
.

By error E of an approximate number p, we mean the difference between the exact number α and its computed approximation p. Thus we define

$$E = \alpha - p. \tag{1.1}$$

If $\alpha > p$, the error E is positive, and if $\alpha < p$, the error E is negative. In many situations, the sign of the error may not be known and might even be irrelevant. Therefore, we define absolute error as

$$|E| = |\alpha - p|. \tag{1.2}$$

The relative error RE of an approximate number p is the ratio of the absolute error of the number to the absolute value of the corresponding exact number α . Thus

$$RE = \frac{|\alpha - p|}{|\alpha|}, \qquad \alpha \neq 0.$$
 (1.3)

If we approximate $\frac{1}{3}$ by 0.333, we have

$$E = \frac{1}{3} \times 10^{-3}$$
 and $RE = 10^{-3}$.

Note that relative error is generally a better measure of the extend of error than the actual error. But one should also note that relative error is undefined if the exact answer is equal to zero. Generally, we shall be interested in E (or sometimes |E|) rather than RE, but when the true value of a quantity is very small or very large, relative errors are more meaningful. For example, if the true value of a quantity is 10^{15} , and error of 10^6 is probably not serious, but this is more meaningfully expressed by saying that $RE = 10^{-9}$. In actual computation of the relative error, we shall often replace the unknown true value by the computed approximate value. Sometimes the quantity

$$\frac{|\alpha - p|}{|\alpha|} \times 100\%,\tag{1.4}$$

is defined as percentage error. From the above example, we have

$$PE = 0.001 \times 100 = 0.1\%$$
.

In investigating the effect of the total error in various methods, we shall often mathematically derive an error, called, error bound and which is a limit on how large the error can be. We shall have the reason to compute error bounds in many situations. This applies to both absolute and relative errors. Note that the error bound can be much larger than the actual error and that this is often the case in practice. Any mathematically derived error bound must account for the worst possible case that can occur and is often based upon certain simplifying assumptions about the problem which in many practical cases cannot be actually tested. For the error bound to be used in any practical way, the user must have a good understanding of how the error bound was derived in order to know how crude it is, that is, how likely it is to over estimate the actual error. Of course, whenever possible, our goal is to eliminate or lesser the effects of errors, rather than trying to estimate them after they occur.

1.3 Sources of Errors

In analysing the accuracy of numerical result, one should be aware of the possible sources of error in each stage of the computational process and of the extend to which these errors can affect the final answer. We will consider that there are three types of errors which occur in a computation. We discuss them step by step as follows.

1.3.1 Human Error

These types of errors arise when the equations of the mathematical model are formed, due to sources such as the idealistic assumptions made to simplify the model, inaccurate measurements of data, miscopying of figures, the inaccurate representation of mathematical constants (for example, if the constant π occurs in an equation, we must replace π by 3.1416 or 3.141593, etc.).

6 1.3 Sources of Errors

1.3.2 Truncation Error

This type of error is caused when we are forced to use mathematical techniques which give approximate, rather than exact, answer. For example, suppose that we use the Maclaurin's series expansion to represent $\sin x$, so that

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \cdots$$

If we want a number that approximates $\sin(\frac{\pi}{2})$, we must terminate the expansion in order to obtain

$$\sin(\frac{\pi}{2}) = \frac{\pi}{2} - \frac{(\pi/2)^3}{3!} + \frac{(\pi/2)^5}{5!} - \frac{(\pi/2)^7}{7!} + E,$$

where E is the truncation error introduced in the calculation. Truncation errors in numerical analysis usually occur because many numerical methods are iterative in nature, with the approximations theoretically becoming more accurate as we take more iterations. As a practical matter, we must stop the iteration after a finite number of steps, thus introducing a truncation error. The Taylor series is the most important means used to derive numerical schemes and analysis truncation errors.

1.3.3 Round-off Error

This type of errors are associated with the limited number of digits numbers in the computer. For example, by rounding off 1.32463672 to six decimal places to give 1.324637. Any further calculation involving such a number will also contain an error. Round-off numbers according to following rules:

- 1. If first discarded digit is less than 5, leave the remaining digits of number unchanged, that is, $48.47263 \approx 48.4726$.
- 2. If the first discarded digit is exceeds 5, add 1 to the retained digit. For example, $48.4726 \approx 48.473$.
- 3. If the first discarded digit is exactly 5 and there are nonzero among those discarded, add, 1 to the last retained digit. For example, $3.0554 \approx 3.06$.
- 4. If the first discarded digit is exactly 5 and all other discarded digits are zero, the last retained digit is left unchanged if it is even, otherwise 1 is added to it. For example,

$$3.05500 \approx 3.06$$

 $3.04500 \approx 3.04$.

with these rules, the error is never larger in magnitude than one-half unit of the place of the nth digit in the rounded number. To understand the nature of round-off errors, it is necessary to learn the ways numbers are stored and additions and subtractions are performed in a computer.

• A solution is correct within k decimal places if the error is less than 0.5×10^{-k} .

If x^* is an approximation to x, then we say that x^* approximates x to k significant digits if k is the largest nonnegative integer for which $\left|\frac{x-x^*}{x}\right| < 5 \times 10^{-k}$.

Chapter 2

Solution of Nonlinear Equations

2.1 Introduction

In this chapter we study one of the fundamental problems of numerical analysis, namely the numerical solution of nonlinear equations. Most equations arising in practice are nonlinear and are rarely of a form which allows the roots to be determined exactly. Consequently, numerical methods are used to solve nonlinear algebraic equations when the equations prove intractable to ordinary mathematical techniques. These numerical methods are all iterative, and they may be used for equations that contain one or several variables. These techniques can be divided into two categories; one-point (need one initial approximation) and two-point (need two initial approximations) methods.

Important Points of the Chapter 2

- I. A nonlinear equation in this chapter may be considered any one of the following types:
 - 1. An equation may be an algebraic equation (a polynomial equation of degree n) expressible in the form:

$$a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0 = 0, \quad a_n \neq 0, \quad n > 1,$$

where $a_n, a_{n-1}, \ldots, a_1$ and a_0 are constants. For example, the following equations are nonlinear.

$$x^{2} + 5x + 6 = 0;$$
 $x^{3} = 2x + 1;$ $x^{200} - 2x + 1 = 0.$

2. The power of the unknown variable (not a positive integer number) involved in the equation must be difficult to manipulate. For example, the following non-polynomial equations are nonlinear

$$x^{-1} + 2x = 1;$$
 $\sqrt{x} + x = 2;$ $x^{2/3} + \frac{2}{x} + 4 = 0.$

3. An equation may be a *transcendental equation*, the equation which involves the trigonometric functions, exponential functions and logarithmic functions. For example, all the following transcendental equations are nonlinear

$$x = cos(x);$$
 $e^x + x - 10 = 0;$ $x + \ln x = 10.$

8 2.2 Method of Bisection

II. Given nonlinear equation must be put in the following form

$$f(x) = 0$$
, where $f(x)$ must be nonlinear function.

III. There may be many roots of the given nonlinear equation but we will seek the approximation of only one of its real root α lies in the given interval [a, b], that is

$$f(\alpha) = 0$$
, where $\alpha \in [a, b]$.

IV. If f(x) is continuous function in a interval [a,b] and f(x) has opposite signs at the end points of the interval, then there must be a root of nonlinear equation f(x) = 0 in [a,b].

V. Root of a nonlinear equation may be **simple** (not repeating) or **multiple** (repeating). Simple root means

$$f(\alpha) = 0$$
 but $f'(\alpha) \neq 0$.

For example, $\alpha_1 = -3$ and $\alpha_2 = -2$ are the simple roots of the nonlinear equation $x^2 + 5x + 6 = 0$. For the multiple root, we mean

$$f(\alpha) = 0$$
 but $f'(\alpha) = 0$.

For example, $\alpha_1 = -2$ and $\alpha_2 = -2$ are the multiple roots of the nonlinear equation $x^2 + 4x + 4 = 0$. **VI.** The methods we will consider in this chapter are iterative methods and they are, bisection method, fixed-point method, Newton method (also called, Newton-Raphson method) and secant method which give us the approximation of single (or simple) root of the nonlinear equation. For the multiple roots of the nonlinear equation we will use other iterative methods, called, the first modified Newton's method (also called the Schroeder's method) and the second modified Newton's method. The iterative methods for the approximation of simple root can be use also for the approximation of the multiple roots but they are very slow. All the numerical methods described in this chapter are applicable to general nonlinear functions. The iterative methods we will discuss in this chapter are basically of two types: one in which the convergence is guaranteed and the other in which the convergence depends on the initial approximation.

VII. Remember that the best method for the approximation of the simple root of nonlinear equation is Newton's method (called quadratic convergent method) and for multiple root of nonlinear equation is modified Newton's method (called quadratic convergent method). Newton's method for multiple root of nonlinear root is called a linear convergent method.

Definition 2.1 (Root of an Nonlinear Equation)

Assume that f(x) is a continuous function. An number α for which $f(\alpha) = 0$ is called a root of the equation f(x) = 0 or a zero of the function f(x).

First, we shall discuss the numerical iterative methods for simple root of nonlinear equations in a *single variable*. The problem here can be simply written down as:

$$f(x) = 0. (2.1)$$

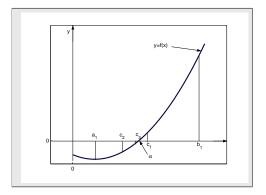


Figure 2.1: Graphical Solution of Bisection Method.

2.2 Method of Bisection

This is one of the simplest iterative technique for determining roots of (2.1) and it needs two initial approximations to start. It is based on the *Intermediate Value Theorem*. This method is also called the *interval-halving method* because the strategy is to bisect or halve the interval from one endpoint of the interval to the other endpoint and then retain the half interval whose end still bracket the root. It is also referred to a *bracketing method* or sometimes called the *Bolzano's method*. The fact that the function is required to change sign only once gives us a way to determine which half interval to retain; we keep the half on which f(x) changes sign or became zero. The basis for this method can be easily illustrated by considering a function y = f(x). Our object is to find an x value for which y is zero. Using this method, we begin by supposing f(x) is a continuous function defined on the interval [a, b] and then by evaluation the function at two x values, say, a and b, such that

The implication is that one of the values is negative and the other is positive. These conditions can be easily satisfied by sketching the function, see Figure 2.1. Obviously, the function is negative at one endpoint a of the interval and positive at other endpoint b and is continuous on $a \le x \le b$. Therefore the root must lies between a and b (by Intermediate Value Theorem) and a new approximation to the root a be calculated as

$$c = \frac{a+b}{2},$$

and, in general

$$c_n = \frac{a_n + b_n}{2}, \qquad n \ge 1. \tag{2.2}$$

The iterative formula (2.2) is known as the *bisection method*.

If $f(c) \approx 0$, then $c \approx \alpha$ is the desired root, and, if not, then there are two possibilities. Firstly, if f(a)f(c) < 0, then f(x) has a zero between point a and point c. The process can then be repeated on the new interval [a, c]. Secondly, if f(a)f(c) > 0 it follows that f(b)f(c) < 0 since it is known that f(b) and f(c) have opposite signs. Hence, f(x) has zero between point c and point c and the process can be repeated with [c, b]. We see that after one step of the process, we have found

10 2.2 Method of Bisection

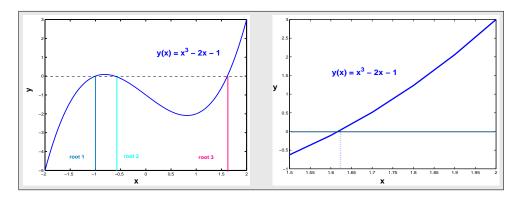


Figure 2.2: Graphical Solution of $x^3 = 2x + 1$ in the intervals [-2, 2] and [-1.5, 2].

either a zero or a new bracketing interval which is precisely half the length of the original one. The process continue until the desired accuracy is achieved. We use the bisection process in the following example.

Example 2.1 Use the bisection method to find the approximation to the root of the equation

$$x^3 = 2x + 1$$
,

that is located in the interval [1.5, 2.0] accurate to within 10^{-2} .

Solution. Since the given function $f(x) = x^3 - 2x - 1$ is a polynomial function and so is continuous on [1.5, 2.0], starting with $a_1 = 1.5$ and $b_1 = 2$, we compute:

$$a_1 = 1.5$$
: $f(a_1) = -0.625$
 $b_1 = 2.0$: $f(b_1) = 3.0$,

and since f(1.5)f(2.0) < 0, so that a root of f(x) = 0 lies in the interval [1.5, 2.0]. Using formula (2.2) (when n = 1), we get:

$$c_1 = \frac{a_1 + b_1}{2} = 1.75;$$
 $f(c_1) = 0.859375.$

Hence the function changes sign on $[a_1, c_1] = [1.5, 1.75]$. To continue, we squeeze from right and set $a_2 = a_1$ and $b_2 = c_1$. Then the midpoint is:

$$c_2 = \frac{a_2 + b_2}{2} = 1.625;$$
 $f(c_2) = 0.041056.$

Continue in this way we obtain a sequence $\{c_k\}$ of approximation shown by Table 2.1. We see that the functional values are approaching zero as the number of iterations is increase. We got the desired approximation to the root of the given equation is $c_6 = 1.617188 \approx \alpha$ after 6 iterations with accuracy $\epsilon = 10^{-2}$.

To use MATLAB command for the bisection method, first we define a function m-file as fn.m for the equation as follows:

Ta	Table 2.1: Solution of $x^3 = 2x + 1$ by disection method					
n	Left		Right	Function Value		
	Endpoint a_n	Midpoint c_n	Endpoint b_n	$f(c_n)$		
01	1.500000	1.750000	2.000000	0.8593750		
02	1.500000	1.625000	1.750000	0.0410156		
03	1.500000	1.562500	1.625000	-0.3103027		
04	1.562500	1.593750	1.625000	-0.1393127		
05	1.593750	1.609375	1.625000	-0.0503273		
06	1.609375	1.617188	1.625000	-0.0049520		

function
$$y = fn(x)$$

 $y = x \cdot 3 - 2 \cdot x - 1;$

then use the single commands:

$$>> s = bisect('fn', 1.5, 2, 1e - 2)$$

We can easily find the roots (1.61803399, -1.00, -0.61803399) of the equation $x^3 = 2x + 1$ by defining the coefficients of the polynomial equation using MATLAB commands as:

$$>> CP = [1\ 0\ -2\ -1];\ Sol = roots(CP);$$

Example 2.2 Find the point of intersection of the graphs $y = x^3 + 2x - 1$ and $y = \sin x$, then use bisection method within accuracy 10^{-3} .

Solution. The graphs in the Figure 2.3 show that there is an intersection at about point (0.66, 0.61). Using the function $f(x) = x^3 + 2x - \sin x - 1$ and the starting interval [0.5, 1.0], we compute:

$$a_1 = 0.5$$
: $f(a_1) = -0.3544$,
 $b_1 = 1.0$: $f(b_1) = 1.1585$.

Since f(x) is continuous on [0.5, 1.0] and f(0.5).f(1.0) < 0, so that a root of f(x) = 0 lies in the interval [0.5, 1.0]. Using formula (2.2) (when n = 1), we get:

$$c_1 = \frac{a_1 + b_1}{2} = 0.75;$$
 $f(c_1) = 0.240236.$

Hence the function changes sign on $[a_1, c_1] = [0.5, 0.75]$. To continue, we squeeze from right and set $a_2 = a_1$ and $b_2 = c_1$. Then the midpoint is:

$$c_2 = \frac{a_2 + b_2}{2} = 0.625;$$
 $f(c_2) = -0.090957.$

Then continue in this manner we obtain a sequence $\{c_k\}$ of approximation shown by Table 2.2.

12 2.2 Method of Bisection

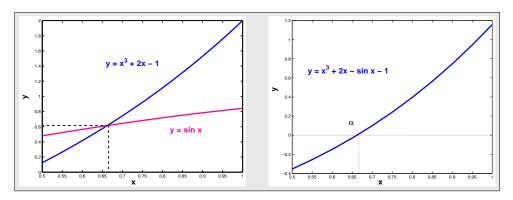


Figure 2.3: Graphical Solution of $\sin x = x^3 + 2x - 1$ and $x^3 + 2x - \sin x = 1$.

Table 2.2: Solution of $x^3 + 2x - \sin x - 1$ by bisection method

n	Left	Right		Function Value
	Endpoint a_n	Endpoint b_n	Midpoint c_n	$f(c_n)$
01	0.5000	1.0000	0.750000	0.240236
02	0.5000	0.7500	0.625000	-0.090957
03	0.6250	0.7500	0.687500	0.065344
:	:	:	:	:
07	0.6563	0.6641	0.660156	-0.005228
08	0.6602	0.6641	0.662109	-0.000302

Program 2.1

MATLAB m-file for the Bisection Method

 $function \ sol = bisect(fn, a, b, tol)$

fa = feval(fn, a); fb = feval(fn, b);

if fa * fb > 0; fprintf('Endpoints have same sign') return end while $abs (b-a) > tol \ c = (a+b)/2$; fc = feval(fn, c);

if fa * fc < 0; b = c; else a = c; end; end; sol=(a + b)/2;

We see that the functional values are approaching zero as the number of iterations is increase. We got the desired approximation to the root of the given equation is $c_8 = 0.662109 \approx \alpha$ after 8 iterations with accuracy $\epsilon = 10^{-3}$.

Theorem 2.1 (Bisection Convergence and Error Theorem)

Let f(x) be continuous function defined on the given initial interval $[a_0, b_0] = [a, b]$ and suppose that f(a)f(b) < 0. Then bisection method (2.2) generates a sequence $\{c_n\}_{n=1}^{\infty}$ approximating $\alpha \in (a, b)$ with the property

$$|\alpha - c_n| \le \frac{b - a}{2^n}, \quad n \ge 1. \tag{2.3}$$

Moreover, to obtain accuracy of

$$|\alpha - c_n| \le \epsilon$$
,

(for $\epsilon = 10^{-k}$) it suffices to take

$$n \ge \frac{\ln\left\{10^k(b-a)\right\}}{\ln 2},\tag{2.4}$$

where k is nonnegative integer.

Proof.

Since both the root α and the midpoint c_1 lie in the interval [a,b], the distance between them cannot be greater than of this width interval. Thus

$$|\alpha - c_n| \le \frac{b_{n-1} - a_{n-1}}{2}$$
, for all n .

Observe that

$$b_1 - a_1 = \frac{b_0 - a_0}{2},$$

then

$$b_2 - a_2 = \frac{b_1 - a_1}{2} = \frac{b_0 - a_0}{2^2}.$$

Finite mathematical induction is used to conclude that

$$b_{n-1} - a_{n-1} = \frac{b_0 - a_0}{2^{n-1}}.$$

Therefore, the error is bounded as follows

$$|\alpha - c_n| \le \frac{b_{n-1} - a_{n-1}}{2} = \frac{b_0 - a_0}{2^n},$$

gives the estimate.

Now to establish the bound on the number of bisections n (or iterations), we simply observe that

$$\frac{b-a}{2^n} \le 10^{-k},$$

together with (2.3) implies that

$$|\alpha - c_n| \le 10^{-k},$$

that is, we wish to calculate a root to within 10^{-k} . Since $2^n \ge 10^k(b-a)$, so by taking logarithms, we get $n \ln 2 \ge \ln \left\{ 10^k(b-a) \right\}$, and solving for n, we get the inequality (2.4).

The above Theorem 2.1 gives us information about bounds for errors in approximation and the number of bisections needed to obtain any given accuracy.

14 2.2 Method of Bisection

Example 2.3 Show that number of iterations of bisection will require to attain an accuracy of 10^{-4} using the starting interval [a, b] is

$$n \ge \frac{\ln(b-a) + 4\ln 5}{\ln 2} + 4.$$

Determine the number of iterations needed to achieve the an approximation with same above given accuracy to the solution of $x^3 - 2x - 1 = 0$ lying in the interval [1.5, 2].

Solution. By using the inequality (2.4), we get

$$n \ge \frac{\ln\left\{10^k(b-a)\right\}}{\ln 2} = \frac{\ln(b-a) + \ln[(5)(2)]^4}{\ln 2},$$

or

$$n \ge \frac{\ln(b-a) + 4[\ln 5 + \ln 2]}{\ln 2} = \frac{\ln(b-a) + 4\ln 5}{\ln 2} + 4.$$

Now by taking a = 1.5 and b = 2 in the above inequality, we get

$$n \ge \frac{\ln(2 - 1.5) + 4\ln 5}{\ln 2} + 4 = 8.2877 + 4 = 12.2877.$$

So no more than thirteen iterations are required to obtain an approximation accurate to within the given accuracy 10^{-4} .

Example 2.4 Find a bound for the number of iterations needed to achieve an approximation with accuracy 10^{-1} to the solution of $xe^x = 1$ lying in the interval [0.5, 1] using the bisection method. Find an approximation to the root with this degree of accuracy.

Solution. Here a = 0.5, b = 1 and k = 1, then by using inequality (2.4), we get

$$n \ge \frac{\ln[10^1(1-0.5)]}{\ln 2} \approx 2.3219.$$

So no more than three iterations are required to obtain an approximation accurate to within 10^{-1} . The given function $f(x) = xe^x - 1$ is continuous on [0.5, 1.0], so starting with $a_1 = 0.5$ and $b_1 = 1$, we compute:

$$a_1 = 0.5$$
: $f(a_1) = -0.1756$,
 $b_1 = 1$: $f(b_1) = 1.7183$,

since f(0.5)f(1) < 0, so that a root of f(x) = 0 lies in the interval [0.5, 1]. Using formula (2.2) (when n = 1), we get:

$$c_1 = \frac{a_1 + b_1}{2} = 0.75;$$
 $f(c_1) = 0.5878.$

Hence the function changes sign on $[a_1, c_1] = [0.5, 0.75]$. To continue, we squeeze from right and set $a_2 = a_1$ and $b_2 = c_1$. Then the bisection formula gives

$$c_2 = \frac{a_2 + b_2}{2} = 0.625;$$
 $f(c_2) = 0.1677.$

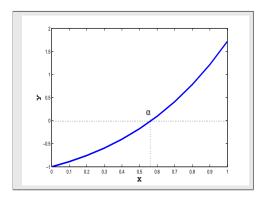


Figure 2.4: Graphical Solution of $xe^x = 1$.

Finally, we have in the similar manner as

$$c_3 = \frac{a_3 + b_3}{2} = 0.5625,$$

the value of the third approximation which is accurate to within 10^{-1} .

It is important to keep in mind that the error analysis gives only a bound for the number of iterations necessary, and in many cases this bound is much larger than the actual number required.

Example 2.5 Use the bisection method to compute the first three approximate values for $\sqrt[4]{18}$. Also, compute an error bound and absolute error for your approximation.

Solution. Consider

$$x = \sqrt[4]{18} = (18)^{1/4}, \quad or \quad x^4 - 18 = 0.$$

Choose the interval [2,2.5] on which the function $f(x) = x^4 - 18$ is continuous and the function f(x) satisfies the sign property, that is

$$f(2)f(2.5) = (-2)(21.0625) = -42.125 < 0.$$

Hence root $\alpha = \sqrt[4]{18} = 2.0598 \in [2, 2.5]$ and we compute its first approximate value by using formula (2.2) (when n = 1) as follows:

$$c_1 = \frac{2.0 + 2.5}{2} = 2.2500$$
 and $f(2.25) = 7.6289$.

Since the function f(x) changes sign on [2.0, 2.25]. To continue, we squeeze from right and use formula (2.2) again to get the following second approximate value of the root α as:

$$c_2 = \frac{2.0 + 2.25}{2} = 2.1250$$
 and $f(2.1250) = 2.3909$.

Then continue in the similar way, the third approximate value of the root α is $c_3 = 2.0625$ with f(2.0625) = 0.0957. Note that the value of the function at each new approximate value is decreasing

16 2.3 Fixed-Point Method

which shows that the approximate values are coming closer to the root α . Now to compute the error bound for the approximation we use the formula (2.3) and get

$$|\alpha - c_3| \le \frac{2.5 - 2.0}{2^3} = 0.0625,$$

which is the possible maximum error in our approximation and

$$|E| = |2.0598 - 2.0625| = 0.0027,$$

be the absolute error in the approximation.

One drawback of the bisection method is the convergence rate is raster slow. However, the rate of convergence is guaranteed. So for this reason it is often used as a started for the more efficient used to find roots of the nonlinear equations. The method may give a false root if f(x) is discontinuous on the given interval [a, b].

Procedure 2.1 (Bisection Method)

- 1. Establish an interval $a \le x \le b$ such that f(a) and f(b) are of opposite sign, that is, f(a).f(b) < 0.
- 2. Choose an error tolerance $(\epsilon > 0)$ value for the function.
- 3. Compute a new approximation for the root: $c_n = \frac{(a_n + b_n)}{2}$; $n = 1, 2, 3, \dots$
- 4. Check tolerance. If $|f(c_n)| \leq \epsilon$, use c_n , $n \geq 1$ for desired root; otherwise continue.
- 5. Check, if $f(a_n)f(c_n) < 0$, then set $b_n = c_n$; otherwise set $a_n = c_n$.
- 6. Go back to step 3, and repeat the process.

2.3 Fixed-Point Method

This is another iterative method to solve the nonlinear equation (2.1) and needs one initial approximation to start. This is a very general method for finding the root of (2.1) and it provides us with a theoretical framework within which the convergence properties of subsequent methods can be evaluated. The basic idea of this method which is also called successive approximation method or function iteration, is to rearrange the original equation

$$f(x) = 0, (2.5)$$

into an equivalent expression of the form

$$x = g(x). (2.6)$$

Any solution of (2.6) is called a fixed-point for the iteration function q(x) and hence a root of (2.5).

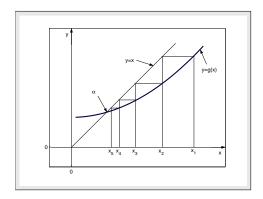


Figure 2.5: Graphical Solution of Fixed-Point Method.

Definition 2.2 (Fixed-Point of a Function)

A fixed-point of a function
$$g(x)$$
 is a real number α such that $\alpha = g(\alpha)$.
For example, $x = 2$ is a fixed-point of the function $g(x) = \frac{x^2 - 4x + 8}{2}$ because $g(2) = 2$.

The task of solving (2.5) is therefore reduced to that of finding a point satisfying the fixed-point condition (2.6). The fixed-point method essentially solves two functions simultaneously; y = x and y = g(x). The point of intersection of these two functions is the solution to x = g(x), and thus to f(x) = 0, see Figure 2.5.

This method is conceptually very simple. Since g(x) is also nonlinear, the solution must be obtained iteratively. An initial approximation to the solution, say, x_0 , must be determined. For choosing the best initial value x_0 for using this iterative method, we have to find an interval [a, b] on which the original function f(x) satisfies the sign property and then use the midpoint $\frac{a+b}{2}$ as the initial approximation x_0 . Then this initial value x_0 is substituted in the function g(x) to determined the next approximation x_1 and so on.

Definition 2.3 (Fixed-Point Method)

The iteration defined in the following

$$x_{n+1} = g(x_n); \qquad n = 0, 1, 2, \dots,$$
 (2.7)

is called the fixed-point method or the fixed-point iteration.

The value of the initial approximation x_0 is chosen arbitrarily and the hope is that the sequence $\{x_n\}_{n=0}^{\infty}$ converges to a number α which will automatically satisfies (2.5). Moreover, since (2.5) is a rearrangement of (2.6), α is guaranteed to be a zero of f(x). In general, there are many different ways of rearranging of (2.6) in (2.5) form. However, only some of these are likely to give rise to successful iterations but sometime we don't have successful iterations. To describe such behaviour, we discuss the following example.

18 2.3 Fixed-Point Method

Example 2.6 Consider the nonlinear equation $x^3 = 2x+1$ which has a root in the interval [1.5, 2.0] using fixed-point method with $x_0 = 1.5$, take three different rearrangements for the equation.

Solution. Let us consider the three possible rearrangement of the given equation as follows:

(i)
$$x_{n+1} = g_1(x_n) = \frac{(x_n^3 - 1)}{2};$$
 $n = 0, 1, 2, ...,$

(ii)
$$x_{n+1} = g_2(x_n) = \frac{1}{(x_n^2 - 2)};$$
 $n = 0, 1, 2, ...,$

(iii)
$$x_{n+1} = g_3(x_n) = \sqrt{\frac{(2x_n+1)}{x_n}}; \qquad n=0,1,2,\ldots,$$

then the numerical results for the corresponding iterations, starting with the initial approximation $x_0 = 1.5$ with accuracy 5×10^{-2} , are given in Table 2.3. We note that the first two considered

Table 2.3: Solution of $x^3 = 2x + 1$ by fixed-point method

n	x_n	$x_{n+1} = g_1(x_n)$	$x_{n+1} = g_2(x_n)$	$x_{n+1} = g_3(x_n)$
		$= (x_n^3 - 1)/2$	$=1/(x_n^2-2)$	$=\sqrt{(2x_n+1)/x_n}$
00	x_0	1.500000	1.500000	1.500000
01	x_1	1.187500	4.000000	1.632993
02	x_2	0.337280	0.071429	1.616284
03	x_3	-0.480816	-0.501279	1.618001
04	x_4	-0.555579	-0.571847	1.618037
05	x_5	-0.585745	-0.597731	1.618034

sequences diverge and the last one converges. This example asks the need for a mathematical analysis of the method. The following theorem gives sufficient conditions for the convergence of the fixed-point iteration.

Theorem 2.2 (Fixed-Point Theorem)

If g is continuously differentiable on the interval [a,b] and $g(x) \in [a,b]$ for all $x \in [a,b]$, then

(a) g has at-least one fixed-point in the given interval [a,b].

Moreover, if the derivative g'(x) of the function g(x) exists on an interval [a,b] which contains the starting value x_0 , with

$$k \equiv \max_{a \le x \le b} |g'(x)| < 1; \quad \text{for all} \quad x \in [a, b].$$
 (2.8)

Then

- (b) The sequence (2.7) will converge to the attractive (unique) fixed-point α in [a,b].
- (c) The iteration (2.7) will converge to α for any initial approximation.

(d) We have the error estimate

$$|\alpha - x_n| \le \frac{k^n}{1 - k} |x_1 - x_0|, \quad \text{for all} \quad n \ge 1.$$
 (2.9)

(e) The limit holds:

$$\lim_{n \to \infty} \frac{\alpha - x_{n+1}}{\alpha - x_n} = g'(\alpha). \tag{2.10}$$

Proof

(a) Suppose g is continuous on [a,b] and $g(x) \in [a,b]$. We need to show it has a fixed point. If g(a) = a and g(b) = b, then the function g has a fixed-point at the endpoints. Suppose that it is not happening, that is, $g(a) \neq a$ and $g(b) \neq b$ and define a function f(x) = g(x) - x which is continuous on [a,b]. Then f(x) has a zero in [a,b] if and only if g(x) has a fixed point in [a,b] but

$$f(a) = g(a) - a > 0,$$

since g(a) is in [a,b] and hence cannot be smaller than a, and we have assumed that g(a) is not equal to a. Similarly,

$$f(b) = g(b) - b < 0,$$

and so by the Intermediate Value Theorem there is a α in the interval (a,b) such that $f(\alpha) = 0$, which implies that $\alpha = g(\alpha)$. Thus the function g(x) has at least one fixed-point in [a,b]. This proves (a).

(b) Suppose now that (2.8) holds, and α and β are two fixed-points of the function g in [a,b]. Then we have

$$\alpha = g(\alpha)$$
 and $\beta = g(\beta)$.

In addition, by the Mean Value Theorem, we have that for any two points α and β in [a,b], there exits a number η such that

$$|\alpha - \beta| = |g(\alpha) - g(\beta)| = |g'(\eta)||\alpha - \beta| \le k|\alpha - \beta|,$$

where $\eta \in (a,b)$. Thus

$$|\alpha - \beta| - k|\alpha - \beta|$$
 or $(1 - k)|\alpha - \beta| \le 0$.

Since k < 1, we must have $\alpha = \beta$; and thus, the function g has a unique fixed-point α in the interval [a, b]. This proves (b).

(c) For the convergence, consider the iteration

$$x_n = g(x_{n-1}), \quad \text{for all} \quad n \ge 1, 2 \dots,$$

and the definition of the fixed-point, that is

$$\alpha = q(\alpha)$$
.

20 2.3 Fixed-Point Method

If we subtract last two equations and take the absolute values, we get

$$|\alpha - x_n| = |g(\alpha) - g(x_{n-1})| \le k|\alpha - x_{n-1}|.$$

The recursion can be solved readily to get

$$|\alpha - x_n| \le k|\alpha - x_{n-1}| \le k^2|\alpha - x - n - 2|\dots \le k^n|\alpha - x_0|,$$
 (2.11)

from which it follows that

as
$$n \to \infty$$
, $k^n \to 0$, (since $k < 1$),

therefore, $x_n \to \alpha$. Hence the iteration converges. This proves (c).

(d) Since we note that

$$|\alpha - x_0| = |\alpha - x_1 + x_1 - x_0| \le |\alpha - x_1| + |x_1 - x_0|$$

 $\le |g(\alpha) - g(x_0)| + |x_1 - x_0| \le k|\alpha - x_0| + |x_1 - x_0|,$

which gives

$$|\alpha - x_0| - k|\alpha - x_0| \le |x_1 - x_0|$$
 or $(1 - k)|\alpha - x_0| \le |x_1 - x_0|$,

and from this it follows that

$$|\alpha - x_0| \le \frac{1}{1-k} |x_1 - x_0|.$$

From (2.11), we can write above equation as follows

$$|\alpha - x_n| \le \frac{k^n}{1 - k} |x_1 - x_0|,$$

which proves (d).

(e) Finally, by subtracting iteration $x_{n+1} = g(x_n)$ and $\alpha = g(\alpha)$, we have

$$\alpha - x_{n+1} = g(\alpha) - g(x_n) = g'(\eta(x))(\alpha - x_n),$$

which implies that

$$\frac{\alpha - x_{n+1}}{\alpha - x_n} = g'(\eta(x)),$$

and by taking limits, we have

$$\lim_{n \to \infty} \frac{\alpha - x_{n+1}}{\alpha - x_n} = \lim_{n \to \infty} g'(\eta(x)) = g'(\alpha),$$

since $\eta(x) \to \alpha$ is forced by the convergence of x_n to α . This proves (e).

Now we come back to our previous Example 2.6 and discuss that why the first two rearrangements we considered, do not converge but on the other hand, last sequence has a fixed-point and converges. Since, we observe that f(1.5)f(2) < 0, then the solution we seek is in the interval [1.5, 2].

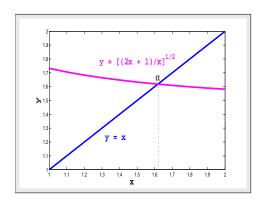


Figure 2.6: Graphical Solution of $x = \sqrt{(2x+1)/x}$.

- (i) For $g_1(x) = \frac{x^3 1}{3}$, we have $g'_1(x) = x^2$, which is greater than unity throughout the interval [1.5, 2]. So by Fixed-Point Theorem 2.2 this iteration will fail to converge.
- (ii) For $g_2(x) = \frac{1}{x^2 2}$, we have $g'_2(x) = \frac{-2x}{(x^2 2)^2}$, and $|g'_2(1.5)| > 1$, so from Fixed-Point Theorem 2.2 this iteration will fail to converge.
- (iii) For $g_3(x) = \sqrt{\frac{2x+1}{x}}$, we have $g_3'(x) = x^{-3/2}/2\sqrt{2x+1} < 1$, for all x in the given interval [1.5,2]. Also, g_3 is decreasing function of x, and $g_3(1.5) = 1.63299$ and $g_3(2) = 1.58114$ both lie in the interval [1.5,2]. Thus $g_3(x) \in [1.5,2]$, for all $x \in [1.5,2]$, so from Fixed-Point Theorem 2.2 the iteration will converge, see Figure 2.6.

Note 2.1 From (2.9) Note that the rate of convergence of the fixed-point method depends on the factor $\frac{k^n}{(1-k)}$; the smaller the value of k, then faster the convergence. The convergence may be very slow if the value of k is very close to 1.

Note 2.2 Assume that g(x) and g'(x) are continuous functions of x for some open interval I, with the fixed-point α contained in this interval. Moreover assume that

$$|g'(\alpha)| < 1, \quad for \quad \alpha \in I,$$

then, there exists an interval [a, b], around the solution α for which all the conditions of Theorem 2.2 are satisfied. But if

$$|g'(\alpha)| > 1$$
, for $\alpha \in I$,

then the sequence (2.7) will not converge to α . In this case α is called a repulsive fixed-point. If

$$|g'(\alpha)| = 0,$$
 for $\alpha \in I$,

then the sequence (2.7) converges very fast to the root α while if

$$|g'(\alpha)| = 1,$$
 for $\alpha \in I$,

22 2.3 Fixed-Point Method

then the convergence the sequence (2.7) is not guaranteed and if the convergence happened, it would be very slow. Thus to get the faster convergence, the value of $|g'(\alpha)|$ should be equal to zero or very close to zero.

Example 2.7 If α and β are roots of the nonlinear equation $x^2 + ax + b = 0$, then show that the following iterative schemes

$$x_{n+1} = -\left(\frac{ax_n + b}{x_n}\right) = g_1(x_n), \qquad n \ge 0,$$

will converge near $x = \alpha$ if $|\alpha| > |\beta|$ and the other scheme

$$x_{n+1} = -\left(\frac{b}{x_n + a}\right) = g_2(x_n), \qquad n \ge 0,$$

will converge near $x = \alpha$ if $|\alpha| < |\beta|$.

Solution. Since α and β are roots of the nonlinear equation $x^2 + ax + b = 0$, then we have

$$\alpha + \beta = -a$$
 and $\alpha\beta = b$.

The first iterative scheme will converge to $x = \alpha$ if

$$|g_1'(\alpha)| < 1$$
, that is, $\left| \frac{b}{\alpha^2} \right| < 1$,

it gives

$$|\alpha|^2 > |b| = |\alpha||\beta|, \quad \textit{gives} \quad |\alpha| > |\beta|.$$

The second iterative scheme will converge to $x = \alpha$ if

$$|g_2'(\alpha)| < 1$$
, that is, $\left| \frac{b}{(\alpha+a)^2} \right| < 1$,

which implies that

$$(\alpha + a)^2 > |b|$$
 or $|\beta|^2 > |b|$, (since $\alpha + a = -\beta$).

Hence

$$|\beta|^2 > |b| = |\alpha||\beta|, \qquad gives \qquad |\beta| > |\alpha|,$$

or, $|\alpha| < |\beta|$.

Example 2.8 Find an interval [a,b] on which fixed-point problem $x = \frac{2 - e^x + x^2}{3}$ will converges. Estimate the number of iterations n within accuracy 10^{-5} .

Solution. Since $x = \frac{2 - e^x + x^2}{3}$ can be written as

$$f(x) = e^x - x^2 + 3x - 2 = 0,$$

and we observe that $f(0)f(1)=(-1)(e^1)<0$, then the solution we seek is in the interval [0,1]. For $g(x)=\frac{2-e^x+x^2}{3}$, we have $g'(x)=\frac{2x-e^x}{3}<1$, for all x in the given interval [0,1]. Also, g is decreasing function of x and g(0)=0.3333 and $g(1)=\frac{3-e}{3}=0.0939$ both lie in the interval [0,1]. Thus $g(x)\in[0,1]$, for all $x\in[0,1]$, so from Fixed-Point Theorem 2.2 the g(x) has a unique fixed-point in [0,1]. Taking $x_0=0.5$, we have

$$x_1 = g(x_0) = \frac{2 - e^{x_0} + x_0^2}{3} = 0.2004.$$

Also, we have

$$k_1 = |g'(0)| = 0.3333$$
 and $k_2 = |g'(1)| = 0.2394$,

which give $k = \max\{k_1, k_2\} = 0.3333$. Thus the error estimate (2.9) within the accuracy 10^{-5} is

$$|\alpha - x_n| \le 10^{-5}$$
, gives $\frac{(0.3333)^n}{1 - 0.3333}(0.2996) \le 10^{-5}$,

and by solving this inequality, we obtain $n \ge 9.7507$. So we need ten approximations to get the desired accuracy for the given problem.

Example 2.9 Convert the equation $x^2 - 5 = 0$ to the fixed-point problem $x = x + c(x^2 - 5)$ with c a nonzero constant. Find a value of c to ensure rapid convergence of the following scheme to $\alpha = \sqrt{5}$

$$x_{n+1} = x_n + c(x_n^2 - 5), n \ge 0.$$

Solution. Given $x^2 - 5 = 0$, and it can be written as for $c \neq 0$

$$c(x^2 - 5) = 0$$
 or $-x + x + c(x^2 - 5) = 0$.

From this we have

$$x = x + c(x^2 - 5) = g(x),$$

and it gives the iterative scheme

$$x_{n+1} = x_n + c(x_n^2 - 5) = g(x_n), \qquad n \ge 0.$$

For guaranteed convergence of this scheme, we mean that

$$|g'(x)| < 1$$
 or $|1 + 2cx| < 1$ or $-1 < 1 + 2cx < 1$.

Moreover, the convergence will be rapid if

$$g'(\alpha) = 1 + 2\alpha c = 0.$$

Since $\alpha = \sqrt{5}$, therefore

$$1 + 2\sqrt{5}c = 0$$
.

Thus, we have $c = -\frac{1}{2\sqrt{5}}$, the required value of c.

2.3 Fixed-Point Method

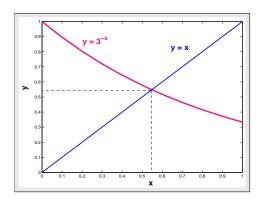


Figure 2.7: Graphical Solution of $x = 3^{-x}$.

Example 2.10 Show that the function $g(x) = 3^{-x}$ on the interval [0,1] has at least one fixed-point but it is not unique.

Solution. Given $x = g(x) = 3^{-x}$, and it can be written as

$$x - 3^{-x} = f(x) = 0.$$

So f(0)(1) = (-1)(2/3) < 0, so f(x) has a root in the interval [0,1], see Figure 2.7. Note that g is decreasing function of x and g(0) = 1 and g(1) = 0.3333 both lie in the interval [0,1]. Thus $g(x) \in [0,1]$, for all $x \in [0,1]$, so from Fixed-Point Theorem 2.2 the function g(x) has at least one fixed-point in [0,1]. Since the derivative of the function g(x) is

$$g'(x) = -3^{-x} \ln 3$$
,

which is less than zero on [0,1], therefore, the function g is decreasing on [0,1]. But $g'(0) = -\ln 3 = -1.0986$, so

$$|g'(x)| > 1$$
 on $(0,1)$.

Thus from Fixed-Point Theorem 2.2 the function g(x) has no unique fixed-point in [0,1].

Example 2.11 Show that the function $g(x) = \sqrt{2x-1}$ on the interval [0,1] that satisfies none of the hypothesis of Theorem 2.2 but still has a unique fixed-point on [0,1].

Solution. Since $x = g(x) = \sqrt{2x - 1}$, it gives

$$x^{2} - 2x + 1 = (x - 1)^{2} = f(x) = 0.$$

Then $x = \alpha = 1 \in [0,1]$ is the root of the nonlinear equation f(x) = 0 and the fixed-point of the function g(x) as g(1) = 1. But notice that the function g(x) is not continuous on the interval [0,1] and the derivative of the function g(x)

$$g'(x) = \frac{1}{\sqrt{2x - 1}},$$

does not exist on the interval (0,1). So all the conditions of Fixed-Point Theorem 2.2 fail.

Example 2.12 Show that the fixed point form of the equation $x = N^{1/3}$ can be written as $x = Nx^{-2}$ and the associated iterative scheme

$$x_{n+1} = Nx_n^{-2}, \qquad n \ge 0,$$

will not successful (diverge) in finding the approximation of cubic root of the positive number N.

Solution. Given $x = N^{1/3}$ and it can be written as

$$x^3 - N = 0$$
 or $x = \frac{N}{x^2} = Nx^{-2}$.

It gives the iterative scheme

$$x_{n+1} = Nx_n^{-2} = g(x_n), \qquad n \ge 0.$$

From this, we have

$$g(x) = Nx^{-2}$$
 and $g'(x) = -2Nx^{-3}$.

Since $\alpha = x = N^{1/3}$, therefore

$$g'(\alpha) = -2N\alpha^{-3}$$
 and $g'(N^{1/3}) = -2N(N^{1/3})^{-3} = -2NN^{-1} = -2$.

Thus

$$|g'(N^{1/3})| = |-2| = 2 > 1,$$

which shows the divergence.

Example 2.13 Which of the following sequences will converge faster to $\sqrt{5}$

(a)
$$x_{n+1} = x_n + 1 - \frac{x_n^2}{5}$$
, (b) $x_{n+1} = \frac{1}{3} \left[3x_n + 1 - \frac{x_n^2}{5} \right]$.

Solution. It can be easily verify by using the Note 2.2. From the first sequence, we have

$$g_1(x) = x + 1 - \frac{x^2}{5}$$
 and $g'_1(x) = 1 - \frac{2x}{5}$,

which implies that

$$|g_1'(\sqrt{5})| = \left|1 - \frac{2\sqrt{5}}{5}\right| = 0.1056 < 1.$$

Similarly, from the second sequence, we have

$$g_2(x) = \frac{1}{3} \left[3x + 1 - \frac{x^2}{5} \right]$$
 and $g_2'(x) = \frac{1}{3} \left[3 - \frac{2x}{5} \right]$, gives, $|g_2'(\sqrt{5})| = 0.701186 < 1$.

We note that both sequences are converging to $\sqrt{5}$ but the sequence (a) will converges faster than the sequence (b) because the value of $|g_1'(\sqrt{5})|$ is closer to zero than by $|g_2'(\sqrt{5})|$.

26 2.3 Fixed-Point Method

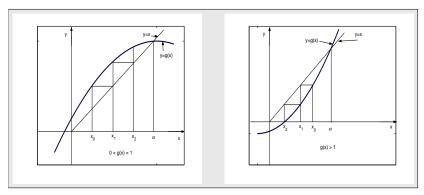


Figure 2.3.1. Convergent function iterations.

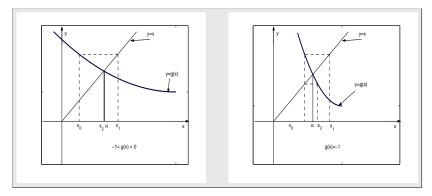


Figure 2.3.2. Divergent function iterations.

The graphical interpretation of the fixed-point method is illustrated in Figure 2.3.1(a) and in Figure 2.3.1(b). The fixed point α is the abscissa of the intersection of the graph of the function g(x) with the line y = x. The ordinate of the function g(x) at x_0 is the value of x_1 . To turn this ordinate into an abscissa, reflect it in the line y = x. We may repeat this process to get x_2, x_3 and so on. It is seen that the iterates see Figure 2.3.1(a) moving (zigzag) towards the fixed-point, while in Figure 2.3.2(b) they going away: the iterations in Figure 2.3.2(a) converge if you start near enough to the fixed point, whereas the other diverge no matter how close you start. The fixed-point in Figure 2.3.2(a) is said to be attractive, and the one in Figure 2.3.2(b) is said to be repulsive.

Example 2.14 Assuming that the iterative schemes generated by

$$x_{n+1} = g_1(x_n) = \frac{1}{3}(x_n^2 + 2);$$
 $n = 0, 1, ...,$
 $x_{n+1} = g_2(x_n) = 3 - \frac{2}{x_n};$ $n = 0, 1, ...,$

both converge, show that they do so to different roots of the same equation $x^2 - 3x + 2 = 0$. Find the third approximation of to the root 1 using suitable iterative scheme, starting with $x_0 = 0.5$.

Solution. If x_n converges to α , then the first iterative scheme gives

$$\alpha = g_1(\alpha) = \frac{1}{3}(\alpha^2 + 2), \text{ that is, } \alpha^2 - 3\alpha + 2 = 0,$$

which has the roots $\alpha = 1$ and $\alpha = 2$. Since the first derivative of the function $g_1(x)$ is

$$g_1'(x) = \frac{2}{3}x,$$

so

$$|g_1'(1)| < 1$$
 and $|g_1'(2)| > 1$.

Hence convergence to $\alpha = 1$. Now if x_n converges to α , then other scheme gives

$$\alpha = g_2(\alpha) = 3 - \frac{2}{\alpha}$$
, that is, $\alpha^2 - 3\alpha + 2 = 0$.

Since the first derivative of the function $g_2(x)$ is

$$g_2'(x) = \frac{2}{x^2},$$

so

$$|g_2'(2)| < 1$$
 and $|g_2'(1)| > 1$.

Hence convergence to $\alpha = 2$. Thus first iterative scheme is the most suitable one. Using this suitable iterative scheme with $x_0 = 0.5$, we get

$$x_1 = g_1(x_0) = 0.7500, \quad x_2 = g_1(x_1) = 0.8542, \quad x_3 = g_1(x_2) = 0.9099,$$

and

$$|\alpha - x_3| = |1 - 0.9099| = 0.0901,$$

the required absolute error.

Example 2.15 Let α_1 and α_2 are two fixed points of the function

$$g(x) = 0.5x^2 - 1.5x + 2.$$

- (a) Find the values of both fixed points.
- (b) For which fixed point in part (a), the fixed-point method will converge.

Solution. Since f(x) = g(x) - x = 0, so we have

$$0.5x^2 - 2.5x + 2 = 0.$$

Solving this quadratic equation, we get

$$x_1 = 1$$
 and $x_2 = 4$.

Then

$$g(1) = 0.5(1^2) - 1.5(1) + 2 = 1$$
 and $g(4) = 0.5(4^2) - 1.5(4) + 2 = 4$,

showing that $\alpha_1 = 1$ and $\alpha_2 = 4$ are the two fixed points of the given function. Since the first derivative of the given function g(x) is

$$g'(x) = x - 1.5,$$

and its absolute value at the both fixed points are

$$|g'(1)| = |1 - 1.5| = 0.5 < 1$$
 and $|g'(4)| = |4 - 1.5| = 2.5 > 1$.

Therefore, from Fixed-Point Theorem 2.2, we conclude that the fixed-point method will converge for the fixed-point $\alpha_1 = 1$ only.

28 2.3 Fixed-Point Method

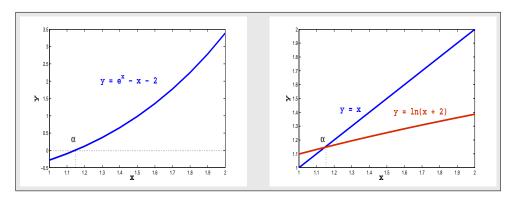


Figure 2.8: Graphical Solution of $e^x = x + 2$ Graphical solution of $x = \ln(x + 2)$.

Example 2.16 One of the possible rearrangement of the nonlinear equation $e^x = x + 2$, which has root in [1,2] is

$$x_{n+1} = g(x_n) = \ln(x_n + 2);$$
 $n = 0, 1, \dots$

- (a) Show that g(x) has a unique fixed-point in [1, 2].
- (b) Use fixed-point iteration formula (2.7) to compute approximation x_3 , using $x_0 = 1.5$.
- (c) Compute an error estimate $|\alpha x_3|$ for your approximation.
- (d) Determine the number of iterations needed to achieve an approximation with accuracy 10^{-2} to the solution of $g(x) = \ln(x+2)$ lying in the interval [1,2] by using the fixed-point iteration method.

Solution. Since, we observe that f(1)f(2) < 0, then the solution we seek is in the interval [1,2].

- (a) For $g(x) = \ln(x+2)$, we have g'(x) = 1/(x+2) < 1, for all x in the given interval [1,2]. Also, g is increasing function of x, and $g(1) = \ln(3) = 1.0986123$ and $g(2) = \ln(4) = 1.3862944$ both lie in the interval [1,2]. Thus $g(x) \in [1,2]$, for all $x \in [1,2]$, so from fixed-point theorem the g(x) has a unique fixed-point, see Figure 2.8.
- (b) using the given initial approximation $x_0 = 1.5$, we have the other approximations as

$$x_1 = g(x_0) = 1.252763, \quad x_2 = g(x_1) = 1.179505, \quad x_3 = g(x_2) = 1.156725.$$

(c) Since a = 1 and b = 2, then the value of k can be found as follows

$$k_1 = |g'(1)| = |1/3| = 0.333$$
 and $k_2 = |g'(2)| = |1/4| = 0.25$,

which give $k = \max\{k_1, k_2\} = 0.333$. Thus using the error formula (2.9), we have

$$|\alpha - x_3| \le \frac{(0.333)^3}{1 - 0.333} |1.252763 - 1.5| = 0.013687.$$

(d) From the error bound formula (2.9), we have

$$\frac{k^n}{1-k}|x_1-x_0| \le 10^{-2}.$$

By using above parts (b) and (c), we have

$$\frac{(0.333)^n}{1 - 0.333} |1.252763 - 1.5| \le 10^{-2}.$$

Solving this inequality, we obtain

$$n \ln(0.333) \le \ln(0.02698), \quad gives, \quad n \ge 3.28539.$$

So we need four approximations to get the desired accuracy for the given problem.

MATLAB command for the above given rearrangement x = g(x) of $f(x) = x^3 - 2x - 1$ by using the initial approximation $x_0 = 1.5$, can be written as follows:

```
function y = fn(x)

y = log(x + 2);

>> x0 = 1.5; tol = 0.01; sol = fixpt('fn', x0, tol);
```

Program 2.2

MATLAB m-file for the Fixed-Point Method function sol=fixpt(fn,x0,tol) old= x0+1; while abs(x0-old) > tol; old=x0; x0 = feval(fn, old); end; sol=x0;

Procedure 2.2 (Fixed-Point Method)

- 1. Choose an initial approximation x_0 such that $x_0 \in [a, b]$.
- 2. Choose a convergence parameter $\epsilon > 0$.
- 3. Compute new approximation x_{new} by using the iterative formula (2.7).
- 4. Check, if $|x_{new} x_0| < \epsilon$ then x_{new} is the desire approximate root; otherwise set $x_0 = x_{new}$ and go to step 3.

2.4 Newton's Method

This is one of the most popular and powerful iterative method for finding roots of the nonlinear equation (2.1). It is also known as the method of tangents because after estimated the actual root, the zero of the tangent to the function at that point is determined. It always converges if the initial approximation is sufficiently close to the exact solution. This method is distinguished from

30 2.4 Newton's Method

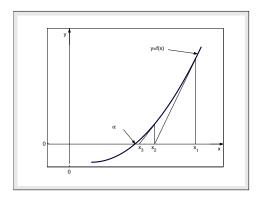


Figure 2.9: Graphical Solution of Newton's Method.

the methods of previous sections by the fact that it requires the evaluation of both the function f(x) and the derivative of the function f'(x), at arbitrary point x. The Newton's method consists geometrically of expanding the tangent line at a current point x_i until it crosses zero, then setting the next guess x_{i+1} to the abscissa of that zero crossing, see Figure 2.9. This method is also called the Newton-Raphson method.

There are many description of the Newton's method. We shall derive the method from the familiar Taylor's series expansion of a function in the neighborhood of a point.

Let $f \in C^2[a, b]$ and let x_n be the *nth* approximation to the root α such that $f'(x_n) \neq 0$ and $|\alpha - x_n|$ is small. Consider the first Taylor polynomial for f(x) expanded about x_n , so we have

$$f(x) = f(x_n) + (x - x_n)f'(x_n) + \frac{(x - x_n)^2}{2}f''(\eta(x)), \tag{2.12}$$

where $\eta(x)$ lies between x and x_n . Since $f(\alpha) = 0$, then (2.12), with $x = \alpha$, gives

$$f(\alpha) = 0 = f(x_n) + (\alpha - x_n)f'(x_n) + \frac{(\alpha - x_n)^2}{2}f''(\eta(\alpha)).$$

Since $|\alpha - x_n|$ is small, then we neglect the term involving $(\alpha - x_n)^2$ and so

$$0 \approx f(x_n) + (\alpha - x_n)f'(x_n).$$

Solving for α , we get

$$\alpha \approx x_n - \frac{f(x_n)}{f'(x_n)},\tag{2.13}$$

which should be better approximation to α than is x_n . We call this approximation as x_{n+1} , then we get

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}, \quad f'(x_n) \neq 0, \quad \text{for all} \quad n \ge 0.$$
 (2.14)

The iterative method (2.14) is called the Newton's method. Usually the Newton's method converges well and quickly but its convergence cannot, however guaranteed and it may sometime converge to a different root from the one expected. In particular, there may be difficulties if initial approximation is not sufficiently close to the actual root. The most serious problem of the Newton's method is that

some functions are difficult to differentiate analytically, and some functions cannot be differentiated analytically at all. The Newton's method is not restricted to one-dimension only. The method readily generalizes to multiple dimensions. It should be noted that this method is suitable for finding real as well as imaginary roots of the polynomials.

Example 2.17 Use the Newton's method to find the root of $x^3 = 2x + 1$ that is located in the interval [1.5, 2.0] accurate to 10^{-2} , take an initial approximation $x_0 = 1.5$.

Solution. Given $f(x) = x^3 - 2x - 1$ and so $f'(x) = 3x^2 - 2$. Now evaluating f(x) and f'(x) at the give approximation $x_0 = 1.5$, gives

$$x_0 = 1.5,$$
 $f(1.5) = -0.625,$ $f'(1.5) = 4.750.$

Using the Newton's iterative formula (2.14), we get

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)} = 1.5 - \frac{(-0.625)}{4.75} = 1.631579.$$

Now evaluating f(x) and f'(x) at the new approximation x_1 , gives

$$x_1 = 1.631579,$$
 $f(1.631579) = 0.0801869,$ $f'(1.631579) = 5.9861501.$

Using the iterative formula (2.14) again to get other new approximation. The successive iterates were shown in the Table 2.4. Just after the third iterations the required root is approximated to

Table 2.4: Solution of $x^3 = 2x + 1$ by Newton's method

n	x_n	$f(x_n)$	$f'(x_n)$	$Error = x - x_n$
00	1.500000	-0.625000	4.750000	0.1180339
01	1.631579	0.0801869	5.9861501	-0.0135451
02	1.618184	0.000878	5.855558	-0.0001501
03	1.618034	0.00000007	5.854102	-0.0000001

be $x_3 = 1.618034$ and the functional value is reduced to 7.0×10^{-8} . Since the exact solution is 1.6180339, so the actual error is 1×10^{-7} . We see that the convergence is quite faster than the methods considered previously.

To get the above results using MATLAB command, firstly the function $x^3 - 2x - 1$ and its derivative $3x^2 - 2$ were saved in m-files called fn.m and dfn.m, respectively written as follows:

function
$$y = fn(x)$$
 function $dy = dfn(x)$
 $y = x.^3 - 2 * x - 1;$ $dy = 3 * x.^2 - 2;$

after which we do the following:

$$>> x0 = 1.5; tol = 0.01; sol = newton('fn', 'dfn', x0, tol);$$

32 2.4 Newton's Method

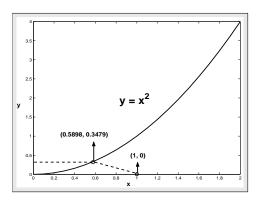


Figure 2.10: Graphical Solution of the function $y = x^2$.

Example 2.18 Use Newton's method to approximate, to within 10^{-4} , the value of x that produces the point on the graph of $y = x^2$ that is closest to (1,0), using initial approximation $x_0 = 1$.

Solution. The distance between an arbitrary point (x, x^2) on the graph of $y = x^2$ and the point (1,0) is

$$d(x) = \sqrt{(x-1)^2 + (x^2 - 0)^2} = \sqrt{x^4 + x^2 - 2x + 1}.$$

Because a derivative is needed to find the critical point of d, it is easier to work with the square of this function

$$F(x) = [d(x)]^2 = x^4 + x^2 - 2x + 1,$$

whose minimum will occur at the same value of x as the minimum of d(x). To minimize F(x), we need x so that

$$F'(x) = 4x^3 + 2x - 2 = 0$$
, gives, $f(x) = 4x^3 + 2x - 2$, $f'(x) = 12x^2 + 2$.

Applying Newton's iterative formula (2.14) to find the approximation of this equation, we have

$$x_{n+1} = x_n - \frac{4x_n^3 + 2x_n - 2}{12x_n^2 + 2}.$$

Finding the approximation to the x within 10^{-4} using the initial approximation $x_0 = 1$, we get

$$x_1 = x_0 - \frac{4x_0^3 + 2x_0 - 2}{12x_0^2 + 2} = 0.7143.$$

Continue in the same manner, we get, $x_2 = 0.6052$, $x_3 = 0.5900$, $x_4 = 0.5898$, $x_5 = 0.5898$. So the point on the graph that is closest to (1,0) has the approximate coordinates (0.5898, 0.3479).

Program 2.3

MATLAB m-file for the Newton's Method function sol=newton(fn,dfn,x0,tol)

old = x0+1; while abs (x0 - old) > tol; old = x0;

x0 = old - feval(fn, old)/feval(dfn, old); end; sol=x0;

Example 2.19 If the difference of two numbers x and y is 6 and the square root of their product is 4, then use Newton's method to approximate, to within 10^{-4} , the largest value of the number x and the corresponding number y using initial approximation $x_0 = 7.5$.

Solution. Given

$$x - y = 6$$
 and $\sqrt{xy} = 4$.

Solving the above equations for x, we have

$$x(x-6) = 16$$
 or $x^2 - 6x - 16 = f(x) = 0$.

Applying Newton's iterative formula (2.14) to find the approximation of this equation, we have

$$x_{n+1} = x_n - \frac{x_n^2 - 6x_n - 16}{2x_n - 6}.$$

Finding the approximation to within 10^{-4} using the initial approximation $x_0 = 7.5$, we get

$$x_1 = x_0 - \frac{x_0^2 - 6x_0 - 16}{2x_0 - 6} = 8.0278,$$

and continue in the same manner, we get the approximations within accuracy 10^{-4} as follows

$$x_2 = 8.0001, \quad x_3 = 8.0000, \quad x_4 = 8.0000.$$

Thus the largest value of number x is 8 and its corresponding y value is 2.

Example 2.20 The graphs of $y = 2 \sin x$ and $y = \ln(x) + k$ touch each other in the neighborhood of point x = 8. Find the value of the constant k and the coordinates of point of contact, use $x_0 = 8$.

Solution. Since we know that the graphs will touch each other if the values of derivatives at their point of contact is same. So for

$$y = 2\sin x$$
, gives, $y' = 2\cos x$,

and

$$y = \ln(x) + k$$
, gives, $y' = \frac{1}{x}$.

Thus

$$2\cos x = \frac{1}{x}$$
, gives, $x\cos x - 0.5 = 0$,

and from this we have the function and its derivative as follows

$$f(x) = x \cos x - 0.5$$
 and $f'(x) = \cos x - x \sin x$.

Using Newton's iterative formula (2.14), we get

$$x_{n+1} = x_n - \frac{x_n \cos x_n - 0.5}{\cos x_n - x_n \sin x_n},$$

and for finding the approximations, starting $x_0 = 8$, we obtain, $x_1 = 7.7936$ and $x_2 = 7.7897$. Taking x = 7.79, we have $y = 2\sin 7.79 = 1.996$. Therefore, the point of contact is (7.79, 1.996). To find the value of k, we solve the equation, $1.996 = \ln(7.79) + k$, and it gives, k = -0.0568, the required value of k.

34 Newton's Method

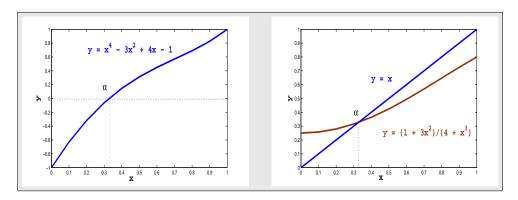


Figure 2.11: Graphical Solution of $x^4 - 3x^2 + 4x = 1$ and $x = (1 + 3x^2)/(4 + x^3)$.

Example 2.21 Successive approximations x_n to the desired root are generated by the scheme

$$x_{n+1} = \frac{1+3x_n^2}{4+x_n^3}, \qquad n \ge 0.$$

Find $f(x_n)$ and $f'(x_n)$ and then use the Newton's method to find the approximation of the root accurate to 10^{-2} , starting with $x_0 = 0.5$.

Solution. Given

$$x = \frac{1+3x^2}{4+x^3} = g(x),$$

and

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

Since

$$f(x) = x - g(x) = 0,$$

therefore, we have

$$f(x_n) = x_n^4 - 3x_n^2 + 4x_n - 1$$
 and $f'(x_n) = 4x_n^3 - 6x_n + 4$.

Using these functions values in the Newton's iterative formula (2.14), we have (see Figure 2.11),

$$x_{n+1} = x_n - \frac{x_n^4 - 3x_n^2 + 4x_n - 1}{4x_n^3 - 6x_n + 4}.$$

Finding the first approximation of the root using the initial approximation $x_0 = 0.5$, we get

$$x_1 = x_0 - \frac{x_0^4 - 3x_0^2 + 4x_0 - 1}{4x_0^3 - 6x_0 + 4} = 0.5 - \frac{0.3125}{1.5} = 0.2917.$$

Similarly, the other approximations can be obtained as

$$x_2 = 0.2917 - \frac{(-0.0813)}{2.3491} = 0.3263;$$
 $x_3 = 0.3263 - \frac{(-0.0029)}{2.1812} = 0.3276.$

Notice that $|x_3 - x_2| = 0.0013$.

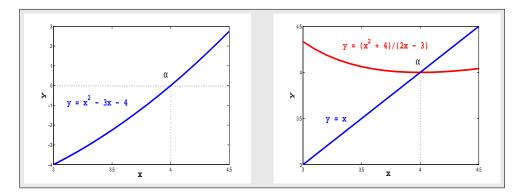


Figure 2.12: Graphical Solution of $x^2 - 3x = 4$ and $x = (x^2 + 4)/(2x - 3)$.

Example 2.22 Develop the iterative formula

$$x_{n+1} = \frac{x_n^2 - b}{2x_n - a}, \qquad n \ge 0,$$

for the approximate roots of the quadratic equation $x^2 - ax + b = 0$ using the Newton's method. Then use the formula to find the third approximation of the positive root of the equation $x^2 - 3x = 4$, starting with $x_0 = 3.5$.

Solution. Given

$$f(x) = x^2 - ax + b,$$

therefore, we have (see Figure 2.12),

$$f(x_n) = x_n^2 - ax_n + b$$
 and $f'(x_n) = 2x_n - a$.

Using these functions values in the Newton's iterative formula (2.14), we have

$$x_{n+1} = x_n - \frac{x_n^2 - ax_n + b}{2x_n - a} = \frac{x_n^2 - b}{2x_n - a}, \quad n \ge 0.$$

Finding the first three approximations of the positive root of $x^2 - 3x = 4$ using the initial approximation $x_0 = 3.5$ and a = 3, b = -4, we use the above formula by taking n = 0, 1, 2 as follows

$$x_1 = \frac{x_0^2 - b}{2x_0 - a} = 4.0625, \quad x_2 = \frac{x_1^2 - b}{2x_1 - a} = 4.0008, \quad x_3 = \frac{x_2^2 - b}{2x_2 - a} = 4.0000,$$

are the possible three approximations. Note that the positive root of $x^2 - 3x - 4 = 0$ is 4, so we have

$$|4 - x_3| = |4 - 4| = 0.0000,$$

the possible absolute error.

The Newton's method is widely used in computers as a basis for the square root and the reciprocal evaluation.

36 2.4 Newton's Method

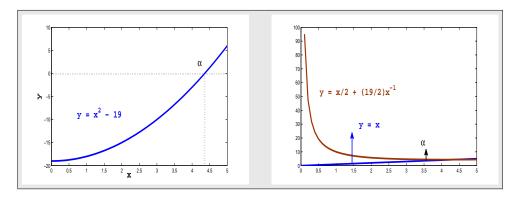


Figure 2.13: Graphical Solution of $x^2 = 19$ and $x = x/2 + (19/2)x^{-1}$.

Example 2.23 Develop an iterative procedure for evaluating pth any root of a positive number N by using Newton's method. Use the developed formula to find second approximation to the square root of 19, taking an initial approximation $x_0 = 5$. Compute absolute error.

Solution. We shall compute $x = N^{1/p}$ by finding a positive root for the nonlinear equation

$$x^p - N = 0.$$

where p is any positive integer and N > 0 is the number whose root is to be found. Therefore, if f(x) = 0, then $x = N^{1/p}$ is the exact root. Let

$$f(x) = x^p - N$$
 and $f'(x) = px^{p-1}$.

Hence, assuming an initial estimate to the root, say, $x = x_0$ and by using iterative formula (2.14), we get

$$x_1 = x_0 - \frac{(x_0^p - N)}{px_0^{p-1}} = x_0 - \frac{x_0^p}{px_0^{p-1}} + \frac{N}{px_0^{p-1}} = (1 - \frac{1}{p})x_0 + \frac{N}{p}x_0^{1-p}.$$

In general, we have

$$x_{n+1} = \left(1 - \frac{1}{p}\right)x_n + \frac{N}{p}x_n^{1-p},\tag{2.15}$$

where p = 2, 3, ..., and n = 0, 1, ...

Since we want the approximations of the square root of number 19, so we take N = 19 and p = 2. Given the initial approximation $x_0 = 5$, then by using the iterative formula (2.15), we get

$$x_1 = 4.4$$
 and $x_2 = 4.3590909$.

After just two iterations the estimated value compares rather favorably with the exact value of $\sqrt{19} \approx 4.3588989$, (see Figure 2.13). Thus the absolute error is

$$|E| = |\sqrt{19} - x_2| = |5.3588989 - 4.3590909| = 0.000192.$$

We can calculate higher roots of a number by using the general iterative formula (2.15).

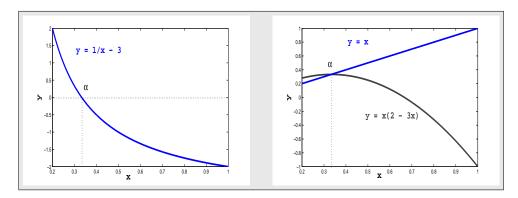


Figure 2.14: Graphical Solution of 1/x = 3 and x = x(2 - 3x).

Example 2.24 Develop an iterative procedure for evaluating the reciprocal of a positive number N by using Newton's method. Use the developed formula to find third approximation to the reciprocal of 3, taking an initial approximation $x_0 = 0.4$. Compute absolute error.

Solution. Consider x = 1/N. This problem can be easily solved by noting that we seek to find a root to the nonlinear equation

$$1/x - N = 0,$$

where N > 0 is the number whose reciprocal is to be found. Therefore, if f(x) = 0, then x = 1/N is the exact root. Let

$$f(x) = 1/x - N$$
 and $f'(x) = -1/x^2$.

Hence, assuming an initial estimate to the root, say, $x = x_0$ and by using iterative formula (2.14), we get

$$x_1 = x_0 - \frac{(1/x_0 - N)}{(-1/x_0^2)} = x_0 + (1/x_0 - N)x_0^2 = x_0 + x_0 - Nx_0^2 = x_0(2 - Nx_0).$$

In general, we have

$$x_{n+1} = x_n(2 - Nx_n), \qquad n = 0, 1, \dots,$$
 (2.16)

We have to find the approximation of the reciprocal of number N=3. Given the initial gauss of say $x_0=0.4$, then by using the iterative formula (2.16), we get

$$x_1 = 0.3200, \qquad x_2 = 0.3328, \qquad x_3 = 0.3333.$$

After just three iterations the estimated value compares rather favorably with the exact value of $1/3 \approx 0.3333$, (see Figure 2.14). Thus the absolute error is

$$|E| = \left| \frac{1}{3} - x_3 \right| = |0.3333 - 0.3333| = 0.0000.$$

We can calculate the other reciprocal of the number in the same way by using the general iterative formula (2.16).

38 2.5 Secant Method

Lemma 2.1 Assume that $f \in C^2[a,b]$ and there exists a number $\alpha \in [a,b]$, where $f(\alpha) = 0$. If $f'(\alpha) \neq 0$, then there exists a number $\delta > 0$ such that the sequence $\{x_n\}_{n=0}^{\infty}$ defined by the iteration

$$x_{n+1} = g(x_n) = x_n - \frac{f(x_n)}{f'(x_n)}, \quad \text{for} \quad n = 0, 1, \dots,$$
 (2.17)

will converges to α for any initial approximation $x_0 \in [\alpha - \delta, \alpha - \delta]$.

The Newton's method uses the iteration function

$$g(x) = x - \frac{f(x)}{f'(x)},$$
 (2.18)

is called the Newton's iteration function. Since $f(\alpha) = 0$, it is easy to see that $g(\alpha) = \alpha$. Thus the Newton's iteration for finding the root of the equation f(x) = 0 is accomplished by finding a fixed-point of the equation g(x) = x.

Procedure 2.3 (Newton's Method)

- 1. Find the initial approximation x_0 for the root by sketching the graph of the function.
- 2. Evaluate function f(x) and the derivative f'(x) at initial approximation. Check: if $f(x_0) = 0$ then x_0 is the desire approximation to a root. But if $f'(x_0) = 0$, then go back to step 1 to choose new approximation.
- 3. Establish Tolerance $(\epsilon > 0)$ value for the function.
- 4. Compute new approximation for the root by using the iterative formula (2.14).
- 5. Check Tolerance. If $|f(x_n)| \le \epsilon$, for $n \ge 0$, then end; otherwise, go back to step 4, and repeat the process.

2.5 Secant Method

Since we known the main obstacle to using the Newton's method is that it may be difficult or impossible to differentiate the function f(x). The calculation of $f'(x_n)$ may be avoided by approximating the slope of the tangent at $x = x_n$ by that of the chord joining the two points $(x_{n-1}, f(x_{n-1}))$ and $(x_n, f(x_n))$, see Figure 2.15.

The slope of the chord (or secant) is

$$f'(x_n) \approx \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}.$$
 (2.19)

Then by using this approximation of the derivative of the function in the Newton's iterative formula (2.14), we get

$$x_{n+1} = x_n - \frac{(x_n - x_{n-1})f(x_n)}{f(x_n) - f(x_{n-1})} = \frac{x_{n-1}f(x_n) - x_nf(x_{n-1})}{f(x_n) - f(x_{n-1})}, \qquad n \ge 1.$$
 (2.20)

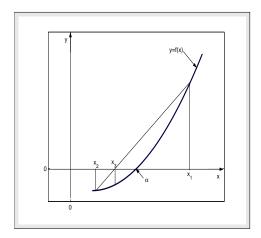


Figure 2.15: Graphical Solution of Secant Method.

Note that when $f(x_n) = f(x_{n-1})$, the calculation of x_{n+1} fails. This is because the chord is horizontal. The iterative formula (2.20) known as the *secant method*. It needs two initial approximations to start. This method is very similar to false position method described in Section 2.3. However, for the secant method it is not necessary for the interval to contain a root and no account is taken of signs of the numbers $f(x_n)$.

This method suffers from the same disadvantages as the Newton's method, that is, convergence to a particular root cannot guaranteed but nevertheless it is a powerful general purpose method. The order of convergence of the secant method is $(1+\sqrt{5})/2\approx 1.618$, so its ultimate convergence is not quite as fast as the Newton's method (order of convergence is quadratically) but the order of convergence of this method is somewhat better than the bisection method, the false position method, and the fixed-point method (all these methods have linear convergence). This is sometimes called superlinear. We will discuss the order of convergence of all these methods in some details later in the chapter.

Example 2.25 Show that the iterative procedure for evaluating the reciprocal of a number N by using the secant method is:

$$x_{n+1} = x_n + (1 - Nx_n)x_{n-1}, \qquad n \ge 1.$$
 (2.21)

Solution. Let N be a positive number and x = 1/N. If f(x) = 0, then $x = \alpha = 1/N$ is the exact zero of the function

$$f(x) = 1/x - N.$$

Since the secant formula is

$$x_{n+1} = x_n - \frac{(x_n - x_{n-1})f(x_n)}{f(x_n) - f(x_{n-1})}, \quad n \ge 1.$$

Hence, assuming the initial estimates to the root, say, $x = x_0, x = x_1$ and by using the secant iterative formula, we have

$$x_2 = x_1 - \frac{(x_1 - x_0)(1/x_1 - N)}{(1/x_1 - N) - (1/x_0 - N)} = x_1 - \frac{(x_1 - x_0)(1/x_1 - N)}{-(x_1 - x_0)/x_1 x_0}.$$

40 2.5 Secant Method

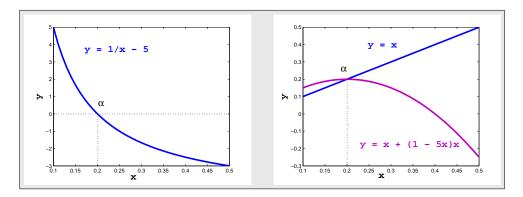


Figure 2.16: Graphical Solution of 1/x = 5 and x = x + (1 - 5x)x.

It gives

$$x_2 = x_1 + (1/x_1 - N)x_1x_0 = x_1 + x_0 - Nx_1x_0 = x_1 + (1 - Nx_1)x_0.$$

In general, this becomes

$$x_{n+1} = x_n + (1 - Nx_n)x_{n-1}, \qquad n = 1, 2, \dots$$

For example, suppose we want the reciprocal of number N=5. Assuming the initial approximations of say $x_0=0$ and $x_1=0.1$, then by using the above iterative formula, we get the first three approximations as follows:

$$x_2 = 0.1, \qquad x_3 = 0.15, \qquad x_4 = 0.175,$$

The estimated value compares rather favorably with exact value of 1/5, (see Figure 2.16).

Example 2.26 Show that the secant method for finding approximation of the square root of a positive number N is

$$x_{n+1} = \frac{x_n x_{n-1} + N}{x_n + x_{n-1}}, \quad n \ge 1.$$
 (2.22)

Carry out the first three approximations for the square root of 9, using $x_0 = 2, x_1 = 2.5$ and also compute absolute error.

Solution. We shall compute $x = N^{1/2}$ by finding a positive root for the nonlinear equation

$$x^2 - N = 0,$$

where N > 0 is the number whose root is to be found. If f(x) = 0, then $x = \alpha = N^{1/2}$ is the exact zero of the function

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

Since the secant formula is

$$x_{n+1} = x_n - \frac{(x_n - x_{n-1})f(x_n)}{f(x_n) - f(x_{n-1})}, \quad n \ge 1.$$

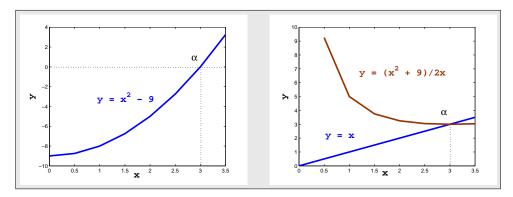


Figure 2.17: Graphical Solution of $x^2 = 9$ and $x = (x^2 + 9)/2x$.

Hence, assuming the initial estimates to the root, say, $x = x_0, x = x_1$, and by using the secant iterative formula, we have

$$x_2 = x_1 - \frac{(x_1 - x_0)(x_1^2 - N)}{(x_1^2 - N) - (x_0^2 - N)} = x_1 - \frac{(x_1 - x_0)(x_1^2 - N)}{(x_1 - x_0)(x_1 + x_0)} = \frac{x_1 x_0 + N}{(x_1 + x_0)}.$$

In general, we have

$$x_{n+1} = \frac{x_n x_{n-1} + N}{x_n + x_{n-1}}, \qquad n = 1, 2, \dots,$$

the secant formula for approximation of the square root of number N. Now using this formula for approximation of the square root of N = 9, taking $x_0 = 2$ and $x_1 = 2.5$, we have

$$x_2 = 3.1111, x_3 = 2.9901, x_4 = 2.9998.$$

Hence

Absolute
$$Error = |9^{1/2} - x_4| = |3 - 2.9998| = 0.0002,$$

is the possible absolute error, (see Figure 2.17).

Example 2.27 Use the secant method to find the approximate root of the following equation within the accuracy 10^{-2} take $x_0 = 1.5$ and $x_1 = 2.0$ as starting values

$$x^3 = 2x + 1.$$

Solution. Since $f(x) = x^3 - 2x - 1$ and

$$x_0 = 1.5,$$
 $f(x_0) = -0.625,$
 $x_1 = 2.0,$ $f(x_1) = 3.0,$

therefore, we see that $f(x_0) \neq f(x_1)$. Hence, one can use the iterative formula (2.20), to get new approximation:

$$x_2 = \frac{x_0 f(x_1) - x_1 f(x_0)}{f(x_1) - f(x_0)} = \frac{(1.5)(3.0) - (2.0)(-0.625)}{3.0 - (-0.625)} = 1.586207,$$

and $f(x_2) = -0.18434$. Similar way, we can find the other possible approximation of the root. A summary of the calculations is given in Table 2.5.

42 2.5 Secant Method

	Table 2.0. Selection of a 2a + 1 sy second meeting a								
	n	x_{n-1}	x_n	x_{n+1}	$f(x_{n+1})$				
	01	1.500000	2.000000	1.586207	-0.1814342				
İ	02	2.000000	1.586207	1.609805	-0.0478446				
	03	1.586207	1.609805	1.618257	0.0013040				

To use MATLAB command for the secant method, the function has been used in the m-file as fn.m, then the first few iterations are easily performed by the following sequence of MATLAB commands:

>>
$$x0 = 1.5; x1 = 2; x2 = x1 - (x1 - x0)/(fn(x1) - fn(x0)) * fn(x1)$$

>> $x0 = x1; x1 = x2; x3 = x2 - (x2 - x1)/(fn(x2) - fn(x1)) * fn(x2)$

The last two commands can be repeated to generate the subsequent iterates shown in Table 2.5.

Program 2.4

MATLAB m-file for the Secant Method

function sol=secant(fn,a,b,tol)

$$x0 = a; x1 = b; fa = feval(fn, x0); fb = feval(fn, x1);$$
 while abs(x1-old)> tol $new = x1 - fb * (x1 - x0)/(fb - fa);$ $x0 = x1; fa = fb; x1 == new; fb = feval(fn, new);$ end; sol=new;

We see that the functional values are approaching zero as the number of iterations is increased a and note that this method converges faster than the methods discussed in Sections 2.1 and 2.2 because we got the desired approximation $c_3 = 1.618032$ after 4 iterations with accuracy $\epsilon = 10^{-2}$.

One can use MATLAB single command to get the same above results by the secant method as

$$>> a = 1.5; b = 2.0; tol = 1e - 2; sol = secant('fn', a, b, tol);$$

Example 2.28 The root of the equation $x^2e^{-x/2} = 1$ is in the interval $[x_0, x_1]$. Apply secant method to find approximation of the root correct to 3 decimal places using $x_0 = 1.42$ and $x_1 = 1.43$.

Solution. Since f(1.42) = -0.0086 and f(1.43) = 0.00034, so using secant iterative formula (2.20), we obtain new approximation as

$$x_2 = \frac{x_0 f(x_1) - x_1 f(x_0)}{f(x_1) - f(x_0)} = \frac{(1.42)(0.00034) - (1.43)(-0.0086)}{0.00034 - (-0.0086)} = 1.4296.$$

As $f(x_2) = -0.00011$, we get another new approximation as

$$x_3 = \frac{x_1 f(x_2) - x_2 f(x_1)}{f(x_2) - f(x_1)} = \frac{(1.43)(-0.00011) - (1.4296)(0.00034)}{(-0.00011 - 0.00034)} = 1.4292.$$

Since x_2 and x_3 agree to 3 decimal places, so the required root is 1.429.

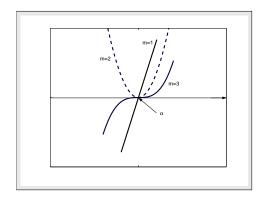


Figure 2.18: Multiple roots of f(x) = 0.

Procedure 2.4 (Secant Method)

- 1. Choose the two initial approximation x_0 and x_1 .
- 2. Check, if $f(x_0) = f(x_1)$, go to step 1 otherwise, continue.
- 3. Establish Tolerance $(\epsilon > 0)$ value for the function.
- 4. Compute new approximation for the root by using the iterative formula (2.20).
- 5. Check tolerance. If $|x_n x_{n-1}| \le \epsilon$, for $n \ge 1$, then end; otherwise, go back to step 4, and repeat the process.

In the preceding discussion the restriction was made that $f'(\alpha) \neq 0$, where α is the solution to f(x) = 0. The rapid convergence of both the Newton's method and the secant method depend because of this restriction. From the definition of the Newton's method, it is clear that difficulties might occur if $f'(x_n)$ goes to zero simultaneously with $f(x_n)$. In particular, the Newton's method and the secant method will generally give problems if $f'(\alpha) = 0$ when $f(\alpha) = 0$. In the following section we investigate this situation and will uncover an interesting fact, namely, how fast the iteration converges. We will see that both the Newton's method and the secant method will continue to converge, but not as rapidly as we expect.

2.6 Multiplicity of a Root

So far we discussed about the function which has simple root. Now we will discuss about the function which has multiple roots. A root is called a *simple root* if it is distinct, otherwise roots that are of the same order of magnitude are called *multiple*.

Definition 2.4 (Order of a Root)

The equation f(x) = 0 has a root α of order m, if there exists a continuous function h(x), and f(x) can be expressed as the product

$$f(x) = (x - \alpha)^m h(x), \quad \text{where} \quad h(\alpha) \neq 0.$$
 (2.23)

44 2.6 Multiplicity of a Root

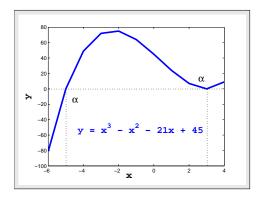


Figure 2.19: Graphical Solution of $x^3 - x^2 - 21x + 45 = 0$.

So h(x) can be used to obtain the remaining roots of f(x) = 0. It is called polynomial deflation. •

A root of order m = 1 is called a *simple root* and if m > 1 it is called *multiple root*. In particular, a root of order m = 2 is sometimes called a *double root*, and so on.

The behavior of the graph of f(x) near a root of multiplicity m (m = 1, 2, 3) is shown in Figure 2.18. It can be seen that when α is a root of odd multiplicity, the graph of f(x) will cross the x-axis at $(\alpha, 0)$; and when α has even multiplicity the graph will be tangent to but will not cross the x-axis at $(\alpha, 0)$. Moreover, the higher the value of m the flatter the graph will be near the point $(\alpha, 0)$. Sometime it is more difficult to deal with the Definition 2.4 concerning about the order of the root. We will use the following Lemma which will illuminate these concepts.

Lemma 2.2 Assume that function f(x) and its derivatives $f'(x), f''(x), \dots, f^{(m)}(x)$ are defined and continuous on an interval about $x = \alpha$. Then f(x) = 0 has a root α of order m if and only if

$$f(\alpha) = f'(\alpha) = f''(\alpha) = \dots = f^{(m-1)}(\alpha) = 0, \qquad f^{(m)}(\alpha) \neq 0.$$
 (2.24)

For example, consider the equation $f(x) = x^3 - x^2 - 21x + 45 = 0$, which has three roots; a simple root at $\alpha = -5$ and a double root at $\alpha = 3$. This can be verified by considering the derivatives of the function as follows

$$f'(x) = 3x^2 - 2x - 21,$$
 $f''(x) = 6x - 2.$

At the value $\alpha = -5$, we have f(5) = 0 and $f'(5) = 64 \neq 0$, so by (2.24), we see that m = 1. Hence $\alpha = -5$ is a simple root of the equation. For the value $\alpha = 3$, we have

$$f(3) = 0,$$
 $f'(3) = 0,$ $f''(3) = 16 \neq 0,$

so that m = 2 by (2.24), hence $\alpha = 3$ is a double root of the equation. Note that this function f(x) has the factorization and can be written in the form of (2.23) as (see Figure 2.19),

$$f(x) = (x-3)^2(x+5).$$

Note that for a simple root α of the nonlinear equation f(x) = 0, means that

$$f(\alpha) = 0$$
 but $f'(\alpha) \neq 0$.

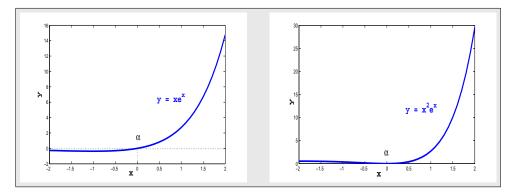


Figure 2.20: Graphical Solution of $xe^x = 0$ and $x^2e^x = 0$.

But for multiple root α of the nonlinear equation, we must have

$$f(\alpha) = 0$$
 and $f'(\alpha) = 0$.

The order of multiplicity of the multiple root can be easily find out by taking the higher derivatives of the function at α unless the higher derivative becomes nonzero at α . Then the order of nonzero higher derivative will be the order of multiplicity of the multiple root.

Example 2.29 Find the multiplicity of the root $\alpha = 1$ of the equation $x \ln x = \ln x$.

Solution. From the given equation, we have

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

Thus the multiplicity of the root $\alpha = 1$ of the given equation is 2.

Usually we don't know in advance that an equation has multiple roots, although we might suspect it from sketching the graph. Many problems which leads to multiple roots, are in fact ill-posed. The methods we discussed so far cannot be guaranteed to converge efficiently for all problems. In particular, when a given function has a multiple root which we require, the methods we have described will either not converge at all or converge more slowly. For example, the Newton's method converges very fast to simple root but converges more slowly when used for functions involving multiple roots.

Example 2.30 Consider the following two nonlinear equations

(1)
$$xe^x = 0$$
 (2) $x^2e^x = 0$.

- (a) Find the Newton's method for the solutions of the given equations.
- (b) Explain why one of the sequences converges much faster than the other to the root $\alpha = 0$.

46 2.6 Multiplicity of a Root

Solution. (a) For the first equation, we have

$$f(x) = xe^x$$
 and $f'(x) = (1+x)e^x$.

Then the Newton's method for the solution of the first equation is

$$x_{n+1} = g_1(x_n) = x_n - \frac{f(x_n)}{f'(x_n)} = \frac{x_n^2}{(1+x_n)}, \quad n \ge 0,$$

which is the first sequence. Similarly, we can find the Newton's method for the solution of the second equation as follows:

$$x_{n+1} = g_2(x_n) = x_n - \frac{x_n^2 e^{x_n}}{(2x_n + x_n^2)e^{x_n}} = \frac{x_n + x_n^2}{(2 + x_n)}, \quad n \ge 0,$$

and it is the second sequence.

(b) From the first sequence, we have

$$g_1(x) = \frac{x^2}{(1+x)}$$
 and $g'_1(x) = \frac{x^2+2x}{(1+x)^2}$.

Then

$$|g_1'(\alpha)| = |g_1'(0)| = \left|\frac{0}{1}\right| = 0,$$

which shows that the first sequence converges to zero. Similarly, from the second sequence, we have

$$g_2(x) = \frac{x+x^2}{(2+x)}$$
 and $g'_2(x) = \frac{x^2+4x+2}{(2+x)^2}$.

Thus

$$|g_2'(0)| = \left|\frac{2}{4}\right| = \frac{1}{2} < 1,$$

which shows that the second sequence is also converges to zero. Since the value of $|g_1'(0)|$ is smaller than $|g_2'(0)|$, therefore, the first sequence converges faster than the second one.

Note that in the above Example 2.30, the root $\alpha = 0$ is the simple root for the first equation (see Figure 2.20) because

$$f(0) = 0$$
 but $f'(0) = 1 \neq 0$,

and for the second equation it is a multiple root because

$$f(0) = 0$$
 and $f'(0) = 0$.

Therefore, the Newton's method converges very fast for the first equation and converges very slow for the second equation. However, in some cases simple modifications can be made to the methods to maintain the rate of convergence. Two such modified methods are considered here, called the Newton modified methods.

First Modified Newton's Method

If we wish to determine a root of known multiplicity m for the equation f(x) = 0, then the first Newton's modified method (also called the Schroeder's method) may be used. It has the form

$$x_{n+1} = x_n - m \frac{f(x_n)}{f'(x_n)}, \qquad n = 0, 1, 2, \dots$$
 (2.25)

It is assumed that we have an initial approximation x_0 . The similarity to the Newton's method is obvious and like the Newton's method it converges very fast for the multiple roots. The major disadvantage of this method is that the multiplicity of the root must be known in advance and this is generally not the case in practice.

Example 2.31 Show that nonlinear equation $\frac{1}{e^{(1-x)}} = x$ has a root at x = 1. Use the first modified Newton's method to find its first three approximations using $x_0 = 0$.

Solution. Since $f(x) = 1 - xe^{(1-x)}$. First we show that $\alpha = 1$ is the zero of the given function as

$$f(\alpha) = f(1) = 1 - 1e^0 = 1 - 1 = 0.$$

To check whether it is simple or multiple zero of f(x), we do the following

$$f'(x) = -e^{1-x} + xe^{(1-x)}$$
 and $f'(\alpha) = f'(1) = -1 + 1 = 0$,

which means that $\alpha = 1$ is the multiple zero of the given function. To find its order of multiplicity, we do

$$f''(x) = 2e^{(1-x)} - xe^{(1-x)}$$
 and $f''(\alpha) = f''(1) = 2 - 1 = 1 \neq 0$,

hence $\alpha=1$ is a zero of multiplicity 2 of the given function. Now we have to find the first three approximations to the multiple zero $\alpha=1$ of the given function by using the first modified Newton's method which can be written as

$$x_{n+1} = x_n - m \frac{f(x_n)}{f'(x_n)} = x_n - m \frac{1 - x_n e^{(1 - x_n)}}{(x_n - 1)e^{(1 - x_n)}}, \qquad n \ge 0,$$

where m is the order of multiplicity of the zero of the function.

For n = 0, 1, 2 and m = 2, with initial approximation $x_0 = 0$, we have

$$x_1 = 0.7358, \qquad x_2 = 0.9782, \qquad x_3 = 0.9998,$$

are the required first three approximations to $\alpha = 1$, (see Figure 2.21).

Second Modified Newton's Method

An alternative approach to this problem that does not require any knowledge of the multiplicity of the root is to replace the function f(x) in the equation by q(x), where

$$q(x) = \frac{f(x)}{f'(x)}.$$

48 2.6 Multiplicity of a Root

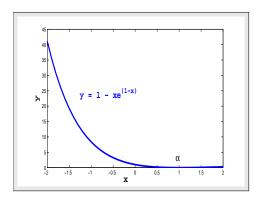


Figure 2.21: Graphical Solution of $1 - xe^{(1-x)} = 0$.

One can show that q(x) has only a simple root at $x = \alpha$. Thus the Newton's method applied to find a root of q(x) will avoid any problems of multiple roots. If

$$f(x) = (x - \alpha)^m h(x),$$

then

$$f'(x) = m(x - \alpha)^{m-1}h(x) + (x - \alpha)^m h'(x).$$

Thus

$$q(x) = \frac{(x - \alpha)h(x)}{[mh(x) + (x - \alpha)h'(x)]}.$$

Obviously we find that q(x) has the root α to multiplicity one. So with this modification, the Newton's method becomes

$$x_{n+1} = x_n - \frac{q(x_n)}{q'(x_n)},$$

which gives

$$x_{n+1} = x_n - \frac{f(x_n)f'(x_n)}{[f'(x_n)]^2 - [f(x_n)][f''(x_n)]}, \quad n = 0, 1, 2, \dots$$
 (2.26)

This iterative formula (2.26) is known as the *second modified Newton's method*. The disadvantage of this method is that we must calculate a further higher derivative. A similar modification can be made to the secant method.

Example 2.32 Use the second modified Newton's method to find the first approximation x_1 to the multiple root of the nonlinear equation $1 - \cos(x) = 0$, using $x_0 = 0.1$.

Solution. Since $f(x) = 1 - \cos x$, we have $f'(x) = \sin x$ and $f''(x) = \cos x$. Now using the second modified Newton's formula (2.26)

$$x_{n+1} = x_n - \frac{f(x_n)f'(x_n)}{[f'(x_n)]^2 - [f(x_n)][f''(x_n)]}, \qquad n \ge 0,$$

we have

$$x_{n+1} = x_n - \frac{(1 - \cos x_n)(\sin x_n)}{[\sin x_n]^2 - (1 - \cos x_n)(\cos x_n)}, \quad n \ge 0.$$

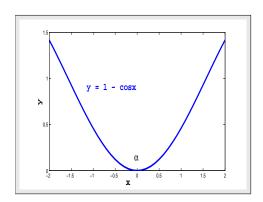


Figure 2.22: Graphical Solution of 1 - cos x = 0.

For n = 0 and the initial approximation $x_0 = 0.1$, we have

$$x_1 = x_0 - \frac{(1 - \cos x_0)(\sin x_0)}{[\sin x_0]^2 - (1 - \cos x_0)(\cos x_0)} = 0.1 - \frac{(1 - \cos 0.1)(\sin 0.1)}{[\sin 0.1]^2 - (1 - \cos 0.1)(\cos 0.1)} = 0.098,$$

which is the required first approximation to $\alpha = 0$, (see Figure 2.22).

Example 2.33 Show that the function $f(x) = e^x - \frac{x^2}{2} - x - 1$ has zero of multiplicity 3 at $\alpha = 0$ and then, find the approximate solution of the zero of the function with the help of the Newton's method, first and second modified Newton's methods, by taking initial approximation $x_0 = 1.5$ within an accuracy of 10^{-4} .

Solution. Since $\alpha = 0$ is a root of f(x), (see Figure 2.23), so

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

$$f'(x) = e^{x} - x - 1, \quad f'(0) = 0,$$

$$f''(x) = e^{x} - 1, \quad f''(0) = 0,$$

$$f'''(x) = e^{x}, \quad f'''(0) = 1 \neq 0,$$

the function has zero of multiplicity 3. In Table 2.6 we showed the comparison of three methods.

To use MATLAB command for the first modified Newton's method (2.25) and the second modified Newton's method (2.26), we define a function m-file as fn1.m and its derivatives m-files as dfn1.m and ddfn1.m for the equation as follows:

$$\begin{array}{ll} function \ y = fn1(x) & function \ dy = dfn1(x) & function \ ddy = ddfn1(x) \\ y = exp(x) - x.^2/2 - x - 1; & dy = exp(x) - x - 1; & ddy = exp(x) - 1; \end{array}$$

then use the following commands:

$$>> x0 = 1.5; m = 3; tol = 1e - 4; \ sol = mnewton1('fn1', 'dfn1', x0, m, tol)$$

50 2.6 Multiplicity of a Root

Table 2.6: Comparison results of three methods for the Example 2.33						
	Newton's Method	1st. M.N. Method	2nd. M.N. Method			
n	x_n	$ x_n $	x_n			
00	1.500000	1.500000	1.500000			
01	1.067698	0.2030926	-0.297704			
02	0.745468	3.482923e-03	-6.757677e-03			
03	0.513126	1.010951e-06	-3.798399e-06			
25	7.331582e-05					

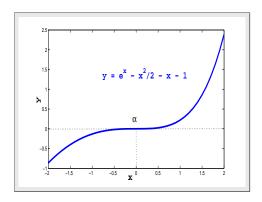


Figure 2.23: Graphical Solution of $e^x - x^2/2 - x - 1 = 0$.

and

$$>> x0 = 1.5; tol = 1e - 4;$$

 $>> sol = mnewton2('fn1', 'dfn1', 'ddfn1', x0, tol)$

We note that for the multiple root the both modified Newton's methods converge very fast as they took 4 iterations to converge while the Newton's method converges very slow and took 25 iterations to converge for the same accuracy.

Program 2.5

MATLAB m-file for First Modified Newton's Method function sol=mnewton1(fn1,dfn1,x0,m,tol) old = x0+1; while abs (x0 - old) > tol; old = x0; fa=feval(fn,old); fb=feval(dfn,old); x0 = old - (m * fa)/fb; end; sol=x0;

Note that when the order of multiplicity of a root of the equation f(x) = 0 is not known, then the second modified Newton's formula (2.26) can be used. MATLAB m-file can be written as follows

Program 2.6

MATLAB m-file for Second Modified Newton's Method function sol=mnewton2(fn1,dfn1,ddfn1,x0,tol)

old = x0+1; while abs (x0 - old) > tol; old = x0;

fa=feval(fn,old); fb=feval(dfn,old); fc=feval(ddfn,old);

 $x0 = old - (fa * fb)/((fb). ^2 - (fa * fc));$ end; sol=x0;

2.7 Convergence of Iterative Methods

Now we define the order of the convergence of functional iteration schemes discussed in the previous sections. This is a measure of how rapidly a sequence converges.

Definition 2.5 (Order of Convergence)

Suppose that the sequence $\{x_n\}_{n=0}^{\infty}$ converges to α , and let $e_n = \alpha - x_n$ define the error of the nth iterate. If two positive constants $\beta \neq 0$ and R > 0 exist, and

$$\lim_{n \to \infty} \frac{|\alpha - x_{n+1}|}{|\alpha - x_n|^R} = \lim_{n \to \infty} \frac{|e_{n+1}|}{|e_n|^R} = \beta, \tag{2.27}$$

then the sequence is said to converge to α with order of convergence R. The number β is called the asymptotic error constant. The cases R=1,2 are given special consideration.

If R = 1, the convergence of the sequence $\{x_n\}_{n=0}^{\infty}$ is called linear.

If R=2, the convergence of the sequence $\{x_n\}_{n=0}^{\infty}$ is called quadratic.

If R is large, the sequence $\{x_n\}$ converges rapidly to α ; that is, (2.27) implies that for large values of n we have the approximation $|e_{n+1}| \approx \beta |e_n|^R$. For example, suppose that R = 2 and $|e_n| \approx 10^{-3}$; then we could expect that $|e_{n+1}| \approx \beta \times 10^{-6}$.

Example 2.34 Show that the following sequence

$$x_{n+1} = \frac{1}{2}x_n\left(1 + \frac{N}{x_n^2}\right), \quad n \ge 0,$$

will converge quadratically to \sqrt{N} .

Solution. Since the sequence is given as

$$x_{n+1} = \frac{1}{2}x_n\left(1 + \frac{N}{x_n^2}\right),$$

and $\alpha = \sqrt{N}$, then we have

$$x_{n+1} - \sqrt{N} = \frac{1}{2} x_n \left(1 + \frac{N}{x_n^2} \right) - \sqrt{N} = \frac{1}{2} \left(x_n + \frac{N}{x_n} - 2\sqrt{N} \right)$$
$$= \frac{1}{2} \left(\sqrt{x_n} - \frac{\sqrt{N}}{\sqrt{x_n}} \right)^2 = \frac{1}{2x_n} (x_n - \sqrt{N})^2.$$

Thus

$$e_{n+1} = \frac{1}{2x_n}e_n^2$$
 or $e_{n+1} \propto e_n^2$,

which shows the quadratic convergence.

Example 2.35 Find the value of a and b so that the rate of convergence of the iterative scheme

$$x_{n+1} = ax_n + b\left(\frac{N}{x_n^2}\right), \quad for \quad n \ge 0,$$

for computing $N^{1/3}$ becomes quadratically or higher.

Solution. Since we have

$$x = N^{1/3} \qquad or \qquad x^3 = N,$$

therefore

$$f(x) = x^3 - N.$$

Let α be the exact root and $e_n = \alpha - x_n$, be the error in nth step, then by substituting

$$\alpha^3 = N$$
, $x_n = \alpha + e_n$, $x_{n+1} = \alpha + e_{n+1}$,

in the given iterative scheme, we get

$$\alpha + e_{n+1} = a(\alpha + e_n) + b\left(\frac{\alpha^3}{(\alpha + e_n)^2}\right)$$

$$= a(\alpha + e_n) + b\left(\frac{\alpha^3}{\alpha^2 (1 + e_n/\alpha)^2}\right)$$

$$= a(\alpha + e_n) + b\alpha\left(1 + \frac{e_n}{\alpha}\right)^{-2}$$

$$= a(\alpha + e_n) + b\alpha\left\{1 - 2\frac{e_n}{\alpha} + 3\left(\frac{e_n}{\alpha}\right)^2 - \dots\right\}$$

$$= a(\alpha + e_n) + b\alpha - 2be_n + 3b\frac{e_n^2}{\alpha} - \dots$$

Thus

$$e_{n+1} = (a+b-1)\alpha + (a-2b)e_n + O(e_n^2) + \dots$$

Now for the method to become of order 2, we must have

$$a + b - 1 = 0$$
 and $a - 2b = 0$,

and by solving this system for a and b, we have, $a = \frac{2}{3}$ and $b = \frac{1}{3}$, the required values.

Quadratically convergent sequences generally converge much faster than those that converge only linearly, but many techniques that generate convergent sequences do so only linearly. The following two lemmas tell us about the conditions of the linear convergence and the quadratic convergence of the sequences.

Lemma 2.3 (Linear Convergence)

Let g is continuously differentiable on the interval [a,b] and suppose that $g(x) \in [a,b]$ for all $x \in [a,b]$. Suppose that g'(x) is continuous on (a,b) with

$$|g'(x)| \le k < 1;$$
 for all $x \in (a, b)$.

If $g'(\alpha) \neq 0$, then for any $x_0 \in [a, b]$, the sequence $x_{n+1} = g(x_n)$, for $n \geq 0$, converges only linearly to the unique fixed-point α in [a, b].

Example 2.36 Consider an iterative scheme

$$x_{n+1} = 0.4 + x_n - 0.1x_n^2, \qquad n \ge 0.$$

Will this scheme converge to the fixed-point $\alpha = 2$? If yes, find its rate of convergence.

Solution. Since

$$g(x) = 0.4 + x - 0.1x^{2}$$
 and $g(2) = 0.4 + 2 - 0.1(2)^{2} = 2$,

which shows that the scheme converges to $\alpha = 2$. Also

$$g'(x) = 1 - 0.2x$$
, gives, $g'(2) = 1 - 0.4 = 0.6 \neq 0$.

Therefore, the scheme converges linearly.

Lemma 2.4 (Quadratic Convergence)

Let α be a solution of the equation x = g(x). Suppose that $g'(\alpha) = 0$ and g'' is continuous on an open interval (a,b) containing α . Then there exists a $\delta > 0$ such that, for $x_0 \in [\alpha - \delta, \alpha + \delta]$, the sequence $\{x_n\}_{n=0}^{\infty}$ defined by the iteration $x_{n+1} = g(x_n)$, for $n \geq 0$, converges at least quadratically to α .

Example 2.37 The iterative scheme

$$x_{n+1} = 2 - (1+a)x_n + ax_n^2, \qquad n \ge 0,$$

converges to $\alpha = 1$ for some values of a. Find the value of a for which the convergence is at least quadratic.

Solution. Given

$$g(x) = 2 - (1+a)x + ax^2$$
 and $g(1) = 2 - (1+a) + a = 1$.

Thus, the given iterative scheme converges to 1. Also

$$g'(x) = -(1+a) + 2ax,$$

and so

$$g'(1) = 0 = -(1+a) + 2a$$
, gives, $a = 1$.

Thus, the convergence of the given iterative scheme is at least quadratic for the value of a = 1.

Note 2.3 The sequence $\{x_n\}_{n=0}^{\infty}$ defined by the iteration

$$x_{n+1} = g(x_n), \quad for \qquad n \ge 0,$$

converges only quadratically to α if

$$g'(\alpha) = 0$$
 but $g''(\alpha) \neq 0$.

and cubically (order three) to α if

$$g'(\alpha) = 0$$
, $g''(\alpha) = 0$ but $g'''(\alpha) \neq 0$.

In the similar manner the higher order of convergence can be achieved.

Example 2.38 (a) Find the values of k_1 and k_2 such that the iterative scheme

$$x_{n+1} = k_1 x_n^2 + \frac{k_2}{x_n} - 5, \qquad n \ge 0,$$

converges quadratically to $\alpha = 1$.

(b) What is the order of convergence of the iteration

$$x_{n+1} = \frac{x_n(x_n^2 + 3k)}{3x_n^2 + k}, \quad k > 0,$$

as it converges to the fixed-point $\alpha = \sqrt{k}$.

Solution. (a) Given

$$g(x) = k_1 x^2 + \frac{k_2}{x} - 5,$$

and at fixed-point $\alpha = 1$, we have

$$g(1) = 1 = k_1 + k_2 - 5$$
, gives, $k_1 + k_2 = 6$.

Also

$$g'(x) = 2k_1x - \frac{k_2}{r^2},$$

and the convergence is quadratic at $\alpha = 1$, we have

$$g'(1) = 0 = 2k_1 - k_2$$
, gives, $2k_1 - k_2 = 0$.

Solving these two equations for unknowns k_1 and k_2 , we obtain

$$k_1 = 2$$
 and $k_2 = 4$.

Note that

$$g''(x) = 2k_1 + \frac{2k_2}{x^3}$$
 and $g''(1) = 12 \neq 0$.

(b) Since the given iteration is

$$x_{n+1} = \frac{x_n(x_n^2 + 3k)}{3x_n^2 + k} = g(x_n), \text{ which gives }, g(x) = \frac{x(x^2 + 3k)}{3x^2 + k}.$$

The first derivative of g(x) can be found as

$$g'(x) = \frac{3(x^2 - k)^2}{(3x^2 + k)^2}.$$

To find the order of convergence of the iteration, we have to check the derivative g'(x) at fixed-point $x = \alpha = \sqrt{k}$, if it is equal to zero, then order is at least quadratic, otherwise linear. So

$$g'(\sqrt{k}) = \frac{3[(\sqrt{k})^2 - k]^2}{[3(\sqrt{k})^2 + k]^2} = 0.$$

Therefore, the order of convergence for the given iteration is at least quadratic. One can find the second derivative of g(x) as

$$g''(x) = \frac{48xk(x^2 - k)}{(3x^2 + k)^3}, \quad g''(\sqrt{k}) = 0, \quad but \ g'''(x) = \frac{-48k(9x^4 - 18kx^2 + b^2)}{(3x^2 + k)^4}, \quad g'''(\sqrt{k}) = \frac{3}{2k} \neq 0.$$

Hence, the order of convergence for the given iteration is exactly cubic.

Now we discuss the rate of the convergence of all the iterative methods for the nonlinear equations which we discussed in the previous sections.

Case 2.1 (Bisection Method)

The convergence of the bisection method is very slow. At each step we gain one binary digit in accuracy. Since $10^{-1} \approx 2^{-3.3}$, we gain on the average one decimal digit per 3.3 steps. Note that the rate of convergence is completely independent of the function f(x). This is because we only make use of the sign of the computed function values. To investigate the rate of convergence of the bisection method, we consider the following example.

Example 2.39 If α is the fixed-point of the equation x = g(x) in [x, b]. Then show that rate of convergence of the bisection method is linear.

Solution. Since the bisection iteration function is define in the interval [x,b], so by using the bisection formula (2.2), we have

$$g(x) = \frac{x+b}{2}$$
 and $g'(x) = \frac{1}{2}$.

So at $x = \alpha$, we have

$$g'(\alpha) = \frac{1}{2} \neq 0,$$

therefore, by the Lemma 2.3, the convergence is linear.

Case 2.2 (Fixed-point Method)

The convergence rate of the fixed-point iteration can be analyzed as follows. The general procedure is given by

$$x_{n+1} = g(x_n), \qquad n = 0, 1, 2, \dots$$
 (2.28)

Let $x = \alpha$ denote the solution to f(x) = 0, so $f(\alpha) = 0$ and $\alpha = g(\alpha)$. Then

$$x_{n+1} - \alpha = e_{n+1} = g(x_n) - g(\alpha), \tag{2.29}$$

where e_{n+1} denote the error of the (n+1)th iterate. Expressing $g(\alpha)$ in the Taylor series about x_n gives:

$$g(\alpha) = g(x_n) + g'(\eta)(\alpha - x_n), \qquad x_n \le \eta \le \alpha.$$
(2.30)

Solving (2.30) for $g(x_n) - g(\alpha)$ and substituting into (2.29), we get

$$e_{n+1} = g'(\eta)e_n,$$
 (2.31)

or

$$|e_{n+1}| = |g'(\eta)||e_n|. \tag{2.32}$$

Now suppose that $|g'(x)| \le k < 1$ for all values of x in an interval. If x_1 is choose in this interval, x_2 will also be in the interval and the fixed-point iteration method will converge, since

$$\left| \frac{e_{n+1}}{e_n} \right| = |g'(\eta)| < 1.$$
 (2.33)

Convergence is linear since e_{n+1} is linearly dependent on e_n . If $|g'(\eta)| > 1$, the procedure diverges. If $|g'(\eta)| < 1$, but close to one, convergence is quite slow.

Example 2.40 (a) Show that $\alpha = 1$ is a unique fixed-point of

$$g(x) = \frac{x^2 - 4x + 7}{4}.$$

(b) Find the rate of convergence of the sequence

$$x_n = \frac{x_{n-1}^2 - 4x_{n-1} + 7}{4}.$$

Solution. (a) Firstly, we show that $\alpha = 1$ is a fixed-point of g(x) by showing that g(1) = 1 and it happened because

$$g(1) = \frac{1 - 4 + 7}{4} = 1.$$

It is unique also because

$$g'(x) = \frac{2x-4}{4}$$
 and $|g'(1)| = 0.5 < 1$.

(b) To find the rate of convergence of the given sequence, we have

$$g(x) = \frac{x^2 - 4x + 7}{4}$$
 and $g'(x) = \frac{2x - 4}{4}$.

Taking $x = \alpha = 1$, gives

$$g'(1) = \frac{2-4}{4} = -\frac{1}{2} \neq 0.$$

Hence the rate of the convergence of the given sequence is linear or slow.

Case 2.3 (Newton's Method)

The convergence rate of the Newton's method can be analyzed as follows. The general procedure

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}, \qquad n = 0, 1, 2, \dots,$$

is of the form $x_{n+1} = g(x_n)$. Consequently, if the method converges, then the absolute value of the derivative of the function g(x) with respect to x must be less than one, that is, |g'(x)| < 1. Since

$$g(x) = x - \frac{f(x)}{f'(x)},$$

then

$$g'(x) = 1 - \frac{f'(x)f'(x) - f(x)f''(x)}{[f'(x)]^2} = \frac{f(x)f''(x)}{[f'(x)]^2}.$$

Hence if

$$\left| \frac{f(x)f''(x)}{[f'(x)]^2} \right| < 1, \quad or \quad |f(x)f''(x)| < [f'(x)]^2,$$
 (2.34)

on an interval about the root α , the method will converge for any initial approximation in the interval. The (2.34) represents a sufficient condition for convergence. It is evident that f'(x) must not be zero. This is an important factor to consider when choosing the initial x value.

Now we show that the Newton's method is quadratically convergent for the simple root. Let $x = \alpha$ denote the solution to f(x) = 0, so $f(\alpha) = 0$ and $\alpha = g(\alpha)$. Since $x_{n+1} = g(x_n)$, we can write

$$x_{n+1} - \alpha = e_{n+1} = g(x_n) - g(\alpha), \tag{2.35}$$

where e_n denote the error of the nth iterate. Let us expand $g(x_n)$ as a Taylor series in terms of $(x_n - \alpha)$ with the second derivative term as the remainder:

$$g(x_n) = g(\alpha) + g'(\alpha)(x_n - \alpha) + \frac{g''(\eta)}{2}(x_n - \alpha)^2, \quad x_n \le \eta \le \alpha.$$

Since

$$g'(\alpha) = \frac{f(\alpha)f''(\alpha)}{[f'(\alpha)]^2} = 0,$$

because $f(\alpha) = 0$, we have

$$g(x_n) = g(\alpha) + \frac{g''(\eta)}{2}(x_n - \alpha)^2.$$

Solving above equation for $(g(x_n) - g(\alpha))$ and substituting into (2.35), we get

$$e_{n+1} = g(\alpha) - g(x_n) = -\frac{g''(\eta)}{2} (e_n)^2.$$
 (2.36)

This implies that each error is (in the limit) proportional to the square of the previous error, that is, the Newton's method is quadratically convergent.

Example 2.41 Show that the sequence $x_{n+1} = \frac{1}{2}x_n\left(3 - \frac{x_n^2}{N}\right)$, $n \ge 0$, has convergence of the second order with the limit \sqrt{N} .

Solution. Let

$$x_{n+1} - \sqrt{N} = \frac{1}{2}x_n\left(3 - \frac{x_n^2}{N}\right) - \sqrt{N} = \frac{1}{2}x_n\left(1 - \frac{x_n^2}{N}\right) + (x_n - \sqrt{N}),$$

or

$$e_{n+1} = \frac{1}{2N} x_n \left(N - x_n^2 \right) + (x_n - \sqrt{N}) = (x_n - \sqrt{N}) \left[1 - \frac{x_n}{2N} (x_n + \sqrt{N}) \right].$$

Thus

$$e_{n+1} = \frac{(x_n - \sqrt{N})}{2N} \left[(N - x_n^2) + (N - x_n \sqrt{N}) \right] = -\left[\frac{(x_n + 2\sqrt{N})}{2N} \right] (x_n - \sqrt{N})^2,$$

which shows that $e_{n+1} \propto e_n^2$, the quadratic convergence.

Since we know that rate of convergence of the Newton's method is linear if the function has multiple root. In the following example we discuss the rate of the convergence of the Newton's method for the multiple roots.

Example 2.42 If $x = \alpha$ is a simple root of f(x) = 0, then show that the rate of convergence of the Newton's method is at least quadratic.

Solution. Consider the Newton's iteration function which is define as follows:

$$g(x) = x - \frac{f(x)}{f'(x)}.$$

Since α is a simple root of nonlinear equation f(x) = 0, so

$$f(\alpha) = 0$$
 but $f'(\alpha) \neq 0$.

Thus taking derivative of g(x), we get

$$g'(x) = 1 - \frac{f'(x)f'(x) - f(x)f''(x)}{[f'(x)]^2} = \frac{f(x)f''(x)}{[f'(x)]^2}.$$

At $x = \alpha$, we know that $f(\alpha) = 0$ and $f'(\alpha) \neq 0$, so we have

$$g'(\alpha) = 1 - \frac{f'(\alpha)f'(\alpha) - f(\alpha)f''(\alpha)}{[f'(\alpha)]^2} = \frac{f(\alpha)f''(\alpha)}{[f'(\alpha)]^2} = 0.$$

Thus from Lemma 2.4, the rate of convergence of Newton's method is at least quadratic.

Example 2.43 If $x = \alpha$ is a double root of f(x) = 0, then show that the rate of convergence of the Newton's method is linear.

Solution. Consider the Newton's iteration function which is define as follows:

$$g(x) = x - \frac{f(x)}{f'(x)}.$$

Since α is a double root of nonlinear equation f(x) = 0, so

$$f(\alpha) = 0$$
, $f'(\alpha) = 0$, $f''(\alpha) \neq 0$.

Thus taking derivative of g(x), we get

$$g'(x) = 1 - \frac{f'(x)f'(x) - f(x)f''(x)}{[f'(x)]^2} = \frac{f(x)f''(x)}{[f'(x)]^2}.$$

At $x = \alpha$, we know that $f(\alpha) = 0$ and $f'(\alpha) = 0$, so we have indeterminate form $\left(\frac{0}{0}\right)$. Using L'Hôpital's rule we have

$$g'(x) = \frac{f'(x)f''(x) + f(x)f'''(x)}{2[f'(x)f''(x)]},$$

once again we get the indeterminate form. So again applying the $L'H\hat{o}pital$'s rule, we obtain

$$g'(x) = \frac{(f''(x))^2 + 2f'(x)f'''(x) + f'(x)f^{(4)}}{2[(f''(x))^2 + f'(x)f'''(x)]},$$

and at $x = \alpha$, we get

$$g'(\alpha) = \frac{(f''(\alpha))^2 + 0 + 0}{2[(f''(x))^2 + 0]} = \frac{1}{2} \neq 0.$$

Thus from Lemma 2.3, the rate of convergence of Newton's method is linear.

In the following example we discuss the rate of convergence of Newton's method for the root of nonlinear equation with multiplicity m without using the $L'H\hat{o}$ pital's rule.

Example 2.44 If $x = \alpha$ is a root of multiplicity m of f(x) = 0, then show that the rate of convergence of the Newton's method is linear.

Solution. Consider the Newton's iteration function which is define as follows:

$$g(x) = x - \frac{f(x)}{f'(x)}.$$

Since the function f(x) has multiple root, so

$$f(x) = (x - \alpha)^m h(x),$$

and its derivative is

$$f'(x) = m(x - \alpha)^{m-1}h(x) + (x - \alpha)^m h'(x).$$

Substituting the values of the f(x) and f'(x) in the above equation, we get

$$g(x) = x - \frac{(x - \alpha)^m h(x)}{(m(x - \alpha)^{m-1} h(x) + (x - \alpha)^m h'(x))},$$

or

$$g(x) = x - \frac{(x - \alpha)h(x)}{(mh(x) + (x - \alpha)h'(x))}.$$

Then

$$g'(x) = 1 - \{([mh(x) + (x - \alpha)][h(x) + (x - \alpha)h'(x)] - [(x - \alpha)h(x)] \\ [mh'(x) + h'(x) + (x - \alpha)h''(x)]\}/([mh(x) + (x - \alpha)h'(x)]^2).$$

At $x = \alpha$, and since $f(\alpha) = 0$, we have

$$g'(\alpha) = 1 - \frac{[mh(\alpha)][h(\alpha)]}{[mh(\alpha)]^2} = 1 - \frac{1}{m} \neq 0,$$
 because $m > 1$.

Therefore, the Newton's method converges to a multiple zero from any sufficiently close approximation and the convergence is linear (by the Lemma 2.3), with ration $(1 - \frac{1}{m})$. In particular for a double root, the ration is $\frac{1}{2}$, which is comparable with the convergence of the bisection method. •

Case 2.4 (Secant Method)

The convergence rate of the secant method can be analyzed as follows. The general procedure is

$$x_{n+1} = x_n - \frac{f(x_n)(x_n - x_{n-1})}{f(x_n) - f(x_{n-1})}$$
(2.37)

As before, let $x_{n-1} = \alpha - e_{n-1}, x_n = \alpha - e_n$, and $x_{n+1} = \alpha - e_{n+1}$. Then

$$e_{n+1} = e_n - \frac{f(\alpha - e_n)(e_n - e_{n-1})}{f(\alpha - e_n) - f(\alpha - e_{n-1})}$$

Since by using the Taylor's theorem

$$y_n = f(\alpha - e_n) = f(\alpha) - e_n f'(\alpha) + \frac{e_n^2}{2!} f''(\alpha) - \cdots$$

and

$$y_{n-1} = f(\alpha - e_{n-1}) = f(\alpha) - e_{n-1}f'(\alpha) + \frac{e_{n-1}^2}{2!}f''(\alpha) - \cdots$$

therefore, we have

$$\begin{split} e_{n+1} &= e_n - \frac{(e_n - e_{n-1})[f(\alpha) - e_n f'(\alpha) + 1/2e_n^2 f''(\alpha) - \cdots]}{[-(e_n - e_{n-1})f'(\alpha) + 1/2(e_n^2 - e_{n-1}^2)f''(\alpha) - \cdots]} \\ &= e_n - \left[\frac{-e_n f'(\alpha) + 1/2e_n^2 f''(\alpha) - \cdots}{-f'(\alpha) + 1/2(e_n + e_{n-1})f''(\alpha) - \cdots} \right] \quad (because f(\alpha) = 0) \\ &= e_n - \frac{1}{f'(\alpha)} [-e_n f'(\alpha) + 1/2e_n^2 f''(\alpha) - \cdots] \\ &\times [-1 + 1/2(e_n + e_{n-1}) \frac{f''(\alpha)}{f'(\alpha)} - \cdots]^{-1} \\ &= e_n - \frac{1}{f'(\alpha)} [e_n f'(\alpha) + 1/2e_n^2 f''(\alpha) + \cdots] \\ &\times [1 - 1/2(e_n + e_{n-1}) \frac{f''(\alpha)}{f'(\alpha)} + \cdots]^{-1} \\ &= e_n - \frac{1}{f'(\alpha)} [e_n f'(\alpha) - 1/2e_n^2 f''(\alpha) + 1/2e_n(e_n + e_{n-1}) f''(\alpha) - \cdots] \\ &= e_n - \frac{1}{f'(\alpha)} [-e_n f'(\alpha) + 1/2e_n e_{n-1} f''(\alpha) - \cdots] \\ &= -\frac{f''(\alpha)}{2f'(\alpha)} e_n e_{n-1} + \cdots \end{split}$$

Hence

$$e_{n+1} \approx Ke_n e_{n-1}, \quad where \quad K = -\frac{f''(\alpha)}{2f'(\alpha)},$$
 (2.38)

so that the each error is proportional to the product of the previous two errors. By comparison with the Newton's method convergence we expect the rate of the convergence of the secant method is inferior to that of the Newton's method. If we put

$$e_n = \beta e_{n-1}^R,$$

where R is the order of convergence and the constant β is called the asymptotic error constant, then we obtain

$$e_{n+1} = \beta e_n^R = \beta (\beta e_{n-1}^R)^R,$$

and

$$\beta(\beta e_{n-1}^R)^R \approx K\beta e_{n-1}^R e_{n-1}.$$

Thus

$$e_{n-1}^{R^2} \approx \lambda e_{n-1}^{R+1},$$

for some constant λ , it follows that

$$R^2 = R + 1.$$

Solving this quadratic equation, we get

$$R = \frac{1 \pm \sqrt{5}}{2} = 1.61803,$$

neglecting the negative value. This formula for R tells us that

$$|e_n| \approx \beta |e_{n-1}|^{1.61803}$$
.

Thus the error of the secant method is of order 1.61803, which is in-between 1 and 2. This shows that the order of the secant method is better than the bisection method and fixed-point method but less than the Newton's method.

Remember, however, that the secant method does not require the derivative of f(x) to be evaluated at each step, so that in many ways the secant method is a very attractive alternative to the standard Newton's method.

Example 2.45 If $x = \alpha$ is a root of multiplicity m of nonlinear equation f(x) = 0, then show that the rate of convergence of the first modified Newton's method is at least quadratic.

Solution. The first modified Newton's iteration function is define as follows:

$$g(x) = x - m \frac{f(x)}{f'(x)}. (2.39)$$

Since the function f(x) has multiple root, so

$$f(x) = (x - \alpha)^m h(x),$$

and its derivative is

$$f'(x) = m(x - \alpha)^{m-1}h(x) + (x - \alpha)^m h'(x).$$

Substituting the values of the f(x) and f'(x) in (2.39), we get

$$g(x) = x - \frac{m(x - \alpha)^m h(x)}{(m(x - \alpha)^{m-1} h(x) + (x - \alpha)^m h'(x))},$$

or

$$g(x) = x - \frac{m(x - \alpha)h(x)}{(mh(x) + (x - \alpha)h'(x))}.$$

Then

$$g'(x) = 1 - m\{([mh(x) + (x - \alpha)][h(x) + (x - \alpha)h'(x)] - [(x - \alpha)h(x)] - [mh'(x) + h'(x) + (x - \alpha)h''(x)]\}/([mh(x) + (x - \alpha)h'(x)]^2).$$

At $x = \alpha$, and since $f(\alpha) = 0$, we have

$$g'(\alpha) = 1 - \frac{[m^2 h^2(\alpha)]}{[mh(\alpha)]^2},$$

it gives

$$g'(\alpha) = 0.$$

Therefore, the modified Newton's method converges to a multiple root α and the convergence is at least quadratically (by Lemma 2.4).

Similarly, if $x = \alpha$ is a root of multiplicity m of f(x) = 0, then by using the Example 2.45 one can easily show that the rate of convergence of the Newton's method is linear. As the Newton iteration function is defined by

$$g(x) = x - m\frac{f(x)}{f'(x)},$$

and proceeding in the same way as we did in the Example 2.45, one can get

$$g'(\alpha) = 1 - \frac{1}{m} \neq 0$$
, because $m > 1$.

Hence the Newton's method converges to a multiple root α from any sufficiently close approximation and the convergence is linear (by Lemma 2.3) with ration $(1-\frac{1}{m})$. In particular for a double root, the ration is $\frac{1}{2}$, which is comparable with the convergence of the bisection method.

2.8 Systems of Nonlinear Equations

A system of nonlinear algebraic equations may arise when one is dealing with problems involving optimization and numerical integration (Gauss quadratures). Generally, the system of equations may not be of the polynomial variety. Therefore a system of n equations in n unknowns is called nonlinear if one or more of the equations in the systems is/are nonlinear.

The numerical methods we discussed so far have been concerned with finding a root of a nonlinear algebraic equation with one independent variable. We now consider methods for solving systems of nonlinear algebraic equations in which each equation is a function of a specified number of variables. Consider the system of two nonlinear equations with two variables

$$f_1(x,y) = 0, (2.40)$$

and

$$f_2(x,y) = 0. (2.41)$$

The problem can be stated as follows:

Given the continuous functions $f_1(x,y)$ and $f_2(x,y)$, find the values $x=\alpha$ and $y=\beta$ such that

$$f_1(\alpha, \beta) = 0$$
 and $f_2(\alpha, \beta) = 0.$ (2.42)

The function $f_1(x,y)$ and $f_2(x,y)$ may be algebraic equations, transcendental or any nonlinear relationship between the input x and y and the output $f_1(x,y)$ and $f_2(x,y)$. The solutions to (2.40) and (2.41) are the intersections of the $f_1(x,y) = f_2(x,y) = 0$, see Figure 2.24. This problem is considerably more complicated then solution of a single nonlinear equation. The one-point iterative method discussed in the previous Section 2.5 for the solution of a single equation may be extended to system. So to solve the system of nonlinear equations we have many methods but we will use the Newton's method.

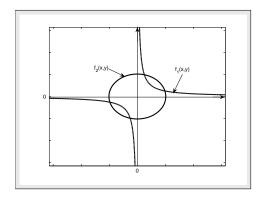


Figure 2.24: Nonlinear equation in two variables.

Newton's Method

Consider the two nonlinear equations specified by the equations (2.40) and (2.41). Suppose that (x_n, y_n) is an approximation to a root (α, β) , then by using the Taylor's theorem for functions of two variables for $f_1(x, y)$ and $f_2(x, y)$ expanding about (x_n, y_n) , we have

$$f_{1}(x,y) = f_{1}(x_{n} + (x - x_{n}), y_{n} + (y - y_{n}))$$

$$= f_{1}(x_{n}, y_{n}) + (x - x_{n}) \frac{\partial f_{1}(x_{n}, y_{n})}{\partial x} + (y - y_{n}) \frac{\partial f_{1}(x_{n}, y_{n})}{\partial y} + \cdots$$

and

$$f_{2}(x,y) = f_{2}(x_{n} + (x - x_{n}), y_{n} + (y - y_{n}))$$

$$= f_{2}(x_{n}, y_{n}) + (x - x_{n}) \frac{\partial f_{2}(x_{n}, y_{n})}{\partial x} + (y - y_{n}) \frac{\partial f_{2}(x_{n}, y_{n})}{\partial y} + \cdots$$

Since $f_1(\alpha, \beta) = 0$ and $f_2(\alpha, \beta) = 0$, these equations, with $x = \alpha$ and $y = \beta$, give

$$0 = f_1(x_n, y_n) + (\alpha - x_n) \frac{\partial f_1(x_n, y_n)}{\partial x} + (\beta - y_n) \frac{\partial f_1(x_n, y_n)}{\partial y} + \cdots$$
$$0 = f_2(x_n, y_n) + (\alpha - x_n) \frac{\partial f_2(x_n, y_n)}{\partial x} + (\beta - y_n) \frac{\partial f_2(x_n, y_n)}{\partial y} + \cdots$$

The Newton's method has a condition that initial approximation (x_n, y_n) should sufficiently close to exact root (α, β) , therefore, the higher order terms may be neglected to obtain

$$0 \approx f_1(x_n, y_n) + (\alpha - x_n) \frac{\partial f_1(x_n, y_n)}{\partial x} + (\beta - y_n) \frac{\partial f_1(x_n, y_n)}{\partial y}$$

$$0 \approx f_2(x_n, y_n) + (\alpha - x_n) \frac{\partial f_2(x_n, y_n)}{\partial x} + (\beta - y_n) \frac{\partial f_2(x_n, y_n)}{\partial y}$$
(2.43)

We see that this represents a system of two linear algebraic equations for α and β . Of course, since the higher order terms are omitted in the derivation of these equations, their solution (α, β) is no longer an exact root of (2.42) and (2.43). However, it will usually be a better approximation than (x_n, y_n) , so replacing (α, β) by (x_{n+1}, y_{n+1}) in (2.42) and (2.43), gives the iterative scheme

$$0 = f_1(x_n, y_n) + (x_{n+1} - x_n) \frac{\partial f_1(x_n, y_n)}{\partial x} + (y_{n+1} - y_n) \frac{\partial f_1(x_n, y_n)}{\partial y}$$

$$0 = f_2(x_n, y_n) + (x_{n+1} - x_n) \frac{\partial f_2(x_n, y_n)}{\partial x} + (y_{n+1} - y_n) \frac{\partial f_2(x_n, y_n)}{\partial y}$$

Then writing in the matrix form, we have

$$\begin{pmatrix}
\frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\
\frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y}
\end{pmatrix}
\begin{pmatrix}
x_{n+1} - x_n \\
y_{n+1} - y_n
\end{pmatrix} = -\begin{pmatrix} f_1 \\
f_2
\end{pmatrix},$$
(2.44)

where f_1, f_2 and their partial derivatives f_{1x}, f_{1y} are evaluated at (x_n, y_n) . Hence

$$\begin{pmatrix} x_{n+1} \\ y_{n+1} \end{pmatrix} = \begin{pmatrix} x_n \\ y_n \end{pmatrix} - \begin{pmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{pmatrix}^{-1} \begin{pmatrix} f_1 \\ f_2 \end{pmatrix}. \tag{2.45}$$

We call the following matrix J a Jacobian matrix

$$J = \begin{pmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{pmatrix}. \tag{2.46}$$

Note that (2.44) can be written in the simplified form as follows

$$\begin{pmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{pmatrix} \begin{pmatrix} h \\ k \end{pmatrix} = -\begin{pmatrix} f_1 \\ f_2 \end{pmatrix},$$

where h and k can be evaluated as

$$h = \frac{\left(-f_1 \frac{\partial f_2}{\partial y} + f_2 \frac{\partial f_1}{\partial y}\right)}{\left(\frac{\partial f_1}{\partial x} \frac{\partial f_2}{\partial y} - \frac{\partial f_1}{\partial y} \frac{\partial f_2}{\partial x}\right)} \quad \text{and} \quad k = \frac{\left(f_1 \frac{\partial f_2}{\partial x} - f_2 \frac{\partial f_1}{\partial x}\right)}{\left(\frac{\partial f_1}{\partial x} \frac{\partial f_2}{\partial y} - \frac{\partial f_1}{\partial y} \frac{\partial f_2}{\partial x}\right)},$$
(2.47)

where all functions are to be evaluated at (x, y). The Newton's method for a pair of equations in two unknowns is therefore

$$\begin{pmatrix} x_{n+1} \\ y_{n+1} \end{pmatrix} = \begin{pmatrix} x_n \\ y_n \end{pmatrix} + \begin{pmatrix} h \\ k \end{pmatrix}, \quad n = 0, 1, 2, \dots$$
 (2.48)

where (h, k) are given by (2.47) evaluated at (x_n, y_n) .

At a starting approximation (x_0, y_0) , the functions $f_1, f_{1x}, f_{1y}, f_2, f_{2x}$ and f_{2y} are evaluated. The linear equations are then solved for (x_1, y_1) and whole process is repeated until convergence is obtained. By comparison of the (2.14) and (2.45) shows that the above procedure is indeed an extension of the Newton's method in one variable, where division by f' generalized to pre-multiplication by J^{-1} .

Example 2.46 For the following system of two equations

$$x^3 + 3y^2 = 21
 x^2 + 2y = -2$$

Find the Jacobian matrix and its inverse using initial approximation (1,-1), then find the first approximation by using the Newton's method.

Solution. Given

$$f_1(x,y) = x^3 + 3y^2 - 21, \quad f_{1x} = 3x^2, \quad f_{1y} = 6y,$$

 $f_2(x,y) = x^2 + 2y + 2, \quad f_{2x} = 2x, \quad f_{2y} = 2.$

At the given initial approximation $x_0 = 1$ and $y_0 = -1$, we have

$$f_1(1,-1) = -17, \quad \frac{\partial f_1}{\partial x} = f_{1x} = 3, \quad \frac{\partial f_1}{\partial y} = f_{1y} = -6,$$

$$f_2(1,-1) = 1, \frac{\partial f_1}{\partial x} = f_{2x} = 2, \frac{\partial f_2}{\partial y} = f_{2y} = 2.$$

The Jacobian matrix J at the given initial approximation can be calculated as

$$J = \begin{pmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{pmatrix} = \begin{pmatrix} 3 & -6 \\ 2 & 2 \end{pmatrix} \quad and \quad J^{-1} = \frac{1}{18} \begin{pmatrix} 2 & 6 \\ -2 & 3 \end{pmatrix},$$

is the inverse of the Jacobian matrix. Now to find the first approximation we have to solve the following equation

$$\begin{pmatrix} x_1 \\ y_1 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \end{pmatrix} - \frac{1}{18} \begin{pmatrix} 2 & 6 \\ -2 & 3 \end{pmatrix} \begin{pmatrix} -17 \\ 1 \end{pmatrix} = \begin{pmatrix} 2.5556 \\ -3.0556 \end{pmatrix},$$

the required first approximation.

Example 2.47 Solve the following system of two equations using the Newton's method with given accuracy $\epsilon = 10^{-5}$.

$$4x^3 + y = 6$$
$$x^2y = 1$$

Assume $x_0 = 1.0$ and $y_0 = 0.5$ as starting values.

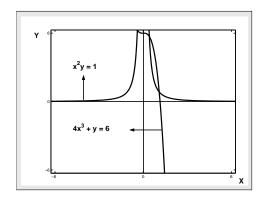


Figure 2.25: Graphical solution of the given nonlinear system.

Solution. Obviously this system of nonlinear equations has an exact solution of x = 1.088282 and y = 0.844340, (see Figure 2.25). Let us look how the Newton's method is used to approximate these roots. The first partial derivatives are as follows:

$$f_1(x,y) = 4x^3 + y - 6, \quad f_{1x} = 12x^2, \quad f_{1y} = 1,$$

 $f_2(x,y) = x^2y - 1, \quad f_{2x} = 2xy, \quad f_{2y} = x^2.$

At the given initial approximation $x_0 = 1.0$ and $y_0 = 0.5$, we get

$$f_1(1.0, 0.5) = -1.5, \quad \frac{\partial f_1}{\partial x} = f_{1x} = 12, \quad \frac{\partial f_1}{\partial y} = f_{1y} = 1.0,$$

$$f_2(1.0, 0.5) = -0.5, \quad \frac{\partial f_1}{\partial x} = f_{2x} = 1.0, \quad \frac{\partial f_2}{\partial y} = f_{2y} = 1.0.$$

The Jacobian matrix J and its inverse J^{-1} at the given initial approximation can be calculated as follows:

$$J = \begin{pmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{pmatrix} = \begin{pmatrix} 12.0 & 1.0 \\ 1.0 & 1.0 \end{pmatrix} \quad and \quad J^{-1} = \frac{1}{11.0} \begin{pmatrix} 1.0 & -1.0 \\ -1.0 & 12.0 \end{pmatrix}.$$

The Jacobian matrix can be find out by using MATLAB commands as follows:

>>
$$syms \ x \ y$$

>> $fun = [4 * x ^3 + y - 6, x ^2 * y - 1];$
>> $var = [x, y]; \ R = jacobian(f, var);$

Substituting all these values in (2.46), we get the first approximation as follows:

$$\left(\begin{array}{c} x_1 \\ y_1 \end{array}\right) = \left(\begin{array}{c} 1.0 \\ 0.5 \end{array}\right) - \frac{1}{11.0} \left(\begin{array}{c} 1.0 & -1.0 \\ -1.0 & 12.0 \end{array}\right) \left(\begin{array}{c} -1.5 \\ -0.5 \end{array}\right) = \left(\begin{array}{c} 1.090909 \\ 0.909091 \end{array}\right).$$

Similarly, the second iteration gives

$$\begin{pmatrix} x_2 \\ y_2 \end{pmatrix} = \begin{pmatrix} 1.090909 \\ 0.909091 \\ = \begin{pmatrix} 1.088264 \\ 0.844686 \end{pmatrix} - \frac{1}{15.012077} \begin{pmatrix} 1.190082 & -1.0 \\ -1.983471 & 14.280989 \end{pmatrix} \begin{pmatrix} 0.102178 \\ 0.081893 \end{pmatrix}$$

The first two and the further steps of the method are listed in Table 2.7.

Table 2.7: Solution of a system of two nonlinear equations									
	n	x-approx.	y-approx.	1st. func.	2nd. func.				
		x_n	y_n	$f_1(x_n, y_n)$	$f_2(x_n,y_n)$				
	00	1.000000	0.500000	-1.50000	-0.500000				
	01	1.090909	0.909091	0.102178	0.081893				
	02	1.088264	0.844686	0.000091	0.000377				
	03	1.088282	0.844340	0.000001	0.000001				

Table 2.7: Solution of a system of two nonlinear equations

A typical iteration of this method for this pair of equations can be implemented in MATLAB command window using:

```
>> f1 = 4 * x0^3 + y0 - 6; f2 = x0^2 * y0 - 1;

>> f1x = 12 * x0^2 2; f1y = 1; f2x = 2 * x_0 * y_0; f2y = x0^2 2;

>> D = f1x * f2y - f1y * f2x;

>> h = (f2 * f1y - f1 * f2y)/D; k = (f1 * f2x - f2 * f1x)/D;

>> x0 = x0 + h; y0 = y0 + k;
```

Using the starting value (1.0, 0.5), we found the possible approximations as shown in Table 2.7. \bullet

We see that the values of both the functional are approaching zero as the number of iterations is increased. We got the desired approximations to the roots after 3 iterations with accuracy $\epsilon = 10^{-5}$. The Newton's method is fairly easy to implement for the case of two equations in two unknowns. We first need the function m-files for the equations and the partial derivatives. For the equations in the Example 2.47, we do the following:

```
\begin{array}{ll} function \ f = fn2(v) & function \ J = dfn2(v) \\ \% f \ and \ v \ are \ vector \ quantities & \% Jacobian \ matrix \ for fn2.m \\ x = v(1); \ y = v(2); & x = v(1); \ y = v(2); \\ f(1) = 4 * x. ^2 + y - 6; & J(1,1) = 12 * x. ^2; \ J(1,2) = 1; \\ f(2) = x. ^2 * y - 1; & J(2,1) = 2 * x * y; \ J(2,2) = x. ^2; \end{array}
```

Then the following MATLAB commands can be used to generate the solution of the Example 2.47:

$$>> s = newton2('fn2', 'dfn2', [1.0, 0.5], 1e - 5)$$

The m-file Newton2.m will need both the function and its partial derivatives as well as starting vector and a tolerance. The following code can be used.

Program 2.7

MATLAB m-file for Newton's Method for a Nonlinear System function sol=newton2(fn2,dfn2,x0,tol)

old=x0+1; while max(abs(x0-old))>tol; old=x0;

f = feval(fn2, old); f1 = f(1); f2 = f(2); J = feval(dfn2, old);

f1x = J(1,1); f1y = J(1,2); f2x = J(2,1); f2y = J(2,2);

 $D = f1x * f2y - f1y * f2x; \ \ h = (f2 * f1y - f1 * f2y)/D;$

 $k = (f1 * f2x - f2 * f1x)/D; \ x0 = old + [h, k]; \text{ end; sol} = x0;$

Similarly, for a large system of equations it is convenient to use vector notation. Consider the system

$$\mathbf{f}(\mathbf{x}) = \mathbf{0},$$

where $\mathbf{f} = (f_1, f_2, \dots, f_n)^T$ and $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$. Denoting the *nth* iterate by $\mathbf{x}^{[n]} = (x_1^{[n]}, x_2^{[n]}, x_3^{[n]}, \dots, x_n^{[n]})^T$, then the Newton's method is defined by

$$\mathbf{x}^{[n+1]} = \mathbf{x}^{[n]} - \left[J(\mathbf{x}^{[n]}) \right]^{-1} \mathbf{f}(\mathbf{x}^{[n]}), \tag{2.49}$$

where the Jacobian matrix J is defined as

$$J = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \dots & \frac{\partial f_n}{\partial x_n} \end{pmatrix}.$$

Since the iterative formula (2.49) involves the inverse of Jacobian J, in practice we do not attempt to find this explicitly. In stead of using the form of (2.49) we use the following form

$$J(\mathbf{x}^{[\mathbf{n}]})\mathbf{Z}^{[\mathbf{n}]} = -\mathbf{f}(\mathbf{x}^{[n]}), \tag{2.50}$$

where $\mathbf{Z}^{[n]} = \mathbf{x}^{[n+1]} - \mathbf{x}^{[n]}$.

This represents a system of linear equations for $\mathbf{Z}^{[n]}$ and can be solved by any methods described in the next Chapter 3. Once $\mathbf{Z}^{[n]}$ has been found, the next iterate is calculated from

$$\mathbf{x}^{[n+1]} = \mathbf{Z}^{[n]} + \mathbf{x}^{[n]}.\tag{2.51}$$

There are two major disadvantages with method:

1. The method may not converges unless the initial approximation is a good one. Unfortunately, there are no general means by which an initial solution can be obtained. One can assume such values for which $\det(J) \neq 0$. This does not guarantee convergence but it does provide some guidance as to the appropriateness of one's initial approximation.

70 2.10 Exercises

2. The method requires the user to provide the derivatives of each function with respect to each variable. Therefore one must evaluate the n functions and the n^2 derivatives at each iteration. So solving systems of nonlinear equations is a difficult task. For systems of nonlinear equations that have analytical partial derivatives, Newton's method can be used; otherwise, multi-dimensional minimization techniques should be used.

Procedure 2.5 (Newton's Method for Two Nonlinear Equations)

- 1. Choose the initial guess for the roots of the system, so that the determinant of the Jacobian matrix is not zero.
- 2. Establish Tolerance $\epsilon(>0)$.
- 3. Evaluate the Jacobian at initial approximations and then find inverse of Jacobian.
- 4. Compute new approximation to the roots by using iterative formula (2.51).
- 5. Check tolerance limit. If $||(x_n, y_n) (x_{n-1}, y_{n-1})|| \le \epsilon$, for $n \ge 0$, then end; otherwise, go back to step 3, and repeat the process.

2.9 Exercises

- 1. Find the root of $f(x) = e^x 2 x$ in the interval [-2.4, -1.6] accurate to 10^{-4} using bisection method.
- 2. Use bisection method to find solutions accurate to within 10^{-4} on the interval [-5,5] of the following functions:

(a)
$$f(x) = x^5 - 10x^3 - 4$$
, (b) $f(x) = 2x^2 + \ln(x) - 3$, (c) $f(x) = \ln(x) + 30e^{-x} - 3$.

- 3. The following equations have a root in the interval [0, 1.6]. Determine these with an error less than 10^{-4} using bisection method.
 - (a) $2x e^{-x} = 0$; (b) $e^{-3x} + 2x 2 = 0$.
- 4. Estimate the number of iterations needed to achieve an approximation with accuracy 10^{-4} to the solution of $f(x) = x^3 + 4x^2 + 4x 4$ lying in the interval [0, 1] using bisection method.
- 5. Use the bisection method for $f(x) = x^3 3x + 1$ in [1, 3] to find:
 - (a) The first eight approximation to the root of the given equation.
 - (b) Find an error estimate $|\alpha x_8|$.
- 6. The cubic equation $x^3 3x 20 = 0$ can be written as

(a)
$$x = \frac{(x^3 - 20)}{3}$$
, (b) $x = \frac{3}{(x^3 - 3)}$, (c) $x = (3x + 20)^{1/3}$.

Choose the form which satisfies the condition |g'(x)| < 1 on [3,4] and then find third approximation x_3 when $x_0 = 3.5$.

- 7. Consider the nonlinear equation $g(x) = \frac{1}{2}e^{0.5x}$ defined on the interval [0, 1]. Then
 - (a) Show that there exists a unique fixed-point for g in [0,1].
 - (b) Use fixed-point iterative method to compute x_3 , set $x_0 = 0$.
 - (c) Compute an error bound for your approximation in part (b).
- 8. An equation $x^3 2 = 0$ can be written in form x = g(x) in two ways:

(a)
$$x = g_1(x) = x^3 + x - 2$$
, (b) $x = g_2(x) = \frac{(2 + 5x - x^3)}{5}$
Generate first four approximations from $x_{n+1} = g_i(x_n)$, $i = 1, 2$ by using $x_0 = 1.2$. Show

which sequence converge to $2^{1/3}$ and why?

- 9. Find value of k such that the iterative scheme $x_{n+1} = \frac{x_n^2 4kx_n + 7}{4}$, $n \ge 0$ converges to 1. Also, find the rate of convergence of the iterative scheme
- 10. Write the equation $x^2 6x + 5 = 0$ in the form x = g(x), where $x \in [0, 2]$, so that the iteration $x_{n+1} = g(x_n)$ will converge to the root of the given equation for any initial approximation $x_0 \in [0, 2].$
- 11. Which of the following iterations

(a)
$$x_{n+1} = \frac{1}{4} \left(x_n^2 + \frac{6}{x_n} \right)$$
, (b) $x_{n+1} = \left(4 - \frac{6}{x_n^2} \right)$

is suitable to find a root of the equation $x^3 = 4x^2 - 6$ in the interval [3,4]? Estimate the number of iterations required to achieve 10^{-3} accuracy, starting from $x_0 = 3$.

- 12. An equation $e^x = 4x^2$ has a root in [4,5]. Show that we cannot find that root using x = $g(x) = \frac{1}{2}e^{x/2}$ for the fixed-point iteration method. Can you find another iterative formula which will locate that root? If yes, then find third iterations with $x_0 = 4.5$. Also find the error bound.
- 13. Let $f(x) = e^x + 3x^2$. Find Newton's formula $g(x_k)$. Start with $x_0 = 4$ and $x_0 = -0.5$, compute
- 14. Use Newton's formula for the reciprocal of square root of a number 15 and then find the 3rd approximation of number, with $x_0 = 0.05$.
- 15. Use Newton's method to find solution accurate to within 10^{-4} of the equation $\tan(x) 7x = 0$, with initial approximation $x_0 = 4$.
- 16. Find Newton's formula for $f(x) = x^3 3x + 1$ in [1, 3] to calculate x_3 , if $x_0 = 1.5$. Also, find the rate of convergence of the method.
- 17. Rewrite the nonlinear equation $g(x) = \frac{1}{2}e^{0.5x}$ which defined in the interval [0, 1] in the equivalent form f(x) = 0 and then use the Newton's method with $x_0 = 0.5$ to find third approximation x_3 .
- 18. Given the iterative scheme $x_{n+1} = x_n \frac{f(x_n)}{f'(x_n)}$, $n \ge 0$ with $f(\alpha) = f'(\alpha) = 0$ and $f''(\alpha) \ne 0$. Find the rate of convergence for this scheme

72 2.10 Exercises

- 19. Find x_4 for $x^3 2x 5 = 0$ by secant method using $x_0 = 2$ and $x_1 = 3$.
- 20. Solve the equation $e^{-x} x = 0$ by secant method, using $x_0 = 0$ and $x_1 = 1$, accurate to 10^{-4} .
- 21. Use secant method to find a solution accurate to within 10^{-4} for $\ln(x) + x 5 = 0$ on [3, 4].
- 22. Find the root of multiplicity of the function $f(x) = (x-1)^2 \ln(x)$ at $\alpha = 1$.
- 23. Show that if f(x) has a root of multiplicity m at $x = \alpha$, then

$$f^{(n)}(x) = 0,$$
 $n = 1, 2, ..., m - 1.$

- 24. Show that the root of multiplicity of the function $f(x) = x^4 x^3 3x^2 + 5x 2$ is 3 at $\alpha = 1$. Estimate the number of iterations required to solve the problem with accuracy 10^{-4} , start with the starting value $x_0 = 0.5$ by using:
 - (a) Newton's method; (b) First modified Newton's method; (c) Second modified Newton's method
- 25. If f(x), f'(x) and f''(x) are continuous and bounded on a certain interval containing $x = \alpha$ and if both $f(\alpha) = 0$ and $f'(\alpha) = 0$ but $f''(\alpha) \neq 0$, show that

$$x_{n+1} = x_n - 2\frac{f(x_n)}{f'(x_n)}$$

will converge quadratically if x_n is in the interval.

- 26. Show that iterative scheme $x_{n+1} = 1 + x_n \frac{x_n^2}{2}$, $n \ge 0$ converges to $\sqrt{2}$. Find the rate of convergence of the sequence.
- 27. Let α be the exact solution of the function f(x) = 0 such that $f'(\alpha) \neq 0$, $f''(\alpha) \neq 0$, then find the conditions on the constant K under which the rate of convergence of the sequence $x_{n+1} = x_n^2 Kf(x_n)$, $n = 0, 1, 2, \ldots$ is quadratic.
- 28. Solve the following system using the Newton's method:

$$4x^3 + y = 6$$
$$x^2y = 1$$

Start with initial approximation $x_0 = y_0 = 1$. Stop when successive iterates differ by less than 10^{-7} .

29. Solve the following system using the Newton's method:

$$\begin{array}{rcl}
x & + e^y & = & 68.1 \\
\sin x & -y & = & -3.6
\end{array}$$

Start with initial approximation $x_0 = 2.5$, $y_0 = 4$, compute the first three approximations.

Chapter 3

Systems of Linear Algebraic Equations

3.1 Introduction

When engineering systems are modeled, the mathematical description is frequently developed in terms of set of algebraic simultaneous equations. Sometimes these equations are non-linear and sometimes linear. In this chapter we discuss systems of simultaneous linear equations and describe the numerical methods for the approximate solutions of such systems. The solution of a system of simultaneous linear algebraic equations is probably one of the most important topics in the engineering computation. Problems involving simultaneous linear equations arise in the areas of elasticity, electric-circuit analysis, heat transfer, vibrations and so on. Also, the numerical integration of some types of ordinary and partial differential equations may be reduced to the solution of such a system of equations. It has been estimated, for example, that about 75% of all scientific problems require the solution of a system of linear equations at one stage or another. It is therefore important to be able to solve linear problems efficiently and accurately.

Important Points of the Chapter 3

- I. In this chapter we look for the solutions of systems of linear equations.
- II. Linear systems may be simultaneous (number of linear equations and unknowns variables are equal) or underdetermined (number of linear equations less than unknowns variables) or overdetermined (number of linear equations more than unknowns variables). Here, we shall discuss only simultaneous systems.
- III. Matrix form of linear system is $A\mathbf{x} = \mathbf{b}$, where A called coefficient matrix, column matrix \mathbf{b} is right hand constant and column matrix \mathbf{x} be the unknowns.
- IV. Linear systems may be nonhomogeneous (right hand vector $\mathbf{b} \neq 0$) or homogeneous ($\mathbf{b} = 0$).
- V. Linear systems may have unique solution or no solution or infinitely many solutions.

74 3.1 Introduction

VI. Linear systems may be nonsingular (determinant of coefficients matrix A not equal to zero) or singular (determinant of coefficients matrix A equal to zero). Nonsingular systems have unique solution while singular systems have either no solution or infinitely many solutions.

VII. Solutions of linear systems can be obtained by both direct and indirect (iterative) methods.

VII. Linear systems may be well-conditioned (small condition number) or ill-conditioned (large condition number) .

Definition 3.1 (Linear equation)

It is an equation in which the highest exponent in a variable term is no more than one. The graph of such equation is a straight line.

A linear equation in two variables x_1 and x_2 is an equation that can be written in the form

$$a_1x_1 + a_2x_2 = b$$
,

where a_1, a_2 , and b are real numbers. Note that this is the equation of a straight line in the plane. For example, the equations

$$5x_1 + 2x_2 = 2$$
, $\frac{4}{5}x_1 + 2x_2 = 1$, $2x_1 - 4x_2 = \pi$,

are all linear equations in two variables.

A linear equation in n variables x_1, x_2, \ldots, x_n is an equation that can be written as

$$a_1x_1 + a_2x_2 + \dots + a_nx_n = b,$$

where a_1, a_2, \ldots, a_n are real numbers and called the *coefficients* of unknown variables x_1, x_2, \ldots, x_n and the real number b, the right-hand side of equation, is called the *constant term* of the equation.

Definition 3.2 (System of Linear Equations)

A system of linear equations (or linear system) is simply a finite set of linear equations.

For example,

$$\begin{array}{rcrrr} 4x_1 & - & 2x_2 & = & 5 \\ 3x_1 & + & 2x_2 & = & 4 \end{array}$$

is the system of two equations in two variables x_1 and x_2 , while

is the system of three equations in the four variables x_1, x_2, x_3 and x_4 .

In order to write a general system of m linear equations in the n variables x_1, \ldots, x_n , we have

or, in compact form the system (3.1) can be written

$$\sum_{j=1}^{n} a_{ij} x_j = b_i, \qquad i = 1, 2, \dots, m.$$
(3.2)

For such system we seek all possible ordered sets of numbers c_1, \ldots, c_n which satisfies all m equations when they are substituted for the variables x_1, x_2, \ldots, x_n . Any such set $\{c_1, c_2, \ldots, c_n\}$, is called a *solution* of the system of linear equations (3.1) or (3.2).

There are three possible types of linear systems arise in engineering problems and they are:

1. If there are more equations than unknown variables (m > n), then the system is usually called *overdetermined*. Typically, an overdetermined system has no solution. For example, the following system has no solution.

$$\begin{array}{rcl}
4x_1 & = & 8 \\
3x_1 & + & 9x_2 & = & 13 \\
& & 3x_2 & = & 9
\end{array}$$

2. If there are more unknown variables than the number of the equations (n > m), then the system is usually called *underdetermined*. Typically, an underdetermined system has infinite number of solutions. For example, the following system has infinitely many solutions.

$$x_1 + 5x_2 = 45$$

 $3x_2 + 4x_3 = 21$

3. If there are same number of equations as the unknown variables (m = n), then the system is usually called *simultaneous* system. It has unique solution if the system satisfies the certain conditions (which we will discuss below). For example, the system

has unique solution $x_1 = 2, x_2 = -4$ and $x_3 = 1$.

Most engineering problems fall into this category. In this chapter we will solve the simultaneous linear systems using many numerical methods.

76 3.1 Introduction

A simultaneous system of linear equations is said to be *linear independent* if no equation in the system can be expressed as a linear combination of the others. Under these circumstances a unique solution exists. For example, the following system of linear equations

$$\begin{array}{rclrcrcr}
2x_1 & + & x_2 & - & x_3 & = & 1 \\
x_1 & - & 2x_2 & + & 3x_3 & = & 4 \\
x_1 & + & x_2 & & = & 1
\end{array}$$

are linear independent and therefore, has unique solution $x_1 = 1, x_2 = 0$ and $x_3 = 1$. However, the system

does not have a unique solution since the equations are not linear independent; the first equation is equal to the second equation plus twice the third equation.

Theorem 3.1 (Solution of a Linear System)

Every system of linear equations has either no solution, exactly one solution, or infinitely many solutions.

For example, in the case of a system of two equations in two variables, we can have these three possibilities for the solutions of the linear system. Firstly, the two lines (since the graph of linear equation is straight line) may be parallel and distinct, in this case there is no solution to the system because the two lines do not intersect each other at any point. For example, consider the following system

$$\begin{array}{rcrrr} x_1 & + & x_2 & = & 1 \\ 2x_1 & + & 2x_2 & = & 3 \end{array}$$

From the graphs (see Figure 3.1(a)) of the given two equations show that lines are parallel so, the given system has no solution. It can be proved algebraically simply by multiplying the first equation of the system by 2, to get a system of the form

$$2x_1 + 2x_2 = 2$$

 $2x_1 + 2x_2 = 3$

which is not possible.

Secondly, the two lines may not be parallel, and they meet exactly one point, so in this case the system has exactly one solution. For example, consider the following system

$$\begin{array}{rcrrr} x_1 & - & x_2 & = & -1 \\ 3x_1 & - & x_2 & = & 3 \end{array}$$

From the graphs (see Figure 3.1(b)) of these two equations, we can see that the lines intersect in exactly one point, namely, the (2,3), and so the system has exactly one solution, $x_1 = 2$, $x_2 = 3$. To show algebraically, if we substitute $x_2 = x_1 + 1$ in the second equation, we have $3x_1 - x_1 - 1 = 3$, or $x_1 = 2$ and using this value of x_1 in $x_2 = x_1 + 1$, gives $x_2 = 3$.

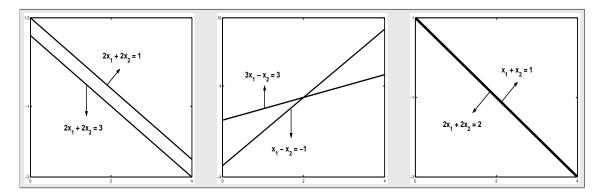


Figure 3.1: Three possible solutions of simultaneous systems.

Finally, the two lines may actually be the same line, and so in this case, every point on the lines gives a solution to the system, and therefore, there are infinitely many solutions. For example, consider the following system

$$\begin{array}{rcl} x_1 & + & x_2 & = & 1 \\ 2x_1 & + & 2x_2 & = & 2 \end{array}$$

Here, both equations have same line for their graph, see Figure 3.1(c). So this system has infinitely many solutions because any point on this line gives a solution to this system. Since any solution of first equation is also solution of the second equation. For example, if we set $x_2 = x_1 - 1$ and choose $x_1 = 0, x_2 = 1$ and $x_1 = 1, x_2 = 0$, and so on.

Note that a system of equations with no solution is said to be *inconsistent system* and if it has at least one solution, it is said to be *consistent system*.

3.1.1 Linear System in Matrix Notation

To write the general simultaneous system of n linear equations in the n unknown variables x_1, x_2, \ldots, x_n , is

The system of linear equations (3.3) can be written as the single matrix equation

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}. \tag{3.4}$$

78 3.1 Introduction

If we compute the product of the two matrices on the left-hand side of (3.9), we have

$$\begin{pmatrix} a_{11}x_1 & + & a_{12}x_2 & + & \cdots & + & a_{1n}x_n \\ a_{21}x_1 & + & a_{22}x_2 & + & \cdots & + & a_{2n}x_n \\ \vdots & & \vdots & & \vdots & & \vdots \\ a_{n1}x_1 & + & a_{n2}x_2 & + & \cdots & + & a_{nn}x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}.$$
(3.5)

But two matrices are equal if and only if their corresponding elements are equal. Hence the single matrix equation (3.9) is equivalent to the system of the linear (3.3). If we define

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix},$$

the coefficient matrix, the column matrix of unknowns, and the column matrix of constants, respectively, then the system (3.3) can be written very compactly as

$$A\mathbf{x} = \mathbf{b},\tag{3.6}$$

which is called the *matrix form* of the system of linear equations (3.3). The column matrices \mathbf{x} and \mathbf{b} are called *vectors*. If right-hand sides of the equal signs of (3.6) are not zero, then the linear system (3.6) is called a *nonhomogeneous system*, and we will find that all the equations must be independent to obtain unique solution. If the constants \mathbf{b} of (3.6) are added to the coefficient matrix A as a column of elements in the position shown below

$$[A|\mathbf{b}] = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} & \vdots & b_1 \\ a_{21} & a_{22} & \cdots & a_{2n} & \vdots & b_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} & \vdots & b_n \end{pmatrix},$$
(3.7)

then the matrix $[A|\mathbf{b}]$ is called *augmented matrix* of the system (3.6). In many instances, it may be found convenient to operate on the augmented matrix instead of manipulating the equations. It is customary to put a bar between the last two columns of the augmented matrix remind us where the last column come from. However, the bar is not absolutely necessary. The coefficient and augmented matrices of a linear system will play key roles in our methods of solving linear systems.

Using MATLAB commands we can define augmented matrix as follows:

$$>> A = [1\ 2\ 3; 4\ 5\ 6; 7\ 8\ 9];\ b = [10; 11; 12];\ Aug = [A\ b];\ Aug = [A\ eye(3)]$$

Homogeneous Linear System

If all of the constant terms b_1, b_2, \ldots, b_n on the right-hand sides of the equal signs of the linear system (3.6) are zero, then the system (3.6) is called a homogeneous system, and it can be written

as

$$a_{11}x_{1} + a_{12}x_{2} + \cdots + a_{1n}x_{n} = 0$$

$$a_{21}x_{1} + a_{22}x_{2} + \cdots + a_{2n}x_{n} = 0$$

$$\vdots \qquad \vdots \qquad \vdots \qquad \vdots$$

$$a_{n1}x_{1} + a_{n2}x_{2} + \cdots + a_{nn}x_{n} = 0$$

$$(3.8)$$

The system of linear equations (3.8) can be written as the single matrix equation

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}. \tag{3.9}$$

Also, it can be written in more compact form as

$$A\mathbf{x} = \mathbf{0},\tag{3.10}$$

where

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}, \quad \mathbf{0} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

It can be seen by inspection of the homogeneous system (3.10) that one of its solution is $\mathbf{x} = \mathbf{0}$, such solution, in which all of the unknowns are zero, is called the *trivial solution* or *zero solution*. For the general nonhomogeneous linear system there are three possibilities: no solution, one solution, or an infinitely many solutions. For the general homogeneous system, there are only two possibilities: either the zero solution is the only solution or there are an infinitely many solutions (called non-trivial solutions). Of course, it is usually non-trivial solutions that are of interest in physical problems. A non-trivial solution to the homogeneous system can occurs with certain conditions on the coefficient matrix A which we will discuss later.

3.2 Properties of Matrices and Determinant

To discuss the solution of the linear systems, it will be necessary to introduce the basic algebraic properties of matrices which make it possible to describe the linear systems in a concise way that makes solving a system of n linear equations as easy as possible.

3.2.1 Introduction of Matrices

A matrix can be described as a rectangular array of elements that can be represented as follows:

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}.$$

$$(3.11)$$

The number $a_{11}, a_{12}, \ldots, a_{mn}$ that make up the array are called the elements of the matrix. The first subscript for the element denotes the row and the second denotes the column in which the element appear. The elements of a matrix may take many forms. They could be all numbers (real or complex), or variables, or functions, or integrals, or derivatives, or even matrices themselves.

The *order* or *size* of a matrix is specified by the number of rows (m) and column (n); thus the matrix A in (3.11) is of order m by n, usually written as $m \times n$.

A vector can be considered as a special case of a matrix having only one row or one column. A row vector containing n elements is a $1 \times n$ matrix, called a (row matrix), and a column vector of n elements is an $n \times 1$ matrix, called a (column matrix). A matrix of order 1×1 is called a scalar.

Definition 3.3 (Matrix Equality)

Two matrices $A = (a_{ij})$ and $B = (b_{ij})$ are equal if they are the same size and the corresponding elements in A and B are equal, that is

$$A = B$$
 if and only if $a_{ij} = b_{ij}$,

for i = 1, 2, ..., m and j = 1, 2, ..., n. For example, the following matrices

$$A = \begin{pmatrix} 1 & -1 & 2 \\ 1 & 3 & 2 \\ 2 & 4 & 3 \end{pmatrix} \quad and \quad B = \begin{pmatrix} 1 & -1 & z \\ 1 & 3 & 2 \\ x & y & w \end{pmatrix},$$

are equal if and only if x = 2, y = 4, z = 2 and w = 3.

Definition 3.4 (Addition of Matrices)

Let $A = (a_{ij})$ and $B = (b_{ij})$ are both $m \times n$ matrices, then the sum A + B of two matrices of same size is a new matrix $C = (c_{ij})$ each of whose elements is the sum of the two corresponding elements in the original matrices, that is

$$c_{ij} = a_{ij} + b_{ij}$$
, for $i = 1, 2, ..., m$ and $j = 1, 2, ..., n$.

For example, let

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \quad and \quad B = \begin{pmatrix} 4 & 1 \\ 5 & 2 \end{pmatrix},$$

then

$$\left(\begin{array}{cc} 1 & 2 \\ 3 & 4 \end{array}\right) + \left(\begin{array}{cc} 4 & 1 \\ 5 & 2 \end{array}\right) = \left(\begin{array}{cc} 5 & 3 \\ 8 & 6 \end{array}\right) = C.$$

Using MATLAB commands adding two matrices A and B of same size, results in the answer C another matrix of same size, are:

$$>> A = [1\ 2; 3\ 4];\ B = [4\ 1; 5\ 2];\ C = A + B$$

Definition 3.5 (Difference of Matrices)

Let A and B are $m \times n$ matrices, we write A + (-1)B as A - B and call this the difference of two matrices of same size is a new matrix C each of whose elements is the difference of the two corresponding elements in the original matrices. For example, let

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \quad and \quad B = \begin{pmatrix} 4 & 1 \\ 5 & 2 \end{pmatrix}.$$

Then

$$\left(\begin{array}{cc} 1 & 2 \\ 3 & 4 \end{array}\right) - \left(\begin{array}{cc} 4 & 1 \\ 5 & 2 \end{array}\right) = \left(\begin{array}{cc} -3 & 1 \\ -2 & 2 \end{array}\right) = C.$$

Note that (-1)B = -B is obtained by multiplying each entries of matrix B by (-1), called the scalar multiple of matrix B by -1. The matrix -B is called negative of the matrix B.

Definition 3.6 (Multiplication of Matrices)

The multiplication of two matrices is defined only when the number of columns in the first matrix is equal to the number of rows in the second. If an $m \times n$ matrix A is multiplied by an $n \times p$ matrix B, then the product matrix C is an $m \times p$ matrix where each term is defined by

$$c_{ij} = \sum_{k=1}^{n} a_{ik} b_{kj},$$

for each i = 1, 2, ..., m, and j = 1, 2, ..., p. For example, let

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \quad and \quad B = \begin{pmatrix} 4 & 1 \\ 5 & 2 \end{pmatrix}.$$

Then

$$\left(\begin{array}{cc} 1 & 2 \\ 3 & 4 \end{array}\right) \left(\begin{array}{cc} 4 & 1 \\ 5 & 2 \end{array}\right) = \left(\begin{array}{cc} 4+10 & 1+4 \\ 12+20 & 3+8 \end{array}\right) = \left(\begin{array}{cc} 14 & 5 \\ 32 & 11 \end{array}\right) = C.$$

Note that even AB is defined, the product BA may not be defined. Moreover, a simple multiplication of two square matrices of same size will show that even BA is defined, it need not equal to AB, that is, they do not commute. For example, if

$$A = \begin{pmatrix} 1 & 2 \\ -1 & 3 \end{pmatrix} \quad and \quad B = \begin{pmatrix} 2 & 1 \\ 0 & 1 \end{pmatrix},$$

then

$$AB = \left(\begin{array}{cc} 2 & 3 \\ -2 & 2 \end{array} \right) \quad while \quad BA = \left(\begin{array}{cc} 1 & 7 \\ -1 & 3 \end{array} \right).$$

Thus $AB \neq BA$.

Using MATLAB commands matrix multiplication has the standard meaning as well. Multiplying two matrices A and B of size $m \times p$ and $p \times n$ respectively, results in the answer C another matrix of size $m \times n$, are:

$$>> A = [1\ 2; -1\ 3];\ B = [2\ 1; 0\ 1];\ C = A*B$$

3.2.2 Some Special Matrix Forms

There are many special types of matrices that are encountered frequently in engineering analysis. We discuss some of the them in the following.

Definition 3.7 (Square Matrix)

A matrix A which has the same number of rows m and columns n, that is, m = n, defined as

$$A = (a_{ij}),$$
 for $i = 1, 2, ..., n$ and $j = 1, 2, ..., n$,

is called a square matrix. For example, the following matrices

$$A = \begin{pmatrix} 1 & 2 \\ -1 & 3 \end{pmatrix} \quad and \quad B = \begin{pmatrix} 2 & 1 & 2 \\ 1 & 2 & 3 \\ 0 & 1 & 5 \end{pmatrix},$$

are square matrices because both have the same numbers of rows and columns.

Definition 3.8 (Null Matrix)

It is a matrix in which all elements are zero, that is

$$A = (a_{ij}) = \mathbf{0},$$
 for $i = 1, 2, ..., n$ and $j = 1, 2, ..., n$.

It is also called zero matrix. It may be either rectangular or square. For example, the following matrices

$$A = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad and \quad B = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix},$$

are the zero matrices.

Definition 3.9 (Identity Matrix)

It is a square matrix in which the main diagonal elements are equal to 1, and is defined as follows

$$\mathbf{I} = (a_{ij}) = \begin{cases} a_{ij} = 0, & if \ i \neq j, \\ a_{ij} = 1, & if \ i = j. \end{cases}$$

An example of 4×4 identity matrix may be written as

$$\mathbf{I}_4 = \left(\begin{array}{cccc} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{array}\right).$$

The Identity matrix (also called unit matrix) serves somewhat the same purpose in matrix algebra as does the number one (unity) in scalar algebra. It is called the identity matrix because multiplication of a matrix by it will result in a same matrix. For a square matrix A of order n, it can be seen that

$$\mathbf{I}_n A = A \mathbf{I}_n = A.$$

Similarly, for rectangular matrix B of order $m \times n$, we have

$$\mathbf{I}_m B = B \mathbf{I}_n = B.$$

The multiplication of an identity matrix by itself results in a same identity matrix.

In MATLAB identity matrices are created with eye function, which can take either one or two input arguments.

$$>> I = eye(n); I = eye(m, n)$$

Definition 3.10 (Transpose Matrix)

The transpose of a matrix A, which is a new matrix formed by interchanging the rows and columns of the original matrix. If the original matrix A is of order $m \times n$, then the transpose matrix, as A^T , will be of order $n \times m$, that is

If
$$A = (a_{ij})$$
, for $i = 1, 2, ..., m$ and $j = 1, 2, ..., n$,

then

$$A^{T} = (a_{ji}), \text{ for } i = 1, 2, ..., n \text{ and } j = 1, 2, ..., m.$$

The transpose of a matrix A can be found by using MATLAB command as follows:

$$>> A = [1 \ 2 \ 3; \ 4 \ 5 \ 6; \ 7 \ 8 \ 9]; B = A'$$

It is to be noted that

1.
$$(A^T)^T = A$$

$$(A^T)^T = A$$
 2. $(A_1 + A_2)^T = A_1^T + A_2^T$

3.
$$(A_1A_2)^T = A_2^T A_1^T$$

$$(A_1A_2)^T = A_2^T A_1^T$$
 4. $(\alpha A)^T = \alpha A^T$, α is a scalar.

Definition 3.11 (Inverse Matrix)

An $n \times n$ matrix A has an inverse or is invertible if there exists an $n \times n$ matrix B such that

$$AB = BA = I_n$$
.

Then the matrix B is called the inverse of A and is denoted by A^{-1} . For example, let

$$A = \begin{pmatrix} 2 & 3 \\ 2 & 2 \end{pmatrix}$$
 and $B = \begin{pmatrix} -1 & \frac{3}{2} \\ 1 & -1 \end{pmatrix}$.

Then we have

$$AB = BA = \mathbf{I}_2,$$

which means that B is an inverse of A. The invertible matrix is also called, nonsingular matrix. •

To find the inverse of the square matrix A using MATLAB commands we do as follows:

$$>> A = [2 \ -1 \ 0 \ 0; -1 \ 2 \ -1 \ 0; \ 0 \ -1 \ 2 \ -1; \ 0 \ 0 \ -1 \ 2]; \ Ainv = INVMAT(A)$$

Program 3.1

MATLAB m-file for finding inverse of a matrix

function [Ainv]=INVMAT(A)

```
[n,n]=size(A); I=zeros(n,n); for i=1:n; I(i,i)=1; end; m(1:n,1:n)=A; m(1:n,n+1:2*n)=I; for i=1:n; m(i,1:2*n)=m(i,1:2*n)/m(i,i); for k=1:n; if i^*=k m(k,1:2*n)=m(k,1:2*n)-m(k,i)*m(i,1:2*n); end; end; end; invrs=m(1:n,n+1:2*n);
```

MATLAB built-in function inv(A) can be also use to calculate the inverse of a square matrix A if A is invertible. If the matrix A is not invertible, then the matrix A is called singular. There are some well-known properties of the invertible matrix which are define as follows.

Theorem 3.2 If the matrix A is invertible, then

- 1. It has exactly one inverse. If B and C are the inverses of A, then B = C.
- 2. Its inverse matrix A^{-1} is also invertible and $(A^{-1})^{-1} = A$.
- 3. Its product with another invertible matrix is invertible, and the inverse of the product is the product of the inverses in the reverse order. If A and B are invertible matrices of the same size, then AB is invertible and $(AB)^{-1} = B^{-1}A^{-1}$.
- 4. Its transpose matrix A^T is invertible and $(A^T)^{-1} = (A^{-1})^T$.
- 5. The kA for any non-zero k is invertible, that is, $(kA)^{-1} = \frac{1}{k}A^{-1}$.
- 6. The A^k for any k is also invertible, that is, $(A^k)^{-1} = (A^{-1})^k$.
- 7. Its size 1×1 is invertible when it is nonzero. If A = (a), then $A^{-1} = (\frac{1}{a})$.
- 8. The formula for A^{-1} when n=2 is

$$A^{-1} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}^{-1} = \frac{1}{a_{11}a_{22} - a_{12}a_{21}} \begin{pmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{pmatrix},$$

provided that $a_{11}a_{22} - a_{12}a_{21} \neq 0$.

Definition 3.12 (Diagonal Matrix)

It is a square matrix having all elements equal to zero except those on main diagonal, that is

$$A = (a_{ij}) = \begin{cases} a_{ij} = 0, & \text{if } i \neq j, \\ a_{ij} \neq 0, & \text{if } i = j. \end{cases}$$

Note that all diagonal matrices are invertible if all diagonal entries are nonzero.

The MATLAB diag function is used to either create a diagonal matrix from a vector or it extract the diagonal entries of a matrix. If the input argument of the diag function is a vector, MATLAB uses the vector to create a diagonal matrix:

>>
$$x = [2, 2, 2]; A = diag(x)$$

>> $B = [2 -4 \ 1; 6 \ 10 -3; 0 \ 5 \ 8]; M = diag(B)$

The matrix A is called the *scalar* matrix because it has all the elements on the main diagonal equal to the same scalars 2. Multiplication of a square matrix and a scalar matrix is commutative, and the product is also a diagonal matrix.

Definition 3.13 (Upper-Triangular Matrix)

It is a square matrix which has zero elements below and to the left of the main diagonal. The diagonal as well as the above diagonal elements can take on any value, that is

$$U = (u_{ij}), \quad where \quad u_{ij} = 0, \quad if \quad i > j.$$

An example of such a matrix is

$$U = \left(\begin{array}{rrr} 1 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \end{array}\right).$$

The upper-triangular matrix is called upper-unit-triangular matrix if the diagonal elements are equal to one. This type of matrix is used in solving linear algebraic equations by LU decomposition with Crout's method. Also, if the main diagonal elements of the upper-triangular matrix are zero, then

$$A = \left(\begin{array}{ccc} 0 & a_{12} & a_{13} \\ 0 & 0 & a_{23} \\ 0 & 0 & 0 \end{array}\right),$$

is called the strictly upper-triangular matrix. This type of matrix will be used in solving linear systems by iterative methods.

Using MATLAB command triu(A) we can create an upper triangular matrix from a matrix A as

$$>> A = [1\ 2\ 3; 4\ 5\ 6; 7\ 8\ 9];\ U = triu(A)$$

Also we can create strictly upper-triangular matrix, that is, an upper-triangular matrix with zero diagonal, from a given matrix A by using MATLAB built-in function $\mathbf{triu}(A,I)$ as follows:

$$>> A = [1\ 2\ 3; 4\ 5\ 6; 7\ 8\ 9];\ U = triu(A, I)$$

Definition 3.14 (Lower-Triangular Matrix)

It is a square matrix which has zero elements above and to the right of the main diagonal and the rest of the elements can take on any value, that is

$$L = (l_{ij}), \quad where \quad l_{ij} = 0, \quad if \quad i < j.$$

An example of such a matrix is

$$L = \left(\begin{array}{ccc} 2 & 0 & 0 \\ 3 & 1 & 0 \\ 4 & 5 & 3 \end{array}\right).$$

The lower-triangular matrix is called lower-unit-triangular matrix if the diagonal elements are equal to one. This type of matrix is used in solving linear algebraic equations by LU Decomposition with Doolittle's method. Also, if the main diagonal elements of the lower-triangular matrix are zero, then the matrix

$$A = \left(\begin{array}{ccc} 0 & 0 & 0\\ a_{21} & 0 & 0\\ a_{31} & a_{32} & 0 \end{array}\right),$$

is called the strictly lower-triangular matrix. We will use this type of matrix in solving the linear systems by using iterative methods.

In similar way like upper-triangular matrices we can create lower-triangular matrix and strictly lower-triangular matrix from a given matrix A by using MATLAB built-in functions $\mathbf{tril}(A)$ and $\mathbf{tril}(A,I)$ respectively.

Note that all the triangular matrices (upper or lower) with nonzero diagonal entries are invertible.

Definition 3.15 (Symmetric Matrix)

A symmetric matrix is one in which the elements a_{ij} of a matrix A, in the ith row and jth column equal to the element a_{ji} in the jth row and ith column which means that

$$A^T = A$$
, that is $a_{ij} = a_{ji}$, for $i \neq j$.

Note that any diagonal matrix, including the identity, is symmetric. A lower- or upper-triangular matrix is symmetric if and only if it is, in fact, a diagonal matrix.

One way to generate a symmetric matrix is to multiply a matrix by its transpose, since A^TA is symmetric for any A. To generate a symmetric matrix using MATLAB commands we do as

$$>> A = [1:4;5:8;9:12]; B = A'*A; C = A*A$$

Example 3.1 Find all the values of a, b and c for which the following matrix is symmetric:

$$A = \begin{pmatrix} 4 & a+b+c & 0 \\ -1 & 3 & b-c \\ -a+2b-2c & 1 & b-2c \end{pmatrix}.$$

Solution. If the given matrix is symmetric, then $A = A^T$, that is

$$A = \begin{pmatrix} 4 & a+b+c & 0 \\ -1 & 3 & b-c \\ -a+2b-2c & 1 & b-2c \end{pmatrix} = \begin{pmatrix} 4 & -1 & -a+2b-2c \\ a+b+c & 3 & 1 \\ 0 & b-c & b-2c \end{pmatrix} = A^{T},$$

which implies that

$$0 = -a + 2b - 2c$$
, $-1 = a + b + c$, $1 = b - c$.

Solving above system, we get, a=2, b=-1, c=-2 and using these values, we have the given matrix of the form

 $A = \left(\begin{array}{rrr} 4 & -1 & 0 \\ -1 & 3 & 1 \\ 0 & 1 & 3 \end{array} \right).$

Theorem 3.3 If A and B are symmetric matrices with same size, and if k is any scalar, then

- 1. A^T is also symmetric.
- 2. A + B and A B are symmetric.
- 3. kA is also symmetric.

Note that product of symmetric matrices is not symmetric in general but the product is symmetric if and only if the matrices commute. Also, note that if A is a square matrix, then the matrices A, AA^T and A^TA are either all nonsingular or all singular.

If for a matrix A, the $a_{ij} = -a_{ji}$ for $i \neq j$ and the main diagonal elements are not all zero, then the matrix A is called skew matrix. If all the elements on the main diagonal of a skew matrix are zero, then the matrix is called skew symmetric, that is

$$A = -A^T$$
, with $a_{ij} = -a_{ji}$, for $i \neq j$ and $a_{ii} = 0$.

Any square matrix may be split into the sum of a symmetric and a skew symmetric matrix. Thus

$$A = \frac{1}{2}(A + A^{T}) + \frac{1}{2}(A - A^{T}),$$

where $\frac{1}{2}(A+A^T)$ is symmetric matrix and $\frac{1}{2}(A-A^T)$ is skew symmetric matrix. The following matrices

$$\left(\begin{array}{ccc} 1 & 2 & 3 \\ 2 & 4 & 5 \\ 3 & 5 & 6 \end{array}\right), \quad \left(\begin{array}{ccc} 1 & 2 & 3 \\ -2 & 4 & -5 \\ -3 & 5 & 6 \end{array}\right), \quad \left(\begin{array}{ccc} 0 & 2 & 3 \\ -2 & 0 & 5 \\ -3 & 5 & 0 \end{array}\right),$$

are the examples of symmetric, skew and skew symmetric matrices respectively.

Definition 3.16 (Band Matrix)

A $n \times n$ square matrix A is called a band matrix if there exists positive integers p and q, with 1 < p and q < n such that

$$a_{ij} = 0$$
 for $p \le j - i$ or $q \le i - j$.

The number p describes the number of diagonals above, and including, the main diagonal on which nonzero entries may lie. The number q describes the number of diagonals below, and including, the

main diagonal on which nonzero entries may lie. The number p + q - 1 is called the bandwidth of the matrix A, which tells us how many of the diagonals can contain nonzero entries. For example, the following matrix

$$A = \left(\begin{array}{cccc} 1 & 2 & 3 \\ 2 & 3 & 4 & 5 \\ 0 & 5 & 6 & 7 \\ 0 & 0 & 7 & 8 \end{array}\right),$$

is banded with p=3 and q=2, and so the bandwidth is equal to 4. An important property of the band matrix is called the tridiagonal matrix, in this case p=q=2, that is, all nonzero elements lie either on or directly above or below the main diagonal. For such type of matrix, the Gaussian elimination is particular simpler. In general, the nonzero elements of a tridiagonal matrix lie in three bands: the superdiagonal, diagonal and subdiagonal. For example, the following matrix

$$A = \begin{pmatrix} 1 & 2 & & & \\ 2 & 3 & 1 & & & \\ & 3 & 2 & 1 & & \\ & 2 & 4 & 3 & & \\ & & 1 & 2 & 3 & \\ & & & 1 & 6 & 4 \\ & & & & 3 & 4 \end{pmatrix},$$

is a tridiagonal matrix. A matrix which is predominantly zero is called a sparse matrix. A band matrix or a tridiagonal matrix is a sparse matrix but the nonzero elements of a sparse matrix are not necessarily near the diagonal.

3.2.3 The Determinant of Matrix

The determinant is a certain kind of a function that associates a real number with a square matrix. We will denote the determinant of a square matrix A by det(A) or |A|.

Definition 3.17 (Determinant of Matrix)

Let $A = (a_{ij})$ be an $n \times n$ square matrix then a determinant of A is given by:

1.
$$\det(A) = a_{11},$$
 if $n = 1.$

2.
$$\det(A) = a_{11}a_{22} - a_{12}a_{21}$$
, if $n = 2$.

For example, if

$$A = \begin{pmatrix} 4 & 2 \\ -3 & 7 \end{pmatrix} \quad \text{and} \quad B = \begin{pmatrix} 6 & 3 \\ 2 & 5 \end{pmatrix},$$

then

$$det(A) = (4)(7) - (-3)(2) = 34$$
 and $det(B) = (6)(5) - (3)(2) = 24$.

Notice that the determinant of a 2×2 matrix is given by the difference of the products of the two diagonals of a matrix. The determinant of a 3×3 matrix is defined in terms of determinants of 2×2 matrices and the determinant of a 4×4 matrix is defined in terms of determinants of 3×3

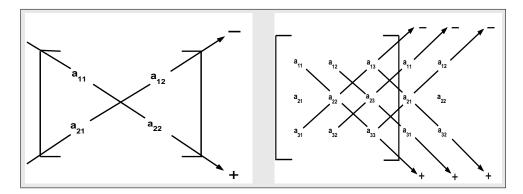


Figure 3.2: Direct evaluation of 2×2 and 3×3 determinants.

matrices and so on.

MATLAB function det(A) calculated the determinant of the square matrix A as:

$$>> A = [2\ 2; 6\ 7];\ B = det(A)$$

Other way to find the determinants of only 2×2 and 3×3 matrices can be found easily and quickly using diagonals (or direct evaluation). For 2×2 matrix, the determinant can be obtained by forming the product of the entries on the line from left to right and subtracting from this number the product of the entries on the line from right to left. For a matrix of size 3×3 , the diagonals of an array consisting of the matrix with the two first columns added to the right are used. Then the determinant can be obtained by forming the sum of the products of the entries on the lines from left to right, and subtract from this number the products of the entries on the lines from right to left, as shown in Figure 3.2.

Thus for 2×2 matrix

$$|A| = a_{11}a_{22} - a_{12}a_{21},$$

and for 3×3 matrix

$$|A| = a_{11}a_{22}a_{33} + a_{12}a_{23}a_{31} + a_{13}a_{21}a_{32} - a_{13}a_{22}a_{31} - a_{11}a_{23}a_{32} - a_{12}a_{21}a_{33}$$
 (diagonal products from left to right) (diagonal products from right to left)

For example, the determinant of 2×2 matrix can be computed as

$$|A| = \begin{vmatrix} 12 & 5 \\ -7 & 6 \end{vmatrix} = (12)(6) - (5)(-7) = 72 + 35 = 107,$$

and the determinant of 3×3 matrix can be obtained as

$$|A| = \begin{vmatrix} 4 & 5 & 6 \\ -3 & 8 & 2 \\ 4 & 9 & 7 \end{vmatrix} = [(4)(8)(7) + (5)(2)(4) + (6)(-3)(9)] - [(6)(8)(4) + (4)(2)(9) + (5)(-3)(7)] = 102 - 159 = -57.$$

For finding the determinants of the higher-order matrices, we will define the following concepts of *minor* and *cofactor* of the matrices.

Definition 3.18 (Minors of a Matrix)

The minor M_{ij} of all elements a_{ij} of a matrix A of order $n \times n$ as the determinant of the sub-matrix of order $(n-1) \times (n-1)$ obtained from A by deleting the ith row and jth column (also called ijth minor of A). For example, let

$$A = \left(\begin{array}{ccc} 2 & 3 & -1 \\ 5 & 3 & 2 \\ 4 & -2 & 4 \end{array}\right),$$

then, the minor M_{11} will be obtained by deleting the first row and the first column of the given matrix A, that is,

$$M_{11} = \begin{vmatrix} 3 & 2 \\ -2 & 4 \end{vmatrix} = (3)(4) - (-2)(2) = 12 + 4 = 16.$$

Similarly, we can find the other possible minors of the given matrix as follows:

$$M_{12} = 12$$
, $M_{13} = -22$, $M_{21} = 10$, $M_{22} = 12$,

and

$$M_{23} = -16$$
, $M_{31} = 9$, $M_{32} = 9$, $M_{33} = -9$,

which are the required minors of the given matrix.

Definition 3.19 (Cofactor of a Matrix)

The cofactor A_{ij} of all elements a_{ij} of a matrix A of order $n \times n$ is given by

$$A_{ij} = (-1)^{i+j} M_{ij},$$

where M_{ij} is the minor of all elements a_{ij} of a matrix A. For example, the cofactor A_{11} of the elements a_{11} of the matrix

$$A = \left(\begin{array}{ccc} 2 & 3 & -1 \\ 5 & 3 & 2 \\ 4 & -2 & 4 \end{array}\right),$$

is computed as follows

$$A_{11} = (-1)^{1+1} M_{11} = M_{11} = 16.$$

Similarly, for other elements, we have

$$A_{12} = -12, \ A_{13} = -22, \ A_{21} = -10, \ A_{22} = 12,$$

$$A_{23} = 16, \ A_{31} = 9, \ A_{32} = -9, \ A_{33} = -9.$$

Program 3.2

MATLAB m-file for finding minors and cofactors of a matrix function CofA = cofactor(A,i,j)

[m,n] = size(A); if m = n error('Matrix must be square') end

$$A1 = A([1:i-1,i+1:n],[1:j-1,j+1:n]); Minor = det(A1); CofA = (-1)^ (i+j)*det(Minor);$$

Definition 3.20 (Cofactor Expansion of Determinant of a Matrix)

Let A be a square matrix, then we define determinant of A is the sum of the products of the elements of the first row and their cofactors. If A is 3×3 matrix, then its determinant is define as

$$\det(A) = |A| = a_{11}A_{11} + a_{12}A_{12} + a_{13}A_{13}.$$

Similarly, more general for $n \times n$ matrix, we define as

$$\det(A) = |A| = \sum_{1}^{n} a_{ij} A_{ij}, \quad n > 2,$$
(3.12)

where summation is on i for any fixed value of jth column $(1 \le j \le n)$, or on j for any fixed value of ith row $(1 \le i \le n)$ and A_{ij} is the cofactor of element a_{ij} .

Example 3.2 Find the minors and cofactors of the matrix A and use it to evaluate the determinant of the matrix

$$A = \left(\begin{array}{ccc} 3 & 1 & -4 \\ 2 & 5 & 6 \\ 1 & 4 & 8 \end{array}\right).$$

Solution. The minors of A are calculated as follows

$$M_{11} = 16, \qquad M_{12} = 10, \qquad M_{13} = 3,$$

and from these values of the minors, we can calculate the cofactors of the elements of the given matrix as follows

$$A_{11} = 16, \qquad A_{12} = -10, \qquad A_{13} = 3.$$

Now by using the cofactor expansion along the first row, we can find the determinant of the matrix as follows

$$\det(A) = a_{11}A_{11} + a_{12}A_{12} + a_{13}A_{13} = (3)(16) + (1)(-10) + (-4)(3) = 26.$$

Note that in the above Example 3.2, we computed the determinant of the matrix by using the cofactor expansion along the first row but it can also be found along the first column of the matrix. To get the results of the Example 3.2, we use the MATLAB command window as follows:

$$>> A = [3\ 1\ -4; 2\ 5\ 6; 1\ 4\ 8]; Det A = Cof Fexp(A);$$

Program 3.3

MATLAB m-file for determinant of a matrix by cofactor expansion function DetA = CofFexp(A)

$$[m,n] = size(A);$$
 if $m = n$ error('Matrix must be square') end; $a = A(1,:); c = [$]; for $i=1:n;$ $c1i = cofactor(A,1,i);$ $c = [c;c1i];$ end; $DetA = a*c;$

Theorem 3.4 (The Laplace Expansion Theorem)

The determinant of an $n \times n$ matrix $A = \{a_{ij}\}$, when $n \geq 2$, can be computed as

$$\det(A) = a_{i1}A_{i1} + a_{i2}A_{i2} + \dots + a_{in}A_{in} = \sum_{j=1}^{n} a_{ij}A_{ij},$$

which is called the cofactor expansion along the ith row and also as

$$\det(A) = a_{1j}A_{1j} + a_{2j}A_{2j} + \dots + a_{nj}A_{nj} = \sum_{i=1}^{n} a_{ij}A_{ij},$$

is called cofactor expansion along jth column. It is called Laplace Expansion Theorem.

Note that the cofactor and minor of an element a_{ij} differs only in sign, that is, $A_{ij} = \pm M_{ij}$. A quick way for determining whether to use the + or - is to use the fact that the sign relating A_{ij} and M_{ij} is in the *ith* row and *jth* column of the *checkerboard* array

$$\begin{pmatrix} + & - & + & - & + & \cdots \\ - & + & - & + & - & \cdots \\ + & - & + & - & + & \cdots \\ - & + & - & + & - & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix}.$$

For example, $A_{11} = M_{11}$, $A_{21} = -M_{21}$, $A_{12} = -M_{12}$, $A_{22} = M_{22}$ and so on.

Definition 3.21 (Cofactor Matrix)

If A is any $n \times n$ matrix and A_{ij} is the cofactor of a_{ij} , then the matrix

$$Cof(A) = \begin{pmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \cdots & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nn} \end{pmatrix},$$

is called the matrix of cofactor from A. For example, the cofactor of the matrix

$$A = \left(\begin{array}{ccc} 3 & 2 & -1 \\ 1 & 6 & 3 \\ 2 & -4 & 0 \end{array}\right),$$

can be calculated as follows:

$$A_{11} = 12$$
, $A_{12} = 6$, $A_{13} = -16$, $A_{21} = 4$, $A_{22} = 2$, $A_{23} = 16$, $A_{31} = 12$, $A_{32} = -10$, $A_{33} = 16$.

So that the matrix

$$Cof(A) = \begin{pmatrix} 12 & 6 & -16 \\ 4 & 2 & 16 \\ 12 & -10 & 16 \end{pmatrix}.$$

is the required of the cofactor matrix.

Definition 3.22 (Adjoint of a Matrix)

If A is any $n \times n$ matrix and A_{ij} is the cofactor of a_{ij} of A, then the transpose of this matrix is called the adjoint of A and is denoted by Adj(A). For example, the cofactor matrix of the following matrix

$$A = \left(\begin{array}{ccc} 3 & 2 & -1 \\ 1 & 6 & 3 \\ 2 & -4 & 0 \end{array}\right),$$

is calculated as

$$Cof(A) = \begin{pmatrix} 12 & 6 & -16 \\ 4 & 2 & 16 \\ 12 & -10 & 16 \end{pmatrix}.$$

So by taking its transpose, we get the matrix

$$\begin{pmatrix} 12 & 6 & -16 \\ 4 & 2 & 16 \\ 12 & -10 & 16 \end{pmatrix}^T = \begin{pmatrix} 12 & 4 & 12 \\ 6 & 2 & -10 \\ -16 & 16 & 16 \end{pmatrix} = Adj(A),$$

which is called the adjoint of the given matrix A.

Example 3.3 Find the determinant of the following matrix using cofactor expansion and show that det(A) = 0 when x = 4

$$A = \left(\begin{array}{ccc} x+2 & x & 2\\ 1 & x-1 & 3\\ 4 & x+1 & x \end{array}\right).$$

Solution. Using the cofactor expansion along the first row, we compute the determinant of the given matrix as follows:

$$|A| = a_{11}A_{11} + a_{12}A_{12} + a_{13}A_{13},$$

where

$$A_{11} = M_{11} = x^2 - 4x - 3, \ A_{12} = -M_{12} = -x + 12, \ A_{13} = -3x + 5.$$

Thus

$$|A| = (x+2)[x^2 - 4x - 3] + x[-x+12] + 2[-3x+5] = x^3 - 3x^2 - 5x + 4.$$

Now taking x = 4, we get

$$|A| = (4)^3 - 3(4)^2 - 5(4) + 4 = 64 - 48 - 20 + 4 = 0,$$

which is the required determinant of the matrix at x = 4.

The following are special properties which will be helpful in reducing the amount of work involved in evaluating determinants.

Theorem 3.5 (Properties of the Determinant)

Let A be an $n \times n$ matrix:

1. The determinant of a matrix A is zero if any row or column is zero or equal to a linear combination of other rows and columns.

For example, if

$$A = \left(\begin{array}{ccc} 3 & 1 & 0 \\ 2 & 1 & 0 \\ 4 & 3 & 0 \end{array}\right),$$

then det(A) = 0.

2. A determinant of a matrix A is changed in sign if the two rows or two columns are interchange. For example, if

$$A = \left(\begin{array}{cc} 3 & 2 \\ 4 & 5 \end{array}\right),$$

then det(A) = 7, but for the matrix

$$B = \left(\begin{array}{cc} 4 & 5 \\ 3 & 2 \end{array}\right),$$

obtained from the matrix A by interchanging its rows, we have det(B) = -7.

3. The determinant of a matrix A is equal to the determinant of its transposed. For example, if

$$A = \left(\begin{array}{cc} 5 & 3 \\ 4 & 4 \end{array}\right),$$

then det(A) = 8, and for the matrix

$$B = \left(\begin{array}{cc} 5 & 4 \\ 3 & 4 \end{array}\right),$$

obtained from the matrix A by taking its transpose, we have

$$det(B) = 8 = det(A)$$
.

5. If the matrix B is obtained from the matrix A by multiplying every element in one row or in one column by k, then determinant of the matrix B is equal to k times the determinant of A. For example, if

$$A = \left(\begin{array}{cc} 6 & 5 \\ 3 & 4 \end{array}\right),$$

then det(A) = 9, but for the matrix

$$B = \left(\begin{array}{cc} 12 & 10\\ 3 & 4 \end{array}\right),$$

obtained from the matrix A by multiplying its first row by 2, we have

$$det(B) = 18 = 2(9) = 2 det(A).$$

6. If the matrix B is obtained from the matrix A by adding to a row (or a column) of a multiple of another row (or another column) of A, then determinant of the matrix B is equal to the determinant of A. For example, if

$$A = \left(\begin{array}{cc} 4 & 3 \\ 5 & 4 \end{array}\right)$$

then det(A) = 1, and for the matrix

$$B = \left(\begin{array}{cc} 4 & 3\\ 13 & 10 \end{array}\right),$$

obtained from the matrix A by adding to its second row 2 times the first row, we have

$$\det(B) = 1 = \det(A).$$

7. If two rows or two columns of a matrix A are identical, then the determinant is zero. For example, if

$$A = \begin{pmatrix} 2 & 3 \\ 2 & 3 \end{pmatrix}$$
, then $\det(A) = 0$.

8. The determinant of a product of matrices is the product of the determinants of all matrices. For example, if

$$A = \begin{pmatrix} 3 & 4 & 5 \\ 3 & 2 & 1 \\ 2 & 1 & 6 \end{pmatrix} \quad and \quad B = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 2 & 3 \\ 1 & 3 & 5 \end{pmatrix},$$

then det(A) = -36 and det(A) = -3. Also,

$$AB = \left(\begin{array}{ccc} 24 & 29 & 46\\ 12 & 13 & 20\\ 12 & 24 & 39 \end{array}\right),$$

then det(AB) = 108. Thus

$$\det(A)\det(B) = (-36)(-3) = 108 = \det(AB).$$

9. The determinant of a triangular matrix (upper-triangular or lower-triangular matrix) is equal to the product of all their main diagonal elements. For example, if

$$A = \begin{pmatrix} 3 & 4 & 5 \\ 0 & 4 & 7 \\ 0 & 0 & 5 \end{pmatrix}, \quad then \quad \det(A) = (3)(4)(5) = 60.$$

10. The determinant of an $n \times n$ matrix A times scalar multiple k equal to k^n times the determinant of the matrix A, that is $\det(kA) = k^n \det(A)$. For example, if

$$A = \left(\begin{array}{rrr} 3 & 4 & 5 \\ 2 & 3 & 6 \\ 1 & 0 & 5 \end{array}\right),$$

then det(A) = 14, and for the matrix

$$B = 2A = \left(\begin{array}{ccc} 6 & 8 & 10 \\ 4 & 6 & 12 \\ 2 & 0 & 10 \end{array}\right),$$

obtained from the matrix A by multiply by 2, we have

$$\det(B) = 112 = 8(14) = 2^3 \det(A).$$

11. The determinant of the kth power of a matrix A equal to the kth power of the determinant of the matrix A, that is $\det(A^k) = (\det(A))^k$. For example, if

$$A = \left(\begin{array}{ccc} 2 & -2 & 0 \\ 2 & 3 & -1 \\ 1 & 0 & 1 \end{array}\right),$$

then det(A) = 12, and for the matrix

$$B = A^3 = \begin{pmatrix} -18 & -30 & 12 \\ 24 & -3 & -9 \\ 3 & -12 & 3 \end{pmatrix},$$

obtained by taking cubic power of the matrix A, we have

$$\det(B) = 1728 = (12)^3 = (\det(A))^3.$$

12. The determinant of a scalar matrix (1×1) is equal to the element itself. For example, if A = (8), then det(A) = 8.

Example 3.4 Find all the values of α for which det(A) = 0, where

$$A = \left(\begin{array}{ccc} \alpha - 3 & 1 & 0 \\ 0 & \alpha - 1 & 1 \\ 0 & 2 & \alpha \end{array}\right).$$

Solution. We find the determinant of the given matrix by using the cofactor expansion along the first row, so we compute

$$|A| = a_{11}A_{11} + a_{12}A_{12} + a_{13}A_{13}$$

$$= (\alpha - 3) \begin{vmatrix} \alpha - 1 & 1 \\ 2 & \alpha \end{vmatrix} - 1 \begin{vmatrix} 0 & 1 \\ 0 & \alpha \end{vmatrix} + 0 \begin{vmatrix} 0 & \alpha - 1 \\ 0 & 2 \end{vmatrix}$$

$$= (\alpha - 3)(\alpha + 1)(\alpha - 2).$$

Given det(A) = 0, implies that, alpha = -1, 2, 3, the required values of α .

Theorem 3.6 If A is an invertible matrix, then

1.
$$\det(A) \neq 0 \qquad 2. \qquad \det(A^{-1}) = \frac{1}{\det(A)} \qquad 3. \qquad A^{-1} = \frac{Adj(A)}{\det(A)}.$$

4.
$$(adj(A))^{-1} = \frac{A}{\det(A)} = adj(A^{-1})$$
 5. $\det(adj(A)) = \det(A)^{n-1}$.

By using the Theorem 3.6 we can find the inverse of a matrix by showing that determinant of a matrix not equal to zero and by using adjoint and determinant of the given matrix A.

Example 3.5 For what values of α the following matrix has an inverse:

$$A = \left(\begin{array}{ccc} 1 & 0 & \alpha \\ 2 & 2 & 1 \\ 0 & 2\alpha & 1 \end{array}\right).$$

Solution. We find the determinant of the given matrix by using the cofactor expansion along the first row as follows:

$$|A| = a_{11}A_{11} + a_{12}A_{12} + a_{13}A_{13},$$

which is equal to

$$|A| = (1)A_{11} + (0)A_{12} + (\alpha)A_{13} = A_{11} + \alpha A_{13} = 2 - 2\alpha + 4\alpha^2.$$

From the Theorem 3.6 we know that the matrix has an inverse if $det(A) \neq 0$, so

$$|A| = 2 - 2\alpha + 4\alpha^2 = 2(2\alpha + 1)(\alpha - 1) \neq 0.$$

Hence the given matrix has an inverse if $\alpha \neq -1/2$ and $\alpha \neq 1$.

Example 3.6 Use the adjoint method to compute the inverse of the the following matrix

$$A = \left(\begin{array}{ccc} 1 & 2 & -1 \\ 2 & -1 & 1 \\ 1 & 2 & 2 \end{array}\right).$$

Also, find the inverse and determinant of the adjoint matrix.

Solution. First we compute the determinant of the given matrix as follows:

$$\det(A) = |A| = a_{11}A_{11} + a_{12}A_{12} + a_{13}A_{13} = (1)(-4) - (2)(3) + (-1)(5) = -15,$$

and the compute the nine cofactors as follows:

$$A_{11} = -4$$
, $A_{12} = -3$, $A_{13} = 5$, $A_{21} = -6$, $A_{22} = 3$, $A_{23} = 0$, $A_{31} = 1$, $A_{32} = -3$, $A_{33} = -5$.

Thus we have the cofactor matrix and the adjoint matrix as follows

$$Cof(A) = \begin{pmatrix} -4 & -3 & 5 \\ -6 & 3 & 0 \\ 1 & -3 & -5 \end{pmatrix}, \ adj(A) = \begin{pmatrix} -4 & -3 & 5 \\ -6 & 3 & 0 \\ 1 & -3 & -5 \end{pmatrix}^{T} = \begin{pmatrix} -4 & -6 & 1 \\ -3 & 3 & -3 \\ 5 & 0 & -5 \end{pmatrix}.$$

To get adjoint of the matrix of the Example 3.6, we use MATLAB command window as:

$$>> A = [1 \ 2 \ -1; 2 \ -1 \ 1; 1 \ 2 \ 2]; \ AdjA = Adjoint(A);$$

Program 3.4

MATLAB m-file for adjoint of a matrix

 $function \ AdjA = Adjoint(A)$

 $[m,n]=size(A); \ \ if \ m\ \ \ =\ n\ \ error(\ Matrix\ must\ be\ square\ ')\ \ end;\ A1=[\]; \ for\ i=1:n;\ for\ j=1:n;\ A1=[A1;cofactor(A,i,j)];end;end;\ AdjA=reshape(A1,n,n);$

Then by using the Theorem 3.6 we can have the inverse of the matrix as follows:

$$A^{-1} = \frac{Adj(A)}{\det(A)} = -\frac{1}{15} \begin{pmatrix} -4 & -6 & 1\\ -3 & 3 & -3\\ 5 & 0 & -5 \end{pmatrix} = \begin{pmatrix} 4/15 & 2/5 & -1/15\\ 1/5 & -1/5 & 1/5\\ -1/3 & 0 & 1/3 \end{pmatrix}.$$

Using the Theorem 3.6 we can compute the inverse of the adjoint matrix as follows:

$$(adj(A))^{-1} = \frac{A}{\det(A)} = \begin{pmatrix} -1/15 & -2/15 & 1/15 \\ -2/15 & 1/15 & -1/15 \\ -1/15 & -2/15 & -2/15 \end{pmatrix},$$

and $det(adj(A)) = (det(A))^{3-1} = (-15)^2 = 225$.

Example 3.7 If det(A) = 3 and det(B) = 4, then show that

$$\det(A^2B^{-1}A^TB^3) = 432.$$

Solution. By using the properties of the determinant of the matrix, we have

$$\det(A^2B^{-1}A^TB^3) = \det(A^2)\det(B^{-1})\det(A^T)\det(B^3),$$

which can be also written as

$$\det(A^2B^{-1}A^TB^3) = (\det(A))^2 \frac{1}{\det(B)} (\det(A)) (\det(B))^3.$$

Now using the given information, we get

$$det(A^2B^{-1}A^TB^3) = (3)^2 \frac{1}{4}(3)(4)^3 = 3^3 4^2 = 432,$$

the required solution.

3.2.4 Matrix Inversion Method

If matrix A is nonsingular, then the linear system (3.6) always has a unique solution for each **b**, since the inverse matrix A^{-1} exists, so the solution of the system (3.6) can formally expressed as

$$A^{-1}A\mathbf{x} = A^{-1}\mathbf{b}$$
, or $\mathbf{I}\mathbf{x} = A^{-1}\mathbf{b}$,

gives

$$\mathbf{x} = A^{-1}\mathbf{b}.\tag{3.13}$$

If A is a square invertible matrix, there exists a sequence of elementary row operations that carry A to the identity matrix **I** of the same size, that is, $A \longrightarrow \mathbf{I}$. This same sequence of row operations carries **I** to A^{-1} , that is, $\mathbf{I} \longrightarrow A^{-1}$. This can be also written as

$$[A|\mathbf{I}] \longrightarrow [\mathbf{I}|A^{-1}].$$

Example 3.8 Use matrix inversion method to find unique solution the linear system $A\mathbf{x} = \mathbf{b}$, where

$$A = \begin{pmatrix} 1 & 2 & 0 \\ -2 & 1 & 2 \\ -1 & 1 & 1 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}.$$

Solution. First we compute the inverse of the given matrix which has the form

$$A^{-1} = \left(\begin{array}{ccc} 1 & 2 & -4 \\ 0 & -1 & 2 \\ 1 & 3 & -5 \end{array}\right),$$

and then we can find unique solution of the given system as

$$\mathbf{x} = A^{-1}\mathbf{b} = \begin{pmatrix} 1 & 2 & -4 \\ 0 & -1 & 2 \\ 1 & 3 & -5 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} -1 \\ 1 \\ -1 \end{pmatrix},$$

the solution of the given system by the matrix inversion method.

Thus, when the matrix inverse A^{-1} of the coefficient matrix A is computed, the solution vector \mathbf{x} of linear system is simply the product of inverse matrix A^{-1} and the right-hand side vector \mathbf{b} .

Using MATLAB commands the linear system of equations defined by the coefficient matrix A and the right hand-side vector \mathbf{b} using matrix inverse method is solved with:

$$>> A = [1\ 2\ 0; -2\ 1\ 2; -1\ 1\ 1];\ b = [1; 1; 1];\ x = A \setminus b$$

Theorem 3.7 For an $n \times n$ matrix A, the following properties are equivalent:

- 1. The inverse of matrix A exists, that is, A is nonsingular.
- 2. The determinant of matrix A is nonzero.
- 3. The homogeneous system $A\mathbf{x} = \mathbf{0}$ has a trivial solution $\mathbf{x} = \mathbf{0}$.
- 4. The nonhomogeneous system $A\mathbf{x} = \mathbf{b}$ has a unique solution.

Not all matrices have inverses. Singular matrices don't have inverse and thus the corresponding system of equations does not have a unique solution. The inverse of a matrix can also be computed by using the following numerical methods for linear systems, called, Gauss-elimination method, Gauss-Jordan method and LU-decomposition method but the best and simplest method for finding the inverse of a matrix is to perform the Gauss-Jordan method on the augmented matrix with identity matrix of same size.

3.3 Solutions of Linear Systems of Equations

Now we shall discuss numerical methods for solving system of linear equations. We shall discuss both direct and indirect (iterative) methods for the solution of given linear systems. In direct method we shall discuss the familiar technique called the *method of elimination* to find the solution of linear systems. This method starts with the augmented matrix of the given linear system and obtain a matrix of a certain form. This new matrix represents a linear system that has exactly the same solutions as the given origin system. In indirect methods we shall discuss Jacobi and Gauss-Seidel methods.

The following basic theorems on the solvability of linear systems are proved in linear algebra.

Theorem 3.8 A homogeneous system of n equations in n unknowns has a solution other than the trivial solution if and only if the determinant of the coefficients matrix A vanishes, that is matrix A is singular.

Theorem 3.9 (Necessary and Sufficient Condition for a unique solution)

A nonhomogeneous system of n equations in n unknowns has a unique solution if and only if the determinant of a coefficients matrix A is not vanishes, that is, A is nonsingular.

Before, we discuss numerical methods for solving linear system, we introduce the most important numerical quantity associated with a matrix.

Definition 3.23 (Rank of a Matrix)

The rank of a matrix A is the number of pivots. An $m \times n$ matrix will, in general, have a rank r, where r is an integer and $r \leq \min\{m,n\}$. If $r = \min\{m,n\}$, then the matrix is said to be full rank. If $r < \min\{m,n\}$, then the matrix is said to be rank deficient.

In principle, the rank of a matrix can be determined by using the Gaussian elimination process in which the coefficient matrix A is reduced to upper-triangular form U. After reducing the matrix to triangular form, we find that the rank is the number of columns with nonzero values on the diagonal of U. In practice, especially for large matrices, round-off errors during the row operation may cause a loss of accuracy in this method of rank computation.

Theorem 3.10 For a system of n equations in n unknowns written in the form $A\mathbf{x} = \mathbf{b}$, then solution \mathbf{x} of a system exists and is unique for any \mathbf{b} if and only if rank(A) = n.

Conversely, if $\operatorname{rank}(A) < n$ for an $n \times n$ matrix A, then the system of equations $A\mathbf{x} = \mathbf{b}$ may or may not be consistent. Such a system may not have solution, or the solution, if it exists, will not be unique. For example, the rank of the following matrix is 3.

$$A = \left(\begin{array}{rrr} 1 & 2 & 4 \\ 1 & 1 & 5 \\ 1 & 1 & 6 \end{array}\right).$$

In MATLAB command, the built-in rank function can be use to estimate the rank of a matrix:

$$>> A = [1\ 2\ 4; 1\ 1\ 5; 1\ 1\ 6];\ rank(A)$$

Note that:

$$\operatorname{rank}(AB) \leq \min(\operatorname{rank}(A), \operatorname{rank}(B))$$

 $\operatorname{rank}(A+B) \leq \operatorname{rank}(A) + \operatorname{rank}(B)$
 $\operatorname{rank}(AA^T) = \operatorname{rank}(A) = \operatorname{rank}(A^TA)$

Although the rank of a matrix is very useful to categorize the behaviour of matrices and systems of equations, the rank of a matrix is usually not computed.

3.4 Direct Numerical Methods for Linear Systems

To solve the systems of linear equations using the numerical methods, there are two types of methods available. Methods of first type are called *direct methods* or *elimination methods*. The other type of the numerical methods are called *iterative methods*. In this chapter we will discuss both type of the numerical methods. The first type of methods find the solution in a finite number of steps. These methods are guaranteed to succeed and are recommended for general purpose. Here, we will consider *Gaussian elimination method* and its variants and *LU decomposition*, by Doolittle's and Crout's methods.

The direct method refers to a procedure for computing a solution from a form that is mathematically exact. We shall begin with simple method, called Gaussian elimination method and its variants and then continue with the methods involving triangular matrices, symmetric and tridiagonal matrices.

3.4.1 Gaussian Elimination Method

It is one of the most popular and widely used direct method for solving linear systems of algebraic equations. No method of solving linear systems requires fewer operations than the Gaussian procedure. The goal of the Gaussian elimination method for solving linear systems is to convert the original system into the equivalent upper-triangular system and from which each unknown is determined by backward substitution.

The Gaussian elimination procedure start with forward elimination, in which the first equation in the linear system is used to eliminate the first variable from the rest of (n-1) equations. Then the new second equation is used to elimination second variable from the rest of (n-2) equations, and so on. If (n-1) such elimination is performed then the resulting system will be the triangular form. Once this forward elimination is completed, we can determine whether the system is overdetermined or underdetermined or has a unique solution. If it has a unique solution, then the backward substitution is used to solve the triangular system easily and one can find the unknown variables involve in the system.

Now we shall describe the method in detail for a system of n linear equations. Consider the following a system of n linear equations:

$$a_{11}x_{1} + a_{12}x_{2} + a_{13}x_{3} + \cdots + a_{1n}x_{n} = b_{1}$$

$$a_{21}x_{1} + a_{22}x_{2} + a_{23}x_{3} + \cdots + a_{2n}x_{n} = b_{2}$$

$$a_{31}x_{1} + a_{32}x_{2} + a_{33}x_{3} + \cdots + a_{3n}x_{n} = b_{3}$$

$$\vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots$$

$$a_{n1}x_{1} + a_{n2}x_{2} \qquad a_{n3}x_{3} + \cdots + a_{nn}x_{n} = b_{n}$$

$$(3.14)$$

Forward Elimination

Consider first equation of the given system (3.14)

$$a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n = b_1,$$
 (3.15)

as first pivotal equation with first pivot element a_{11} . Then the first equation times multiples $m_{i1} = (a_{i1}/a_{11}), i = 2, 3, ..., n$, is subtracted from the *ith* equation to eliminate first variable x_1 , producing an equivalent system

$$a_{11}x_{1} + a_{12}x_{2} + a_{13}x_{3} + \cdots + a_{1n}x_{n} = b_{1}$$

$$a_{22}^{(1)}x_{2} + a_{23}^{(1)}x_{3} + \cdots + a_{2n}^{(1)}x_{n} = b_{2}^{(1)}$$

$$a_{32}^{(1)}x_{2} + a_{33}^{(1)}x_{3} + \cdots + a_{3n}^{(1)}x_{n} = b_{3}^{(1)}$$

$$\vdots \qquad \vdots \qquad \vdots \qquad \vdots \qquad \vdots$$

$$a_{n2}^{(1)}x_{2} + a_{n3}^{(1)}x_{3} + \cdots + a_{nn}^{(1)}x_{n} = b_{n}^{(1)}$$

$$(3.16)$$

Now consider a second equation of the system (3.16), which is

$$a_{22}^{(1)}x_2 + a_{23}^{(1)}x_3 + \dots + a_{2n}^{(1)}x_n = b_2^{(1)},$$
 (3.17)

as second pivotal equation with second pivot element $a_{22}^{(1)}$. Then the second equation times multiples $m_{i2} = (a_{i2}^{(1)}/a_{22}^{(1)}), i = 3, ..., n$, is subtracted from the *ith* equation to eliminate second variable x_2 , producing an equivalent system

$$a_{11}x_{1} + a_{12}x_{2} + a_{13}x_{3} + \cdots + a_{1n}x_{n} = b_{1}$$

$$a_{22}^{(1)}x_{2} + a_{23}^{(1)}x_{3} + \cdots + a_{2n}^{(1)}x_{n} = b_{2}^{(1)}$$

$$a_{33}^{(2)}x_{3} + \cdots + a_{3n}^{(2)}x_{n} = b_{3}^{(2)}$$

$$\vdots \qquad \vdots \qquad \vdots \qquad \vdots$$

$$a_{n2}^{(2)}x_{3} + \cdots + a_{nn}^{(2)}x_{n} = b_{n}^{(2)}$$

$$(3.18)$$

Now consider a third equation of the system (3.18), which is

$$a_{33}^{(2)}x_3 + \dots + a_{3n}^{(2)}x_n = b_3^{(2)},$$
 (3.19)

as the third pivotal equation with third pivot element $a_{33}^{(2)}$. Then the third equation times multiples $m_{i3} = (a_{i3}^{(2)}/a_{33}^{(2)}), i = 4, ..., n$, is subtracted from the *ith* equation to eliminate third variable x_3 . Similarly, after (n-1)th steps, we have the nth pivotal equation which have only one unknown variable x_n , that is

$$a_{11}x_{1} + a_{12}x_{2} + a_{13}x_{3} + \cdots + a_{1n}x_{n} = b_{1} + a_{22}^{(1)}x_{2} + a_{23}^{(1)}x_{3} + \cdots + a_{2n}^{(1)}x_{n} = b_{2}^{(1)} + a_{33}^{(2)}x_{3} + \cdots + a_{3n}^{(2)}x_{n} = b_{3}^{(2)} \vdots \vdots \vdots a_{nn}^{(n-1)}x_{n} = b_{n}^{(n-1)}$$

$$(3.20)$$

with nth pivotal element $a_{nn}^{(n-1)}$. After getting the upper-triangular system which is equivalent to the original system, the forward elimination is completed.

Backward Substitution

After the triangular set of equations has been obtained, the last equation of the system (3.20) yields the value of x_n directly. The value is then substituted into the equation next to the last one of the system (3.20) to obtain a value of x_{n-1} , which is, in turn, used along with the value of x_n in the second to the last equation to obtain a value of x_{n-2} , and so on. Mathematical formula can be obtain for the backward substitution

$$x_{n} = \frac{b_{n}^{(n-1)}}{a_{nn}^{(n-1)}}$$

$$x_{n-1} = \frac{1}{a_{n-1}^{(n-2)}} \left(b_{n-1}^{(n-2)} - a_{n-1n}^{(n-2)} x_{n} \right)$$

$$\vdots$$

$$x_{1} = \frac{1}{a_{11}} \left(b_{1} - \sum_{j=2}^{n} a_{1j} x_{j} \right)$$
(3.21)

The Gaussian elimination can be carried out by writing only the coefficients and the right-hand side terms in a matrix form, which means the augmented matrix form. Indeed, this is exactly what a computer program for the Gaussian elimination does. Even for hand calculation, the augmented matrix form is more convenient than writing all set of equations. The augmented matrix is formed as follows

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} & | & b_1 \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2n} & | & b_2 \\ a_{31} & a_{32} & a_{33} & \cdots & a_{3n} & | & b_3 \\ \vdots & \vdots & \vdots & \vdots & \vdots & | & & \\ a_{n1} & a_{n2} & a_{n3} & \cdots & a_{nn} & | & b_n \end{pmatrix}.$$

$$(3.22)$$

The operations used in the Gaussian elimination method can now be applied to the augmented matrix. Consequently system (3.20) is now written directly as follows:

$$\begin{pmatrix}
a_{11} & a_{12} & a_{13} & \cdots & a_{1n} & | & b_1 \\
& a_{21}^{(1)} & a_{23}^{(1)} & \cdots & a_{2n}^{(1)} & | & b_2^{(1)} \\
& & a_{33}^{(2)} & \cdots & a_{3n}^{(2)} & | & b_3^{(2)} \\
& & \vdots & \vdots & | & & & \\
& & & a_{nn}^{(n-1)} & | & b_n^{(n-1)}
\end{pmatrix},$$
(3.23)

from which the unknowns are determined as before by using backward substitution. The number of multiplications and divisions for the Gaussian elimination method for one \mathbf{b} vector is approximately

$$N = \left(\frac{n^3}{3}\right) + n^2 - \left(\frac{n}{3}\right). \tag{3.24}$$

Simple Gaussian Elimination Method

Firstly, we will solve the linear system using the simplest variation of the Gaussian elimination method, called the *simple* Gaussian elimination or the Gaussian elimination without pivoting. The basic of this variation is that all the possible diagonal elements (called pivot elements) should be nonzero. If at any stage it becomes zero, then interchange that row with any below row with nonzero element at that position. After getting upper-triangular matrix, we use backward substitution to get the solution of the given linear system.

Example 3.9 Solve the following linear system using the simple Gaussian elimination method

$$x_1 + 2x_2 + x_3 = 2$$

 $2x_1 + 5x_2 + 3x_3 = 1$
 $x_1 + 3x_2 + 4x_3 = 5$

solution. The process begins with the augmented matrix form

$$\left(\begin{array}{ccccc} 1 & 2 & 1 & \vdots & 2 \\ 2 & 5 & 3 & \vdots & 1 \\ 1 & 3 & 4 & \vdots & 5 \end{array}\right).$$

Since $a_{11} = 1 \neq 0$, so we wish to eliminate the elements a_{21} and a_{31} by subtracting from the second and third rows the appropriate multiples of the first row. In this case the multiples are $m_{21} = \frac{2}{1} = 2$ and $m_{31} = \frac{1}{1} = 1$. Hence

$$\left(\begin{array}{cccccc}
1 & 2 & 1 & \vdots & 2 \\
0 & 1 & 1 & \vdots & -3 \\
0 & 1 & 3 & \vdots & 3
\end{array}\right).$$

As $a_{22}^{(1)} = 1 \neq 0$, therefore, we eliminate entry in $a_{32}^{(1)}$ position by subtracting the multiple $m_{32} = \frac{1}{1} = 1$ of the second row from the third row, to get

$$\left(\begin{array}{ccccc}
1 & 2 & 1 & \vdots & 2 \\
0 & 1 & 1 & \vdots & -3 \\
0 & 0 & 2 & \vdots & 6
\end{array}\right).$$

Obviously, the original set of equations has been transformed to an upper-triangular form. Since all the diagonal elements of the obtaining upper-triangular matrix are nonzero, which means that the coefficient matrix of the given system is nonsingular and therefore, the given system has a unique solution. Now expressing the set in algebraic form yields

Now using backward substitution, we get

$$2x_3 = 6,$$
 $gives$ $x_3 = 3,$ $x_2 = -x_3 - 3 = -(3) - 3 = -6,$ $gives$ $x_2 = -6,$ $x_1 = 2 - 2x_2 - x_3 = 2 - 2(-6) - 3,$ $gives$ $x_1 = 11,$

which is the required solution of the given system.

The above results can be obtained using MATLAB commands as follows:

$$>> B = [1\ 2\ 1\ 2; 2\ 5\ 3\ 1; 1\ 3\ 4\ 5];\ x = WP(B);\ disp(x)$$

Program 3.5

MATLAB m-file for the Simple Gaussian Elimination Method function x=WP(B)

```
[n,t]=size(B); U=B; for k=1:n-1; for i=k:n-1; m=U(i+1,k)/U(k,k);for j=1:t; U(i+1,j)=U(i+1,j)-m*U(k,j);end;end end; i=n; x(i,1)=U(i,t)/U(i,i); for i=n-1:-1:1; s=0; for k=n:-1:i+1; s=s+U(i,k)*x(k,1); end; x(i,1)=(U(i,t)-s)/U(i,i); end; B; U; x; end
```

In the simple description of Gaussian elimination without pivoting just given, we used the kth equation to eliminate variable x_k from equations $k+1,\ldots,n$ during the kth step of the procedure. This is possible only if at the beginning of the kth step, the coefficient $a_{kk}^{(k-1)}$ of x_k in equation k is not zero. Since these coefficients are used as denominators both in the multipliers m_{ij} and in the backward substitution equations. But this does not necessarily mean that the linear system is not solvable, but that the procedure of solution must be altered.

Example 3.10 Solve the following linear system using the simple Gaussian elimination method

Solution. Writing the given system in the augmented matrix form

$$\left(\begin{array}{ccccc} 0 & 1 & 1 & \vdots & 1 \\ 1 & 2 & 2 & \vdots & 1 \\ 2 & 1 & 2 & \vdots & 3 \end{array}\right).$$

To solve this system, the simple Gaussian elimination method will fail immediately because the element in the first row on the leading diagonal, the pivot, is zero. Thus it is impossible to divide that row by the pivot value. Clearly, this difficulty can be overcome by rearranging the order of the rows; for example by making the first row the second, gives

$$\left(\begin{array}{ccccc} 1 & 2 & 2 & \vdots & 1 \\ 0 & 1 & 1 & \vdots & 1 \\ 2 & 1 & 2 & \vdots & 3 \end{array}\right).$$

Now we use the usual elimination process. The first elimination step is to eliminate the element $a_{31} = 2$ from the third row by subtracting a multiple $m_{31} = \frac{2}{1} = 2$ of row 1 from row 3, gives

$$\left(\begin{array}{ccccc} 1 & 2 & 2 & \vdots & 1 \\ 0 & 1 & 1 & \vdots & 1 \\ 0 & -3 & -2 & \vdots & 1 \end{array}\right).$$

We finished with the first elimination step since the element a_{21} is already eliminated from second row. The second elimination step is to eliminate the element $a_{32}^{(1)} = -3$ from the third row by subtracting a multiple $m_{32} = \frac{-3}{1}$ of row 2 from row 3, gives

$$\left(\begin{array}{ccccc} 1 & 2 & 2 & \vdots & 1 \\ 0 & 1 & 1 & \vdots & 1 \\ 0 & 0 & 1 & \vdots & 4 \end{array}\right).$$

Obviously, the original set of equations has been transformed to an upper-triangular form. Now expressing the set in algebraic form yields

Using backward substitution, we get, $x_1 = -1$, $x_2 = -3$, $x_3 = 4$, the solution of the system.

Example 3.11 Solve the linear system using the simple Gaussian elimination method

$$x_1 + x_2 + x_3 = 3$$

 $2x_1 + 2x_2 + 3x_3 = 7$
 $x_1 + 2x_2 + 3x_3 = 6$

Solution. Writing the given system in the augmented matrix form

$$\left(\begin{array}{ccccc} 1 & 1 & 1 & \vdots & 3 \\ 2 & 2 & 3 & \vdots & 7 \\ 1 & 2 & 3 & \vdots & 6 \end{array}\right).$$

First elimination step is to eliminate the elements $a_{21} = 2$ and $a_{31} = 1$ from second and third rows by subtracting multiples $m_{21} = \frac{2}{1} = 2$ and $m_{31} = \frac{1}{1} = 1$ of row 1 from row 2 and row 3 respectively, gives

$$\left(\begin{array}{ccccc} 1 & 1 & 1 & \vdots & 3 \\ 0 & 0 & 1 & \vdots & 1 \\ 0 & 1 & 2 & \vdots & 3 \end{array}\right).$$

We finished with the first elimination step. To start the second elimination step, since we note that the element $a_{22}^{(1)} = 0$, called the second pivot element, so the simple Gaussian elimination cannot continue in its present form. Therefore, we interchange the rows 2 and 3, to get

$$\left(\begin{array}{ccccc}
1 & 1 & 1 & \vdots & 3 \\
0 & 1 & 2 & \vdots & 3 \\
0 & 0 & 1 & \vdots & 1
\end{array}\right).$$

We finished with the second elimination step since the element $a_{32}^{(1)}$ is already eliminated from third row. Obviously, the original set of equations has been transformed to an upper-triangular form. Now expressing the set in algebraic form yields

Using backward substitution, we get, $x_1 = 1$, $x_2 = 1$, $x_3 = 1$, the solution of the system.

Example 3.12 Use the simple Gaussian elimination method, find all values of a and b for which the following linear system is consistent or inconsistent.

Solution. Writing the given system in the augmented matrix form

$$\left(\begin{array}{cccc} 2 & -1 & 3 & 1 \\ 4 & 2 & 2 & 2a \\ 2 & 1 & 1 & b \end{array}\right),$$

in which we wish to eliminate the elements a_{21} and a_{31} by subtracting from the second and third rows the appropriate multiples of the first row. In this case the multiples are $m_{21} = 2$ and $m_{31} = 1$. Hence

$$\left(\begin{array}{cccc} 2 & -1 & 3 & 1 \\ 0 & 4 & -4 & 2a - 2 \\ 0 & 2 & -2 & b - 1 \end{array}\right).$$

We finished with the first elimination step. The second elimination step is to eliminate element $a_{32}^{(1)}=2$ by subtracting a multiple $m_{32}=\frac{2}{4}=\frac{1}{2}$ of row 2 from row 3, gives

$$\left(\begin{array}{cccc} 2 & -1 & 3 & 1 \\ 0 & 4 & -4 & 2a - 2 \\ 0 & 0 & 0 & b - a \end{array}\right).$$

We finished with the second column. So third row of the equivalent upper-triangular system is

$$0x_1 + 0x_2 + 0x_3 = b - a. (3.25)$$

Firstly, if (3.25) has no constraint on unknowns x_1, x_2 , and x_3 , then the upper-triangular system represents only two non-trivial equations, namely

in three unknowns. As a result, one of the unknowns can be chosen arbitrarily, say $x_3 = x_3^*$, then x_2^* and x_1^* can be obtained by using backward substitution

$$x_2^* = a/2 - 1/2 + x_3^*; \quad x_1^* = \frac{1}{2}(1 + a/2 - 1/2 - 2x_3^*).$$

Hence

$$\mathbf{x}^* = \left[\frac{1}{2}(1/2 + a/2 - 2x_3^*), 1/2a - 1/2 + x_3^*, x_3^*\right]^T,$$

is an approximation solution of given system for any value of x_3^* for any real value of a. Hence the given linear system is consistent (infinitely many solutions).

Secondly, when $b-a \neq 0$, in this case (3.25) puts a restriction on unknowns x_1, x_2 and x_3 that is impossible to satisfy. So the system cannot have any solutions and therefore, it is inconsistent. •

Example 3.13 For what values of α the following linear system has (i) Unique solution, (ii) No solution, (iii) Infinitely many solutions, by using the simple Gaussian elimination method. Use smallest positive integer value of α to get the unique solution of the system.

Solution. Writing the given system in the augmented matrix form

$$\left(\begin{array}{cccc}
1 & 3 & \alpha & 4 \\
2 & -1 & 2\alpha & 1 \\
\alpha & 5 & 1 & 6
\end{array}\right),$$

and by using the following multiples $m_{21} = 2$ and $m_{31} = \alpha$, we get

$$\begin{pmatrix}
1 & 3 & \alpha & 4 \\
0 & -7 & 0 & -7 \\
0 & 5 - 3\alpha & 1 - \alpha^2 & 6 - 4\alpha
\end{pmatrix}.$$

Now using the multiple $m_{32} = \frac{5-3\alpha}{-7}$, gives

$$\left(\begin{array}{cccc}
1 & 3 & \alpha & 4 \\
0 & -7 & 0 & -7 \\
0 & 0 & 1 - \alpha^2 & 1 - \alpha
\end{array}\right).$$

So if $1-\alpha^2 \neq 0$, then we have the unique solution of the given system while for $\alpha = \pm 1$, we have no unique solution. If $\alpha = 1$, then we have infinitely many solution because third row of above matrix gives

$$0x_1 + 0x_2 + 0x_3 = 0$$

and when $\alpha = -1$, we have

$$0x_1 + 0x_2 + 0x_3 = 2$$
,

which is not possible, so no solution.

Since we can not take $\alpha = 1$ for the unique solution, so can take next positive integer $\alpha = 2$, which gives us upper-triangular system of the form

Solving this system using backward substitution, we get, $x_1 = 1/3$, $x_2 = 1$, $x_3 = 1/3$, the required unique solution of the given system using smallest positive integer value of α .

Theorem 3.11 An upper-triangular matrix A is nonsingular if and only if all its diagonal elements are different from zero.

Example 3.14 Use the simple Gaussian elimination method to find all the values of α which make the following matrix singular.

$$A = \left(\begin{array}{ccc} 1 & -1 & \alpha \\ 2 & 2 & 1 \\ 0 & \alpha & -1.5 \end{array}\right).$$

Solution. Applying the forward elimination step of the simple Gaussian elimination on the given matrix A and eliminate the element a_{21} by subtracting from the second row the appropriate multiple of the first row. In this case the multiple is given as

$$\begin{pmatrix} 1 & -1 & \alpha \\ 0 & 4 & 1 - 2\alpha \\ 0 & \alpha & -1.5 \end{pmatrix}.$$

We finished with the first elimination step. The second elimination step is to eliminate element $a_{32}^{(1)} = \alpha$ by subtracting a multiple $m_{32} = \frac{\alpha}{4}$ of row 2 from row 3, gives

$$\begin{pmatrix} 1 & -1 & \alpha \\ 0 & 4 & 1 - 2\alpha \\ 0 & 0 & -1.5 - \frac{\alpha(1 - 2\alpha)}{4} \end{pmatrix}.$$

To show that the given matrix is singular, we have to set the third diagonal element equal to zero (by Theorem 3.11), that is

$$-1.5 - \frac{\alpha(1-2\alpha)}{4} = 0$$
, or $2\alpha^2 - \alpha - 6 = 0$.

Solving the above quadratic equation, we get, $\alpha = -\frac{3}{2}$ and $\alpha = 2$, the possible values of α which make the given matrix singular.

Example 3.15 Use the smallest positive integer value of α to find the unique solution of the linear system $A\mathbf{x} = [1, 6, -4]^T$ by simple Gaussian elimination method, where

$$A = \left(\begin{array}{ccc} 1 & -1 & \alpha \\ 2 & 2 & 1 \\ 0 & \alpha & -1.5 \end{array} \right).$$

Solution. Since we know from the previous Example 3.14 that the given matrix A is singular when $\alpha = -\frac{3}{2}$ and $\alpha = 2$. To find the unique solution we take the smallest positive integer value $\alpha = 1$ and consider the augmented matrix as follows:

$$\begin{pmatrix} 1 & -1 & 1 & \vdots & 1 \\ 2 & 2 & 1 & \vdots & 6 \\ 0 & 1 & -1.5 & \vdots & -4 \end{pmatrix}.$$

Applying the forward elimination step of the simple Gaussian elimination on the given matrix A and eliminate the element a_{21} by subtracting from the second row the appropriate multiple $m_{21} = 2$ of the first row, gives

$$\begin{pmatrix} 1 & -1 & 1 & \vdots & 1 \\ 0 & 4 & -1 & \vdots & 4 \\ 0 & 1 & -1.5 & \vdots & -4 \end{pmatrix}.$$

The second elimination step is to eliminate element $a_{32}^{(1)} = 1$ by subtracting a multiple $m_{32} = \frac{1}{4}$ of row 2 from row 3, gives

$$\begin{pmatrix} 1 & -1 & 1 & \vdots & 1 \\ 0 & 4 & -1 & \vdots & 4 \\ 0 & 0 & -5/4 & \vdots & -5 \end{pmatrix}.$$

Now expressing the set in algebraic form yields

Using backward substitution, we get, $x_1 = -1$, $x_2 = 2$, $x_3 = 4$, the unique solution.

Note that the inverse of the nonsingular matrix A can be easily determined by using the simple Gaussian elimination method. Here, we have to consider the augmented matrix as a combination of the given matrix A and the identity matrix \mathbf{I} (same size as of A). To find the inverse matrix $B = A^{-1}$ we must solve the linear system in which the jth column of the matrix B is the solution of the linear system with right-hand side the jth column of the matrix \mathbf{I} .

Example 3.16 Use the simple Gaussian elimination method to find the inverse of the following matrix

$$A = \left(\begin{array}{ccc} 2 & -1 & 3 \\ 4 & -1 & 6 \\ 2 & -3 & 4 \end{array}\right).$$

Solution. Suppose that the inverse $A^{-1} = B$ of the given matrix exists and let

$$AB = \begin{pmatrix} 2 & -1 & 3 \\ 4 & -1 & 6 \\ 2 & -3 & 4 \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \mathbf{I}.$$

Now to find the elements of the matrix B, we apply the simple Gaussian elimination on the augmented matrix

$$[A|\mathbf{I}] = \begin{pmatrix} 2 & -1 & 3 & \vdots & 1 & 0 & 0 \\ 4 & -1 & 6 & \vdots & 0 & 1 & 0 \\ 2 & -3 & 4 & \vdots & 0 & 0 & 1 \end{pmatrix}.$$

Applying the forward elimination step of the simple Gaussian elimination on the given matrix A and eliminate the elements $a_{21} = 4$ and $a_{31} = 2$ by subtracting from the second and the third rows the appropriate multiples $m_{21} = \frac{4}{2} = 2$ and $m_{31} = \frac{2}{2} = 1$ of the first row. It gives

$$\left(\begin{array}{cccccccc} 2 & -1 & 3 & \vdots & 1 & 0 & 0 \\ 0 & 1 & 0 & \vdots & -2 & 1 & 0 \\ 0 & -2 & 1 & \vdots & -1 & 0 & 1 \end{array}\right).$$

We finished with the first elimination step. The second elimination step is to eliminate element $a_{32}^{(1)} = -2$ by subtracting a multiple $m_{32} = \frac{-2}{1} = -2$ of row 2 from row 3, gives

We solve the first system

$$\begin{pmatrix} 2 & -1 & 3 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} b_{11} \\ b_{21} \\ b_{31} \end{pmatrix} = \begin{pmatrix} 1 \\ -2 \\ -5 \end{pmatrix},$$

by using backward substitution, we get

which gives $b_{11}=7,\ b_{21}=-2,\ b_{31}=-5.$ Similarly, the solution of the second linear system

$$\begin{pmatrix} 2 & -1 & 3 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} b_{12} \\ b_{22} \\ b_{32} \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 2 \end{pmatrix},$$

can be obtained as follows:

$$2b_{12} - b_{22} + 3b_{32} = 0$$

 $b_{22} = 1$
 $b_{32} = 2$

which gives $b_{12} = -5/2$, $b_{22} = 1$, $b_{32} = 2$. Finally, the solution of the third linear system

$$\begin{pmatrix} 2 & -1 & 3 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} b_{13} \\ b_{23} \\ b_{33} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix},$$

can be obtained as follows:

$$2b_{13} - b_{23} + 3b_{33} = 0$$

 $b_{23} = 0$
 $b_{33} = 1$

and it gives $b_{13}=-3/2,\ b_{23}=0,\ b_{33}=1.$ Hence the elements of the inverse matrix B are

$$B = A^{-1} = \begin{pmatrix} 7 & -\frac{5}{2} & -\frac{3}{2} \\ -2 & 1 & 0 \\ -5 & 2 & 1 \end{pmatrix},$$

which is the required inverse of the given matrix A.

Procedure 3.1 [Gaussian Elimination Method]

- 1. Form the augmented matrix, $B = [A|\mathbf{b}]$.
- 2. Check first pivot element $a_{11} \neq 0$, then move to the next step; otherwise, interchange rows so that $a_{11} \neq 0$.
- 3. Multiply row one by multiplier $m_{i1} = \frac{a_{i1}}{a_{11}}$ and subtract to the ith row for i = 2, 3, ..., n.
- 4. Repeat the steps 2 and 3 for the remaining pivots elements unless coefficient matrix A becomes upper-triangular matrix U.
- 5. Use backward substitution to solve x_n from the nth equation $x_n = \frac{b_n^{n-1}}{a_{nn}}$ and solve the other (n-1) unknowns variables by using (3.21).

The use of non-zero pivots is sufficient for the theoretical correctness of the simple Gaussian elimination, but more care must be taken if one is to obtain reliable results.

Example 3.17 Consider a linear system

which has exact solution $\mathbf{x} = [1.00010, 0.99990]^T$. Now we solve this system by the simple Gaussian elimination. The first elimination step is to eliminate first variable x_1 from second equation by subtracting multiple $m_{21} = 10000$ of first equation from second equation, gives

Using backward substitution to get the solution $\mathbf{x}^* = [0, 1]^T$. Thus a computational disaster has occurred. But if we interchange the equations, we obtain

$$\begin{array}{rcl} x_1 & + & x_2 & = & 2 \\ 0.000100x_1 & + & x_2 & = & 1 \end{array}$$

Apply the Gaussian elimination again, and we got the solution $\mathbf{x}^* = [1,1]^T$. This solution is as good as one would hope. So, we conclude from this example that it is not enough just to avoid zero pivot, one must also avoid relatively small one.

Here we need some pivoting strategies which help us to over come these difficulties facing during the process of simple Gaussian elimination.

3.4.2 Pivoting Strategies

Since we know that simple Gaussian elimination is applied to a problem with no pivotal elements zero. However, the method does not work if the first coefficient of the first equation or if a diagonal coefficients becomes zero in the process of solution because they are used as denominators in a forward elimination.

Pivoting is used to change sequential order of the equations for two purposes, first to prevent diagonal coefficients from becoming zero, and second, to make each diagonal coefficient larger in magnitude than any other coefficient below it, that is, to decrease the round-off errors. The equations are not mathematical affected by changes of the sequential order, but changing the order makes coefficient become non-zero. Even when all diagonal coefficients are non-zero, the changes of order increases accuracy of the computations. The standard pivoting strategy which handled these difficulties easily are explained below.

Partial Pivoting

Here we develop an implementation of the Gaussian elimination which utilizes the pivoting strategy discussed above. In using the Gaussian elimination by partial pivoting(or row pivoting), the basic approach is to use the largest (in absolute value) element on or below the diagonal in the column of current interest as the pivotal element for elimination in the rest of that column.

One immediate effect of this will be to force all the multiples used to be not greater than 1 in absolute value. This will inhibit the growth of error in the rest of elimination phase and in subsequent backward substitution.

At stage k of forward elimination, it is necessary, therefore, to be able to identify the largest element from $|a_{kk}|, |a_{k+1,k}|, \ldots, |a_{nk}|$, where these a_{ik} 's are the elements in the current partially triangularized coefficient matrix. If this maximum occurs in row p, then pth and kth rows of the augmented matrix are interchange and the elimination proceed as usual. In solving n linear equations, a total of $N = \frac{n(n+1)}{2}$ coefficients must be examined.

Example 3.18 Solve the following linear system using the Gaussian elimination with partial pivoting

$$x_1 + x_2 + x_3 = 1$$

 $2x_1 + 3x_2 + 4x_3 = 3$
 $4x_1 + 9x_2 + 16x_3 = 11$

Solution. For the first elimination step, since 4 is the largest absolute coefficient of first variable x_1 , therefore, the first row and the third row are interchange, giving us

$$4x_1 + 9x_2 + 16x_3 = 11$$

 $2x_1 + 3x_2 + 4x_3 = 3$
 $x_1 + x_2 + x_3 = 1$

Eliminate first variable x_1 from the second and third rows by subtracting the multiples $m_{21} = \frac{2}{4}$ and $m_{31} = \frac{1}{4}$ of row 1 from row 2 and row 3 respectively, gives

$$4x_1 + 9x_2 + 16x_3 = 11$$

$$- 3/2x_2 - 4x_3 = -5/2$$

$$- 5/4x_2 - x_3 = -7/5$$

For the second elimination step, -3/2 is the largest absolute coefficient of second variable x_2 , so eliminate second variable x_2 from the third row by subtracting the multiple $m_{32} = \frac{5}{6}$ of row 2 from row 3, gives

Obviously, the original set of equations has been transformed to an equivalent upper-triangular form. Now using backward substitution, gives, $x_1 = 1$, $x_2 = -1$, $x_3 = 1$, the required solution.

The following MATLAB commands will give the same results as we obtained in the preceding Example 3.18 of the Gaussian elimination method with partial pivoting:

$$>> B = [1 \ 1 \ 1 \ 1; 2 \ 3 \ 4 \ 3; 4 \ 9 \ 16 \ 11]; \ x = PP(B); \ disp(x)$$

Program 3.6MATLAB m-file for Gaussian Elimination by Partial Pivoting function x=PP(B)% $B = input('input \ matrix \ in \ form[A/b]');$ [n,t] = size(B); U = B; for $M = 1:n-1; \ mx(M) = abs(U(M,M)); \ r = M;$ for i = M+1:n; if mx(M) < abs(U(i,M)); mx(M)=abs(U(i,M)); r = i; end; end rw1(1,1:t)=U(r,1:t); rw2(1,1:t)=U(M,1:t); U(M,1:t)=rw1; U(r,1:t)=rw2; for k=M+1:n m=U(k,M)/U(M,M); for j=M:t; U(k,j) = U(k,j) - m*U(M,j); end;end; i=n; x(i)=U(i,t)/U(i,i); for i=n-1:-1:1; s=0; for k=n:-1:i+1; s=s+U(i,k)*x(k); end; x(i)=(U(i,t)-s)/U(i,i); end; B; U; x; end

Procedure 3.2 [Partial Pivoting]

- 1. Suppose we are about to work on the ith column of the matrix. Then we search that portion of the ith column below and including the diagonal, and find the element that has the largest absolute value. Let p denote the index of the row that contains this element.
- 2. Interchange row i and p.
- 3. Proceed with the elimination Procedure 3.1.

3.4.3 Gauss-Jordan Method

This method is a modification of the Gaussian elimination method. The Gauss-Jordan method is although inefficient for practical calculation but is often useful for theoretical purposes. The basic of this method is to convert the given matrix into a diagonal form. The forward elimination of the Gauss-Jordan method is identical to that of Gaussian elimination method. However, Gauss-Jordan elimination uses backward elimination rather than backward substitution. In the Gauss-Jordan method the forward elimination and backward elimination need not be separated. This is possible because a pivot element can be used to eliminate the coefficients not only below but also above at the same time. If this approach is taken, the form of the coefficients matrix become diagonal when elimination by the last pivot are completed. The Gauss-Jordan method simply yields a transformation of the augmented matrix of the form

$$[A|\mathbf{b}] \to [\mathbf{I}|\mathbf{c}],$$

where I is the identity matrix and c is the column matrix, which represents the possible solution of the given linear system.

The Gauss-Jordan method particularly well suited to compute the inverse of a matrix through the transformation

$$[A|\mathbf{I}] \to [\mathbf{I}|A^{-1}].$$

Note if the inverse of the matrix can be found, then the solution of the linear system can be computed easily from the product of matrix A^{-1} and column matrix **b**, that is

$$\mathbf{x} = A^{-1}\mathbf{b}.\tag{3.26}$$

Note that one can get easily the solution of linear system $A\mathbf{x} = \mathbf{b}$ and inverse of the coefficient matrix A together by the Gauss-Jordan method using the augmented matrix of the form

$$[A|\mathbf{b}|\mathbf{I}] \to [\mathbf{I}|\mathbf{x}|A^{-1}].$$

Example 3.19 Apply the Gauss-Jordan method to find the inverse of the coefficient matrix and also the solution of the linear system $A\mathbf{x} = \mathbf{b}$, where

$$A = \begin{pmatrix} 1 & 2 \\ 1 & 3 \end{pmatrix} \quad and \quad \mathbf{b} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}.$$

Solution. Consider the following augmented matrix

$$[A|\mathbf{b}|\mathbf{I} = \begin{pmatrix} 1 & 2 & \vdots & 1 & \vdots & 1 & 0 \\ 1 & 3 & \vdots & 2 & \vdots & 0 & 1 \end{pmatrix}.$$

Then we have

$$\equiv \left(\begin{array}{ccccc} 1 & 2 & \vdots & 1 & \vdots & 1 & 0 \\ 0 & 1 & \vdots & 1 & \vdots & -1 & 1 \end{array}\right) \equiv \left(\begin{array}{cccccc} 1 & 0 & \vdots & -1 & \vdots & 3 & -2 \\ 0 & 1 & \vdots & 1 & \vdots & -1 & 1 \end{array}\right).$$

Thus we obtain the inverse of the matrix A

$$A^{-1} = \left(\begin{array}{cc} 3 & -2 \\ -1 & 1 \end{array} \right),$$

and $x_1 = -1$, $x_2 = 1$, the solution of the given system.

The above results can be obtained using MATLAB commands, we do the following:

$$>> Ab = [A|b|I] = [1\ 2\ 1\ 1\ 0; 1\ 3\ 2\ 0\ 1];\ [I|inv(A)] = GaussJ(Ab);$$

3.4.4 LU Decomposition Method

This is another direct method to find the solution of the system of linear equations. The LU decomposition (or factorization method) is a modification of the elimination method. Here we decompose or factorize the coefficient matrix A into the product of two triangular matrices in the form $A = LU, \tag{3.27}$

where L is a lower-triangular matrix and U is the upper-triangular matrix. Both are of same size as the coefficients matrix A. To solve a number of linear equations sets in which the coefficients matrices are all identical but the right-hand side are different, then the LU decomposition is more efficient than elimination method. Specifying the diagonal elements of either L and U makes the factoring unique. The procedure based on unity elements on the diagonal of matrix L is called *Doolittle's method* (or Gauss factorization), while the procedure based on unity elements on the diagonal of matrix U is called C method.

The general forms of L and U are written as

$$L = \begin{pmatrix} l_{11} & 0 & \cdots & 0 \\ l_{21} & l_{22} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ l_{n1} & l_{n2} & \cdots & l_{nn} \end{pmatrix} \quad \text{and} \quad U = \begin{pmatrix} u_{11} & u_{12} & \cdots & u_{1n} \\ 0 & u_{22} & \cdots & u_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & u_{nn} \end{pmatrix}, \tag{3.28}$$

such that $l_{ij} = 0$ for i < j and $u_{ij} = 0$ for i > j.

Consider a linear system

$$A\mathbf{x} = \mathbf{b},\tag{3.29}$$

and let A be factored into the product of L and U, as shown by (3.28). Then the linear system (3.29) becomes

$$LU\mathbf{x} = \mathbf{b}$$
,

which can be written as

$$L\mathbf{y} = \mathbf{b}$$
, where $\mathbf{y} = U\mathbf{x}$.

The unknown elements of matrix L and matrix U are computed by equating corresponding elements in matrices A and LU in a systematic way. Once the matrices L and U have been constructed, the solution of system (3.29) can be computed in the following two steps:

1. Solve the lower-triangular system $L\mathbf{y} = \mathbf{b}$.

By using the forward elimination, we will find the components of the unknown vector \mathbf{y} , by using the following steps:

$$y_1 = b_1,$$

 $y_i = b_i - \sum_{j=1}^{i-1} l_{ij} y_j, \qquad i = 2, 3, \dots, n$ (3.30)

2. Solve the upper-triangular system $U\mathbf{x} = \mathbf{y}$.

By using the *backward substitution*, we will find the components of the unknown vector \mathbf{x} , by using the following steps:

$$x_{n} = \frac{y_{n}}{u_{nn}},$$

$$x_{i} = \frac{1}{u_{ii}} \left[y_{i} - \sum_{j=i+1}^{n} u_{ij} x_{j} \right], \quad i = n-1, n-2, \dots, 1$$

$$(3.31)$$

Thus the relationship of the matrices L and U to the original matrix A is given by the following theorem.

Theorem 3.12 If the Gaussian elimination can be performed on the linear system $A\mathbf{x} = \mathbf{b}$ without row interchanges, then the matrix A can be factored into the product of a lower-triangular matrix L and an upper-triangular matrix U, that is

$$A = LU$$

where the matrices L and U are of the same size as A.

Theorem 3.13 Let A be an $n \times n$ matrix that has an LU factorization, that is

$$A = LU$$
.

If A has rank n (that is, all pivots are non-zeros), then L and U are uniquely determined by A. •

Now we discuss the two possible variations of the LU decomposition to find the solution of the nonsingular linear system in the following.

Doolittle's Method

In Doolittle's method(which is also called the Gauss factorization), the upper-triangular matrix U is obtained by forward elimination of the Gaussian elimination method and the lower-triangular matrix L containing the multiples used in the Gaussian elimination process as the elements below the diagonal with unity elements on the main diagonal, that is,

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ m_{21} & 1 & 0 \\ m_{31} & m_{32} & 1 \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{pmatrix}.$$

Example 3.20 Construct the LU decomposition of the following matrix A by using the Gauss factorization (that is, the LU decomposition by Doolittle's method).

Solution. Applying the forward elimination step of Simple Gauss-elimination to the given matrix

$$A = \left(\begin{array}{rrr} 1 & 2 & 1 \\ 2 & 5 & 3 \\ 1 & 3 & 4 \end{array}\right),$$

using the multiples $m_{21} = 2$ and $m_{31} = 1$, we get

$$\equiv \left(\begin{array}{ccc} 1 & 2 & 1 \\ 0 & 1 & 1 \\ 0 & 1 & 3 \end{array}\right).$$

Similarly, by using the multiple $m_{32} = 1$, we obtain

$$\equiv \left(\begin{array}{ccc} 1 & 2 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 2 \end{array} \right) = U.$$

Hence we obtained the LU-decomposition of the given matrix as follows

$$\begin{pmatrix} 1 & 2 & 1 \\ 2 & 5 & 3 \\ 1 & 3 & 4 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 2 \end{pmatrix},$$

where the unknown elements of matrix L are the used multiples and the matrix U is same as we obtained in forward elimination process.

Example 3.21 Construct the LU decomposition of the following matrix A by using the Gauss factorization (that is, the LU decomposition by Doolittle's method). Find the value(s) of α for which the following matrix

$$A = \left(\begin{array}{rrr} 1 & -1 & \alpha \\ -1 & 2 & -\alpha \\ \alpha & 1 & 1 \end{array}\right),$$

is singular. Also, find the unique solution of the linear system $A\mathbf{x} = [1, 1, 2]^T$ by using the smallest positive integer value of α .

Solution. Since we know that

$$A = \begin{pmatrix} 1 & -1 & \alpha \\ -1 & 2 & -\alpha \\ \alpha & 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ m_{21} & 1 & 0 \\ m_{31} & m_{32} & 1 \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{pmatrix} = LU.$$

Now we will use only the forward elimination step of the simple Gaussian elimination method to convert the given matrix A into the upper-triangular matrix U. Since $a_{11} = 1 \neq 0$, so we wish to eliminate the elements $a_{21} = -1$ and $a_{31} = \alpha$ by subtracting from the second and third rows the appropriate multiples of the first row. The multiples are $m_{21} = -1 = l_{21}$ and $m_{31} = \alpha = l_{31}$, so

$$\left(\begin{array}{cccc}
1 & -1 & \alpha \\
0 & 1 & 0 \\
0 & 1 + \alpha & 1 - \alpha^2
\end{array}\right).$$

As $a_{22}^{(1)}=1\neq 0$, therefore, we eliminate entry in $a_{32}^{(1)}=1+\alpha$ position by subtracting the multiple $m_{32}=\frac{1+\alpha}{1}=l_{32}$ of the second row from the third row, to get

$$\left(\begin{array}{cccc}
1 & -1 & \alpha \\
0 & 1 & 0 \\
0 & 0 & 1 - \alpha^2
\end{array}\right).$$

Obviously, the original set of equations has been transformed to an upper-triangular form. Thus

$$\begin{pmatrix} 1 & -1 & \alpha \\ -1 & 2 & -\alpha \\ \alpha & 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ \alpha & 1 + \alpha & 1 \end{pmatrix} \begin{pmatrix} 1 & -1 & \alpha \\ 0 & 1 & 0 \\ 0 & 0 & 1 - \alpha^2 \end{pmatrix}.$$

which is the required decomposition of A. The matrix will be singular if the third diagonal element $1-\alpha^2$ of the upper-triangular U is equal to zero (Theorem 3.11), gives, $\alpha=\pm 1$. To find the unique solution of the given system we take $\alpha=2$ and it gives

$$\begin{pmatrix} 1 & -1 & 2 \\ -1 & 2 & -2 \\ 2 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 2 & 3 & 1 \end{pmatrix} \begin{pmatrix} 1 & -1 & 2 \\ 0 & 1 & 0 \\ 0 & 0 & -3 \end{pmatrix}.$$

Write MATLAB m-file **Decompd**.m to factored a nonsingular matrix A into a unit lower triangular matrix L and an upper triangular matrix U and using the following MATLAB commands as:

$$>> A = [1 \ 1 \ -2; -1 \ 2 \ -2; 2 \ 1 \ 1]; Sol = Decomp(A);$$

Program 3.7

MATLAB m-file for Decomposition of Matrix function Sol = Decomp(A)

 $[n,n] = size(A); \ U=A; \ L=eye(n); \ for \ i=1:n; \ for \ k=i=1:n; \ L(k,i)=U(k,i)/U(i,i); \ U(k,i:n)=U(k,i:n)-U(i,i:n)*L(k,i); end; end; \ Sol=[L,U]; \ det A=prod(diag(U))$

Now solving the lower-triangular system $L\mathbf{y} = \mathbf{b}$ for unknown vector \mathbf{y} , that is

$$\begin{pmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 2 & 3 & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix}.$$

Performing forward substitution yields

Using the m-file ForwardSubs.m

```
\begin{aligned} &functiony = ForwardSubs(L,b)\\ &[n,n] = size(L); \ y = zeros(n,1);\\ &fork = 1:n; \ y(k) = (b(k)-L(k,1:k-1)*y(1:k-1))/L(k,k) \ end \end{aligned}
```

and the following MATLAB commands to generate the solution of lower-triangular system as:

$$>> L = [1\ 0\ 0; -1\ 1\ 0; 2\ 3\ 1];\ b = [1;\ 1;\ 2];\ sol = ForwardSubs(L,b);$$

Then solving the upper-triangular system $U\mathbf{x} = \mathbf{y}$ for unknown vector \mathbf{x} , that is

$$\begin{pmatrix} 1 & -1 & 2 \\ 0 & 1 & 0 \\ 0 & 0 & -3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ -6 \end{pmatrix}.$$

Performing backward substitution yields

the approximate solution of the given system.

Using the m-file BackwardSubs.m

```
function x = Backward Subs(U, y)
[n, n] = size(U); \ x = zeros(n, 1); \ x(n) = y(n)/U(n, n);
for k = n - 1: -1: 1; x(k) = (y(k) - U(k, k + 1: n) * x(k + 1: n))/U(k, k); \ end
```

and the following MATLAB commands to generate the solution of upper-triangular system as:

$$>> U = [1 \ -1 \ 2; 0 \ 1 \ 0; 0 \ 0 \ -3]; \ y = [1; \ 1; \ 2]; \ sol = BackwardSubs(U,y);$$

There is an other way to find the values of the unknown elements of the matrices L and U, which we describe in the following example.

Example 3.22 Construct the LU decomposition of the following matrix using Doolittle's method

$$A = \left(\begin{array}{rrr} 1 & 2 & 4 \\ 1 & 3 & 3 \\ 2 & 2 & 2 \end{array}\right).$$

Solution. Since

$$A = LU = \begin{pmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{pmatrix}.$$

Performing the multiplication on the right-hand side, gives

$$\begin{pmatrix} 1 & 2 & 4 \\ 1 & 3 & 3 \\ 2 & 2 & 2 \end{pmatrix} = \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ l_{21}u_{11} & l_{21}u_{12} + u_{22} & l_{21}u_{13} + u_{23} \\ l_{31}u_{11} & l_{31}u_{12} + l_{32}u_{22} & l_{31}u_{13} + l_{32}u_{23} + u_{33} \end{pmatrix}.$$

Then equating elements of first column to obtain

Now equating elements of second column to obtain

Finally, equating elements of third column to obtain

Thus we obtain

$$\begin{pmatrix} 1 & 2 & 4 \\ 1 & 3 & 3 \\ 2 & 2 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 2 & -2 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 & 4 \\ 0 & 1 & -1 \\ 0 & 0 & -8 \end{pmatrix},$$

the factorization of the given matrix.

The general formula for getting elements of L and U corresponding to the coefficient matrix A for a set of n linear equations can be written as

$$u_{ij} = a_{ij} - \sum_{k=1}^{i-1} l_{ik} u_{kj}, \qquad 2 \leq i \leq j$$

$$l_{ij} = \frac{1}{u_{ii}} \left[a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \right], \qquad i > j \geq 2$$

$$u_{ij} = a_{1j}, \qquad i = 1$$

$$l_{ij} = \frac{a_{i1}}{u_{11}} = \frac{a_{i1}}{a_{11}}, \qquad j = 1$$

$$(3.32)$$

Example 3.23 Solve the following linear system by LU decomposition using Doolittle's method

$$A = \begin{pmatrix} 1 & 2 & 4 \\ 1 & 3 & 3 \\ 2 & 2 & 2 \end{pmatrix} \quad and \quad \mathbf{b} = \begin{pmatrix} -2 \\ 3 \\ -6 \end{pmatrix}.$$

Solution. Since the factorization of the coefficient matrix A has been already constructed in the Example 3.22 as

$$\begin{pmatrix} 1 & 2 & 4 \\ 1 & 3 & 3 \\ 2 & 2 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 2 & -2 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 & 4 \\ 0 & 1 & -1 \\ 0 & 0 & -8 \end{pmatrix}.$$

Then solving the first system $L\mathbf{y} = \mathbf{b}$ for unknown vector \mathbf{y} , that is

$$\begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 2 & -2 & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} -2 \\ 3 \\ -6 \end{pmatrix}.$$

Performing forward substitution yields

Then solving the second system $U\mathbf{x} = \mathbf{y}$ for unknown vector \mathbf{x} , that is

$$\begin{pmatrix} 1 & 2 & 4 \\ 0 & 1 & -1 \\ 0 & 0 & -8 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} -2 \\ 5 \\ 8 \end{pmatrix}.$$

Performing backward substitution yields

the approximate solution of the given system.

We can also write MATLAB m-file called, Doolittle.m to get the solution of the linear system by LU decomposition by using Doolittle's method. In order to reproduce above results using MATLAB commands, we do the following:

$$>> A = [1\ 2\ 4; 1\ 3\ 3; 2\ 2\ 2];\ b = [-2\ 3\ -6];\ sol = Doolittle(A, b);$$

Program 3.8

MATLAB m-file for using the Doolittle's Method

function sol = Doolittle(A,b)

[n,n]=size(A); u=A;l=zeros(n,n);

for i=1:n-1; if abs(u(i,i)) > 0; for i1=i+1:n; m(i1,i)=u(i1,i)/u(i,i);

for j=1:n; u(i1, j) = u(i1, j) - m(i1, i) * u(i, j); end; end; end;

for i=1:n; l(i,1)=A(i,1)/u(1,1); end; for j=2:n; for i=2:n; s=0;

for k=1:j-1; s = s + l(i, k) * u(k, j); end

l(i,j)=(A(i,j)-s)/u(j,j); end; end y(1)=b(1)/l(1,1);

for k=2:n; sum=b(k); for i=1:k-1; sum = sum - l(k, i) * y(i); end

y(k)=sum/l(k,k); end; x(n)=y(n)/u(n,n); for k=n-1:-1:1; sum=y(k);

for i=k+1:n; sum = sum - u(k,i) * x(i); end; x(k)=sum/u(k,k); end; l; u; y; x

Procedure 3.3 [LU Decomposition by Doolittle's Method]

- 1. Take the nonsingular matrix, A.
- 2. If possible, decompose the matrix A = LU using (3.32).
- 3. Solve linear system $L\mathbf{y} = \mathbf{b}$ using (3.30).
- 4. Solve linear system $U\mathbf{x} = \mathbf{y}$ using (3.31).

Example 3.24 Use LU-factorization method with Doolittle's method $(l_{ii} = 1)$ to find values of α for which the following linear system has unique solution and infinitely many solutions. Write down the solution for both cases.

Solution. We use Simple Gauss-elimination method to convert the following matrix of the given system by using the multiples $m_{21} = 2$, $m_{31} = -1$ and $m_{32} = 1/4$,

$$A = \left(\begin{array}{rrr} 1 & 0.5 & \alpha \\ 2 & -3 & 1 \\ -1 & -1.5 & 2.5 \end{array}\right),$$

into equivalent an upper-triangular matrix form

$$\begin{pmatrix} 1 & 0.5 & \alpha \\ 0 & -4 & 1 - 2\alpha \\ 0 & 0 & 0.5\alpha - 0.25 \end{pmatrix},$$

to get LU-factorization of A in the following form

$$A = \begin{pmatrix} 1 & 0.5 & \alpha \\ 2 & -3 & 1 \\ -1 & -1.5 & 2.5 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & 0.25 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0.5 & \alpha \\ 0 & -4 & 1 - 2\alpha \\ 0 & 0 & 1.5\alpha + 2.25 \end{pmatrix} = LU.$$

Then by solving the lower-triangular system of the form $L\mathbf{y} = [0.5, -1, -1]^T$ and obtained the solution $\mathbf{y} = [0.5, -2, 0]^T$. Now solving the upper-triangular system $U\mathbf{x} = \mathbf{y}$ of the form

$$\begin{pmatrix} 1 & 0.5 & \alpha \\ 0 & -4 & 1 - 2\alpha \\ 0 & 0 & 1.5\alpha + 2.25 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0.5 \\ -2 \\ 0 \end{pmatrix}.$$

From last equation we have

$$(1.5\alpha + 2.25)x_3 = 0$$
,

so for unique solution of the given system $(1.5\alpha + 2.25) \neq 0$ (nonsingular), which implies that $x_3 = 0$. Using backward substitution, we have $x_2 = 0.5$ and $x_1 = 0.25$. Thus, $[0.25, 0.5, 0]^T$ is the unique solution of the given system.

If $(1.5\alpha + 2.25) = 0$ (singular), that is, $\alpha = -1.5$, then for this we must have infinitely many solutions. So to get the infinitely many solutions, we have to solve the following resulting system

$$x_1 + 0.5x_2 + \alpha x_3 = 0.5$$

 $- 4x_2 + (1-2\alpha)x_3 = -2$
 $0x_3 = 0$

By taking $\alpha = -1.5$ and If we choose $x_3 = t \in R$, $t \neq 0$, then we have $x_2 = 0.5 + t$ and $x_1 = 0.25 + t$, so $\mathbf{x}^* = [0.25 + t, 0.5 + t, t]^T$ is the infinitely many solutions of the given system.

Example 3.25 Use LU-factorization method with Doolittle's method $(l_{ii} = 1)$ to find the constant α such that the following homogeneous linear system has non-trivial solutions. Find these solutions.

Solution. Using Simple Gauss-elimination method, we can easily find fatorization of A as

$$A = \begin{pmatrix} 1 & 1 & 0 \\ 3 & \alpha & 5 \\ 0 & 7 & 3 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ 0 & 7/(\alpha - 3) & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 \\ 0 & (\alpha - 3) & 5 \\ 0 & 0 & 3 - 35/(\alpha - 3) \end{pmatrix}.$$

Since by one of the property of the determinant

$$det(A) = det(LU) = det(L) det(U).$$

So when using LU decomposition by Doolittle's method, then

$$\det(A) = \det(U) = \prod_{i=1}^{n} u_{ii} = (u_{11}u_{22}\cdots u_{nn}),$$

where det(L) = 1 because L is lower-triangular matrix and all its diagonal elements are unity. Thus the determinant of the given matrix A is

$$|A| = |U| = (\alpha - 3)(3 - 35/(\alpha - 3)) = 3\alpha - 44, \quad \alpha \neq 3.$$

So |A| = 0, gives, $\alpha = 44/3$ and for this value of α we have non-trivial solutions. By solving the lower-triangular system of the form $L\mathbf{y} = [0,0,0]^T$, we obtained the solution $\mathbf{y} = [0,0,0]^T$. Now solving the upper-triangular system $U\mathbf{x} = \mathbf{y}$ of the form

$$\begin{pmatrix} 1 & 1 & 0 \\ 0 & 35/3 & 5 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}.$$

If we choose $x_3 = t \in R$, $t \neq 0$, then, $x_2 = (-3/7)t$ and $x_1 = (3/7)t$, then the non-trivial solutions of the given system is $\mathbf{x}^* = [3t/7, -3t/7, t]^T$.

Crout's Method

The Crout's method, in which matrix U has unity on the main diagonal, is similar to Doolittle's method in all other aspects. The L and U matrices are obtained by expanding the matrix equation A = LU term by term to determine the elements of the L and U matrices.

Example 3.26 Construct the LU decomposition of the following matrix using Crout's method

$$A = \left(\begin{array}{ccc} 1 & 2 & 3 \\ 6 & 5 & 4 \\ 2 & 5 & 6 \end{array}\right).$$

Solution. Since

$$A = LU = \begin{pmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{pmatrix} \begin{pmatrix} 1 & u_{12} & u_{13} \\ 0 & 1 & u_{23} \\ 0 & 0 & 1 \end{pmatrix}.$$

Performing the multiplication on the right-hand side, gives

$$\begin{pmatrix} 1 & 2 & 3 \\ 6 & 5 & 4 \\ 2 & 5 & 6 \end{pmatrix} = \begin{pmatrix} l_{11} & l_{11}u_{12} & l_{11}u_{13} \\ l_{21} & l_{21}u_{12} + l_{22} & l_{21}u_{13} + l_{22}u_{23} \\ l_{31} & l_{31}u_{12} + l_{32} & l_{31}u_{13} + l_{32}u_{23} + l_{33} \end{pmatrix}.$$

Then equating elements of first column to obtain

$$\begin{array}{rcl}
1 & = & l_{11}, \\
6 & = & l_{21}, \\
2 & = & l_{21}
\end{array}$$

Then equating elements of second column to obtain

$$2 = l_{11}u_{12},$$
 $u_{12} = 2,$ $5 = l_{21}u_{12} + l_{22},$ $l_{22} = 5 - 12 = -7,$ $5 = l_{31}u_{12} + l_{32},$ $l_{32} = 5 - 4 = 1.$

Finally, then equating elements of third column to obtain

$$3 = l_{11}u_{13},$$
 $u_{13} = 3,$ $u_{23} = (4-18)/-7 = 2,$ $u_{23} = (6-6-2) = -2.$

Thus we get

$$\left(\begin{array}{ccc} 1 & 2 & 3 \\ 6 & 5 & 4 \\ 2 & 5 & 6 \end{array}\right) = \left(\begin{array}{ccc} 1 & 0 & 0 \\ 6 & -7 & 0 \\ 2 & 1 & -2 \end{array}\right) \left(\begin{array}{ccc} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{array}\right),$$

the factorization of the given matrix.

The general formula for getting elements of L and U corresponding to the coefficient matrix A for a set of n linear equations can be written as

$$l_{ij} = a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj}, \quad i \ge j, \quad i = 1, 2, \dots, n$$

$$u_{ij} = \frac{1}{l_{ii}} [a_{ij} - \sum_{k=1}^{i-1} l_{ik} u_{kj}], \quad i < j, \quad j = 2, 3, \dots, n$$

$$l_{ij} = a_{i1}, \quad j = 1$$

$$u_{ij} = \frac{a_{ij}}{a_{11}}, \quad i = 1$$

$$(3.33)$$

Example 3.27 Solve the following linear system by LU decomposition using Crout's method

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 6 & 5 & 4 \\ 2 & 5 & 6 \end{pmatrix} \quad and \quad \mathbf{b} = \begin{pmatrix} 1 \\ -1 \\ 5 \end{pmatrix}.$$

Solution. Since the factorization of the coefficient matrix A has been already constructed in the Example (3.26) as

$$\begin{pmatrix} 1 & 2 & 3 \\ 6 & 5 & 4 \\ 2 & 5 & 6 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 6 & -7 & 0 \\ 2 & 1 & -2 \end{pmatrix} \begin{pmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{pmatrix}.$$

Then solving the first system $L\mathbf{y} = \mathbf{b}$ for unknown vector \mathbf{y} , that is

$$\begin{pmatrix} 1 & 0 & 0 \\ 6 & -7 & 0 \\ 2 & 1 & -2 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ 5 \end{pmatrix}.$$

Performing forward substitution yields

Then solving the second system $U\mathbf{x} = \mathbf{y}$ for unknown vector \mathbf{x} , that is

$$\begin{pmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}.$$

Performing backward substitution yields

and we obtained the approximate solution $\mathbf{x}^* = [-2, 3, -1]^T$ of the given system.

The above results can be reproduced by using MATLAB command as follows:

$$>> A = [1\ 2\ 3; 6\ 5\ 4; 2\ 5\ 6];\ b = [1\ -1\ 5];\ sol = Crout(A,b);$$

Program 3.9

MATLAB m-file for the Crout's Method

function sol = Crout(A, b)

```
function soi = Crout(A, b)  [n,n] = \operatorname{size}(A); \ u = \operatorname{zeros}(n,n); \ l = u; \ \text{for } i = 1:n; \ u(i,i) = 1; \ \text{end}; \ l(1,1) = A(1,1); \\ \text{for } i = 2:n; \ u(1,i) = A(1,i)/l(1,1); \ l(i,1) = A(i,1); \ \text{end}; \ \text{for } i = 2:n; \ \text{for } j = 2:n; \ s = 0; \\ \text{if } i <= j; \ K = i - 1; \ \text{else}; \ K = j - 1; \ \text{end}; \ \text{for } k = 1:K; \ s = s + l(i,k) * u(k,j); \ \text{end} \\ \text{if } j > i; \ u(i,j) = (A(i,j) - s)/l(i,i); \ \text{else} \ l(i,j) = A(i,j) - s; \ \text{end; end; end} \\ y(1) = b(1)/l(1,1); \ \text{for } k = 2:n; \ \text{sum} = b(k); \ \text{for } i = 1:k - 1; \ sum = sum - l(k,i) * y(i); \ \text{end} \\ y(k) = \operatorname{sum}/l(k,k); \ \text{end} \ x(n) = y(n)/u(n,n); \ \text{for } k = n - 1: -1:1; \ \text{sum} = y(k); \\ \text{for } i = k + 1:n; \ sum = sum - u(k,i) * x(i); \ \text{end}; \ x(k) = \operatorname{sum}/u(k,k); \ \text{end}; \ l; \ u; \ y; \ x; \\ \end{cases}
```

Note that we can also find the LU-decomposition of a matrix A by using simple Gauss-elimination method. We start with the product matrices of the form $\mathbf{I}A$ and convert it to the equivalent form LU, that is, we have to convert right matrix A to unit upper-triangular matrix U. We describe the procedure in the following example.

Example 3.28 Solve the following system using LU-decomposition by Crout's method

solution. The Crout's method makes LU factorization a byproduct of Gaussian elimination. To illustrate, let the given matrix of the system is

$$\left(\begin{array}{ccc}
1 & 2 & 0 \\
-1 & 0 & -2 \\
-3 & -5 & 1
\end{array}\right).$$

The process begins with the product matrices form

$$\mathbf{I}A = \left(\begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array}\right) \left(\begin{array}{ccc} 1 & 2 & 0 \\ -1 & 0 & -2 \\ -3 & -5 & 1 \end{array}\right).$$

In each of the steps below, we arrange so that the product of the two matrices is always equal to the original matrix A. Now the first step of Gaussian elimination on the right factor is to divide the first row by the pivot element. Then the Crout's rule copies the pivot element to the matching element of the left factor at the same time we divide. The next step in Gaussian elimination is to eliminate all the elements below the pivot element. This is done by multiplying the first row by below (n-1) eliminating elements, subtracting the product from the (n-1) rows, and putting the result in the (n-1) rows. The Crout's rule copies all those eliminating elements into the matching elements of the left factor. We repeat the same procedure for the remaining pivot elements. Thus we obtain

$$A = \left(\begin{array}{rrr} 1 & 0 & 0 \\ -1 & 1 & 0 \\ -3 & 0 & 1 \end{array}\right) \left(\begin{array}{rrr} 1 & 2 & 0 \\ 0 & 2 & -2 \\ 0 & 1 & 1 \end{array}\right).$$

The product of matrices is still equal to A.

$$A = \left(\begin{array}{rrr} 1 & 0 & 0 \\ -1 & 2 & 0 \\ -3 & 1 & 1 \end{array}\right) \left(\begin{array}{rrr} 1 & 2 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 2 \end{array}\right).$$

The product of matrices is still equal to A.

$$A = \begin{pmatrix} 1 & 0 & 0 \\ -1 & 2 & 0 \\ -3 & 1 & 2 \end{pmatrix} \begin{pmatrix} 1 & 2 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{pmatrix} = LU.$$

The product is still A and we have achieved the desired factorization. Now solving the first system $L\mathbf{y} = \mathbf{b}$ for unknown vector \mathbf{y} , that is

$$\begin{pmatrix} 1 & 0 & 0 \\ -1 & 2 & 0 \\ -3 & 1 & 2 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix}.$$

Performing forward substitution yields

Then solving the second system $U\mathbf{x} = \mathbf{y}$ for unknown vector \mathbf{x} , that is

$$\begin{pmatrix} 1 & 2 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 3 \\ 2 \\ 4 \end{pmatrix}.$$

 $Performing\ backward\ substitution\ yields$

and we obtained $\mathbf{x}^* = [-9, 6, 4]^T$, the approximate solution of the given system.

Note that determinant of A for LU decomposition by Crout's method is

$$\det(A) = \det(L) = \prod_{i=1}^{n} l_{ii} = (l_{11}l_{22} \cdots l_{nn}),$$

where det(U) = 1 because U is unit upper-triangular matrix.

Procedure 3.4 [LU Decomposition by the Crout's Method]

- 1. Take the nonsingular matrix, A.
- 2. If possible, decompose the matrix A = LU using (3.33).
- 3. Solve linear system $L\mathbf{y} = \mathbf{b}$ using (3.30).
- 4. Solve linear system $U\mathbf{x} = \mathbf{y}$ using (3.31).

Note that factorization method is also used to invert matrices. Their usefulness for this purpose is based on the fact that triangular matrices are easily inverted. Once the factorization has been effected, the inverse of a matrix A is found from the formula

$$A^{-1} = (LU)^{-1} = U^{-1}L^{-1}. (3.34)$$

Then

$$UA^{-1} = L^{-1}$$
, where $LL^{-1} = I$.

A practical way of calculating the determinant is to use the forward elimination process of the Gaussian elimination or, alternatively, the LU decomposition.

If no pivoting is used, calculation of the determinant using the LU decomposition is very easy.

Example 3.29 Find determinant and inverse of the following matrix using LU decomposition by Doolittle's method.

$$A = \left(\begin{array}{rrr} 1 & -2 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 2 \end{array}\right).$$

Solution. Since we know that

$$A = \begin{pmatrix} 1 & -2 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ m_{21} & 1 & 0 \\ m_{31} & m_{32} & 1 \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{pmatrix} = LU.$$

Now we will use only the forward elimination step of the simple Gaussian elimination method to convert the given matrix A into the upper-triangular matrix U. Since $a_{11} = 1 \neq 0$, so we wish to eliminate the elements $a_{21} = 1$ and $a_{31} = 1$ by subtracting from the second and third rows the appropriate multiples of the first row. The multiples are $m_{21} = 1$ and $m_{31} = 1$, so

$$\left(\begin{array}{ccc} 1 & -2 & 1 \\ 0 & 1 & 0 \\ 0 & 3 & 1 \end{array}\right).$$

As $a_{22}^{(1)}=1\neq 0$, therefore, we eliminate entry in $a_{32}^{(1)}=3$ position by subtracting the multiple $m_{32}=3$ of the second row from the third row, to get

$$\left(\begin{array}{ccc} 1 & -2 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array}\right).$$

Obviously, the original set of equations has been transformed to an upper-triangular form. Thus

$$\left(\begin{array}{ccc} 1 & -2 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & 2 \end{array}\right) = \left(\begin{array}{ccc} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 3 & 1 \end{array}\right) \left(\begin{array}{ccc} 1 & -2 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array}\right),$$

which is the required decomposition of A.

Now we find the determinant of the matrix A as follows:

$$\det(A) = \det(U) = u_{11}u_{22}u_{33} = (1)(1)(1) = 1.$$

To find the inverse of the matrix A, first we will compute the inverse of the lower-triangular matrix L^{-1} from

$$LL^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 3 & 1 \end{pmatrix} \begin{pmatrix} l'_{11} & 0 & 0 \\ l'_{21} & l'_{22} & 0 \\ l'_{31} & l'_{32} & l'_{33} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = I,$$

by using the forward substitution.

To solve the first system

$$\begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 3 & 1 \end{pmatrix} \begin{pmatrix} l'_{11} \\ l'_{21} \\ l'_{31} \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix},$$

by using forward substitution, we get, $l'_{11} = 1$, $l'_{21} = -1$, $l'_{31} = 2$. Similarly, the solution of the second linear system

$$\begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 3 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ l'_{22} \\ l'_{32} \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix},$$

can be obtained as, $l'_{22} = 1$, $l'_{32} = -3$. Finally, the solution of the third linear system

$$\begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 3 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ l'_{33} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix},$$

gives $l'_{33} = 1$. Hence the elements of the matrix L^{-1} are

$$L^{-1} = \left(\begin{array}{rrr} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 2 & -3 & 1 \end{array} \right),$$

which is the required inverse of the lower-triangular matrix L. To find the inverse of the given matrix A, we will solve the system

$$UA^{-1} = \begin{pmatrix} 1 & -2 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} a'_{11} & a'_{12} & a'_{13} \\ a'_{21} & a'_{22} & a'_{23} \\ a'_{31} & a'_{32} & a'_{33} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 2 & -3 & 1 \end{pmatrix} = L^{-1},$$

by using backward substitution.

We solve the first system

$$\begin{pmatrix} 1 & -2 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} a'_{11} \\ a'_{21} \\ a'_{31} \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix},$$

by using backward substitution, we get, $a'_{11} = -3$, $a'_{21} = -1$, $a'_{31} = 2$. Similarly, the solution of the second linear system

$$\begin{pmatrix} 1 & -2 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} a'_{12} \\ a'_{22} \\ a'_{32} \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ -3 \end{pmatrix},$$

can be obtained as, $a'_{12} = 5$, $a'_{22} = 1$, $a'_{32} = -3$. Finally, the solution of the third linear system

$$\begin{pmatrix} 1 & -2 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} a'_{13} \\ a'_{23} \\ a'_{33} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix},$$

can be obtained as, $a'_{13} = -1$, $a'_{23} = 0$, $a'_{33} = 1$. Hence the elements of the inverse matrix A^{-1} are

$$A^{-1} = \begin{pmatrix} -3 & 5 & -1 \\ -1 & 1 & 0 \\ 2 & -3 & 1 \end{pmatrix}.$$

For the LU decomposition we have not used pivoting for the sake of simplicity. However, pivoting is important for the same reason as in the Gaussian elimination. We know that pivoting in the Gaussian elimination is equivalent to interchanging the rows of coefficients matrix together with the terms on the right-hand side. This indicates that pivoting may be applied to the LU decomposition as long as the interchanging is applied to the left and right terms in the same way. When performing pivoting in the LU decomposition, the changes in the order of the rows are recorded. The same reordering is then applied to the right-hand side terms before starting the solution in accordance with the forward elimination and backward substitution steps.

Since we know that not every matrix has a direct LU decomposition. We define the following matrix which gives the sufficient condition for the LU decomposition of the matrix. It also, helps us for the convergence of the iterative methods for solving linear systems.

Definition 3.24 (Strictly Diagonally Dominant Matrix)

A square matrix is said to be strictly diagonally dominant (SDD) if the absolute value of each element on the main diagonal is greater than the sum of the absolute values of all the other elements in that row. Thus, strictly diagonally dominant matrix is defined as

$$|a_{ii}| > \sum_{\substack{j=1\\j\neq i}}^{n} |a_{ij}|, \quad for \quad i = 1, 2, \dots, n.$$
 (3.35)

Example 3.30 The matrix

$$A = \left(\begin{array}{rrr} 7 & 3 & 1\\ 1 & 6 & 3\\ -2 & 4 & 8 \end{array}\right),$$

is strictly diagonally dominant since

but the following matrix

$$B = \left(\begin{array}{ccc} 6 & -3 & 4\\ 3 & 7 & 3\\ 5 & -4 & 10 \end{array}\right),$$

is not strictly diagonally dominant since

$$|6| > |-3| + |4|$$
, that is, $6 > 7$,

which is not true. Strictly diagonally dominant matrix occurs naturally in a wide variety of practical applications, and when solving a strictly diagonally dominant system by Gauss-elimination method, partial pivoting is never required.

•

Theorem 3.14 If a matrix A is strictly diagonally dominant, then:

1. Matrix A is nonsingular.

•

- 2. Gaussian elimination without row interchange can be performed on the linear system $A\mathbf{x} = \mathbf{b}$.
- 3. Matrix A has LU factorization.

Example 3.31 Solve the following linear system using the simple Gaussian elimination method and also, find the LU decomposition of the matrix using Doolittle's method and Crout's method

$$5x_1 + x_2 + x_3 = 7$$

 $2x_1 + 6x_2 + x_3 = 9$
 $x_1 + 2x_2 + 9x_3 = 12$

solution. Start with the augmented matrix form

$$\left(\begin{array}{ccccc} 5 & 1 & 1 & \vdots & 7 \\ 2 & 6 & 1 & \vdots & 9 \\ 1 & 2 & 9 & \vdots & 12 \end{array}\right),$$

and since $a_{11} = 5 \neq 0$, so we can eliminate the elements a_{21} and a_{31} by subtracting from the second and third rows the appropriate multiples of the first row. In this case the multiples are given

$$m_{21} = \frac{2}{5}$$
 and $m_{31} = \frac{1}{5}$.

Hence

$$\begin{pmatrix}
5 & 1 & 1 & \vdots & 7 \\
0 & 28/5 & 3/5 & \vdots & 31/5 \\
0 & 9/5 & 44/5 & \vdots & 53/5
\end{pmatrix}.$$

As $a_{22}^{(1)} = 28/5 \neq 0$, therefore, we eliminate entry in $a_{32}^{(1)}$ position by subtracting the multiple $m_{32} = \frac{1.8}{5.6} = 9/28$ of the second row from the third row, to get

$$\begin{pmatrix}
5 & 1 & 1 & \vdots & 7 \\
0 & 28/5 & 3/5 & \vdots & 31/5 \\
0 & 0 & 43/5 & \vdots & 43/5
\end{pmatrix}.$$

Obviously, the original set of equations has been transformed to an upper-triangular form. Since all the diagonal elements of the obtaining upper-triangular matrix are nonzero, which means that the coefficient matrix of the given system is nonsingular and therefore, the given system has a unique solution. Now expressing the set in algebraic form yields

$$5x_1 + x_2 + x_3 = 7$$

 $(28/5)x_2 + (3/5)x_3 = 31/5$
 $(43/5)x_3 = 43/5$

Now using backward substitution to get the solution of the system as

$$\begin{array}{rclrcl} (43/5)x_3 & = & 43/5, & gives & x_3 = 1, \\ (28/5)x_2 & = & -(3/5)x_3 + 31/5, & gives & x_2 = 1, \\ 5x_1 & = & 7 - x_2 - x_3, & gives & x_1 = 1. \end{array}$$

Since we know that in using LU decomposition by Doolittle's method the unknown elements of matrix L are the used multiples and the matrix U is same as we obtained in forward elimination process of the simple Gauss elimination. Thus the LU decomposition of the matrix A can be obtained by using Doolittle's method as follows:

$$A = \begin{pmatrix} 5 & 1 & 1 \\ 2 & 6 & 1 \\ 1 & 2 & 9 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 2/5 & 1 & 0 \\ 1/5 & 9/28 & 1 \end{pmatrix} \begin{pmatrix} 5 & 1 & 1 \\ 0 & 28/5 & 3/5 \\ 0 & 0 & 43/5 \end{pmatrix} = LU.$$

Similarly, the LU decomposition of the matrix A by using Crout's method can be obtained as

$$A = \begin{pmatrix} 5 & 1 & 1 \\ 2 & 6 & 1 \\ 1 & 2 & 9 \end{pmatrix} = \begin{pmatrix} 5 & 0 & 0 \\ 2 & 28/5 & 0 \\ 1 & 9/5 & 43/5 \end{pmatrix} \begin{pmatrix} 1 & 1/5 & 1/5 \\ 0 & 1 & 1/10 \\ 0 & 0 & 1 \end{pmatrix} = LU.$$

Thus the conditions of the Theorem 3.14 are satisfied.

3.5 Norms of Vectors and Matrices

For solving linear systems, we discuss a method for quantitatively measuring the distance between vectors in \mathbb{R}^n , the set of all column vectors with real components, to determine whether the sequence of vectors that results from using an direct method converges to a solution of the system. To define a distance in \mathbb{R}^n , we use the notation of the *norm* of a vector.

Vector Norms

It is sometimes useful to have a scalar measure of the magnitude of a vector. Such a measure is called a *vector norm* and for a vector \mathbf{x} is written as $\|\mathbf{x}\|$.

A vector norm on \mathbb{R}^n is a function, from \mathbb{R}^n to \mathbb{R} satisfying:

- 1. $\|\mathbf{x}\| > 0$ for all $\mathbf{x} \in \mathbf{R}^n$.
- 2. $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = \mathbf{0}$.
- 3. $\|\alpha \mathbf{x}\| = |\alpha| \|\mathbf{x}\|$, for all $\alpha \in \mathbf{R}$, $\mathbf{x} \in \mathbf{R}^n$.
- 4. $\|\mathbf{x} + \mathbf{y}\| < \|\mathbf{x}\| + \|\mathbf{y}\|$, for all $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$.

There are three norms in \mathbb{R}^n that are most commonly used in applications, called l_1 -norm, l_2 -norm, and l_{∞} -norm, and are defined for the given vectors $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$ as

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|, \quad \|\mathbf{x}\|_2 = \left(\sum_{i=1}^n x_i^2\right)^{1/2}, \quad \|\mathbf{x}\|_{\infty} = \max_{1 \le i \le n} |x_i|.$$

The l_1 -norm is called the *absolute norm*, the l_2 -norm is frequently called the *Euclidean norm* as it is just the formula for distance in ordinary three-dimensional Euclidean space extended to dimension n. Finally, the l_{∞} -norm is called the *maximum norm* or occasionally the *uniform norm*. All these three norms are also called the *natural norms*.

Example 3.32 Compute l_p -norms $(p = 1, 2, \infty)$ of the vector $\mathbf{x} = [-5, 3, -2]^T$ in \mathbf{R}^3 .

Solution. These l_p -norms $(p = 1, 2, \infty)$ of the given vector are:

$$\|\mathbf{x}\|_{1} = |x_{1}| + |x_{2}| + |x_{3}| = |-5| + |3| + |-2| = 10.$$

$$\|\mathbf{x}\|_{2} = (x_{1}^{2} + x_{2}^{2} + x_{3}^{2})^{1/2} = \left[(-5)^{2} + (3)^{2} + (-2)^{2} \right]^{1/2} \approx 6.1644.$$

$$\|\mathbf{x}\|_{\infty} = \max\{|x_{1}|, |x_{2}|, |x_{3}|\} = \max\{|-5|, |3|, |-2|\} = 5.$$

In MATLAB command the built-in **norm** function computes l_p -norms of vectors. If only one argument is passed to norm, the l_2 -norm is returned and for two arguments, the second one is used to specify the value of p. For example,

$$>> x = [-5 \ 3 \ -2]; \ v = norm(x); \ v = norm(x,2); \ v = norm(x,1), \ v = norm(x,inf)$$

The internal MATLAB constant inf is used to select the l_{∞} -norm.

Matrix Norms

A matrix norm is a measure of how well one matrix approximates another, or, more accurately, of how well their difference approximates the zero matrix. An iterative procedure for inverting a matrix produces a sequence of approximate inverses. Since in practices such a process must be terminated, it is desirable to have some measure of the error of approximate inverse.

So a matrix norm on the set of all $n \times n$ matrices is a real-valued function, $\|.\|$, defined on this set, satisfying for all $n \times n$ matrices A and B and all real number α as follows:

- 1. ||A|| > 0, $A \neq \mathbf{0}$.
- 2. ||A|| = 0, $A = \mathbf{0}$.
- 3. ||I|| = 1, I is the identity matrix.
- 4. $\|\alpha A\| = |\alpha| \|A\|$, for some scalar $\alpha \in \mathbf{R}$.
- 5. $||A + B|| \le ||A|| + ||B||$.
- 6. $||AB|| \le ||A|| ||B||$.
- 7. $||A B|| \ge |||A|| ||B|||$.

Several norms for matrices have been defined, we shall use the following three natural norms l_1, l_2 , and l_{∞} for a square matrix of order n:

$$||A||_1 = \max_j \left(\sum_{i=1}^n |a_{ij}|\right) = \text{maximum column-sum.}$$

$$||A||_2 = \max_{||x||_2=1} ||A\mathbf{x}||_2 = \text{spectral norm.}$$

$$||A||_{\infty} = \max_{i} \left(\sum_{j=1}^{n} |a_{ij}| \right) = \text{row-sum norm.}$$

The l_1 -norm and l_{∞} -norm are widely used because they are easy to calculate. The matrix norm $||A||_2$ that corresponds to the l_2 -norm is related the eigenvalues of the matrix. It sometimes has special utility because no other norm is smaller than this norm. It therefore, provides the best measure of the size of a matrix, but is also the most difficult to compute. We will discuss this natural norm later in the chapter.

For $m \times n$ matrix, we can paraphrase the *Frobenius norm* (or *Euclidean norm*), which is not a natural norm and is define as

$$||A||_F = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2\right)^{1/2}.$$

It can be shown that

$$||A||_F = \sqrt{tr(A^T A)},$$

where $tr(A^TA)$ is the trace of a matrix A^TA , that is, the sum of the diagonal entries of A^TA . The Frobenius norm of a matrix is a good measure of the magnitude of a matrix. It is to be noted that $||A||_F \neq ||A||_2$. For a diagonal matrix, all norms have the same values.

Example 3.33 Compute l_p -norms $(p = 1, \infty, F)$ of the following matrix

$$A = \left(\begin{array}{ccc} 4 & 2 & -1 \\ 3 & 5 & -2 \\ 1 & -2 & 7 \end{array}\right).$$

Solution. The l_1 -norm is defined as

$$\sum_{i=1}^{3} |a_{i1}| = |4| + |3| + |1| = 8,$$

$$\sum_{i=1}^{3} |a_{i2}| = |2| + |5| + |-2| = 9,$$

$$\sum_{i=1}^{3} |a_{i3}| = |-1| + |-2| + |7| = 10,$$

so

$$||A||_1 = max\{8, 9, 10\} = 10.$$

Also, l_{∞} -norm is defined as

$$\sum_{j=1}^{3} |a_{1j}| = |4| + |2| + |-1| = 7,$$

$$\sum_{j=1}^{3} |a_{2j}| = |3| + |5| + |-2| = 10,$$

$$\sum_{j=1}^{3} |a_{3j}| = |1| + |-2| + |7| = 10,$$

so

$$||A||_{\infty} = max\{7, 10, 10\} = 10.$$

Finally, we have the l_F -norm of the matrix as

$$||A||_F = (16 + 4 + 1 + 9 + 25 + 4 + 1 + 4 + 49)^{1/2} \approx 10.6301,$$

the Frobenius norm of the given matrix.

Like l_p -norms of vectors, in MATLAB command the built-in **norm** function can be used to compute l_p -norms of matrices. The l_1 -norm, l_{∞} -norm and Frobenius norm of a matrix can be find as

$$>> A = [4\ 2\ -1; 3\ 5\ -2; 1\ -2\ -7];\ norm(A,1);\ norm(A,inf);\ norm(A,inf)$$

3.6 Iterative Methods for Solving Linear Systems

The methods discussed in the previous section for the solution of the system of linear equations have been direct, which required a finite number of arithmetic operations. The elimination methods of solving such systems usually yield sufficiently accurate solutions for approximately 20 to 25 simultaneous equations, where most of the unknowns are present in all of the equations. When the coefficients matrix is sparse (has many zeros), a considerably large number of equations can be handled by the elimination methods. But these methods are generally impractical when many hundreds or thousands of equations must be solved simultaneously.

There are, however, several methods which can be used to solve large numbers of simultaneous equations. These methods are, called *iterative methods* by which an approximation to the solution of a system of linear equations may be obtained. The iterative methods are used most often for large sparse systems of linear equations and efficient in terms of computer storage and time requirement. Systems of this type arise frequently in the numerical solution of boundary value problems and partial differential equations. Unlike the direct methods, the iterative methods may not always yield a solution, even if the determinant of the coefficients matrix is not zero. Here, we consider just two of these iterative methods. These two forms the basis of a family of methods which are designed either to accelerate the convergence or to suit some particular computer architecture.

3.6.1 Jacobi Iterative Method

This is one of the easiest iterative method to find the approximate solution of the system of linear equations (3.42). To explain its procedure, consider a system of three linear equations as follows:

The solution process starts by solving for the first variable x_1 from first equation, second variable x_2 from second equation and third variable x_3 from third equation, gives

Divide both sides of the above three equations by their diagonal elements, a_{11} , a_{22} and a_{33} respectively, to have

$$x_{1} = \frac{1}{a_{11}} \begin{bmatrix} b_{1} - a_{12}x_{2} - a_{13}x_{3} \end{bmatrix}$$

$$x_{2} = \frac{1}{a_{22}} \begin{bmatrix} b_{2} - a_{21}x_{1} - a_{23}x_{3} \end{bmatrix}$$

$$x_{3} = \frac{1}{a_{33}} \begin{bmatrix} b_{3} - a_{31}x_{1} - a_{32}x_{2} \end{bmatrix}$$

Let $\mathbf{x}^{(k)} = \left[x_1^{(k)}, x_2^{(k)}, x_3^{(k)}\right]^T$ be an initial solution of the exact solution \mathbf{x} of the linear system (3.42), then we define an iterative sequence

$$x_{1}^{(k+1)} = \frac{1}{a_{11}} \left[b_{1} - a_{12} x_{2}^{(k)} - a_{13} x_{3}^{(k)} \right]$$

$$x_{2}^{(k+1)} = \frac{1}{a_{22}} \left[b_{2} - a_{21} x_{1}^{(k)} - a_{23} x_{3}^{(k)} \right]$$

$$x_{3}^{(k+1)} = \frac{1}{a_{33}} \left[b_{3} - a_{31} x_{1}^{(k)} - a_{32} x_{2}^{(k)} \right]$$

$$(3.36)$$

where k is the number of iterative steps. Then the form (3.36) is called the Jacobi formula for system of three equations. For a general system of n linear equations, the Jacobi method is defined by

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left[b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k)} - \sum_{j=i+1}^{n} a_{ij} x_j^{(k)} \right]$$

$$i = 1, 2, \dots, n, \quad k = 0, 1, 2, \dots,$$
(3.37)

provided that the diagonal elements $a_{ii} \neq 0$ for each i = 1, 2, ..., n. If the diagonal elements equal to zero, then reordering of the equations can be performed so that no element in the diagonal

position equal to zero. As usual with iterative methods, an initial approximation $x_i^{(0)}$ must be supplied. If we don't have knowledge about the exact solution, it is conventional to start with $x_i^{(0)} = \mathbf{0}$ for all i. The iterations defined by (3.37) are stopped when

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < \epsilon,\tag{3.38}$$

or by using other possible stopping criteria

$$\frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|}{\|\mathbf{x}^{(k+1)}\|} < \epsilon, \tag{3.39}$$

where ϵ is a preassigned small positive number. For this purpose, any convenient norm can be used, the most usual being the l_{∞} -norm.

Example 3.34 Solve the following system of equations using the Jacobi iterative method, using $\epsilon = 10^{-6}$ in the l_{∞} -norm.

$$5x_1 - x_2 + x_3 = 10$$

 $2x_1 + 8x_2 - x_3 = 11$
 $-x_1 + x_2 + 4x_3 = 3$

Start with the initial solution $\mathbf{x}^{(0)} = [0, 0, 0]^T$.

Solution. The Jacobi iterative method for the given system has the form

$$x_1^{(k+1)} = \frac{1}{5} \begin{bmatrix} 10 + x_2^{(k)} - x_3^{(k)} \end{bmatrix}$$

$$x_2^{(k+1)} = \frac{1}{8} \begin{bmatrix} 11 - 2x_1^{(k)} + x_3^{(k)} \end{bmatrix}$$

$$x_3^{(k+1)} = \frac{1}{4} \begin{bmatrix} 3 + x_1^{(k)} - x_2^{(k)} \end{bmatrix}$$

and starting with initial approximation $x_1^{(0)} = 0, x_2^{(0)} = 0, x_3^{(0)} = 0$, then for k = 0, we obtain

$$x_1^{(1)} = \frac{1}{5} \begin{bmatrix} 10 + x_2^{(0)} - x_3^{(0)} \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 10 + 0 - 0 \end{bmatrix} = 2,$$

$$x_2^{(1)} = \frac{1}{8} \begin{bmatrix} 11 - 2x_1^{(0)} + x_3^{(0)} \end{bmatrix} = \frac{1}{8} \begin{bmatrix} 11 - 0 + 0 \end{bmatrix} = 1.375,$$

$$x_3^{(1)} = \frac{1}{4} \begin{bmatrix} 3 + x_1^{(0)} - x_2^{(0)} \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 3 + 0 - 0 \end{bmatrix} = 0.75.$$

The first and subsequent iterations are listed in Table 3.1.

Note that the Jacobi method converges and after 16 iterations we obtained what is obviously the exact solution. Ideally the iteration should stop automatically when we obtained the required accuracy using one of the stopping criteria mentioned by (3.38) or (3.39).

To get the above results using MATLAB command, we do the following:

$$>> Ab = [A|b] = [5 -1 \ 1 \ 10; 2 \ 8 -1 \ 11; -1 \ 1 \ 4 \ 3];$$

 $>> x = [0 \ 0 \ 0]; \ acc = 0.5e - 6; \ JacobiM(Ab, x, acc);$

Table 3.1: Solution of the Example 3.34							
	k	$x_1^{(k)}$	$x_2^{(k)}$	$x_3^{(k)}$			
	0	0.000000	0.000000	0.000000			
	1	2.000000	1.375000	0.750000			
	2	2.125000	0.968750	0.906250			
	÷	:	:	:			
	15	2.000000	0.999999	1.000000			
	16	2.000000	1.000000	1.000000			

Example 3.35 Solve the following system of equations using the Jacobi iterative method.

$$\begin{array}{rclrcrcr}
2x_1 & + & 8x_2 & - & x_3 & = & 11 \\
5x_1 & - & x_2 & + & x_3 & = & 10 \\
-x_1 & + & x_2 & + & 4x_3 & = & 3
\end{array}$$

Start with the initial solution $\mathbf{x}^{(0)} = [0, 0, 0]^T$.

Solution. Results for this linear system are listed in Table 3.2.

Notice that Jacobi method diverges rapidly. Although the given linear system is same as the linear

Table 3.2: Solution of the Example 3.35

Table 3.2. Solderon of the Endingle 3.33							
	k	$x_1^{(k)}$	$x_2^{(k)}$	$x_3^{(k)}$			
	0	0.000000	0.000000	0.000000			
	1	5.500000	-10.0000	0.750000			
	2	45.87500	18.25000	4.625000			
	3	-65.1875	224.0000	7.656250			

system of the previous Example 3.34 except the first and second equations are interchanged. From this example we concluded that Jacobi iterative method is not always convergent.

Program 3.10

MATLAB m-file for the Jacobi Iterative Method for Linear System function x=JacobiM(Ab,x,acc)

[n,t]=size(Ab); b=Ab(1:n,t); R=1; k=1; d(1,1:n+1)=[0 x]; while R > acc

for i=1:n; sum=0; for j=1:n; if j = i

sum = sum + Ab(i, j) * d(k, j + 1); end; x(1, i) = (1/Ab(i, i)) * (b(i, 1) - sum);end;end

k=k+1; d(k,1:n+1)=[k-1 x]; $R=\max(abs((d(k,2:n+1)-d(k-1,2:n+1))))$;

if k > 10 & R > 100 ('Jacobi Method is diverges') break; end; end; x=d;

Procedure 3.5 [Jacobi Method]

1. Check the coefficient matrix A is strictly diagonally dominant (for quaranteed convergence).

- 2. Initialize the first approximation $\mathbf{x}^{(0)}$ and pre-assigned accuracy ϵ .
- 3. Compute the constant $\mathbf{c} = D^{-1}\mathbf{b} = \frac{b_i}{a_{ii}}$, for $i = 1, 2, \dots, n$.
- 4. Compute the Jacobi iteration matrix $T_J = -D^{-1}(L+U)$.
- 5. Solve for the approximate solutions $\mathbf{x}_i^{(k+1)} = T_J \mathbf{x}_i^{(k)} + \mathbf{c}, \quad i = 1, 2, \dots, n$ and $k = 0, 1, \dots$
- 6. Repeat step 5 until $\|\mathbf{x}_i^{(k+1)} \mathbf{x}_i^{(k)}\| < \epsilon$.

3.6.2 Gauss-Seidel Iterative Method

This is one of the most popular and widely used iterative method to find the approximate solution of the system of linear equations. This iterative method is a modification of the Jacobi iterative method and give us good accuracy by using the most recently calculated values.

From the Jacobi iterative formula (3.37), it is seen that the new estimates for solution \mathbf{x} are computed from the old estimates and only when all the new estimates have been determined are then used in the right-hand side of the equation to perform the next iteration. But the Gauss-Seidel method is to make use of the new estimates in the right-hand side of the equation as soon as they become available. For example, the Gauss-Seidel formula for the system of three equations can be define an iterative sequence

$$x_1^{(k+1)} = \frac{1}{a_{11}} \left[b_1 - a_{12} x_2^{(k)} - a_{13} x_3^{(k)} \right]$$

$$x_2^{(k+1)} = \frac{1}{a_{22}} \left[b_2 - a_{21} x_1^{(k+1)} - a_{23} x_3^{(k)} \right]$$

$$x_3^{(k+1)} = \frac{1}{a_{22}} \left[b_3 - a_{31} x_1^{(k+1)} - a_{32} x_2^{(k+1)} \right]$$
(3.40)

For a general system of n linear equations, the Gauss-Seidel iterative method defined as

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left[b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^{n} a_{ij} x_j^{(k)} \right]$$

$$i = 1, 2, \dots, n, \quad k = 0, 1, 2, \dots$$
(3.41)

The Gauss-Seidel iterative method is sometimes called the method of *successive iteration*, because the most recent values of all \mathbf{x}_i are used in the calculation.

Example 3.36 Solve the following system of equations using the Gauss-Seidel iterative method, with $\epsilon = 10^{-6}$ in l_{∞} -norm.

Start with the initial solution $\mathbf{x}^{(0)} = [0, 0, 0]^T$.

Solution. The Gauss-Seidel iteration for the given system is

$$x_1^{(k+1)} = \frac{1}{5} \Big[10 + x_2^{(k)} - x_3^{(k)} \Big]$$

$$x_2^{(k+1)} = \frac{1}{8} \Big[11 - 2x_1^{(k+1)} + x_3^{(k)} \Big]$$

$$x_3^{(k+1)} = \frac{1}{4} \Big[3 + x_1^{(k+1)} - x_2^{(k+1)} \Big]$$

and starting with initial approximation $x_1^{(0)} = 0, x_2^{(0)} = 0, x_3^{(0)} = 0$, then for k = 0, we obtain

$$x_1^{(1)} = \frac{1}{5} \begin{bmatrix} 10 + x_2^{(0)} - x_3^{(0)} \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 10 + 0 - 0 \end{bmatrix} = 2,$$

$$x_2^{(1)} = \frac{1}{8} \begin{bmatrix} 11 - 2x_1^{(1)} + x_3^{(0)} \end{bmatrix} = \frac{1}{8} \begin{bmatrix} 11 - 4 + 0 \end{bmatrix} = 0.875,$$

$$x_3^{(1)} = \frac{1}{4} \begin{bmatrix} 3 + x_1^{(1)} - x_2^{(1)} \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 3 + 2 - 0.875 \end{bmatrix} = 1.03125.$$

The first and subsequent iterations are listed in Table 3.3.

The above results can be obtained using MATLAB command as follows:

>>
$$Ab = [A|b] = [5 -1 \ 1 \ 10; 2 \ 8 -1 \ 11; -1 \ 1 \ 4 \ 3];$$

>> $x = [0 \ 0 \ 0]; \ acc = 0.5e - 6; \ GaussSM(Ab, x, acc);$

Note that the Gauss-Seidel method converged and required 10 iterations to obtain the correct solution for the given system, which is 6 iterations less than required by the Jacobi method for the same Example 3.34.

Table 3.3: Solution of the Example 3.36

k	$x_1^{(k)}$	$x_{2}^{(k)}$	$x_3^{(k)}$
0	0.000000	0.000000	0.000000
1	2.000000	0.875000	1.031250
2	1.968750	1.011719	0.989258
3	2.004492	0.997534	1.001740
:	:	:	:
9	2.000000	0.999999	1.000000
10	2.000000	1.000000	1.000000

Program 3.11

MATLAB m-file for the Gauss-Seidel Iterative Method for Linear System function x=GaussSM(Ab,x,acc)

 $[n,t] = \text{size}(Ab); \ b = Ab(1:n,t); R = 1; \ k = 1; \ d(1,1:n+1) = [0 \ x]; \ k = k+1; \ \text{while} \ R > acc$ for i=1:n; sum=0; for j=1:n if j <= i-1; sum = sum + Ab(i,j) * d(k,j+1); elseif j >= i+1 sum = sum + Ab(i,j) * d(k - 1, j + 1); end; end; x(1,i) = (1/Ab(i,i)) * (b(i,1) - sum); d(k,1)=k-1; d(k,i+1)=x(1,i); end; R=max(abs((d(k,2:n+1)-d(k-1,2:n+1)))); k=k+1; if $R > 100 \ \& k > 10;$ ('Gauss-Seidel method is Diverges') break ;end;end;x=d;

Example 3.37 Solve the following system of equations using the Gauss-Seidel iterative method.

$$2x_1 + 8x_2 - x_3 = 11$$

 $5x_1 - x_2 + x_3 = 10$
 $-x_1 + x_2 + 4x_3 = 3$

Start with the initial solution $\mathbf{x}^{(0)} = [0, 0, 0]^T$.

Solution. Results for this linear system are listed in Table 3.4. Note that in this case the Gauss-Seidel method diverges rapidly. Although the given linear system is same as the linear system of the previous Example 3.36 except the first and second equations are interchanged. From this example we concluded that the Gauss-Seidel iterative method is not always convergent.

Table 3.4: Solution of the Example 3.37

k	$x_1^{(k)}$	$x_{2}^{(k)}$	$x_3^{(k)}$		
0	0.000000	0.000000	0.000000		
1	5.500000	17.5000	-2.25000		
2	-65.6250	-340.375	69.43750		
3	1401.719	7068.031	-1415.83		

Procedure 3.6 [Gauss-Seidel Method]

- 1. Check the coefficient matrix A is strictly diagonally dominant (for guaranteed convergence).
- 2. Initialize the first approximation $\mathbf{x}^{(0)} \in \mathbf{R}$ and pre-assigned accuracy ϵ .
- 3. Compute the constant $\mathbf{c} = (D+L)^{-1}\mathbf{b}$.
- 4. Compute the Gauss-Seidel iteration matrix $T_G = -(D+L)^{-1}U$.
- 5. Solve for the approximate solutions $x_i^{(k+1)} = T_G x_i^{(k)} + \mathbf{c}, \quad i = 1, 2, \dots, n$ and $k = 0, 1, \dots$
- 6. Repeat step 5 until $\|\mathbf{x}_i^{(k+1)} \mathbf{x}_i^{(k)}\| < \epsilon$.

From the Examples 3.34 and (3.36), we noted that the solution by the Gauss-Seidel method converges more quickly than the Jacobi method. In general, we may state that **if both the Jacobi method and the Gauss-Seidel method are converge, then the Gauss-Seidel method converges more quickly**. This is generally the case but not always true. In fact, there are some linear systems for which the Jacobi method converges but the Gauss-Seidel method does not, and others for which the Gauss-Seidel method converges but the Jacobi method does not.

3.6.3 Matrix Forms of Iterative Methods for Linear System

The iterative methods to solve the system of linear equations

$$A\mathbf{x} = \mathbf{b},\tag{3.42}$$

start with an initial approximation $\mathbf{x}^{(0)} \in \mathbf{R}$ to the solution \mathbf{x} of the linear system (3.42), and generates a sequence of vectors $\{\mathbf{x}^{(k)}\}_{k=0}^{\infty}$ that converges to \mathbf{x} . Most of these iterative methods involve a process that converts the system (3.42) into an equivalent system of the form

$$\mathbf{x} = T\mathbf{x} + \mathbf{c},\tag{3.43}$$

for some square matrix T and vector \mathbf{c} . After the initial vector $\mathbf{x}^{(0)}$ is selected, the sequence of approximate solutions vector is generated by computing

$$\mathbf{x}^{(k+1)} = T\mathbf{x}^{(k)} + \mathbf{c}, \quad \text{for} \quad k = 0, 1, 2, \dots$$
 (3.44)

The sequence is terminated when the error is sufficiently small, that is

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < \epsilon$$
, for small positive ϵ . (3.45)

Note that a matrix T is called iteration matrix and a vector \mathbf{c} is a column matrix. We can find the forms of these matrices easily for both iterative methods as follows. Let a matrix A can be written as

$$A = L + D + U, (3.46)$$

where L is strictly lower-triangular, U is strictly upper-triangular, and D is the diagonal parts of the coefficients matrix A, that is

$$L = \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 \\ a_{21} & 0 & 0 & \cdots & 0 \\ a_{31} & a_{32} & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \cdots & 0 \end{pmatrix}, \qquad U = \begin{pmatrix} 0 & a_{12} & a_{13} & \cdots & a_{1n} \\ 0 & 0 & a_{23} & \cdots & a_{2n} \\ 0 & 0 & 0 & \cdots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 \end{pmatrix},$$

and

$$D = \begin{pmatrix} a_{11} & 0 & 0 & \cdots & 0 \\ 0 & a_{22} & 0 & \cdots & 0 \\ 0 & 0 & a_{33} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & a_{nn} \end{pmatrix}.$$

Then the linear system (3.42) can be written as

$$(L+D+U)\mathbf{x} = \mathbf{b}. (3.47)$$

Now we find forms of both matrices T and c which help us to solve the linear system.

Jacobi Iterative Method

The equation (3.47) can be written as

$$D\mathbf{x} = -(L+U)\mathbf{x} + \mathbf{b}.$$

Since matrix D is nonsingular, so we can write above equation as

$$\mathbf{x} = -D^{-1}(L+U)\mathbf{x} + D^{-1}\mathbf{b},$$

which can be put in the form

$$\mathbf{x}^{(k+1)} = T_J \mathbf{x}^{(k)} + \mathbf{c}_j, \quad \text{for} \quad k = 0, 1, 2, \dots,$$
 (3.48)

which is the matrix form of Jacobi iterative method and where

$$T_J = -D^{-1}(L+U)$$
 and $\mathbf{c}_i = D^{-1}\mathbf{b}$, (3.49)

are called Jacobi iteration matrix and Jacobi constant column matrix, respectively and and their elements are defined by

$$t_{ij} = \frac{a_{ij}}{a_{ii}},$$
 $i, j = 1, 2, \dots, n, \quad i \neq j,$ $t_{ij} = 0,$ $i = j,$ $c_i = \frac{b_i}{a_{ii}},$ $i = 1, 2, \dots, n.$

Note that the diagonal elements of Jacobi iteration matrix T_J are always zero.

Gauss-Seidel Iterative Method

The equation (3.47) can also be written as

$$(L+D)\mathbf{x} = -U\mathbf{x} + \mathbf{b}.$$

Since lower-triangular matrix (L+D) is nonsingular, so we can write above equation as

$$\mathbf{x} = -(L+D)^{-1}U\mathbf{x} + (L+D)^{-1}\mathbf{b},$$

which can be put in the form

$$\mathbf{x}^{(k+1)} = T_G \mathbf{x}^{(k)} + \mathbf{c}_G, \quad \text{for} \quad k = 0, 1, 2, \dots,$$
 (3.50)

which is the matrix form of Gauss-Seidel iterative method and where

$$T_G = -(L+D)^{-1}U$$
 and $\mathbf{c}_G = (L+D)^{-1}\mathbf{b}$, (3.51)

are called Gauss-Seidel iteration matrix and Gauss-seidel constant column matrix, respectively.

Example 3.38 Consider the following system

$$6x_1 + 2x_2 = 1$$

 $x_1 + 7x_2 - 2x_3 = 2$
 $3x_1 - 2x_2 + 9x_3 = -1$

- (a) Find the matrix form of iterative (Jacobi and Gauss-Seidel) methods.
- (b) If $\mathbf{x}^{(k)} = [x_1^{(k)}, x_2^{(k)}, x_3^{(k)}]^T$, then writing the iterative forms of part(a) in the component forms and find the exact solution of the given system.
- (c) Find formulas for the error $e^{(k+1)}$ in the (n+1)th step.
- (d) Find the second approximation of the error $e^{(2)}$ using part (c) if $\mathbf{x}^{(0)} = [0, 0, 0]^T$.

Solution. Since the given matrix A is

$$A = \left(\begin{array}{ccc} 6 & 2 & 0 \\ 1 & 7 & -2 \\ 3 & -2 & 9 \end{array}\right),$$

and so

$$A = L + U + D = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 3 & -2 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 2 & 0 \\ 0 & 0 & -2 \\ 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} 6 & 0 & 0 \\ 0 & 7 & 0 \\ 0 & 0 & 9 \end{pmatrix}.$$

Jacobi Iterative Method

(a) Since the matrix form of the Jacobi iterative method can be written as

$$\mathbf{x}^{(k+1)} = T_J \mathbf{x}^{(k)} + \mathbf{c}_J, \qquad k = 0, 1, 2, \dots,$$

where

$$T_J = -D^{-1}(L+U)$$
 and $\mathbf{c}_J = D^{-1}\mathbf{b}$.

One can easily compute the Jacobi iteration matrix T_J and the vector \mathbf{c}_J as follows:

$$T_{J} = -\begin{pmatrix} \frac{1}{6} & 0 & 0 \\ 0 & \frac{1}{7} & 0 \\ 0 & 0 & \frac{1}{9} \end{pmatrix} \begin{pmatrix} 0 & 2 & 0 \\ 1 & 0 & -2 \\ 3 & -2 & 0 \end{pmatrix} = \begin{pmatrix} 0 & -\frac{2}{6} & 0 \\ -\frac{1}{7} & 0 & \frac{2}{7} \\ -\frac{3}{9} & \frac{2}{9} & 0 \end{pmatrix} \quad and \quad \mathbf{c} = \begin{pmatrix} \frac{1}{6} \\ \frac{2}{7} \\ -\frac{1}{9} \end{pmatrix}.$$

Thus the matrix form of Jacobi iterative method is

$$\mathbf{x}^{(k+1)} = \begin{pmatrix} 0 & -\frac{2}{6} & 0 \\ -\frac{1}{7} & 0 & \frac{2}{7} \\ -\frac{3}{9} & \frac{2}{9} & 0 \end{pmatrix} \mathbf{x}^{(k)} + \begin{pmatrix} \frac{1}{6} \\ \frac{2}{7} \\ -\frac{1}{9} \end{pmatrix}, \quad k = 0, 1, 2.$$

(b) Now by writing the above iterative matrix form of in the component form, we have

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 & -\frac{1}{3} & 0 \\ -\frac{1}{7} & 0 & \frac{2}{7} \\ -\frac{1}{3} & \frac{2}{9} & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} + \begin{pmatrix} \frac{1}{6} \\ \frac{2}{7} \\ -\frac{1}{9} \end{pmatrix},$$

and it is equivalent to

Now solving for x_1, x_2 and x_3 , we get

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 1/12 \\ 1/4 \\ -1/12 \end{pmatrix},$$

which is the exact solution of the given system.

(c) Since the error in the (n+1)th step is defined as

$$\mathbf{e}^{(k+1)} = \mathbf{x} - \mathbf{x}^{(k+1)}.$$

therefore, we have

$$\mathbf{e}^{(k+1)} = \begin{pmatrix} 1/12 \\ 1/4 \\ -1/12 \end{pmatrix} - \begin{pmatrix} 0 & -\frac{2}{6} & 0 \\ -\frac{1}{7} & 0 & \frac{2}{7} \\ -\frac{3}{9} & \frac{2}{9} & 0 \end{pmatrix} \mathbf{x}^{(k)} - \begin{pmatrix} \frac{1}{6} \\ \frac{2}{7} \\ -\frac{1}{9} \end{pmatrix}.$$

This can be also written as

$$\mathbf{e}^{(k+1)} = \begin{pmatrix} 1/12 \\ 1/4 \\ -1/12 \end{pmatrix} - \begin{pmatrix} 0 & -\frac{2}{6} & 0 \\ -\frac{1}{7} & 0 & \frac{2}{7} \\ -\frac{3}{9} & \frac{2}{9} & 0 \end{pmatrix} \begin{pmatrix} 1/12 \\ 1/4 \\ -1/12 \end{pmatrix} + \begin{pmatrix} 0 & -\frac{2}{6} & 0 \\ -\frac{1}{7} & 0 & \frac{2}{7} \\ -\frac{3}{9} & \frac{2}{9} & 0 \end{pmatrix} \mathbf{e}^{(k)} - \begin{pmatrix} \frac{1}{6} \\ \frac{2}{7} \\ -\frac{1}{9} \end{pmatrix},$$

or

$$\mathbf{e}^{(k+1)} = \begin{pmatrix} 0 & -\frac{2}{6} & 0 \\ -\frac{1}{7} & 0 & \frac{2}{7} \\ -\frac{3}{9} & \frac{2}{9} & 0 \end{pmatrix} \mathbf{e}^{(k)},$$

(because $\mathbf{x}^{(k)} = \mathbf{x} - \mathbf{e}^{(k)}$) which is the required error in the (n+1)th step.

(d) Now finding the first approximation of the error, we have to compute the following

$$\mathbf{e}^{(1)} = \begin{pmatrix} 0 & -\frac{2}{6} & 0 \\ -\frac{1}{7} & 0 & \frac{2}{7} \\ -\frac{3}{9} & \frac{2}{9} & 0 \end{pmatrix} \mathbf{e}^{(0)},$$

where

$$\mathbf{e}^{(0)} = \mathbf{x} - \mathbf{x}^{(0)}.$$

Using $\mathbf{x}^{(0)} = [0, 0, 0]^T$, we have

$$\mathbf{e}^{(0)} = \begin{pmatrix} 1/12 \\ 1/4 \\ -1/12 \end{pmatrix} - \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 1/12 \\ 1/4 \\ -1/12 \end{pmatrix}.$$

Thus

$$\mathbf{e}^{(1)} = \begin{pmatrix} 0 & -\frac{2}{6} & 0 \\ -\frac{1}{7} & 0 & \frac{2}{7} \\ -\frac{3}{9} & \frac{2}{9} & 0 \end{pmatrix} \begin{pmatrix} 1/12 \\ 1/4 \\ -1/12 \end{pmatrix} = \begin{pmatrix} -1/12 \\ -1/28 \\ 1/36 \end{pmatrix}.$$

Similarly, for the second approximation of the error, we have to compute the following

$$\mathbf{e}^{(2)} = \begin{pmatrix} 0 & -\frac{2}{6} & 0 \\ -\frac{1}{7} & 0 & \frac{2}{7} \\ -\frac{3}{9} & \frac{2}{9} & 0 \end{pmatrix} \mathbf{e}^{(1)},$$

or

$$\mathbf{e}^{(2)} = \begin{pmatrix} 0 & -\frac{2}{6} & 0 \\ -\frac{1}{7} & 0 & \frac{2}{7} \\ -\frac{3}{6} & \frac{2}{6} & 0 \end{pmatrix} \begin{pmatrix} -1/12 \\ -1/28 \\ 1/36 \end{pmatrix} = \begin{pmatrix} 1/84 \\ 5/252 \\ 5/252 \end{pmatrix},$$

which is the required second approximation of the error.

Gauss-Seidel Iterative Method

(a) Now by using Gauss-Seidel method, first we compute the Gauss-Seidel iteration matrix T_G and the vector \mathbf{c}_G as follows:

$$T_G = \begin{pmatrix} 0 & -\frac{1}{3} & 0 \\ 0 & \frac{1}{21} & \frac{2}{7} \\ 0 & \frac{23}{189} & \frac{4}{63} \end{pmatrix} \quad and \quad \mathbf{c}_G = \begin{pmatrix} \frac{1}{6} \\ \frac{11}{42} \\ -\frac{41}{378} \end{pmatrix}.$$

Thus the matrix form of Gauss-Seidel iterative method is

$$\mathbf{x}^{(k+1)} = \begin{pmatrix} 0 & -\frac{1}{3} & 0 \\ 0 & \frac{1}{21} & \frac{2}{7} \\ 0 & \frac{23}{189} & \frac{4}{63} \end{pmatrix} \mathbf{x}^{(k)} + \begin{pmatrix} \frac{1}{6} \\ \frac{11}{42} \\ -\frac{41}{378} \end{pmatrix}, \quad k = 0, 1, 2.$$

(b) Writing the above iterative form in component form, we get

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 & -\frac{1}{3} & 0 \\ 0 & \frac{1}{21} & \frac{2}{7} \\ 0 & \frac{23}{189} & \frac{4}{63} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} + \begin{pmatrix} \frac{1}{6} \\ \frac{11}{42} \\ -\frac{41}{378} \end{pmatrix},$$

and it is equivalent to

$$x_1 = -1/3x_2 + 1/6$$

 $x_2 = 1/21x_2 + 2/7x_3 + 11/42$
 $x_3 = 23/189x_2 + 4/63x_3 - 41/378$

Now solving for x_1, x_2 and x_3 , we get

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 1/12 \\ 1/4 \\ -1/12 \end{pmatrix},$$

which is the exact solution of the given system.

(c) The error in the (n+1)th step can be easily computed as

$$\mathbf{e}^{(k+1)} = \begin{pmatrix} 0 & -\frac{1}{3} & 0 \\ 0 & \frac{1}{21} & \frac{2}{7} \\ 0 & \frac{23}{189} & \frac{4}{63} \end{pmatrix} \mathbf{e}^{(k)}.$$

(d) The first and second approximations of the error can be calculated as follows:

$$\mathbf{e}^{(1)} = \begin{pmatrix} 0 & -\frac{1}{3} & 0 \\ 0 & \frac{1}{21} & \frac{2}{7} \\ 0 & \frac{23}{189} & \frac{4}{63} \end{pmatrix} \mathbf{e}^{(0)} = [-1/12, -1/84, 19/756]^T,$$

and

$$\mathbf{e}^{(2)} = \begin{pmatrix} 0 & -\frac{1}{3} & 0 \\ 0 & \frac{1}{21} & \frac{2}{7} \\ 0 & \frac{23}{189} & \frac{4}{63} \end{pmatrix} \mathbf{e}^{(1)} = [1/252, 5/756, 1/6804]^{T}.$$

3.6.4 Convergence Criteria of Iterative Methods

Since we noted that the Jacobi method and the Gauss-Seidel method do not always converge to the solution of the given system of linear equations. Here we need some conditions which make the both methods converge. The sufficient conditions for the convergence of both iterative methods are discussed in the following theorems.

Theorem 3.15 (First Sufficient Condition for Convergence)

If the matrix A is strictly diagonally dominant (SDD), then for any choice of initial approximation $\mathbf{x}^{(0)} \in \mathbf{R}$ both the Jacobi method and the Gauss-Seidel method give sequence $\{x^{(k)}\}_{k=0}^{\infty}$ of approximations that converge to the solution of the linear system.

-	0.0. Dolumon of the Example					
	k	$x_1^{(k)}$	$x_2^{(k)}$	$x_3^{(k)}$		
	0	0.000	0.000	0.000		
	1	1.200	0.667	1.167		
	2	1.167	1.050	0.544		
	3	0.889	0.744	0.428		
	:	:	:			
	7	0.994	0.777	0.572		
	8	1.004	0.787	0.576		

Table 3.5: Solution of the Example 3.39

Example 3.39 Rearrange the following linear system of equations

$$x_1 + 6x_2 - 3x_3 = 4$$

 $2x_1 + 2x_2 + 6x_3 = 7$
 $5x_1 + 2x_2 - x_3 = 6$

such that the convergence of both iterative methods (Jacobi and Gauss-Seidel) is guaranteed. Use initial solution $\mathbf{x}^{(0)} = [0,0,0]^T$, compute approximation solution within accuracy 10^{-2} .

Solution. For the guarantee convergence of iterative methods, the system must be SDD form, so rearrange the given system in the following form

$$5x_1 + 2x_2 - x_3 = 6$$

 $x_1 + 6x_2 - 3x_3 = 4$
 $2x_1 + 2x_2 + 6x_3 = 7$

Jacobi Iterative Method:

$$\begin{array}{rclcrcl} x_1^{(k+1)} & = & \frac{1}{5} \Big[6 & - & 2x_2^{(k)} & + & x_3^{(k)} \Big] \\ \\ x_2^{(k+1)} & = & \frac{1}{6} \Big[4 & - & x_1^{(k)} & + & 3x_3^{(k)} \Big] \\ \\ x_3^{(k+1)} & = & \frac{1}{6} \Big[7 & - & 2x_1^{(k)} & - & 2x_2^{(k)} \Big] \end{array}$$

Starting with $\mathbf{x}^{(0)} = [0, 0, 0]^T$, the first and subsequent iterations are listed in Table 3.5.

Gauss-Seidel Iterative Method:

$$x_1^{(k+1)} = \frac{1}{5} \begin{bmatrix} 6 - 2x_2^{(k)} + x_3^{(k)} \end{bmatrix}$$

$$x_2^{(k+1)} = \frac{1}{6} \begin{bmatrix} 4 - x_1^{(k+1)} + 3x_3^{(k)} \end{bmatrix}$$

$$x_3^{(k+1)} = \frac{1}{6} \begin{bmatrix} 7 - 2x_1^{(k+1)} - 2x_2^{(k+1)} \end{bmatrix}$$

5.0. Solution of the Example					
k	$x_1^{(k)}$	$x_2^{(k)}$	$\begin{bmatrix} x_3^{(k)} \end{bmatrix}$		
0	0.000	0.000	0.000		
1	1.200	0.467	0.611		
2	1.136	0.783	0.527		
3	0.992	0.765	0.581		
4	1.010	0.789	0.567		
5	0.998	0.784	0.573		
6	1.001	0.786	0.571		

Table 3.6: Solution of the Example 3.39

Starting with initial approximation $\mathbf{x}^{(0)} = [0,0,0]^T$, the first and subsequent iterations are listed in Table 3.6. Note that Gauss-Seidel iterative method converges faster than Jacobi iterative method.

There is another sufficient condition for the convergence of both iterative methods which is defined in the following theorem.

Theorem 3.16 (Second Sufficient Condition for Convergence)

For any initial approximation $\mathbf{x}^{(0)} \in \mathbf{R}$, the sequence $\{x^{(k)}\}_{k=0}^{\infty}$ of approximations defined by

$$\mathbf{x}^{(k+1)} = T\mathbf{x}^{(k)} + \mathbf{c}, \quad \text{for each} \quad k \ge 0, \quad \text{and} \quad \mathbf{c} \ne 0, \tag{3.52}$$

converges to the unique solution of $\mathbf{x} = T\mathbf{x} + \mathbf{c}$ if ||T|| < 1 for any natural matrix norm, and the following **error bounds** hold:

$$\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|T\|^k \|\mathbf{x}^{(0)} - \mathbf{x}\|,$$

$$\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \frac{\|T\|^k}{1 - \|T\|} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|.$$
(3.53)

Note that smaller the value of the ||T||, faster the convergence of the iterative methods.

Example 3.40 Consider the following nonhomogeneous linear system $A\mathbf{x} = \mathbf{b}$, where

$$A = \begin{pmatrix} 5 & 0 & -1 \\ -1 & 3 & 0 \\ 0 & -1 & 4 \end{pmatrix} \quad and \quad \mathbf{b} = \begin{pmatrix} 1 \\ 2 \\ 4 \end{pmatrix}.$$

Find the matrix form of iterative (Jacobi and Gauss-Seidel) methods and show that Gauss-Seidel iterative method converges faster than Jacobi iterative method for the given system.

Solution. Here we will show that the l_{∞} -norm of the Gauss-Seidel iteration matrix T_G is less than the l_{∞} -norm of the Jacobi iteration matrix T_J , that is

$$||T_G||_{\infty} < ||T_J||_{\infty}.$$

The Jacobi iteration matrix T_J can be obtained from the given matrix A as follows

$$T_{J} = -D^{-1}(L+U) = -\begin{pmatrix} 5 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 4 \end{pmatrix}^{-1} \begin{pmatrix} 0 & 0 & -1 \\ -1 & 0 & 0 \\ 0 & -1 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0 & \frac{1}{5} \\ \frac{1}{3} & 0 & 0 \\ 0 & \frac{1}{4} & 0 \end{pmatrix}.$$

Thus the matrix form of Jacobi iterative method is

$$\mathbf{x}^{(k+1)} = \begin{pmatrix} 0 & 0 & \frac{1}{5} \\ \frac{1}{3} & 0 & 0 \\ 0 & \frac{1}{4} & 0 \end{pmatrix} \mathbf{x}^{(k)} + \begin{pmatrix} \frac{1}{5} \\ \frac{2}{3} \\ 1 \end{pmatrix}, \quad k \ge 0.$$

Similarly, Gauss-Seidel iteration matrix T_G is defined as

$$T_G = -(D+L)^{-1}U = -\begin{pmatrix} 5 & 0 & 0 \\ -1 & 3 & 0 \\ 0 & -1 & 4 \end{pmatrix}^{-1} \begin{pmatrix} 0 & 0 & -1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix},$$

and it gives

$$T_G = -\begin{pmatrix} \frac{1}{5} & 0 & 0\\ \frac{1}{15} & \frac{1}{3} & 0\\ \frac{1}{60} & \frac{1}{15} & \frac{1}{4} \end{pmatrix} \begin{pmatrix} 0 & 0 & -1\\ 0 & 0 & 0\\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0 & \frac{1}{5}\\ 0 & 0 & \frac{1}{15}\\ 0 & 0 & \frac{1}{60} \end{pmatrix}.$$

So the matrix form of Gauss-Seidel iterative method is

$$\mathbf{x}^{(k+1)} = \begin{pmatrix} 0 & 0 & \frac{1}{5} \\ 0 & 0 & \frac{1}{15} \\ 0 & 0 & \frac{1}{60} \end{pmatrix} \mathbf{x}^{(k)} + \begin{pmatrix} \frac{1}{5} \\ \frac{11}{15} \\ \frac{71}{60} \end{pmatrix}, \quad k \ge 0.$$

Since the l_{∞} -norm of the matrix T_J is

$$||T_J||_{\infty} = \max\left\{\frac{1}{5}, \frac{1}{3}, \frac{1}{4}\right\} = \frac{1}{3} = 0.3333 < 1,$$

and the l_{∞} -norm of the matrix T_G is

$$||T_G||_{\infty} = \max\left\{\frac{1}{5}, \frac{1}{15}, \frac{1}{60}\right\} = \frac{1}{5} = 0.2000 < 1.$$

Since $||T_G||_{\infty} < ||T_J||_{\infty}$, which shows that Gauss-Seidel method will converge faster than Jacobi method for the given linear system.

Note that the condition ||T|| < 1 is equivalent to the condition that a matrix A is to be strictly diagonally dominant.

For Jacobi method for a general matrix A, the norm of Jacobi iteration matrix is defined as

$$||T_J|| = \max_{1 \le i \le n} \sum_{\substack{j=1 \ j \ne i}}^n \left| \frac{a_{ij}}{a_{ii}} \right|.$$

Thus for $||T_J|| < 1$ is equivalent to requiring

$$\sum_{\substack{j=1\\j\neq i}}^{n} |a_{ij}| < |a_{ii}|,$$

that is, a matrix A is strictly diagonally dominant.

Example 3.41 Consider the following linear system of equations

$$4x_1 - x_2 + x_3 = 12
-x_1 + 3x_2 + x_3 = 1
x_1 + x_2 + 5x_3 = -14$$

- (a) Show that both iterative methods (Jacobi and Gauss-Seidel) will converge by using $||T||_{\infty} < 1$.
- (b) Find second approximation $\mathbf{x}^{(2)}$ when the initial solution is $\mathbf{x}^{(0)} = [4, 3, -3]^T$.
- (c) Compute the error bounds for your approximations.
- (d) How many iterations needed to get an accuracy within 10^{-4} .

Solution. From (3.46), we have

$$A = \begin{pmatrix} 4 & -1 & 1 \\ -1 & 3 & 1 \\ 1 & 1 & 5 \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 \\ -1 & 0 & 0 \\ 1 & 1 & 0 \end{pmatrix} + \begin{pmatrix} 0 & -1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} 4 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 5 \end{pmatrix}$$
$$= L + U + D.$$

Jacobi Method:

(a) Since the Jacobi iteration matrix is defined as

$$T_J = -D^{-1}(L+U),$$

and by using the given information, we have

$$T_{J} = -\begin{pmatrix} \frac{1}{4} & 0 & 0 \\ 0 & \frac{1}{3} & 0 \\ 0 & 0 & \frac{1}{5} \end{pmatrix} \begin{pmatrix} 0 & -1 & 1 \\ -1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix} = \begin{pmatrix} 0 & \frac{1}{4} & -\frac{1}{4} \\ \frac{1}{3} & 0 & -\frac{1}{3} \\ -\frac{1}{5} & -\frac{1}{5} & 0 \end{pmatrix}.$$

Then the l_{∞} norm of the matrix T_J is

$$||T_J||_{\infty} = \max\left\{\frac{2}{4}, \frac{2}{3}, \frac{2}{5}\right\} = \frac{2}{3} < 1.$$

Thus the Jacobi method will converge for the given linear system.

(b) The Jacobi method for the given system is

$$\begin{array}{rclcrcl} x_1^{(k+1)} & = & \frac{1}{4} \Big[12 & + & x_2^{(k)} & - & x_3^{(k)} \Big] \\ \\ x_2^{(k+1)} & = & \frac{1}{3} \Big[1 & + & x_1^{(k)} & - & x_3^{(k)} \Big] \\ \\ x_3^{(k+1)} & = & \frac{1}{5} \Big[-14 & - & x_1^{(k)} & - & x_2^{(k)} \Big] \end{array}$$

Starting with initial approximation $x_1^{(0)} = 4, x_2^{(0)} = 3, x_3^{(0)} = -3$, and for k = 0, 1, we obtain the first and the second approximations as

$$\mathbf{x}^{(1)} = [4.5, 2.6667, -4.2]^T \quad and \quad \mathbf{x}^{(2)} = [4.7167, 3.2333, -4.2333]^T.$$

(c) Using the error bound formula (3.53), we obtain

$$\|\mathbf{x} - \mathbf{x}^{(2)}\| \le \frac{(2/3)^2}{1 - 2/3} \left\| \begin{pmatrix} 4.5 \\ 2.6667 \\ -4.2 \end{pmatrix} - \begin{pmatrix} 4 \\ 3 \\ -3 \end{pmatrix} \right\| \le \frac{4}{3}(1.2) = 1.6.$$

(d) To find the number of iterations, we use the formula (3.53) as

$$\|\mathbf{x} - \mathbf{x}^{(k)}\| \le \frac{\|T_J\|^k}{1 - \|T_J\|} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \le 10^{-4}.$$

It gives

$$\frac{(2/3)^k}{1/3}(1.2) \le 10^{-4}, \quad or \quad (2/3)^k \le \frac{10^{-4}}{3.6}.$$

Taking ln on both sides, we obtain

$$k \ln(2/3) \le \ln\left(\frac{10^{-4}}{3.6}\right)$$
, gives $k \ge 25.8789$, or $k = 26$,

which is the required number of iterations.

Gauss-Seidel Method:

(a) Since the Gauss-Seidel iteration matrix is defined as

$$T_G = -(D+L)^{-1}U$$

and by using the given information, we have

$$T_G = -\begin{pmatrix} \frac{1}{4} & 0 & 0\\ \frac{1}{12} & \frac{1}{3} & 0\\ -\frac{4}{60} & -\frac{1}{15} & \frac{1}{5} \end{pmatrix} \begin{pmatrix} 0 & -1 & 1\\ 0 & 0 & 1\\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & \frac{1}{4} & -\frac{1}{4}\\ 0 & \frac{1}{12} & -\frac{5}{12}\\ 0 & -\frac{4}{60} & \frac{8}{60} \end{pmatrix}.$$

Then the l_{∞} norm of the matrix T_G is

$$||T_G||_{\infty} = \max\left\{\frac{2}{4}, \frac{6}{12}, \frac{12}{60}\right\} = \frac{1}{2} < 1.$$

Thus the Gauss-Seidel method will converge for the given linear system.

(b) The Gauss-Seidel method for the given system is

$$x_1^{(k+1)} = \frac{1}{4} \begin{bmatrix} 12 & + & x_2^{(k)} & - & x_3^{(k)} \end{bmatrix}$$

$$x_2^{(k+1)} = \frac{1}{3} \begin{bmatrix} 1 & + & x_1^{(k+1)} & - & x_3^{(k)} \end{bmatrix}$$

$$x_3^{(k+1)} = \frac{1}{5} \begin{bmatrix} -12 & - & x_1^{(k+1)} & - & x_2^{(k+1)} \end{bmatrix}$$

Starting with initial approximation $x_1^{(0)} = 4, x_2^{(0)} = 3, x_3^{(0)} = -3$, and for k = 0, 1, we obtain the first and the second approximations as

$$\mathbf{x}^{(1)} = [4.5, 2.8333, -4.2667]^T \quad and \quad \mathbf{x}^{(2)} = [4.775, 3.3472, -4.4244]^T.$$

(c) Using the error bound formula (3.53), we obtain

$$\|\mathbf{x} - \mathbf{x}^{(2)}\| \le \frac{(1/2)^2}{1 - 1/2} \left\| \begin{pmatrix} 4.5 \\ 2.8333 \\ -4.2667 \end{pmatrix} - \begin{pmatrix} 4 \\ 3 \\ -3 \end{pmatrix} \right\| \le \frac{1}{2} (1.2667) = 0.6334.$$

(d) To find the number of iterations, we use the formula (3.53) as

$$\|\mathbf{x} - \mathbf{x}^{(k)}\| \le \frac{\|T_J\|^k}{1 - \|T_J\|} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\| \le 10^{-4}.$$

It gives

$$\frac{(1/2)^k}{1/2}(1.2667) \le 10^{-4}, \quad or \quad (1/2)^k \le \frac{10^{-4}}{2.5334}.$$

Taking ln on both sides, we obtain

$$k \ln(1/2) \le \ln\left(\frac{10^{-4}}{2.5334}\right), \quad gives \quad k \ge 14.6084 \quad or \quad k = 15,$$

which is the required number of iterations.

Example 3.42 Consider the following system

$$\begin{array}{rclrcl}
4x_1 & + & x_2 & = & 7 \\
x_1 & + & 2x_2 & = & 0 \\
& & 2x_2 & + & 3x_3 & = & 1
\end{array}$$

If $x^{(0)} = [0,0,0]^T$, then compute an error bound $||x - x^{(10)}||$ for the approximation using Gauss-Seidel method.

Solution. Since we know that error bound formula for the gauss-Seidel method is

$$\|\mathbf{x} - \mathbf{x}^{(k)}\| \le \frac{\|T_G\|^k}{1 - \|T_G\|} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|,$$

and given k = 10, we have

$$\|\mathbf{x} - \mathbf{x}^{(10)}\| \le \frac{\|T_G\|^{10}}{1 - \|T_G\|} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|.$$

So we have to find $||T_G||$ and the first approximation $\mathbf{x}^{(1)}$. Since the Gauss-Seidel iteration matrix is defined as

$$T_G = -(D+L)^{-1}U,$$

and by using the given information, we have

$$T_G = -\left(\begin{array}{ccc} 4 & 1 & 0 \\ 1 & 2 & 0 \\ 0 & 2 & 3 \end{array}\right)^{-1} \left(\begin{array}{ccc} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{array}\right).$$

To find the inverse of the matrix (D+L), we will use Gauss-Jordan method as follows:

$$[(D+L)|I] = \begin{pmatrix} 4 & 0 & 0 & \vdots & 1 & 0 & 0 \\ 1 & 2 & 0 & \vdots & 0 & 1 & 0 \\ 0 & 2 & 3 & \vdots & 0 & 0 & 1 \end{pmatrix} \equiv \begin{pmatrix} 1 & 0 & 0 & \vdots & 1/4 & 0 & 0 \\ 0 & 2 & 0 & \vdots & -1/4 & 1 & 0 \\ 0 & 2 & 3 & \vdots & 0 & 0 & 1 \end{pmatrix}$$

$$\equiv \left(\begin{array}{ccccc} 1 & 0 & 0 & \vdots & 1/4 & 0 & 0 \\ 0 & 1 & 0 & \vdots & -1/8 & 1/2 & 0 \\ 0 & 0 & 3 & \vdots & 1/4 & -1 & 1 \end{array}\right) \equiv \left(\begin{array}{cccccc} 1 & 0 & 0 & \vdots & 1/4 & 0 & 0 \\ 0 & 1 & 0 & \vdots & -1/8 & 1/2 & 0 \\ 0 & 0 & 1 & \vdots & 1/12 & -1/3 & 1/3 \end{array}\right).$$

Thus

$$T_G = -\begin{pmatrix} \frac{1}{4} & 0 & 0 \\ -\frac{1}{8} & \frac{1}{2} & 0 \\ \frac{1}{12} & -\frac{1}{3} & \frac{1}{3} \end{pmatrix} \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & -\frac{1}{4} & 0 \\ 0 & \frac{1}{8} & 0 \\ 0 & -\frac{1}{12} & 0 \end{pmatrix}.$$

Then the l_{∞} norm of the matrix T_G is

$$||T_G||_{\infty} = \max\left\{\frac{1}{4}, \frac{1}{8}, \frac{1}{12}\right\} = \frac{1}{4} = 0.25 < 1.$$

Now to find the first approximation using Gauss-Seidel method, we will the following formula

$$x_1^{(k+1)} = \frac{1}{4} \begin{bmatrix} 7 - x_2^{(k)} \\ x_2^{(k+1)} = \frac{1}{2} \begin{bmatrix} 0 - x_1^{(k+1)} \\ \end{bmatrix}$$

$$x_3^{(k+1)} = \frac{1}{3} \begin{bmatrix} 1 - 2x_2^{(k+1)} \end{bmatrix}$$

Starting with initial approximation $x_1^{(0)} = 0, x_2^{(0)} = 0, x_3^{(0)} = 0$ and for k = 0, we obtain the first approximation as

$$\mathbf{x}^{(1)} = [1.7500, 0.8750, 0.9167]^T.$$

Thus

$$\|\mathbf{x} - \mathbf{x}^{(10)}\| \le \frac{(0.25)^{10}}{0.75} (1.75) = 2.2252 \times 10^{-6},$$

the required an error bound.

Theorem 3.17 If A is a symmetric positive definite matrix with positive diagonal entries, then the Gauss-Seidel method converges to unique solution of the linear system $A\mathbf{x} = \mathbf{b}$.

Example 3.43 Solve the following system of linear equations using Gauss-Seidel iterative methods, using $\epsilon = 10^{-5}$ in the l_{∞} -norm, taking the initial solution $\mathbf{x}^{(0)} = [0, 0, 0, 0]^T$.

$$5x_1 - x_3 - x_3 = 1$$

$$-x_1 - x_2 + 13x_3 - x_4 = 1$$

$$-x_2 + 13x_3 - x_4 = 3$$

5

0.267505

table 3.7. Solution by Gauss-Seider Method				Method
k	$x_1^{(k)}$	$x_2^{(k)}$	$x_3^{(k)}$	$x_4^{(k)}$
0	0.000000	0.000000	0.000000	0.000000
1	0.200000	0.071429	0.328571	0.341270
2	0.265714	0.119274	0.337307	0.346586
3	0.267461	0.120278	0.337518	0.346698
4	0.267504	0.120301	0.337524	0.346700

0.337524

0.346700

0.120302

Table 3.7: Solution by Gauss-Seidel Method

Solution. The matrix

$$A = \left(\begin{array}{rrrr} 5 & 0 & -1 & 0 \\ 0 & 14 & -1 & -1 \\ -1 & -1 & 13 & 0 \\ 0 & -1 & 0 & 9 \end{array}\right),$$

of the given system is symmetric positive definite with positive diagonal entries and Gauss-Seidel formula for the system is

$$\begin{array}{rclcrcl} x_1^{(k+1)} & = & \frac{1}{5} \Big[1 & & + & x_3^{(k)} & & \Big] \\ \\ x_2^{(k+1)} & = & \frac{1}{14} \Big[1 & & + & x_3^{(k)} & + & x_4^{(k)} \Big] \\ \\ x_3^{(k+1)} & = & \frac{1}{13} \Big[4 & + & x_1^{(k+1)} & + & x_2^{(k+1)} & & \Big] \\ \\ x_4^{(k+1)} & = & \frac{1}{9} \Big[3 & & + & x_2^{(k+1)} & & \Big] \end{array}$$

Using initial approximation $\mathbf{x}^{(0)} = [0,0,0,0]^T$, the first and subsequent iterations are listed in Table 3.7.

Notice that Gauss-Seidel method converged very fast (only five iterations) and the approximate solution of the given system $[0.267505, 0.120302, 0.337524, 0.346700]^T$ is equal to the exact solution $[0.267505, 0.120302, 0.337524, 0.346700]^T$ upto six decimal places.

3.7 Errors in Solving Linear Systems

Any computed solution of a linear system must, because of round-off and other errors, be considered an approximate solution. Here we shall consider the most natural method for determining the accuracy of a solution of the linear system. One obvious way of estimating the accuracy of the computed solution \mathbf{x}^* is to compute $A\mathbf{x}^*$ and to see how close $A\mathbf{x}^*$ comes to \mathbf{b} . Thus if \mathbf{x}^* is an approximate solution of the given system $A\mathbf{x} = \mathbf{b}$, we compute a vector

$$\mathbf{r} = \mathbf{b} - A\mathbf{x}^*,\tag{3.54}$$

which is called the *residual vector* and can be easily calculated. The quantity

$$\frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} = \frac{\|\mathbf{b} - A\mathbf{x}^*\|}{\|\mathbf{b}\|},$$

is called the relative residual.

```
Program 3.12

MATLAB m-file for finding Residual Vector function r=RES(A,b,x0)
[n,n]=size(A); for i=1:n; R(i) = b(i); for j=1:n
R(i)=R(i)-A(i,j)*x0(j);end; RES(i)=R(i); end; r=RES'
```

The smallness of the residual then provides a measure of the goodness of the approximate solution \mathbf{x}^* . If every component of vector \mathbf{r} vanishes, then \mathbf{x}^* is the exact solution. If \mathbf{x}^* is a good approximation then we would expect each component of \mathbf{r} to be small, at least in a relative sense. For example, the following linear system

$$x_1 + 2x_2 = 3$$

 $1.0001x_1 + 2x_2 = 3.0001$

has the exact solution $\mathbf{x} = [1, 1]^T$ but has a poor approximate solution $\mathbf{x}^* = [3, 0]^T$. To see how good this solution is, we compute the residual, $\mathbf{r} = [0, -0.0002]^T$, and so $\|\mathbf{r}\|_{\infty} = 0.0002$. Although the norm of the residual vector is small, the approximate solution $\mathbf{x}^* = [3, 0]^T$ is obviously quite poor; in fact $\|\mathbf{x} - \mathbf{x}^*\|_{\infty} = 2$.

To get above results using MATLAB command, we do the following:

```
>> A = [1 \ 2; 1.0001 \ 2]; \ b = [3 \ 3.0001]; \ x0 = [3 \ 0];
>> RESID(A, b, x0); \ x = [1 \ 1]; \ Error = norm((x - x0), inf);
```

We can conclude from the residual that the approximate solution is correct to at most three decimal places. Also, the following linear system

has exact solution $\mathbf{x} = [1, 1, 1, 1]^T$ and having the approximate solution due to the Gaussian elimination without pivoting is

$$\mathbf{x}^* = [1.0000322, 0.99996948, 0.99998748, 1.0000113]^T,$$

and the residual is

$$\mathbf{r} = [0.6 \times 10^{-7}, 0.6 \times 10^{-7}, -0.53 \times 10^{-5}, -0.21 \times 10^{-4}]^T$$

The approximate solution due to the Gaussian elimination with partial pivoting is

$$\mathbf{x}^* = [0.9999997, 0.99999997, 0.99999996, 1.0000000]^T,$$

and the residual is

$$\mathbf{r} = [0.3 \times 10^{-7}, 0.3 \times 10^{-7}, 0.6 \times 10^{-7}, 0.1 \times 10^{-8}]^T.$$

We found that all the elements of the residual for second case (with pivoting) are less than 0.6×10^{-7} , whereas for first case (without pivoting) they are as large as 0.2×10^{-4} . Even without knowing the exact solution, it is clear that the solution obtained in second case is much better than that of first case. The residual provides a reasonable measure of the accuracy of a solution in those cases where the error is primarily due to the accumulation of round-off errors.

Intuitively it would seem reasonable to assume that when $\|\mathbf{r}\|$ is small for a given vector norm, then the error $\|\mathbf{x} - \mathbf{x}^*\|$ would be small as well. In fact this is true for some systems. However, there are systems of equations which do not satisfy this property. Such systems are said to be *ill-conditioned*.

3.7.1 Conditioning of Linear Systems

In solving the linear system numerically we have to see the problem conditioning, algorithm stability, and cost. Above we discussed efficient elimination schemes to solve a linear system and these schemes are stable when pivoting is employed. But there are some ill-conditioned systems which are tough to solve by any method. These types of linear systems are identified in this chapter.

Here, we will present a parameter, the condition number, which quantitatively measures the conditioning of a linear system. The condition number is greater and equal to one and as a linear system becomes more ill-conditioned, the condition number increases. After factoring a matrix, the condition number can be estimated in roughly the same time as it takes to solve a few factored systems $(LU)\mathbf{x} = \mathbf{b}$. Hence, after factoring a matrix, the extra computer time needed to estimate the condition number is usually insignificant. These are systems in which small changes in the coefficients of the system lead to large changes in the solution. For example, consider a linear system

The exact solution is easily verified to be $x_1 = x_2 = 1$. On the other hand, the system

$$\begin{array}{rcl} x_1 & + & x_2 & = & 2 \\ 1.001x_1 & + & x_2 & = & 2.01 \end{array}$$

has the solution $x_1 = 10, x_2 = -8$. Thus change of 1 percent in the coefficients has changed the solution by a factor of 10. If in the above given system, we substitute $x_1 = 10, x_2 = 8$, we find that the residual are $r_1 = 0, r_2 = 0.09$, so that this solution looks a reasonable although it is grossly in error. In practical problems we can expect the coefficients in the system to be subject to small errors, either because of round-off or because of physical measurement. If the system is ill-conditioned, the resulting solution may be grossly in error. Errors of this type, unlike those caused by round-off error accumulation, can not be avoided by careful programming.

We have seen that for ill-conditioned systems the residual is not necessarily a good measure of the accuracy of a solution. How then can we tell when a system is ill-conditioned? In the following we discuss the some possible indicators of ill-conditioned system.

Definition 3.25 (Condition Number of a Matrix)

The number $||A|| ||A^{-1}||$ is called the condition number of a nonsingular matrix A and is denoted by K(A), that is

$$cond(A) = K(A) = ||A|| ||A^{-1}||.$$
(3.55)

Note that the condition number K(A) for A depends on the matrix norm used and can, for some matrices, vary considerably as the matrix norm is changed. Since

$$1 = ||I|| = ||AA^{-1}|| \le ||A|| ||A^{-1}|| = K(A),$$

therefore, the condition number is always in the range $1 \leq K(A) \leq \infty$ regardless of any natural norm. The lower limit is attained for identity matrices and $K(A) = \infty$ if A is singular. So the matrix A is well-behaved (well-conditioned) if K(A) is close to 1 and is increasingly ill-conditioned when K(A) is significantly greater than 1, that is, $K(A) \to \infty$.

The condition numbers provide bounds for the sensitivity of the solution of a set of equations to changes in the coefficient matrix. Unfortunately, the evaluation of any of the condition numbers of a matrix A is not a trivial task since it is necessary first to obtain its inverse.

So if the condition number of a matrix is very large number then this is one of the indicator of the ill-conditioned system. An other indicator of ill-conditioning is when the pivots during the process of elimination suffer a loss of one or more significant figures. Small changes in the right-hand side terms of the system lead to large changes in the solution, gives another indicator of ill-conditioned systems. Also, when the elements of the inverse of the coefficient matrix are large compared to the elements of the coefficients matrix, shows the ill-conditioned system.

Example 3.44 Compute the condition number of the following matrix using the l_{∞} -norm

$$A = \left(\begin{array}{rrr} 2 & -1 & 0 \\ 2 & -4 & -1 \\ -1 & 0 & 2 \end{array}\right).$$

Solution. Since the condition number of a matrix is defined as

$$K(A) = ||A||_{\infty} ||A^{-1}||_{\infty}.$$

First we calculate the inverse of the given matrix which is

$$A^{-1} = \begin{pmatrix} \frac{8}{13} & -\frac{2}{13} & -\frac{1}{13} \\ \frac{3}{13} & -\frac{4}{13} & -\frac{2}{13} \\ \frac{4}{13} & -\frac{1}{13} & \frac{6}{13} \end{pmatrix}.$$

Now we calculate the l_{∞} -norm of both the matrices A and A^{-1} . Since the l_{∞} -norm of a matrix is the maximum of the absolute row sums, we have

$$||A||_{\infty} = max\{|2| + |-1| + |0|, |2| + |-4| + |-1|, |-1| + |0| + |2|\} = 7,$$

and

$$||A^{-1}||_{\infty} = \max\left\{\left|\frac{8}{13}\right| + \left|\frac{-2}{13}\right| + \left|\frac{-1}{13}\right|, \left|\frac{3}{13}\right| + \left|\frac{-4}{13}\right| + \left|\frac{-2}{13}\right|, \left|\frac{4}{13}\right| + \left|\frac{-1}{13}\right| + \left|\frac{6}{13}\right|\right\},$$

which gives

$$||A^{-1}||_{\infty} = \frac{11}{13}.$$

Therefore,

$$K(A) = ||A||_{\infty} ||A^{-1}||_{\infty} = (7) \left(\frac{11}{13}\right) \approx 5.9231.$$

Depending on the application, we might consider this number to be reasonably small and conclude that the given matrix A is reasonably well-conditioned.

To get above results using MATLAB commands, we do the following:

>>
$$A = [2 -1 \ 0; 2 -4 -1; -1 \ 0 \ 2]; Ainv = inv(A)$$

>> $K(A) = norm(A, inf) * norm(Ainv, inf)$

Example 3.45 If the condition number of following matrix A is 8.8671, then find the l_{∞} -norm of its inverse matrix, that is, $||A^{-1}||_{\infty}$

$$A = \begin{pmatrix} 10.2 & 2.4 & 4.5 \\ -2.3 & 7.7 & 11.1 \\ -5.5 & -3.2 & 0.9 \end{pmatrix}.$$

Solution. Since the condition number of a matrix is defined as

$$K(A) = ||A||_{\infty} ||A^{-1}||_{\infty}.$$

First we calculate the l_{∞} -norm of the given matrix A which is the maximum of the absolute row sums, we have

$$||A||_{\infty} = max\{17.1000, 21.1000, 9.6\} = 21.1000,$$

and as it is given K(A) = 8.8671, so we have

$$8.8671 = (21.1000) ||A^{-1}||_{\infty}.$$

Simplifying this, we get $||A^{-1}||_{\infty} = 0.4202$.

Some matrices are notoriously ill-conditioned. For example, consider the 4×4 Hilbert matrix

$$H = \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} \\ \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \frac{1}{7} \end{pmatrix},$$

whose entries are defined by

$$h_{ij} = \frac{1}{(i+j-1)}, \text{ for } i, j = 1, 2, \dots, n.$$

The inverse of the matrix H can be obtained as

$$H^{-1} = \begin{pmatrix} 16 & -120 & 240 & -140 \\ -120 & 1200 & -2700 & 1680 \\ 240 & -2700 & 6480 & -4200 \\ -140 & 1680 & -4200 & 2800 \end{pmatrix}.$$

Then the condition number of the Hilbert matrix is

$$K(H) = ||H||_{\infty} ||H^{-1}||_{\infty} = (2.0833)(13620) \approx 28375,$$

which is quite large. Note that the condition number of Hilbert matrices increase rapidly as the size of the matrices increases. Therefore, large Hilbert matrices are considered to be extremely ill-conditioned.

We might think that if the determinant of a matrix is close to zero, then the matrix is ill-conditioned. However, this is false. Consider the following matrix

$$A = \left(\begin{array}{cc} 10^{-7} & 0\\ 0 & 10^{-7} \end{array} \right),$$

for which det $A = 10^{-14} \approx 0$. One can easily find the condition number of the given matrix as

$$K(A) = ||A||_{\infty} ||A^{-1}||_{\infty} = (10^{-7})(10^{7}) = 1.$$

The matrix A is therefore perfectly conditioned. Thus a small determinant is necessary but not sufficient for a matrix to be ill-conditioned.

The condition number of a matrix K(A) using l_2 -norm can be computed by the built-in function **cond** command in the MATLAB as follows:

$$>> A = [1 -1 2; 3 1 -1; 2 0 1]; K(A) = cond(A)$$

Theorem 3.18 (Error in Linear Systems)

Suppose that \mathbf{x}^* is an approximation to the solution \mathbf{x} of the linear system $A\mathbf{x} = \mathbf{b}$ and A is a nonsingular matrix and \mathbf{r} is the residual vector for \mathbf{x}^* . Then for any natural norm, the error is

$$\|\mathbf{x} - \mathbf{x}^*\| \le \|\mathbf{r}\| \|A^{-1}\|,$$
 (3.56)

and the relative error is

$$\frac{\|\mathbf{x} - \mathbf{x}^*\|}{\|\mathbf{x}\|} \le K(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}, \quad \text{provided that} \quad \mathbf{x} \ne 0, \ \mathbf{b} \ne 0.$$
 (3.57)

Proof. Since $\mathbf{r} = \mathbf{b} - A\mathbf{x}^*$ and A is nonsingular, then

$$A\mathbf{x} - A\mathbf{x}^* = \mathbf{b} - (\mathbf{b} - \mathbf{r}) = \mathbf{r},$$

which implies that

$$A(\mathbf{x} - \mathbf{x}^*) = \mathbf{r}, \quad \text{or} \quad \mathbf{x} - \mathbf{x}^* = A^{-1}\mathbf{r}.$$
 (3.58)

Taking norm on both side, gives

$$\|\mathbf{x} - \mathbf{x}^*\| = \|A^{-1}\mathbf{r}\| \le \|A^{-1}\|\|\mathbf{r}\|.$$

Moreover, since $\mathbf{b} = A\mathbf{x}$, then

$$\|\mathbf{b}\| \le \|A\| \|\mathbf{x}\|, \quad \text{or,} \quad \|\mathbf{x}\| \ge \frac{\|\mathbf{b}\|}{\|A\|}.$$

Hence

$$\frac{\|\mathbf{x} - \mathbf{x}^*\|}{\|\mathbf{x}\|} \le \frac{\|A^{-1}\| \|\mathbf{r}\|}{\|\mathbf{b}\| / \|A\|} \le K(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}.$$

The inequalities (3.56) and (3.57) imply that the quantities $||A^{-1}||$ and K(A) can be used to give an indication of the connection between the residual vector and the accuracy of the approximation. If the quantity $K(A) \approx 1$, the relative error will be fairly close to the relative residual. But if K(A) >> 1, then the relative error could be many times larger than the relative residual.

Example 3.46 Find the condition number of the following matrix (for n = 2, 3, ...)

$$A_n = \left[\begin{array}{cc} 1 & 1 \\ 1 & 1 - 1/n \end{array} \right].$$

If n = 2 and $x^* = [-1.99, 2.99]^T$ be the approximate solution of the linear system $A\mathbf{x} = [1, -0.5]^T$, then find the relative error.

Solution. We can easily find the inverse of the given matrix as

$$A_n^{-1} = \frac{1}{(1-1/n)-1} \left[\begin{array}{cc} 1-1/n & -1 \\ -1 & 1 \end{array} \right] = -n \left[\begin{array}{cc} 1-1/n & -1 \\ -1 & 1 \end{array} \right] = \left[\begin{array}{cc} 1-n & n \\ n & -n \end{array} \right].$$

Then the l_{∞} -norm of both matrices A_n and A_n^{-1} are

$$||A_n||_{\infty} = 2$$
 and $||A_n^{-1}||_{\infty} = 2n$,

and so the condition number of the matrix can be computed as follows:

$$K(A) = ||A_n||_{\infty} ||A_n^{-1}||_{\infty} = (2)(2n) = 4n$$
 and $\lim_{n \to \infty} K(A) = \infty$,

which shows that the matrix A_n is obviously ill-conditioned. Here we expect that the relative error in the calculated solution to a linear system of the form $A_n \mathbf{x} = \mathbf{b}$ could be as much as 4n times the

relative residual.

The residual vector (by taking n = 2) can be calculated as

$$\mathbf{r} = \mathbf{b} - A_2 \mathbf{x}^* = \begin{pmatrix} 1 \\ -0.5 \end{pmatrix} - \begin{pmatrix} 1 & 1 \\ 1 & 0.5 \end{pmatrix} \begin{pmatrix} -1.99 \\ 2.99 \end{pmatrix} = \begin{pmatrix} 0.000 \\ -0.005 \end{pmatrix},$$

and it gives $\|\mathbf{r}\|_{\infty} = 0.005$. Now using (3.57), we obtain

$$\frac{\|\mathbf{x} - \mathbf{x}^*\|}{\|\mathbf{x}\|} \le K(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} = (8) \frac{0.005}{1} = 0.0400,$$

which is the required relative error.

Example 3.47 Consider a following linear system

- (a) Discuss the ill-conditioning of the given linear system.
- (b) If $\mathbf{x}^* = [2.01, 1.01, 1.98]^T$ be an approximate solution of the given system, then find the residual vector \mathbf{r} and its norm $\|\mathbf{r}\|_{\infty}$.
- (c) Estimate the relative error using (3.57).
- (d) Use the simple Gaussian elimination method to find approximate error using (3.58).

Solution. (a) Given the matrix

$$A = \left(\begin{array}{rrr} 1 & 1 & -1 \\ 1 & 2 & -2 \\ -2 & 1 & 1 \end{array}\right),$$

and whose inverse can be computed as

$$A^{-1} = \left(\begin{array}{ccc} 2 & -1 & 0\\ 1.5 & -0.5 & 0.5\\ 2.5 & -1.5 & 0.5 \end{array}\right).$$

Then the l_{∞} -norm of both matrices are

$$||A||_{\infty} = 5$$
 and $||A^{-1}||_{\infty} = 4.5$.

Using the values of both matrices norms, we can find the value of the condition number of A as follows:

$$K(A) = ||A||_{\infty} |||A^{-1}||_{\infty} = (5)(4.5) = 22.5 >> 1,$$

which shows that the matrix is ill-conditioned. Thus the given system is ill-conditioned.

$$>> A = [1\ 1\ -1; 1\ 2\ -2; -2\ 1\ 1];\ K(A) = norm(A, inf)*norm(inv(A), inf)$$

(b) The residual vector can be calculated as

$$\mathbf{r} = \mathbf{b} - A\mathbf{x}^* = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} - \begin{pmatrix} 1 & 1 & -1 \\ 1 & 2 & -2 \\ -2 & 1 & 1 \end{pmatrix} \begin{pmatrix} 2.01 \\ 1.01 \\ 1.98 \end{pmatrix} = \begin{pmatrix} -0.04 \\ -0.07 \\ 0.03 \end{pmatrix},$$

and it gives

$$\|\mathbf{r}\|_{\infty} = 0.07.$$

>>
$$A = [1 \ 1 \ -1; 1 \ 2 \ -2; -2 \ 1 \ 1]; \ b = [1 \ 0 \ -1]';$$

>> $x0 = [2.01 \ 1.01 \ 1.98]'; \ r = RES(A, b, x0); \ rnorm = norm(r, inf);$

(c) From (3.57), we have

$$\frac{\|\mathbf{x} - \mathbf{x}^*\|}{\|\mathbf{x}\|} \le K(A) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}.$$

By using above parts (a) and (b) and the value $\|\mathbf{b}\|_{\infty} = 1$, we obtain

$$\frac{\|\mathbf{x} - \mathbf{x}^*\|}{\|\mathbf{x}\|} \le (22.5) \frac{(0.07)}{1} = 1.575.$$

$$>> RelErr = (K(A) * rnorm)/norm(b, inf);$$

(d) To solve the linear system $A\mathbf{e} = \mathbf{r}$, where

$$A = \begin{pmatrix} 1 & 1 & -1 \\ 1 & 2 & -2 \\ -2 & 1 & 1 \end{pmatrix} \quad and \quad \mathbf{r} = \begin{pmatrix} -0.04 \\ -0.07 \\ 0.03 \end{pmatrix},$$

and $e = x - x^*$. Writing the above system in the augmented matrix form

$$\begin{pmatrix} 1 & 1 & -1 & \vdots & -0.04 \\ 1 & 2 & -2 & \vdots & -0.07 \\ -2 & 1 & 1 & \vdots & 0.03 \end{pmatrix}.$$

After applying forward elimination step of the simple Gauss elimination method, we obtain

$$\left(\begin{array}{ccccc}
1 & 1 & -1 & \vdots & -0.04 \\
0 & 1 & -1 & \vdots & -0.03 \\
0 & 0 & 2 & \vdots & 0.04
\end{array}\right).$$

Now by using the backward substitution, we obtain the solution

$$\mathbf{e}^* = [-0.01, -0.01, 0.02]^T,$$

which is the required approximation of the exact error.

168 3.8 Exercises

$$>> A = [1 \ 1 \ -1 \ -0.04; 1 \ 2 \ -2 \ -0.07; -2 \ 1 \ 1 \ 0.03]; WP(A);$$

Example 3.48 Consider a linear system $A\mathbf{x} = \mathbf{b}$, where

$$A = \begin{pmatrix} 2 & 1 & 2 \\ 1 & 4 & 0 \\ 1 & 2 & 1 \end{pmatrix} \quad and \quad \mathbf{b} = \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix}.$$

- (a) Discuss the conditioning of the given linear system.
- (b) Suppose that **b** is changed to $\mathbf{b}^* = [1, 1, 1.99]^T$. How large a relative change can this change produce in the solution to $A\mathbf{x} = \mathbf{b}$?

Solution. (a) Since the given matrix is

$$A = \left(\begin{array}{ccc} 2 & 1 & 2 \\ 1 & 4 & 0 \\ 1 & 2 & 1 \end{array}\right),$$

and its inverse can be computed as

$$A^{-1} = \begin{pmatrix} 4/3 & 1 & -8/3 \\ -1/3 & 0 & 2/3 \\ -2/3 & -1 & 7/3 \end{pmatrix}.$$

Then the l_{∞} -norm of both matrices are

$$||A||_{\infty} = 5$$
 and $||A^{-1}||_{\infty} = 5$.

Using the values of both matrices norms, we can find the value of the condition number of A as follows:

$$K(A) = ||A||_{\infty} |||A^{-1}||_{\infty} = (5)(5) = 25.$$

(b) Since the change from **b** to **b*** is an error δ **b**, that is, **b*** = **b** + δ **b**, so

$$\delta \mathbf{b} = \begin{pmatrix} -0.01 \\ 0 \\ 0 \end{pmatrix} = -\mathbf{r},$$

and the l_{∞} -norm of this column matrix is, $\|\delta \mathbf{b}\|_{\infty} = 0.01$. From the equation (3.57), we get

$$\frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} \le \frac{25(0.01)}{2} = 0.1250,$$

the possible relative change in the solution to the given linear system.

3.8 Exercises

1. Determine the matrix C given by the following expression

$$C = 2A - 3B,$$

if the matrices A and B are

$$A = \begin{pmatrix} 2 & -1 & 1 \\ -1 & 2 & 3 \\ 2 & 1 & 2 \end{pmatrix}, \qquad B = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & 3 \\ 2 & 1 & 4 \end{pmatrix}.$$

- 2. Find the product AB and BA for the matrices of the Problem 1.
- 3. Show that the product AB of the following rectangular matrices is a singular matrix.

$$A = \begin{pmatrix} 6 & -3 \\ 1 & 4 \\ -2 & 1 \end{pmatrix}, \qquad B = \begin{pmatrix} 2 & -1 & -2 \\ 3 & -4 & -1 \end{pmatrix}.$$

4. Let

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 0 & -1 & 2 \\ 2 & 0 & 2 \end{pmatrix}, \quad B = \begin{pmatrix} 1 & 1 & 2 \\ -1 & 1 & -1 \\ 1 & 0 & 2 \end{pmatrix}, \quad C = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 2 \\ 2 & 0 & 1 \end{pmatrix}.$$

- (a) Compute AB and BA and show that $AB \neq BA$.
- (b) Find (A + B) + C and A + (B + C). (c) Show that $(AB)^T = B^T A^T$.
- 5. Find a value of x and y such that $AB^T = C^T$, where

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 2 & 0 \\ 2 & 1 & 3 \end{pmatrix}, \quad B = \begin{bmatrix} 1 & x & 1 \end{bmatrix}, \quad C = \begin{bmatrix} -2 & -2 & y \end{bmatrix}.$$

6. Find the values of a and b such that each of the following matrix is symmetric:

(a)
$$A = \begin{pmatrix} 1 & 3 & 5 \\ a+2 & 5 & 6 \\ b+1 & 6 & 7 \end{pmatrix}$$
, (b) $B = \begin{pmatrix} -2 & a+b & 2 \\ 3 & 4 & 2a+b \\ 2 & 5 & -3 \end{pmatrix}$,

(c)
$$C = \begin{pmatrix} 1 & 4 & a - b \\ 4 & 2 & a + 3b \\ 7 & 3 & 4 \end{pmatrix}$$
, (d) $D = \begin{pmatrix} 1 & a - 4b & 2 \\ 2 & 8 & 6 \\ 7 & a - 7b & 8 \end{pmatrix}$.

7. Which of the following matrices are skew-symmetric?

(a)

$$A = \begin{pmatrix} 1 & -5 \\ 5 & 0 \end{pmatrix}, \qquad B = \begin{pmatrix} 0 & -4 \\ 4 & 0 \end{pmatrix},$$

170 3.8 Exercises

(b)
$$C = \begin{pmatrix} 1 & 9 \\ -9 & 7 \end{pmatrix}, \qquad D = \begin{pmatrix} 1 & 6 \\ -6 & 2 \end{pmatrix},$$
(c)
$$E = \begin{pmatrix} 0 & 2 & -2 \\ -2 & 0 & 4 \\ 2 & -4 & 0 \end{pmatrix}, \quad F = \begin{pmatrix} 3 & -3 & -3 \\ 3 & 3 & -3 \\ 3 & 3 & 3 \end{pmatrix}.$$

8. Compute the determinant of each of the following matrix using cofactor expansion along any row or column:

$$A = \begin{pmatrix} \cos x & \sin x & 1 \\ 0 & 3\cos x & -3\sin x \\ 0 & 2\sin x & 2\cos x \end{pmatrix}, \quad B = \begin{pmatrix} x & y & z \\ 0 & x^2 & y \\ 0 & y^2 & x \end{pmatrix}, \quad C = \begin{pmatrix} 2x & 0 & z \\ 0 & 2y & -z \\ z & -z & 2z \end{pmatrix}.$$

9. Compute the determinant of each of the following matrix using cofactor expansion along any row or column:

$$A = \begin{pmatrix} 3 & 7 & 6 \\ 0 & 3 & 5 \\ 7 & 4 & 3 \end{pmatrix}, \quad B = \begin{pmatrix} 11 & -6 & 4 \\ -16 & 8 & 6 \\ 5 & 7 & 12 \end{pmatrix}, \quad C = \begin{pmatrix} 4 & -8 & 11 \\ 10 & 1 & 4 \\ 7 & 10 & 8 \end{pmatrix}.$$

10. Find all zeros (values of x such that f(x) = 0) of polynomial $f(x) = \det(A)$ where

$$A = \left(\begin{array}{ccc} x - 1 & 3 & 2 \\ 3 & x & 1 \\ 2 & 1 & x - 2 \end{array}\right).$$

11. Let

$$A = \left(\begin{array}{cc} 1 & 1 \\ 0 & 1 \end{array}\right), \quad B = \left(\begin{array}{cc} 1 & 0 \\ 1 & 1 \end{array}\right),$$

then show that $(AB)^{-1} = B^{-1}A^{-1}$.

12. Find all zeros (values of x such that f(x) = 0) of polynomial $f(x) = \det(A)$ where

$$A = \left(\begin{array}{ccc} x & 0 & 1 \\ 2 & 1 & 3 \\ 0 & x & 2 \end{array} \right).$$

13. Compute the adjoint of each matrix A, and find the inverse of it, if it exists:

(a)
$$A = \begin{pmatrix} 1 & 2 \\ -3 & 4 \end{pmatrix}$$
, (b) $A = \begin{pmatrix} 1 & 2 & -1 \\ 2 & 1 & 4 \\ 1 & 5 & -8 \end{pmatrix}$, (c) $A = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix}$.

14. Find all zeros (values of x such that f(x) = 0) of polynomial $f(x) = \det(A)$ where

$$A = \left(\begin{array}{cccc} x & -8 & 5 & 2 \\ -3 & x & 2 & 1 \\ 3 & 4 & x & 1 \\ 3 & 6 & -5 & 17 \end{array}\right).$$

15. Show that $A(Adj A) = (Adj A)A = \det(A)\mathbf{I}_3$, if

$$A = \left(\begin{array}{rrr} 2 & 1 & 3 \\ -1 & 2 & 0 \\ 3 & -2 & 1 \end{array}\right).$$

16. Find the inverse and determinant of the adjoint matrix of each following matrix:

$$A = \begin{pmatrix} 4 & 1 & 5 \\ 5 & 6 & 3 \\ 5 & 4 & 4 \end{pmatrix}, \quad B = \begin{pmatrix} 3 & 4 & -2 \\ 2 & 5 & 4 \\ 7 & -3 & 4 \end{pmatrix}, \quad C = \begin{pmatrix} 1 & 2 & 4 \\ 1 & 4 & 0 \\ 3 & 1 & 1 \end{pmatrix}.$$

17. Find the inverse and determinant of the adjoint matrix of each following matrix:

$$A = \begin{pmatrix} 3 & 2 & 5 \\ 2 & 5 & 4 \\ 5 & 4 & 6 \end{pmatrix}, \quad B = \begin{pmatrix} 5 & 3 & -2 \\ 3 & 5 & 6 \\ -2 & 6 & 5 \end{pmatrix}, \quad C = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 8 \end{pmatrix}.$$

18. Find inverse of each of the following matrix using determinant:

$$A = \begin{pmatrix} 0 & 1 & 5 \\ 3 & 1 & 2 \\ 2 & 3 & 4 \end{pmatrix}, \quad B = \begin{pmatrix} 2 & 4 & -2 \\ -4 & 7 & 5 \\ 5 & -4 & 4 \end{pmatrix}, \quad C = \begin{pmatrix} 0 & 4 & 2 & -4 \\ 6 & 1 & 4 & -3 \\ 4 & 3 & 1 & 3 \\ 8 & 4 & -3 & 2 \end{pmatrix}.$$

19. Use matrices in Problem 15, solve the following systems using matrix inversion method:

(a)
$$A\mathbf{x} = [1, 1, -3]^T$$
, (b) $B\mathbf{x} = [2, 1, 3]^T$, (c) $C\mathbf{x} = [1, 0, 1]^T$.

20. Solve the following systems using the matrix inversion method: (a)

3.8 Exercises

21. Solve the following systems using the matrix inversion method: (a)

$$3x_1 - 2x_2 - 4x_3 = 7$$

$$5x_1 - 2x_2 - 3x_3 = 8$$

$$7x_1 + 4x_2 + 2x_3 = 9$$
(b)

$$\begin{array}{rclrcrcr}
-3x_1 & + & 4x_2 & + & 3x_3 & = & 11 \\
5x_1 & + & 3x_2 & + & x_3 & = & 12 \\
x_1 & + & x_2 & + & 5x_3 & = & 10
\end{array}$$

(c)
$$x_1 + 4_2 - 8x_3 = 7$$

$$2x_1 + 7x_2 - 5x_3 = -5$$

$$3x_1 - 6x_2 + 6x_3 = 4$$

22. Use the simple Gaussian elimination method to show that the following system does not have a solution

$$3x_1 + x_2 = 1.5$$

 $2x_1 - x_2 - x_3 = 2$
 $4x_1 + 3x_2 + x_3 = 0$

- 23. Solve the Problem 21 using the simple Gaussian elimination method.
- 24. Solve the following systems using the simple Gaussian elimination method:

(a)
$$x_1 - x_2 = -2$$

$$-x_1 + 2x_2 - x_3 = 5$$

$$4x_1 - x_2 + 4x_3 = 1$$
(b)
$$3x_1 + x_2 - x_3 = 5$$

$$3x_1 + x_2 - x_3 = 5$$

$$5x_1 - 3x_2 + 2x_3 = 7$$

$$2x_1 - x_2 + x_3 = 3$$

(c)
$$3x_1 + x_2 + x_3 = 2$$
$$2x_1 + 2x_2 + 4x_3 = 3$$
$$4x_1 + 9x_2 + 16x_3 = 1$$

25. Solve the following systems using the simple Gaussian elimination method

 $x_1 + 3x_2 - 5x_3 = 1$

(a)
$$2x_1 + 5x_2 - 4x_3 = 3 \\ 2x_1 + 2x_2 - x_3 = 1 \\ 3x_1 + 2x_2 - 3x_3 = -5$$
 (b)
$$2x_2 - x_3 = 1 \\ 3x_1 - x_2 + 2x_3 = 4$$

(c)
$$x_1 + 2x_2 = 3 \\ -x_1 - 2x_3 = -5 \\ -3x_1 - 5x_2 + x_3 = -4$$

26. For what values of a and b the following linear system has no solution or infinitely many solutions.

(a)
$$2x_1 + x_2 + x_3 = 2 \\ -2x_1 + x_2 + 3x_3 = a \\ 2x_1 - x_3 = b$$

(b)
$$2x_1 + 3x_2 - x_3 = 1$$

$$x_1 - x_2 + 3x_3 = a$$

$$3x_1 + 7x_2 - 5x_3 = b$$

(c)
$$2x_1 - x_2 + 3x_3 = 3 3x_1 + x_2 - 5x_3 = a -5x_1 - 5x_2 + 21x_3 = b$$

27. Find the value(s) of α so that each of the following linear system has a non-trivial solution: (a)

$$\begin{array}{rclrcrcr}
2x_1 & + & 2x_2 & + & 3x_3 & = & 1 \\
3x_1 & + & \alpha x_2 & + & 5x_3 & = & 3 \\
x_1 & + & 7x_2 & + & 3x_3 & = & 2
\end{array}$$

(b)
$$x_1 + 2x_2 + x_3 = 2$$

$$x_1 + 3x_2 + 6x_3 = 5$$

$$2x_1 + 3x_2 + \alpha x_3 = 6$$

(c)
$$\alpha x_1 + x_2 + x_3 = 7$$

$$x_1 + x_2 - x_3 = 2$$

$$x_1 + x_2 + \alpha x_3 = 1$$

28. Find the inverse of each of the following matrix by using simple Gauss elimination method:

$$A = \begin{pmatrix} 3 & 3 & 3 \\ 0 & 2 & 2 \\ 2 & 4 & 5 \end{pmatrix}, \quad B = \begin{pmatrix} 5 & 3 & 2 \\ 3 & 2 & 2 \\ 2 & 6 & 5 \end{pmatrix}, \quad C = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 5 & 2 \\ 3 & 4 & 3 \end{pmatrix}.$$

29. Find the inverse of each of the following matrix by using simple Gauss elimination method:

$$A = \begin{pmatrix} 3 & 2 & 3 \\ 4 & 2 & 2 \\ 2 & 4 & 3 \end{pmatrix}, \quad B = \begin{pmatrix} 1 & -3 & 2 \\ 3 & 2 & 6 \\ 2 & -6 & 5 \end{pmatrix}, \quad C = \begin{pmatrix} 5 & 2 & 3 \\ 2 & 5 & 5 \\ 3 & 2 & 4 \end{pmatrix}.$$

174 3.8 Exercises

30. Determine the rank of each of the following matrix:

$$A = \begin{pmatrix} 3 & 1 & -1 \\ 2 & 0 & 4 \\ 1 & -5 & 1 \end{pmatrix}, \quad B = \begin{pmatrix} 4 & 1 & 6 \\ -3 & 6 & 4 \\ 5 & 0 & 9 \end{pmatrix}, \quad C = \begin{pmatrix} 17 & 46 & 7 \\ 20 & 49 & 8 \\ 23 & 52 & 9 \end{pmatrix}.$$

31. Solve Problem 25 using the Gaussian elimination with partial pivoting.

32. Let A be an $m \times n$ matrix and B be an $n \times p$ matrix. Show that the rank of AB is less than or equal to the rank of A.

33. Solve the following linear systems using the Gaussian elimination with partial and without pivoting

$$\begin{array}{rcl}
1.001x_1 & + & 1.5x_2 & = & 0 \\
2x_1 & + & 3x_2 & = & 1
\end{array}$$

$$x_1 + 1.001x_2 = 2.001$$

 $x_1 + x_2 = 2$

$$6.122x_1 + 1500.5x_2 = 1506.622$$

 $2000x_1 + 3x_2 = 2003$

34. The elements of the matrix A, the Hilbert matrix, are defined by

$$a_{ij} = 1/(i+j-1)$$
, for $i, j = 1, 2, \dots, n$

Find the solution of the system $A\mathbf{x} = \mathbf{b}$ for n = 4 and $\mathbf{b} = [1, 2, 3, 4]^T$, using the Gaussian elimination by partial pivoting.

35. Find the LU decomposition of each matrix A using Doolittle's method, and then solve the systems.

$$A = \begin{pmatrix} 3 & -2 & 1 & 1 \\ -3 & 7 & 4 & -3 \\ 2 & -5 & 3 & 4 \\ 7 & -3 & 2 & 4 \end{pmatrix}, \qquad \mathbf{b} = \begin{pmatrix} 3 \\ 2 \\ 1 \\ 2 \end{pmatrix}.$$

$$A = \begin{pmatrix} 2 & -4 & 5 & 3 \\ 3 & 5 & -4 & 3 \\ 1 & 6 & 2 & 6 \\ 7 & 2 & 5 & 1 \end{pmatrix}, \qquad \mathbf{b} = \begin{pmatrix} 6 \\ 5 \\ 2 \\ 4 \end{pmatrix}.$$

$$A = \begin{pmatrix} 2 & 2 & 3 & -2 \\ 10 & 2 & 13 & 11 \\ 2 & 5 & 4 & 6 \\ 1 & -4 & -2 & 7 \end{pmatrix}, \qquad \mathbf{b} = \begin{pmatrix} 10 \\ 14 \\ 11 \\ 9 \end{pmatrix}.$$

36. For what value(s) of α each of the following matrix A is singular using Doolittle's method.

$$(\mathbf{a}) \ A = \begin{pmatrix} 1 & -1 & 2 \\ -1 & 3 & -1 \\ \alpha & -2 & 3 \end{pmatrix}, \ (\mathbf{b}) \ A = \begin{pmatrix} 1 & 5 & 7 \\ 4 & 4 & \alpha \\ -2 & \alpha & 9 \end{pmatrix}, \ \ (\mathbf{c}) \ A = \begin{pmatrix} 2 & -4 & \alpha \\ 2 & 4 & 3 \\ 4 & -2 & 5 \end{pmatrix},$$

$$(\mathbf{d}) \ A = \begin{pmatrix} 2 & \alpha & 1 - \alpha \\ 2 & 5 & -2 \\ 2 & 5 & 4 \end{pmatrix}, \ (\mathbf{e}) \ A = \begin{pmatrix} 1 & -1 & 3 \\ 3 & 2 & 3 \\ 4 & \alpha - 2 & 7 \end{pmatrix}, \ (\mathbf{f}) \ A = \begin{pmatrix} 1 & 5 & \alpha \\ 1 & 4 & \alpha - 2 \\ 1 & -2 & 8 \end{pmatrix}.$$

37. Find the determinant of each of the following matrix using the LU decomposition by Doolittle's method:

(a)
$$A = \begin{pmatrix} 2 & 3 & -1 \\ 1 & 2 & 1 \\ 2 & 1 & -6 \end{pmatrix}$$
, (b) $A = \begin{pmatrix} 1 & -2 & 2 \\ 2 & 1 & 1 \\ 1 & 0 & 1 \end{pmatrix}$, (c) $A = \begin{pmatrix} 2 & 4 & 1 \\ 3 & 3 & 2 \\ 4 & 1 & 4 \end{pmatrix}$.

- 38. Use the smallest positive integer to find the unique solution each of the linear system of the problem 38 using LU decomposition by Doolittle's method.
 - (b)

- (d) $A\mathbf{x} = [-8, 11, 8]^T$
- (e)
- $A\mathbf{x} = [5, -6, 2]^T$ (c) $A\mathbf{x} = [11, 13, 10]^T$. $A\mathbf{x} = [32, 23, 12]^T$ (f) $A\mathbf{x} = [-11, 43, 22]^T$.
- 39. Find the determinant of each of the following matrix using the LU decomposition by Crout's method:

(a)
$$A = \begin{pmatrix} 2 & 2 & -1 \\ 1 & 2 & 1 \\ 2 & 1 & -4 \end{pmatrix}$$
, (b) $A = \begin{pmatrix} 2 & -1 & 1 \\ 1 & 2 & 2 \\ 2 & 0 & 2 \end{pmatrix}$,

(c)
$$A = \begin{pmatrix} 4 & 4 & 1 \\ 5 & 4 & 2 \\ 1 & 4 & 4 \end{pmatrix}$$
, (d) $A = \begin{pmatrix} 2 & 4 & 5 \\ 3 & 5 & 3 \\ 4 & 3 & 2 \end{pmatrix}$.

- 40. Find $\|\mathbf{x}\|_1$, $\|\mathbf{x}\|_2$ and $\|\mathbf{x}\|_{\infty}$ for the following vectors.
- (b) $[\sin k, \cos k, 3^k]^T$, for a fixed integer k.
- 41. Find $\|.\|_1$, $\|.\|_{\infty}$ and $\|.\|_e$ for the following matrices.

$$A = \begin{pmatrix} 3 & 1 & -1 \\ 2 & 0 & 4 \\ 1 & -5 & 1 \end{pmatrix}, \quad B = \begin{pmatrix} 4 & 1 & 6 \\ -3 & 6 & 4 \\ 5 & 0 & 9 \end{pmatrix},$$

$$C = \begin{pmatrix} 17 & 46 & 7 \\ 20 & 49 & 8 \\ 23 & 52 & 9 \end{pmatrix}, \quad D = \begin{pmatrix} 3 & 11 & -5 & 2 \\ 6 & 8 & -11 & 6 \\ -4 & -8 & 10 & 14 \\ 13 & 14 & -12 & 9 \end{pmatrix}.$$

176 3.8 Exercises

42. Consider the following matrices

$$A = \begin{pmatrix} -11 & 7 & -8 \\ 5 & 9 & 6 \\ 6 & 3 & 7 \end{pmatrix}, \quad B = \begin{pmatrix} 6 & 2 & 7 \\ -12 & 10 & 8 \\ 3 & -15 & 14 \end{pmatrix},$$

$$C = \begin{pmatrix} 5 & -6 & 4 \\ -7 & 8 & 5 \\ 3 & -9 & 12 \end{pmatrix}, \quad D = \begin{pmatrix} 2 & 1 & -1 & 1 \\ 1 & 3 & 5 & 2 \\ -2 & -3 & 4 & 5 \\ 3 & 4 & -2 & 4 \end{pmatrix}.$$

Find $\|.\|_1$ and $\|.\|_{\infty}$ for (a) A^3 , (b) $A^2 + B^2 + C^2 + D^2$, (c) BC and (d) $C^2 + D^2$.

43. Compute the condition numbers of the following matrices relative to $\|.\|_{\infty}$

(a)
$$\begin{pmatrix} \frac{1}{3} & \frac{1}{2} & \frac{1}{5} \\ \frac{1}{2} & \frac{1}{5} & \frac{1}{3} \\ \frac{1}{5} & \frac{1}{3} & \frac{1}{2} \end{pmatrix}$$
, (b)
$$\begin{pmatrix} 0.03 & 0.01 & -0.02 \\ 0.15 & 0.51 & -0.11 \\ 1.11 & 2.22 & 3.33 \end{pmatrix}$$
, (c)
$$\begin{pmatrix} 1.11 & 1.98 & 2.01 \\ 1.01 & 1.05 & 2.05 \\ 0.85 & 0.45 & 1.25 \end{pmatrix}$$
.

44. The $n \times n$ Hilbert matrix $H^{(n)}$ defined by

$$H_{ij}^{(n)} = \frac{1}{i+j-1}, \qquad 1 \le i, \ j \le n.$$

Find the l_{∞} -norm of the 10×10 Hilbert matrix.

45. The following linear systems have \mathbf{x} as the exact solution and \mathbf{x}^* is an approximate solution. Compute $\|\mathbf{x} - \mathbf{x}^*\|_{\infty}$ and $K(A) \frac{\|\mathbf{r}\|_{\infty}}{\|\mathbf{b}\|_{\infty}}$, where $\mathbf{r} = \mathbf{b} - A\mathbf{x}^*$ is the residual vector.

(a)
$$0.89x_1 + 0.53x_2 = 0.36$$

$$0.47x_1 + 0.28x_2 = 0.19$$

$$\mathbf{x} = [1, -1]^T$$

$$\mathbf{x}^* = [0.702, -0.500]^T$$
(b)
$$0.986x_1 + 0.579x_2 = 0.235$$

$$0.409x_1 + 0.237x_2 = 0.107$$

$$\mathbf{x} = [2, -3]^T$$

$$\mathbf{x}^* = [2.110, -3.170]^T$$
(c)
$$1.003x_1 + 58.090x_2 = 68.12$$

$$5.550x_1 + 321.8x_2 = 377.3$$

$$\mathbf{x} = [10, 1]^T$$

$$\mathbf{x}^* = [-10, 1]^T$$

46. Discuss the ill-conditioning (stability) of the linear system

$$\begin{array}{rcl}
1.01x_1 & + & 0.99x_2 & = & 2 \\
0.99x_1 & + & 1.01x_2 & = & 2
\end{array}$$

If $\mathbf{x}^* = [2, 0]^t$ be an approximate solution of the system, then find the residual vector \mathbf{r} and estimate the relative error.

47. The exact solution of the following linear system

is $\mathbf{x} = [-99, 100]^T$. Change the coefficient matrix slightly to

$$\delta A = \left(\begin{array}{cc} 1 & 1\\ 1 & 0.99 \end{array}\right),$$

and consider the linear system

$$\begin{array}{rcl}
x_1 & + & x_2 & = & 1 \\
x_1 & + & 0.99x_2 & = & 2
\end{array}$$

Compute the change solution $\delta \mathbf{x}$ of the system. Is the matrix A ill-conditioned?

48. The exact solution of the following linear system

$$x_1 + 3x_2 = 4$$

 $1.0001x_1 + 3x_2 = 4.0001$

is $\mathbf{x} = [1, 1]^T$. Change the right-hand vector \mathbf{b} slightly to $\delta \mathbf{b} = [4.0001, 4.0003]^T$ and consider the linear system

$$x_1 + 3x_2 = 4.0001$$

 $1.0001x_1 + 3x_2 = 4.0003$

Compute the change solution $\delta \mathbf{x}$ of the system. Is the matrix A ill-conditioned?

49. The exact solution of the following linear system

$$\begin{array}{rclcrcl} x_1 & + & x_2 & = & 3 \\ x_1 & + & 1.0005x_2 & = & 3.0010 \end{array}$$

is $\mathbf{x} = [1, 2]^T$. Change the coefficient matrix and the right-hand vector **b** slightly to

$$\delta A = \begin{pmatrix} 1 & 1 \\ 1 & 1.001 \end{pmatrix}$$
 and $\delta \mathbf{b} = \begin{pmatrix} 2.99 \\ 3.01 \end{pmatrix}$,

and consider the linear system

$$x_1 + x_2 = 2.99$$

 $x_1 + 1.001x_2 = 3.01$

Compute the change solution $\delta \mathbf{x}$ of the system. Is the matrix A ill-conditioned?

178 3.8 Exercises

50. Find the condition number of the following matrix

$$A_n = \left(\begin{array}{cc} 1 & 1 \\ 1 & 1 - \frac{1}{n} \end{array}\right).$$

Solve the linear system $A_4\mathbf{x} = [2, 2]^T$ and compute the relative residual.

51. Find the Jacobi iteration matrix and its l_{∞} -norm for each of the following matrix.

(a)
$$\begin{pmatrix} 11 & -3 & 2 \\ 4 & 10 & 3 \\ -2 & 5 & 9 \end{pmatrix}$$
, (b) $\begin{pmatrix} 7 & 1 & 1 \\ 3 & 13 & 2 \\ -4 & 3 & 14 \end{pmatrix}$,
(c) $\begin{pmatrix} 8 & 1 & -1 & 0 \\ 2 & 13 & -2 & 1 \\ -1 & 3 & 15 & 2 \\ 1 & 4 & 5 & 20 \end{pmatrix}$, (d) $\begin{pmatrix} 7 & 1 & -3 & 1 \\ 1 & 10 & 2 & -3 \\ 1 & -5 & 25 & 4 \\ 1 & 2 & 3 & 17 \end{pmatrix}$.

52. Find the Gauss-Seidel iteration matrix and its l_{∞} -norm for each of the following matrix.

(a)
$$\begin{pmatrix} 3 & 0 & 1 \\ 1 & 4 & 0 \\ 0 & 2 & 5 \end{pmatrix}$$
, (b) $\begin{pmatrix} 5 & 2 & 1 \\ 4 & 9 & 2 \\ 3 & 1 & 6 \end{pmatrix}$.

53. Solve the following linear systems using the Jacobi method, start with initial approximation $\mathbf{x}^{(0)} = 0$ and iterate until $\|\mathbf{x}^{(k+1)} - x^{(k)}\|_{\infty} \le 10^{-5}$ for each system. (a)

(a)
$$4x_1 - x_2 + x_3 = 7$$

$$4x_1 - 8x_2 + x_3 = -21$$

$$-2x_1 + x_2 + 5x_3 = 15$$
(b)
$$3x_1 + x_2 + x_3 = 5$$

$$2x_1 + 6x_2 + x_3 = 9$$

$$x_1 + x_2 + 4x_3 = 6$$
(c)
$$4x_1 + 2x_2 + x_3 = 1$$

$$x_1 + 7x_2 + x_3 = 4$$

$$x_1 + x_2 + 20x_3 = 7$$

54. Consider the linear system $A\mathbf{x} = \mathbf{b}$, where

$$A = \begin{pmatrix} -5 & 1 & 0 \\ 1 & 5 & -1 \\ 0 & 1 & 2 \end{pmatrix} \quad \text{and} \quad \mathbf{b} = \begin{pmatrix} 4 \\ -2 \\ -5 \end{pmatrix}.$$

Find the Jacobi iteration matrix T_J and show that $||T_J|| < 1$. Use Jacobi method to find first approximate solution $\mathbf{x}^{(1)}$ of the linear system by using $\mathbf{x}^{(0)} = [0, 0, 0]^T$. Also, compute the error bound $||\mathbf{x} - \mathbf{x}^{(10)}||$. Compute the number of steps needed to get the accuracy within 10^{-5} .

55. Consider the following system of equations

- (a) Show that the Jacobi method converges by using $||T_J||_{\infty} < 1$.
- (b) If the first approximate solution of the system by Jacobi method is $\mathbf{x}^{(1)} = [0.25, 0.57, 0.35]^T$, starting with $\mathbf{x}^{(0)} = [0, 0, 0]^T$, then compute an error estimate $\|\mathbf{x} \mathbf{x}^{(20)}\|_{\infty}$.
- 56. If

$$A = \begin{pmatrix} 4 & 1 & 0 \\ 1 & 3 & -1 \\ 0 & -1 & 4 \end{pmatrix} \quad \text{and} \quad \mathbf{b} = \begin{pmatrix} 3 \\ 4 \\ 5 \end{pmatrix}.$$

Find the Jacobi iteration matrix T_J . If the first approximate solution of the given linear system by the Jacobi method is $[3/4, 4/3, 5/4]^T$, using $\mathbf{x}^{(0)} = [0, 0, 0]^T$, then estimate the number of iterations necessary to obtain approximations accurate to within 10^{-6} .

57. Rearrange the following system such that convergence of Gauss-seidel method is guaranteed. Then use $x^{(0)} = [0, 0, 0]^T$ to find first approximation by Gauss-Seidel method. Also, compute an error bound $\|\mathbf{x} - \mathbf{x}^{(10)}\|$.

$$A = \begin{pmatrix} 1 & 2 & 4 \\ 5 & -1 & 1 \\ 0 & 3 & -1 \end{pmatrix} \quad textrmand \quad \mathbf{b} = \begin{pmatrix} 2 \\ 1 \\ -1 \end{pmatrix}.$$

58. Consider the following system of equations

$$4x_1 + 2x_2 + x_3 = 11$$

 $-x_1 + 2x_2 = 3$
 $2x_1 + x_2 + 4x_3 = 16$

- (a) Show that the Gauss-Seidel method converges by using $||T_G||_{\infty} < 1$.
- (b) Compute the second approximation $\mathbf{x}^{(2)}$, starting with $\mathbf{x}^{(0)} = [1, 1, 1]^T$.
- (c) Compute an error estimate $\|\mathbf{x} \mathbf{x}^{(2)}\|_{\infty}$ for your approximation.
- 59. Consider the linear system $A\mathbf{x} = \mathbf{b}$, where

$$A = \begin{pmatrix} -5 & 2 & 1\\ 1 & -10 & 1\\ 1 & 1 & -4 \end{pmatrix}, \qquad \mathbf{b} = \begin{pmatrix} -3\\ 27\\ 4 \end{pmatrix}.$$

Find the Gauss-Seidel iteration matrix T_G and show that $||T_G|| < 1$. Use Gauss-Seidel method to find second approximate solution $\mathbf{x}^{(2)}$ of the linear system using $\mathbf{x}^{(0)} = [-0.5, -2.5, -1.5]^T$. Also, compute the error bound.

180 3.8 Exercises

60. Consider linear system $A\mathbf{x} = \mathbf{b}$, where the coefficient matrix is

$$A = \left(\begin{array}{ccc} 6 & -1 & 1\\ 1 & 5 & -1\\ 1 & 2 & 9 \end{array}\right).$$

Show that Gauss-Seidel method converges faster than Jacobi method. If the first approximate solution of the given linear system by the Gauss-Seidel method is $[0.5, 0.5, 0.5]^T$, using $\mathbf{x}^{(0)} = [0, 0, 0]^T$, then estimate the number of iterations necessary to obtain approximations accurate to within 10^{-5} .

61. Consider the following system

- (a) Find the matrix form of both iterative (Jacobi and Gauss-Seidel) methods.
- (b) If $\mathbf{x}^{(k)} = [x_1^{(k)}, x_2^{(k)}, x_3^{(k)}]^T$, then writing the iterative forms of part(a) in the component forms and find the exact solution of the given system.
- (c) Find formulas for the error $e^{(k+1)}$ in the (n+1)th step.
- (d) Find the second approximation of the error $e^{(2)}$ using the part (c) if $\mathbf{x}^{(0)} = [0, 0, 0]^T$.
- 62. Consider the following linear system $A\mathbf{x} = \mathbf{b}$, where

$$A = \begin{pmatrix} 5 & -1 & 3 \\ 4 & 7 & -2 \\ 6 & -3 & 9 \end{pmatrix} \quad \text{and} \quad \mathbf{b} = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}.$$

Show that Gauss-Seidel method converges for the given linear system. If the first approximate solution of the given linear system by the Gauss-Seidel method is $\mathbf{x}^{(1)} = [0.2, 0.17, 0.26]^T$, by using initial approximation $\mathbf{x}^{(0)} = [0, 0, 0]^T$, then compute an upper bound $\|\mathbf{x} - \mathbf{x}^{(2)}\|_{\infty}$. Also, compute number of steps needed to get accuracy within 10^{-4} .

63. Consider the following system

- (a) Find the matrix form of both iterative (Jacobi and Gauss-Seidel) methods.
- (b) If $\mathbf{x}^{(k)} = [x_1^{(k)}, x_2^{(k)}, x_3^{(k)}]^T$, then writing the iterative forms of part(a) in the component forms and find the exact solution of the given system.
- (c) Find formulas for the error $e^{(k+1)}$ in the (n+1)th step.
- (d) Find the second approximation of the error $e^{(2)}$ using the part (c) if $\mathbf{x}^{(0)} = [0, 0, 0]^T$.

Chapter 4

Polynomial Interpolation and Approximation

4.1 Introduction

In this chapter we describe the numerical methods for the approximation of functions other than the elementary functions. The main purpose of these techniques is to replace a complicated function by one which is simpler and more manageable. We sometimes know the value of a function f(x) at a set of points (say, $x_0 < x_1 < x_2 \cdots < x_n$) but we do not have an analytic expression for f(x) that let us calculate its value at an arbitrary point. We concentrate on techniques which may be adapted if, for example, we have a table of values of function may have been obtained from some physical measurement or some experiments or long numerical calculation that can not be cast into a simple functional form. The task now is to estimate f(x) for an arbitrary point x by, in some sense, drawing a smooth curve through (and perhaps beyond) the data points x_i . If the desired x is in between the largest and smallest of the data point, then the problem is called *interpolation*; and if x is outside that range, it is called *extrapolation*. In this chapter we shall restrict our attention to interpolation. It is a rational process generally used in estimating a missing functional value by taking a weighted average of known functional values at neighbouring data points.

Interpolation scheme must model the function, in between or beyond the known data point, by some plausible functional form. The form should be sufficiently general so as to be able to approximate large classes of functions which might arise in practice. The functional form are polynomials, trigonometric functions, rational functions and exponential functions. However, we shall restrict our attention to polynomials. The polynomial functions are widely used in practice, since they are easy to determine, evaluate, differentiate, and integrable. Polynomial interpolation provides some mathematical tools that can be used in developing methods for approximation theory, numerical differentiation, numerical integration, and numerical solutions of ordinary differential equations and partial differential equations. A set of data points we consider here may be equally spaced or unequally spaced in the independent variable x. Several procedures can be used to fit approximation polynomials for both cases. For example, Lagrange interpolatory polynomial and Newton divided-difference interpolatory polynomial. These methods are quite easy to apply. The general form of a nth-degree polynomial is

$$f(x) = p_n(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n,$$
(4.1)

where n denotes the degree of the polynomial; and a_0, a_1, \ldots, a_n are constants coefficients. Since there are (n+1) coefficients, so (n+1) data points are required to obtain unique value for the coefficients. The important property of polynomials that makes them suitable for approximating functions is due to the following theorem called, the Weierstrass approximation theorem.

Important Points of the Chapter 4

- I. In this chapter we look for the approximate solution of a function at the given arbitrary point.
- II. We shall use polynomial interpolation (approximation of a function at a point $x \in [a, b]$).
- III. Higher the degree of interpolating polynomial better the approximate solution.
- IV. We shall use interpolating polynomial at equally and unequally spaced data points.
- V. Lagrange and Newton's polynomials may be used to find approximation of a function.
- VI. Newton's polynomial needs a table of divided differences of a function.
- VII. Piecewise linear interpolation can be used to obtain approximate solution of a function.

Theorem 4.1 (Weierstrass Approximation Theorem)

If f(x) is a continuous function in the closed interval [a,b] then for every $\epsilon > 0$ there exists a polynomial $p_n(x)$, where the value of n depends on the value of ϵ , such that for all x in [a,b],

$$|f(x) - p_n(x)| < \epsilon. \tag{4.2}$$

Consequently, any continuous function can be approximated to any accuracy by a polynomial of high enough degree.

4.2 Polynomial Interpolation

Suppose we have given a set of (n + 1) data points relating a dependent variables f(x) to an independent variable x as follows

Generally the data points x_0, x_1, \ldots, x_n are arbitrary and assume the interval between two adjacent points is not the same (unequally spaced) and assume that the data points are organized in such a way that $x_0 < x_1 < x_2 < \cdots < x_{n-1} < x_n$. But some times this is not happen.

When the data points in a given functional relationship are not equally spaced, the interpolation problem becomes more difficult to solve. The basis for this assertion lies in the fact that the interpolating polynomial coefficient will depend on the functional values as well as on the data points given in the table.

4.2.1 Lagrange Interpolating Polynomials

It is one of the popular and well known interpolation method to approximate the functions at an arbitrary point x. The Lagrange interpolation method provides a direct approach for determining

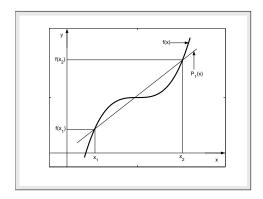


Figure 4.1: Linear Lagrange interpolation.

interpolated values regardless of the data points spacing, that is, it can be fitted to unequally spaced or equally spaced data. To discuss about the Lagrange interpolation method, we start with a simplest form of interpolation, that is, *linear interpolation*. The interpolated value is obtained from the equation of straight line that passes through two tabulated values, one each side of required value. This straight line is a first-degree polynomial. The problem of determining a polynomial of degree one that passes through the distinct points (x_0, y_0) and (x_1, y_1) is the same as approximating the function f(x) for which $f(x_0) = y_0$ and $f(x_1) = y_1$ by means of first degree polynomial interpolation.

Linear Lagrange Interpolating Polynomial

Let us consider the construction of a linear polynomial $p_1(x)$ passing through two data points $(x_0, f(x_0))$ and $(x_1, f(x_1))$, see Figure 4.1. Consider a linear polynomial of the form

$$f(x) = p_1(x) = a_0 + a_1 x. (4.3)$$

Since a polynomial of degree one has two coefficients, so one might expect to be able to choose two conditions, which satisfy

$$p_1(x_k) = f(x_k); \qquad k = 0, 1.$$

When $p_1(x)$ passes through point $(x_0, f(x_0))$, we have

$$p_1(x_0) = a_0 + a_1 x_0 = y_0 = f(x_0),$$

and if it passes through point $(x_1, f(x_1))$, we have

$$p_1(x_1) = a_0 + a_1 x_1 = y_1 = f(x_1).$$

Solving last two equations, gives a unique solution

$$a_0 = \frac{x_0 y_1 - x_1 y_0}{x_0 - x_1}$$
 and $a_1 = \frac{y_1 - y_0}{x_1 - x_0}$. (4.4)

Putting these values in (4.3), we have

$$f(x) = p_1(x) = \left(\frac{x - x_1}{x_0 - x_1}\right) y_0 + \left(\frac{x - x_0}{x_1 - x_0}\right) y_1,$$

Which can also be written as

$$f(x) = p_1(x) = L_0(x)f(x_0) + L_1(x)f(x_1), (4.5)$$

where

$$L_0(x) = \frac{x - x_1}{x_0 - x_1}$$
 and $L_1(x) = \frac{x - x_0}{x_1 - x_0}$. (4.6)

Note that when $x = x_0$, then $L_0(x_0) = 1$ and $L_1(x_0) = 0$. Similarly, when $x = x_1$, then $L_0(x_1) = 0$ and $L_1(x_1) = 1$. The polynomial (4.5) is known as linear Lagrange interpolating polynomial and (4.6) is called the Lagrange coefficient polynomials.

Quadratic Lagrange Interpolating Polynomial

When $p_2(x)$ passes through three points $(x_0, f(x_0)), (x_1, f(x_1))$ and $(x_2, f(x_2))$, we have quadratic Lagrange polynomial as follows

$$f(x) = p_2(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2), \tag{4.7}$$

where the Lagrange coefficients are define as follows:

$$L_0(x) = \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)},$$

$$L_1(x) = \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)},$$

$$L_2(x) = \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)}.$$

$$(4.8)$$

Cubic Lagrange Interpolating Polynomial

Similarly, when $p_3(x)$ passes through 4 points $(x_0, f(x_0)), (x_1, f(x_1)), (x_2, f(x_2))$ and $(x_3, f(x_3)),$ we have the following cubic Lagrange polynomial as follows

$$f(x) = p_3(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2) + L_3(x)f(x_3), \tag{4.9}$$

where the Lagrange coefficients are define as follows:

$$L_{0}(x) = \frac{(x-x_{1})(x-x_{2})(x-x_{3})}{(x_{0}-x_{1})(x_{0}-x_{2})(x_{0}-x_{3})},$$

$$L_{1}(x) = \frac{(x-x_{0})(x-x_{2})(x-x_{3})}{(x_{1}-x_{0})(x_{1}-x_{2})(x_{1}-x_{3})},$$

$$L_{2}(x) = \frac{(x-x_{0})(x-x_{1})(x-x_{3})}{(x_{2}-x_{0})(x_{2}-x_{1})(x_{2}-x_{3})},$$

$$L_{3}(x) = \frac{(x-x_{0})(x-x_{1})(x-x_{2})}{(x_{3}-x_{0})(x_{3}-x_{1})(x_{3}-x_{2})}.$$

$$(4.10)$$

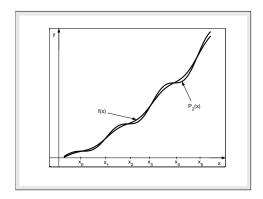


Figure 4.2: General Lagrange interpolation.

Nth Degree Lagrange Interpolating Polynomial

To generalize the concept of the Lagrange interpolation, consider the construction of a polynomial $p_n(x)$ of degree at most n that passes through (n+1) distinct points $(x_0, f(x_0)), \ldots, (x_n, f(x_n))$ (see Figure 4.2) and satisfy the interpolation conditions

$$p_n(x_k) = f(x_k); \qquad k = 0, 1, 2, \dots, n.$$
 (4.11)

Assume that there exists polynomial $L_k(x)$ (k = 0, 1, 2, ..., n) of degree n having the property

$$L_k(x_j) = \begin{cases} 0 & \text{for } k \neq j, \\ 1 & \text{for } k = j, \end{cases}$$
 (4.12)

and

$$\sum_{k=0}^{n} L_k(x) = 1. (4.13)$$

The polynomial $p_n(x)$ is given by

$$f(x) = p_n(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + \dots + L_{i-1}(x)f(x_{i-1}) + L_i(x)f(x_i) + \dots + L_n(x)f(x_n) = \sum_{k=0}^n L_k(x)f(x_k).$$
(4.14)

It is clearly a polynomial of degree at most n and satisfy the conditions (4.11) since

$$p_n(x_i) = L_0(x_i)f(x_0) + L_1(x_i)f(x_1) + \dots + L_{i-1}(x_i)f(x_{i-1}) + L_i(x_i)f(x_i) + \dots + L_n(x_i)f(x_n),$$

which implies that

$$p_n(x_i) = f(x_i).$$

It remains to be shown how the polynomial $L_i(x)$ can be constructed so that they satisfy (4.12). If $L_i(x)$ is to satisfies (4.12), then it must contain a factor

$$(x-x_0)(x-x_1)\cdots(x-x_{i-1})(x-x_{i+1})\cdots(x-x_n).$$
 (4.15)

Since this expression has exactly n terms and $L_i(x)$ is a polynomial of degree n, we can deduce that

$$L_i(x) = A_i(x - x_0)(x - x_1) \cdots (x - x_{i-1})(x - x_{i+1}) \cdots (x - x_n), \tag{4.16}$$

for some multiplicative constant A_i . Let $x = x_i$, then the value of A_i is chosen so that

$$A_i = \frac{1}{(x_i - x_0)(x_i - x_1)\cdots(x_i - x_{i-1})(x_i - x_{i+1})\cdots(x_i - x_n)},$$
(4.17)

where none of the terms in denominator can be zero from the assumption of distinct points. Hence

$$L_i(x) = \prod_{k=0}^n \left(\frac{x - x_k}{x_i - x_k}\right), \quad i \neq k.$$

$$(4.18)$$

The interpolating polynomial can now be readily evaluated by substituting (4.18) into (4.14) to give

$$f(x) = p_n(x) = \sum_{i=0}^n \prod_{k=0}^n \left(\frac{x - x_k}{x_i - x_k}\right) f(x_i), \qquad i \neq k.$$
 (4.19)

This formula is called the Lagrange interpolation formula of degree n and the terms in (4.18) are called the Lagrange coefficient polynomials.

Uniqueness of Lagrange Interpolating Polynomial

To show the uniqueness of the interpolating polynomial $p_n(x)$, we suppose that in addition to the polynomial $p_n(x)$ the interpolation problem has another solution $q_n(x)$ of degree $\leq n$ whose graph passes through (x_i, y_i) , $i = 0, 1, \ldots, n$. Then define

$$r_n(x) = p_n(x) - q_n(x),$$

of the degree not greater than n. Since

$$r_n(x_i) = p_n(x_i) - q_n(x_i) = f(x_i) - f(x_i) = 0,$$

the polynomial $r_n(x)$ vanishes at n+1 point. But by using the following well known result from the theory of equations: "If a polynomial of degree n vanishes at n+1 distinct points, then the polynomial is identically zero". Hence $r_n(x)$ vanishes identically, or equivalently, $p_n(x) = q_n(x)$.

Example 4.1 Let f(x) = 0 be defined on the three numbers -h, 0, h, where $h \neq 0$. Use Lagrange interpolating polynomial to construct the polynomial p(x) which interpolate f(x) at the given numbers. Then show that this polynomial can be written in the following form

$$f(x) = p(x) = \frac{1}{2h^2} [f(-h) - 2f(0) + f(h)]x^2 + \frac{1}{2h} [f(h) - f(-h)]x + f(0).$$

Solution. Given three distinct points $x_0 = -h$, $x_1 = 0$ and $x_2 = h$ and using the quadratic Lagrange interpolating polynomial as follows:

$$f(x) = p_2(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} f(x_0) + \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} f(x_1) + \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} f(x_2),$$

at these data points, we get

$$f(x) = p_2(x) = \frac{(x^2 - xh)}{(-h)(-2h)}f(-h) + \frac{(x^2 - h^2)}{(h)(-h)}f(0) + \frac{(x^2 + xh)}{(2h)(h)}f(h).$$

Separating the coefficients of x^2 , x and constant term, we get

$$f(x) = p_2(x) = \left(\frac{f(-h)}{2h^2} + \frac{f(0)}{-h^2} + \frac{f(h)}{2h^2}\right)x^2 + \left(\frac{-hf(-h)}{2h^2} + \frac{hf(h)}{2h^2}\right)x + \left(\frac{-h^2f(0)}{-h^2}\right)x$$

after simplifying, we obtain

$$f(x) = p(x) = \frac{1}{2h^2} [f(-h) - 2f(0) + f(h)]x^2 + \frac{1}{2h} [f(h) - f(-h)]x + f(0).$$

Example 4.2 Let $p_2(x)$ be the quadratic Lagrange interpolating polynomial for the data: $(1,2), (2,3), (3,\alpha)$. Find α if the constant term in $p_2(x)$ is 5. Find the approximation of f(2.5).

Solution. Consider the quadratic Lagrange interpolating polynomial as follows:

$$f(x) = p_2(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2),$$

using the given data points, we get

$$f(x) = p_2(x) = L_0(x)(2) + L_1(x)(3) + L_2(x)(\alpha),$$

where the Lagrange coefficients can be calculate as follows:

$$L_0(x) = \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} = \frac{(x-2)(x-3)}{(1-2)(1-3)} = \frac{1}{2}(x^2-5x+6),$$

$$L_1(x) = \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} = \frac{(x-1)(x-3)}{(2-1)(2-3)} = -(x^2-4x+3),$$

$$L_2(x) = \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)} = \frac{(x-1)(x-2)}{(3-1)(3-2)} = \frac{1}{2}(x^2-3x+2).$$

Thus

$$f(x) = p_2(x) = \frac{1}{2}(x^2 - 5x + 6)(2) - (x^2 - 4x + 3)(3) + \frac{1}{2}(x^2 - 3x + 2)(\alpha).$$

Separating the coefficients of x^2 , x and constant term, we get

$$f(x) = p_2(x) = \left(-2 + \frac{\alpha}{2}\right)x^2 + \left(7 - \frac{3\alpha}{2}\right)x + (-3 + \alpha).$$

Since the given value of the constant term is 5, using this, we get

$$(-3+\alpha)=5$$
, gives $\alpha=8$.

Thus by using $\alpha = 8$ and x = 2.5, we have

$$f(2.5) \approx p_2(2.5) = 12.50 - 12.50 + 5 = 5$$

the required approximation of the function.

Example 4.3 Let $f(x) = x + \frac{1}{x}$, with points $x_0 = 1, x_1 = 1.5, x_2 = 2.5$ and $x_3 = 3$. Find the quadratic Lagrange polynomial for the approximation of f(2.7). Also, find the relative error.

Solution. Consider the quadratic Lagrange interpolating polynomial as follows:

$$f(x) = p_2(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2).$$

Since the given interpolating point is x = 2.7, therefore, the best three points for the quadratic polynomial should be as follows:

$$x_0 = 1.5, \quad x_1 = 2.5, \quad x_2 = 3,$$

and the function values at these points are

$$f(x_0) = 2.1667$$
, $f(x_1) = 2.9$, $f(x_2) = 3.3333$.

So using these values, we have

$$f(x) = p_2(x) = 2.1667L_0(x) + 2.9L_1(x) + 3.3333L_2(x),$$

where

$$L_0(x) = \frac{(x-2.5)(x-3)}{(1.5-2.5)(1.5-3)} = \frac{1}{1.5}(x^2-5.5x+7.5),$$

$$L_1(x) = \frac{(x-1.5)(x-3)}{(2.5-1.5)(2.5-3)} = \frac{1}{-0.5}(x^2-4.5x+4.5),$$

$$L_2(x) = \frac{(x-1.5)(x-2.5)}{(3-1.5)(3-2.5)} = \frac{1}{0.75}(x^2-4x+3.75).$$

Using these Lagrange coefficients in the polynomial and after simplifying, gives

$$f(x) = p_2(x) = 0.0889x^2 + 0.3776x + 1.4003,$$

which is the required quadratic polynomial. At x = 2.7, we have

$$f(2.7) \approx p_2(2.7) = 3.0679.$$

The relative error is

$$\frac{|f(2.7) - p_2(2.7)|}{|f(2.7)|} = \frac{|3.0704 - 3.0679|}{|3.0704|} = 0.0008.$$

Note that the sum of the Lagrange coefficients is equal to 1 as it should be

$$L_0(2.7) + L_1(2.7) + L_2(2.7) = -0.0400 + 0.7200 + 0.3200 = 1.$$

Using MATLAB command the above results can be reproduce as follows:

>>
$$x = [1.5 \ 2.5 \ 3];$$

>> $y = x + 1/x; \ x0 = 2.7;$
>> $sol = lint(x, y, x0);$

Example 4.4 Using the quadratic Lagrange interpolation formula to find the numbers A, B and C such that $p_2(1.4) = Af(0) + Bf(1) + Cf(2)$. If f(0) = 1, f(1) = 2 and f(2) = 3, then find the approximation of f(1.4).

Solution. Consider the quadratic Lagrange interpolating polynomial as follows:

$$f(x) = p_2(x) = Af(x_0) + Bf(x_1) + Cf(x_2).$$

Using the given values as $x_0 = 0, x_1 = 1, x_2 = 2$ and the interpolating point x = 1.4, we obtain

$$A = L_0(1.4) = \frac{(1.4-1)(1.4-2)}{(0-1)(0-2)} = -0.12,$$

$$B = L_1(1.4) = \frac{(1.4-0)(1.4-2)}{(1-0)(1-2)} = 0.84,$$

$$C = L_2(1.4) = \frac{(1.4-0)(1.4-1)}{(2-0)(2-1)} = 0.28.$$

Thus

$$f(1.4) \approx p_2(1.4) = (-0.12)(1) + (0.84)(2) + (0.28)(3) = 2.4,$$

the required approximation of the function at the point x = 1.4.

Example 4.5 Consider the following table:

$$\begin{array}{c|ccccc} x & 0 & 3 & 7 \\ \hline f(x) & 2 & 4 & 19 \end{array}$$

- (a) Construct the quadratic Lagrange polynomial $p_2(x) = ax^2 + bx + c$ to approximate f(x).
- (b) Use the polynomial in part (a) to interpolate the function f(x) at x = 4.

Solution. (a) Obviously, a quadratic polynomial can be determined so that it passes through the three points. Consider the quadratic Lagrange interpolating polynomial as follows:

$$f(x) = p_2(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2), \tag{4.20}$$

or

$$f(x) = p_2(x) = 2L_0(x) + 4L_1(x) + 19L_2(x). (4.21)$$

The Lagrange coefficients can be calculate as follows:

$$L_0(x) = \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} = \frac{1}{21}(x^2-10x+21),$$

$$L_1(x) = \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} = -\frac{1}{12}(x^2-7x),$$

$$L_2(x) = \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)} = \frac{1}{28}(x^2-3x).$$

Putting these values of the Lagrange coefficients in (4.21), we have

$$f(x) = p_2(x) = \frac{1}{84}(37x^2 - 55x + 168),$$

(with a = 37/84, b = -55/84, c = 2) which is the required quadratic Lagrange polynomial. (b) Now take x = 4 in the above polynomial, we obtain

$$f(4) \approx p_2(4) = \frac{1}{84} [37(4)^2 - 55(4) + 168] = 6.4286,$$

which is the required estimate value of f(4).

By using the following MATLAB commands we can easily find the value of the polynomial at the given point as

$$>> CP = [37/84 - 55/84 \ 168/84];$$

 $>> Sol = polyval(CP, 4);$

Note that the sum of the Lagrange coefficients is equal to 1 as it should be

$$L_0(4) + L_1(4) + L_2(4) = -\frac{1}{7} + 1 + \frac{1}{7} = 1.$$

Using MATLAB command the above results can be reproduce as follows:

$$>> x = [0 \ 3 \ 7]; \ y = [2 \ 4 \ 19];$$

 $>> x0 = 4; \ sol = lint(x, y, x0);$

Program 4.1

MATLAB m-file for the Lagrange Interpolation Method function fi=lint(x,y,x0)

dxi=x0-x; m=length(x); L=zeros(size(y));

L(1) = prod(dxi(2:m))/prod(x(1) - x(2:m));

L(m) = prod(dxi(1:m-1))/prod(x(m) - x(1:m-1));

for j=2:m-1

num = prod(dxi(1:j-1)) * prod(dxi(j+1:m));

dem = prod(x(j) - x(1:j-1)) * prod(x(j) - x(j+1:m));

L(j)=num/dem; end; fi = sum(y.*L);

Example 4.6 (a) Construct the table for $(\alpha, M(\alpha))$ by evaluating the integral

$$M(\alpha) = \int_0^1 (\alpha - e^x) dx,$$

at $\alpha = 1, 3, 5, 6$.

(b) Use the constructed table to find the best approximation of M(4) by using quadratic Lagrange polynomial. Compute the absolute error.

Solution. (a) Since

$$M(\alpha) = \int_0^1 (\alpha - e^x) dx = (\alpha x - e^x) \Big|_0^1 = \alpha - e + 1,$$

so by using the given values of α , we get

$$M(1) = -0.7183$$
, $M(3) = 1.2817$, $M(5) = 3.2817$, $M(6) = 4.2817$.

Thus we have the following table

(b) Since a quadratic polynomial can be determined so that it passes through the three points, let us consider the best form of the constructed table for the quadratic Lagrange interpolating polynomial to approximate M(4) as

$$\begin{array}{c|cccc} \alpha & 3.00 & 5.00 & 6.00 \\ \hline M(\alpha) & 1.2817 & 3.2817 & 4.2817 \end{array}$$

So using the quadratic Lagrange interpolating polynomial

$$M(\alpha) = p_2(\alpha) = L_0(\alpha)f(\alpha_0) + L_1(\alpha)f(\alpha_1) + L_2(\alpha)f(\alpha_2), \tag{4.22}$$

to get the approximation of M(4), we have

$$M(4) \approx p_2(4) = 1.2817L_0(4) + 3.2817L_1(4) + 4.2817L_2(4).$$
 (4.23)

The Lagrange coefficients can be calculate as follows:

$$L_0(4) = \frac{(4-5)(4-6)}{(3-5)(3-6)} = \frac{1}{3},$$

$$L_1(4) = \frac{(4-3)(4-6)}{(5-3)(5-6)} = 1,$$

$$L_2(4) = \frac{(4-3)(4-5)}{(6-3)(6-5)} = -\frac{1}{3}.$$

Putting these values of the Lagrange coefficients in (4.23), we obtain

$$M(4) \approx p_2(4) = \frac{1}{3}(1.2817) + 1(3.2817) - \frac{1}{3}(4.2817) = 2.2817,$$

which is the required approximation of M(4) by the quadratic interpolating polynomial. From the given integral we can obtained the exact value as follows

$$M(4) = \int_0^1 (4 - e^x) dx = (4x - e^x) \Big|_0^1 = 5 - e = 2.2817,$$

so

$$|M(4) - p_2(4)| = |2.2817 - 2.2817| = 0.0000,$$

is the required absolute error in our approximation.

Example 4.7 Consider the following table:

$$\begin{array}{c|ccccc} x & 1 & 2 & 3 \\ \hline f(x) & 2 & 3 & 5 \\ \end{array}$$

If $f(x) = p_2(x) = a_0 + a_1x + a_2x^2$, then find the approximation of f(1.5).

Solution. Using values of f(x) at x = 1, 2 and 3, we get

$$2 = a_0 + a_1 + a_2,$$

$$3 = a_0 + 2a_1 + 4a_2,$$

$$4 = a_0 + 3a_1 + 9a_2.$$

Now solving this linear system for a_0, a_1, a_2 , we obtain

$$a_0 = 2$$
, $a_1 = -0.5$, $a_2 = 0.5$.

Thus

$$f(x) = p_2(x) = 2 - 0.5x + 0.5x^2,$$

and at x = 1.5, we get

$$f(1.5) \approx p_2(1.5) = 2 - 0.75 + 1.125 = 2.375,$$

the required approximation of f(1.5).

Example 4.8 Find the missing term in the following table using best Lagrange polynomial:

Solution. Let $x_0 = 0, x_1 = 1$ and $x_2 = 6$, then obviously, a quadratic polynomial can be determined so that it passes through these three points. Consider the quadratic Lagrange interpolating polynomial as follows:

$$f(x) = p_2(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2), \tag{4.24}$$

and using function values, gives

$$f(x) = p_2(x) = L_0(x) + 2L_1(x) + 3L_2(x). (4.25)$$

The Lagrange coefficients can be calculate as follows:

$$L_0(x) = \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} = \frac{(x-1)(x-6)}{(0-1)(0-6)} = \frac{1}{6}(x^2-7x+6),$$

$$L_1(x) = \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} = \frac{(x-0)(x-6)}{(1-0)(1-6)} = -\frac{1}{5}(x^2-6x),$$

$$L_2(x) = \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)} = \frac{(x-0)(x-1)}{(6-0)(6-1)} = \frac{1}{30}(x^2-x).$$

Putting these values of the Lagrange coefficients in (4.25), we have

$$f(x) = p_2(x) = \frac{1}{6}(x^2 - 7x + 6)(1) - \frac{1}{5}(x^2 - 6x)(2) + \frac{1}{30}(x^2 - x)(3),$$

or

$$f(x) = p_2(x) = \frac{1}{15}(-2x^2 + 17x + 15),$$

which is the required quadratic interpolating polynomial.

Now to find the missing term when x = 2 by using the above polynomial, we have

$$f(2) \approx p_2(2) = \frac{1}{15} \left[-2(2)^2 + 17(2) + 15 \right] = \frac{41}{15} = 2.7333,$$

which is the missing term in the table. We can also find the missing number by considering a quadratic polynomial of the form

$$f(x) = p_2(x) = a + bx + cx^2$$
.

If the curve passes through the points x = 0, 1, 6, then we have

$$1 = a$$

$$2 = a+b+c$$

$$3 = a+6b+36c$$

Solving this system for a, b and c, we obtain

$$a = 1, \quad b = \frac{17}{15}, \quad c = -\frac{2}{15}.$$

Thus

$$f(x) = p_2(x) = 1 + \frac{17}{15}x - \frac{2}{15}x^2$$

and at x = 2, it gives

$$f(2) \approx p_2(2) = 1 + \frac{17}{15}(2) - \frac{2}{15}(2)^2 = 2.7333,$$

the required value of the missing term.

Example 4.9 The equation $x - 9^{-x} = 0$ has a solution in [0, 1]. Compute the Lagrange polynomial on $x_0 = 0, x_1 = 0.5$ and $x_2 = 1$. By setting the interpolating polynomial equal to zero and solving the equation, find an approximate solution to the equation in the given interval [0, 1].

Solution. Let us consider the form of the constructed table for the given function $f(x) = x - 9^{-x}$ at the given points as

$$\begin{array}{c|ccccc} x & 0 & 0.5 & 1 \\ \hline f(x) & -1 & 1/6 & 8/9 \end{array}$$

So using the quadratic Lagrange interpolating polynomial

$$f(x) = p_2(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2) = -L_0(x) + \frac{1}{6}L_1(x) + \frac{8}{9}L_2(x), \quad (4.26)$$

where the values of the Lagrange coefficients can be calculate as follows:

$$L_0(x) = \frac{(x-0.5)(x-1)}{(0-0.5)(0-1)} = 2x^2 - 3x + 1,$$

$$L_1(x) = \frac{(x-0)(x-1)}{(0.5-0)(0.5-1)} = -4x^2 + 4x,$$

$$L_2(x) = \frac{(x-0)(x-0.5)}{(1-0)(1-0.5)} = 2x^2 - x.$$

Putting these values of the Lagrange coefficients in (4.26), we have

$$f(x) = p_2(x) = \frac{1}{18}(-16x^2 + 50x - 18),$$

which is the required quadratic interpolating polynomial. Now setting this polynomial equal to zero, we get

$$0 = p_2(x) = \frac{1}{18}(-16x^2 + 50x - 18),$$

which gives

$$-16x^2 + 50x - 18 = 0$$
, or $8x^2 - 25x + 9 = 0$.

Now solving this quadratic equation, one can get

$$x_1 = 2.70985$$
 and $x_2 = 0.41515$.

Thus the approximate solution to the given equation in [0,1] is $x_2 = 0.41515$.

Note that we can use Lagrange interpolating polynomial to express the given rational function $R(x) = \frac{f(x)}{g(x)}$ as sums of partial fractions. For this we have to construct the table for the function f(x) at x values (the zeros of the denominator function g(x)).

Example 4.10 Use quadratic Lagrange polynomial, express the rational function $\frac{2x^2 - 9x - 9}{x^3 - 9x}$ as sums of partial fractions.

Solution. Given rational function is of the form

$$\frac{f(x)}{g(x)} = \frac{2x^2 - 9x - 9}{x^3 - 9x} = \frac{2x^2 - 9x - 9}{(x+3)(x-0)(x-3)},$$

and it gives the zeros of the function g(x) as x = -3, 0, 3. So the table for the function f(x) at x values is

$$\begin{array}{c|cccc} x & -3 & 0 & 3 \\ \hline f(x) & 36 & -9 & -18 \end{array}$$

The Lagrange coefficients can be calculate as follows:

$$L_0(x) = \frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} = \frac{(x-0)(x-3)}{(-3-0)(-3-3)} = \frac{1}{18}(x-0)(x-3),$$

$$L_1(x) = \frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} = \frac{(x+3)(x-3)}{(0+3)(0-3)} = -\frac{1}{9}(x+3)(x-3),$$

$$L_2(x) = \frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)} = \frac{(x+3)(x-0)}{(3+3)(3-0)} = \frac{1}{18}(x+3)(x-0).$$

By using the quadratic Lagrange interpolating formula

$$f(x) = p_2(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2) = 36L_0(x) - 9L_1(x) - 18L_2(x_2),$$

the polynomial f(x) is given by

$$f(x) = 2(x-0)(x-3) + (x+3)(x-3) - (x+3)(x-0).$$

Thus

$$\frac{2x^2 - 9x - 9}{x^3 - 9x} = \frac{f(x)}{(x+3)(x-0)(x-3)} = \frac{2}{(x+3)} + \frac{1}{(x-0)} - \frac{1}{(x-3)},$$

is the sums of partial fractions of the given rational function.

We can easily check this results by using MATLAB command residue. The residue function takes the ratio of polynomials $\frac{f(x)}{g(x)}$ and return the corresponding partial fraction representation of it. It takes as input the coefficients of the polynomials f(x) and g(x) (in descending powers of x) and returns the coefficients and the pole values (the zeros of the denominator). If we want to find the partial fraction expansion $R(x) = \frac{2x^2 - 9x - 9}{x^3 - 9x}$, the following MATLAB program gives us the coefficients in the expansion. We write the coefficients of the numerator and the denominator in separate vectors and MATLAB gives us the coefficients (c) with the corresponding poles (z) in the expansion. The following MATLAB command lines give the use of the residue function:

$$>> num = [0\ 2\ -9\ -9];\ den = [1\ 0\ -9\ 0];$$

 $>> [c, z] = residue(num, den);$

Error Formula of Lagrange Polynomial

As with any numerical technique, it is important to obtain bounds for the errors involved. Now we discuss the error term when the Lagrange polynomial is used to approximate continuous function f(x). It is not possible, in general, to say how accurately the interpolating polynomial p_n approximates given function f. All can be said with certainty is that $f(x) - p_n(x) = 0$ at $x = x_0, x_1, \ldots, x_n$. However, it is sometimes possible to obtain a bound on the error $f(x) - p_n(x)$ at an intermediate point x using the following theorem.

Error Formulas of Linear, Quadratic and Cubic Lagrange Polynomials

If f(x) has second, third and fourth derivatives on interval I and if it is approximated by the polynomials $p_1(x), p_2(x), p_3(x)$ passing respectively, through 2, 3, 4 data points on I, then the errors E_1, E_2, E_3 are given by

$$E_1 = f(x) - p_1(x) = \frac{f''(\eta(x))}{2!}(x - x_0)(x - x_1), \quad \eta(x) \in I,$$
(4.27)

where $p_1(x)$ is the linear Lagrange polynomial (4.5) and a unknown point $\eta(x) \in (x_0, x_1)$.

$$E_2 = f(x) - p_2(x) = \frac{f'''(\eta(x))}{3!}(x - x_0)(x - x_1)(x - x_2), \quad \eta(x) \in I, \tag{4.28}$$

where $p_2(x)$ is the quadratic Lagrange polynomial (4.7) and a unknown point $\eta(x) \in (x_0, x_2)$.

$$E_3 = f(x) - p_3(x) = \frac{f(4)(\eta(x))}{4!}(x - x_0)(x - x_1)(x - x_2)(x - x_3), \quad \eta(x) \in I, \tag{4.29}$$

where $p_3(x)$ is the cubic Lagrange polynomial (4.9) and a unknown point $\eta(x) \in (x_0, x_3)$.

Continuing in the similar manner, in the following theorem we define the error formula for the nth degree Lagrange polynomial (4.14).

Theorem 4.2 (Error Formula of nth Degree Lagrange Polynomial)

If f(x) has (n+1) derivatives on interval I and if it is approximated by a polynomial $p_n(x)$ passing through (n+1) data points on I, then the error E_n is given by

$$E_n = f(x) - p_n(x) = \frac{f^{(n+1)}(\eta(x))}{(n+1)!} (x - x_0)(x - x_1) \cdots (x - x_n), \quad \eta(x) \in I, \tag{4.30}$$

where $p_n(x)$ is Lagrange interpolating polynomial (4.14) and a unknown point $\eta(x) \in (x_0, x_n)$.

The error formula (4.30) is an important theoretical results because the Lagrange polynomials are used extensively for deriving numerical differentiation and integration methods. Error bounds for these techniques are obtained from Lagrange error formula.

Example 4.11 Let $f(x) = \sqrt{x - x^2}$ and $p_2(x)$ be the quadratic Lagrange interpolating polynomial on $x_0 = 0, x_1$ and $x_2 = 1$. Find the largest value of x_1 in (0,1) for which

$$f(0.5) - p_2(0.5) = -0.25.$$

Solution. Consider the quadratic Lagrange interpolating polynomial as follows:

$$f(x) = p_2(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2).$$

At the given values of $x_0 = 0$, $x_1 = x_1$, $x_2 = 1$, we have, f(0) = 0, $f(x_1) = \sqrt{x_1 - x_1^2}$ and f(1) = 0, gives

$$f(x) = p_2(x) = L_0(x)(0) + L_1(x)(f(x_1\sqrt{x_1 - x_1^2})) + L_2(x)(0),$$

where

$$L_1(x) = \frac{(x-0)(x-1)}{(x_1-0)(x_1-1)} = \frac{x^2-x}{x_1^2-x_1}.$$

Thus

$$f(x) = p_2(x) = \frac{x^2 - x}{x_1^2 - x_1} \sqrt{x_1 - x_1^2}$$
 and $p_2(0.5) = \frac{-0.25}{x_1^2 - x_1} \sqrt{x_1 - x_1^2}$.

Given

$$f(0.5) - p_2(0.5) = -0.25$$
, gives $p_2(0.5) = f(0.5) + 0.25 = 0.5 + 0.25 = 0.75$,

so

$$-0.25 \frac{\sqrt{x_1 - x_1^2}}{(x_1 - x_1^2)} = 0.75, \quad or \quad \sqrt{x_1 - x_1^2} = -3(x_1 - x_1^2).$$

Thus, taking square on both sides, we get

$$x_1 - x_1^2 = 9(x_1 - x_1^2)^2$$
, or $(x_1 - x_1^2)[1 - 9(x_1 - x_1^2)] = 0$,

which can be also written as

$$x_1(1-x_1)[9x_1^2-9x_1+1]=0.$$

Solving this equation for x_1 , we get

$$x_1 = 0$$
, or $x_1 = 1$, or $x_1 = 0.127322$, or $x_1 = 0.872678$.

Thus $x_1 = 0.872678$, the required largest value in the given interval (0,1).

Example 4.12 Show that a bound for the error in the linear interpolation is

$$|f(x) - p_1(x)| \le \frac{h^2}{8}M$$
, where $M = \max_{x_0 \le x \le x_1} |f''(x)|$ and $h = x_1 - x_0$. (4.31)

Solution. Consider two points x_0 and x_1 , then the linear polynomial $p_1(x)$ interpolating f(x) at these points is

$$f(x) = p_1(x) = \frac{(x - x_1)}{(x_0 - x_1)} f(x_0) + \frac{(x - x_0)}{(x_1 - x_0)} f(x_1).$$

By using the given data point, the error formula (4.30) becomes

$$f(x) - p_1(x) = \frac{(x - x_0)(x - x_1)}{2!} f''(\eta(x)),$$

where $\eta(x)$ is a unknown point between x_0 and x_1 . Hence

$$|f(x) - p_1(x)| = \left| \frac{(x - x_0)(x - x_1)}{2!} \right| |f''(\eta(x))|.$$

The value of $f''(\eta(x))$ can not be computed exactly because $\eta(x)$ is not known. But we can bound the error by computing the largest possible value for $|f''(\eta(x))|$. So bound |f''(x)| on $[x_0, x_1]$ can be obtain

$$M = \max_{x_0 \le x \le x_1} |f''(x)|.$$

and so for $|f''(\eta(x))| \leq M$, we have

$$|f(x) - p_1(x)| \le \frac{M}{2} |(x - x_0)(x - x_1)|.$$

Since the maximum of function $g(x) = (x - x_0)(x - x_1)$ in $[x_0, x_1]$ occurs at the critical point $x = \frac{(x_0 + x_1)}{2}$ (g'(x) = 0) and so that maximum is $|(x - x_0)(x - x_1)| = \frac{(x_1 - x_0)^2}{4}$.

This follows easily by noting that the function $(x-x_0)(x-x_1)$ is a quadratic and has two roots x_0 and x_1 , hence its maximum value occurs midway between these roots. Thus, for any $x \in [x_0, x_1]$, we have

$$|f(x) - p_1(x)| \le \frac{(x_1 - x_0)^2}{8}M, \quad or \quad |f(x) - p_1(x)| \le \frac{h^2}{8}M,$$

where $h = x_1 - x_0$.

Example 4.13 Find the linear Lagrange polynomial that passes through the points (0, f(0)) and $(\pi, f(\pi))$ and then use it to approximate the function $f(x) = 2 \cos x$ at $\frac{\pi}{2}$. Find absolute error and an bound for the error in the linear interpolation of f(x).

Solution. Given two points $x_0 = 0$ and $x_1 = \pi$, then the linear Lagrange polynomial $p_1(x)$

$$f(x) = p_1(x) = \frac{(x - x_1)}{(x_0 - x_1)} f(x_0) + \frac{(x - x_0)}{(x_1 - x_0)} f(x_1),$$

interpolating f(x) at these points is

$$f(x) = p_1(x) = \frac{(x-\pi)}{(0-\pi)}f(0) + \frac{(x-0)}{(\pi-0)}f(\pi).$$

By using the function values at the given data point, we get

$$f(x) = p_1(x) = \frac{(x-\pi)}{(0-\pi)}(2) + \frac{(x-0)}{(\pi-0)}(-2) = 2 - \frac{4x}{\pi}$$
 and $f(\pi/2) \approx p_1(\pi/2) = 0$.

Thus absolute error, $|2\cos(\pi/2) - p_1(\pi/2)| = 0$. Since $M = \max_{0 \le x \le \pi} |f''(x)| = \max_{0 \le x \le \pi} |-2\cos x| = 2$ and $h = \pi$, so by using the linear Lagrange error formula (4.31), we get

$$|f(x) - p_1(x)| \le \frac{(\pi - 0)^2}{4} = \frac{\pi^2}{4},$$

which is the required bound of error in the linear interpolation of f(x).

Example 4.14 Use the quadratic Lagrange interpolating polynomial by selecting the best three points from $\{-2,0,1,2,2.5\}$ on the function defined by $f(x) = (x+1)^{1/3}$ to estimate the cube root of $\frac{3}{2}$ and compute an error bound and absolute error.

Solution. Since the given function is a cube root of (x+1), so by taking $x+1=\frac{3}{2}$, we have $x=\frac{1}{2}$, therefore, the best points for the quadratic polynomial are $x_0=0, x_1=1$, and $x_2=2$. Consider a quadratic Lagrange interpolating polynomial as

$$f(x) = p_2(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2), \tag{4.32}$$

and at x = 0.5, gives

$$f(0.5) \approx p_2(0.5) = (1)^{1/3} L_0(0.5) + (2)^{1/3} L_1(0.5) + (3)^{1/3} L_2(0.5).$$
 (4.33)

The Lagrange coefficients can be calculate as follows:

$$L_0(0.5) = \frac{(0.5-1)(0.5-2)}{(0-1)(0-2)} = 0.375,$$

$$L_1(0.5) = \frac{(0.5-0)(0.5-2)}{(1-0)(1-2)} = 0.75,$$

$$L_2(0.5) = \frac{(0.5-0)(0.5-1)}{(2-0)(2-1)} = -0.125.$$

Putting these values of the Lagrange coefficients in (4.33), we have

$$f(0.5) \approx p_2(0.5) = (1)^{1/3}(0.375) + (2)^{1/3}(0.75) + (3)^{1/3}(-0.125) = 1.1396,$$

which is the required approximation of the $\left(\frac{3}{2}\right)^{1/3}$.

To compute an error bound for the approximation of the given function in the interval [0,2], we use the following quadratic error formula

$$|f(x) - p_2(x)| = \frac{|f^{(3)}(\eta(x))|}{3!} |(x - x_0)(x - x_1)(x - x_2)|.$$

As

$$|f^{(3)}(\eta(x))| \le M = \max_{0 \le x \le 2} |f^{(3)}(x)|,$$

and the first three derivatives are

$$f'(x) = \frac{1}{3}(x+1)^{-2/3}, \quad f''(x) = -\frac{2}{9}(x+1)^{-5/3}, \quad f^{(3)}(x) = \frac{10}{27}(x+1)^{-8/3},$$

so

$$M = \max_{0 \le x \le 2} \left| \frac{10}{27} (x+1)^{-8/3} \right| = \frac{10}{27}.$$

Hence

$$|f(0.5) - p_2(0.5)| \le \frac{10/27}{6} |(0.5 - 0)(0.5 - 1)(0.5 - 2)|,$$

and it gives

$$|f(0.5) - p_2(0.5)| \le \frac{10(0.375)}{162} = 0.0232,$$

which is desired error bound. Also, we have the absolute error is given as

$$|f(0.5) - p_2(0.5)| = |(1.5)^{1/3} - 1.1396| = |1.1447 - 1.1396| = 0.0051,$$

which is desired absolute error.

Example 4.15 Consider $f(x) = \sin x$ and its values are known at five points $\{0, 0.2, 0.4, 0.6, 0.8\}$. If the approximation of $\sin 0.28$ by four degree Lagrange interpolating polynomial is 0.2763591, then compute the error bound and the absolute error for the approximation.

Solution. To compute an error bound for the approximation of the given function in the interval [0,0.8], we use the following error formula for Lagrange polynomial degree four

$$|f(x) - p_4(x)| = \frac{|f^{(5)}(\eta(x))|}{5!} |(x - x_0)(x - x_1)(x - x_2)(x - x_3)(x - x_4)|,$$

or

$$|f(x) - p_4(x)| \le \frac{M}{5!} |(x - x_0)(x - x_1)(x - x_2)(x - x_3)(x - x_4)|.$$

Since

$$|f^{(5)}(\eta(x))| \le M = \max_{0 \le x \le 0.8} |f^{(5)}(x)| = \max_{0 \le x \le 0.8} |\cos x| = 1,$$

so

$$|f(0.28) - p_4(0.28)| \le \frac{1}{120} |0.28(0.28 - 0.2)(0.28 - 0.4)(0.28 - 0.6)(0.28 - 0.6)(0.28 - 0.8)| \le 3.7 \times 10^{-6},$$

which is desired error bound. Also, we have to compute absolute error as

$$|f(0.28) - p_4(0.28)| = |\sin 0.28 - p_4(0.28)| = |0.2763556 - 0.2763591| = 3.5 \times 10^{-6},$$

which is desired result.

Example 4.16 Let $p_2(x)$ be the Lagrange polynomial which interpolating $f(x) = x^3 + x + 1$ at the points $x_i = \alpha + (i+1)h$, i = 0, 1, 2, where α is constant and h > 0. Find h such that the error at $x = \alpha$ is bounded above by 10^{-3} .

Solution. To compute an error bound for the approximation of the given function, we use the following error formula for the quadratic Lagrange polynomial as

$$|E| = |f(x) - p_2(x)| = \frac{|f^{(3)}(\eta(x))|}{3!} |(x - x_0)(x - x_1)(x - x_2)|,$$

where

$$x = \alpha$$
, $x_0 = (\alpha + h)$, $x_1 = (\alpha + 2h)$, $x_2 = (\alpha + 3h)$.

Since

$$|f^{(3)}(\eta(x))| \le M = \max_{x_0 \le x \le x_2} |f^{(3)}(x)| = \max_{x_0 \le x \le x_2} |6| = 6,$$

so

$$|E| \le \frac{6}{6} |(\alpha - (\alpha + h))(\alpha - (\alpha + 2h))(\alpha - (\alpha + 3h))| = 6h^3.$$

Given

$$|E| < 10^{-3}, \quad so \quad 6h^3 < 10^{-3},$$

from this we have h < 0.055 and so we can take h = 0.05.

Example 4.17 Use the best Lagrange interpolating polynomial to find the approximation of f(1.5) if f(-2) = 2, f(-1) = 1.5, f(1) = 3.5 and f(2) = 5. Estimate the error bound if the maximum value of $|f^{(4)}(x)|$ is 0.025 in the interval [-2, 2].

Solution. Since the given number of points are, $x_0 = -2$, $x_1 = -1$, $x_2 = 1$, $x_3 = 2$, therefore the best Lagrange interpolating polynomial to find the approximation of f(1.5) will be the cubic. The cubic Lagrange interpolating polynomial for the approximation of the given function is:

$$f(x) = p_3(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2) + L_3(x)f(x_3),$$

and taking f(-2) = 2, f(-1) = 1.5, f(1) = 3.5, f(2) = 5 and the interpolating point x = 1.5, we have

$$f(1.5) \approx p_3(1.5) = L_0(1.5)f(-2) + L_1(1.5)f(-1) + L_2(1.5)f(1) + L_3(1.5)f(2),$$

or

$$f(1.5) \approx p_3(1.5) = 2L_0(1.5) + 1.5L_1(1.5) + 3.5L_2(1.5) + 5L_3(1.5).$$

The Lagrange coefficients can be calculate as follows:

$$L_0(1.5) = \frac{(1.5+1)(1.5-1)(1.5-2)}{(-2+1)(-2-1)(-2-2)} = 0.0521,$$

$$L_1(1.5) = \frac{(1.5+2)(1.5-1)(1.5-2)}{(-1+2)(-1-1)(-1-2)} = -0.1458,$$

$$L_2(1.5) = \frac{(1.5+2)(1.5+1)(1.5-2)}{(1+2)(1+1)(1-2)} = 0.7292,$$

$$L_3(1.5) = \frac{(1.5+2)(1.5+1)(1.5-1)}{(2+2)(2+1)(2-1)} = 0.3646.$$

Putting these values of the Lagrange coefficients in the above equation, we get

$$f(1.5) \approx p_3(1.5) = 2(0.0521) + 1.5(-0.1458) + 3.5(0.7292) + 5(0.3646) = 4.2607,$$

which is the required cubic interpolating polynomial approximation of the function at the given point x = 1.5. Note that $L_0(1.5) + L_1(1.5) + L_2(1.5) + L_3(1.5) = 1$.

To compute an error bound for the approximation of the given function in the interval [-2,2], we use the following cubic error formula

$$|f(x) - p_3(x)| = \frac{|f^{(4)}(\eta(x))|}{4!} |(x - x_0)(x - x_1)(x - x_2)(x - x_3)|.$$

As

$$|f^{(4)}(\eta(x))| \le M = \max_{-2 \le x \le 2} |f^{(4)}(x)| = 0.025,$$

so

$$|f(1.5) - p_3(1.5)| \le \frac{M}{4!} |(1.5+2)(1.5+1)(1.5-1)(1.5-2)|,$$

and it gives

$$|f(1.5) - p_3(1.5)| \le \frac{(0.025)(2.1875)}{24} = 0.0023,$$

which is desired error bound.

Example 4.18 Consider the following table having the data for $f(x) = e^{3x} + \cos 2x$:

$$\begin{array}{c|ccccc} x & 0.1 & 0.2 & 0.4 & 0.5 \\ \hline f(x) & 2.3300 & 2.7432 & 4.0168 & 5.0220 \end{array}$$

Find the approximation of f(0.45) using the best quadratic Lagrange interpolation formula and also estimate an error bound and absolute error for the approximation.

Solution. Using the data points 0.2, 0.4, 0.5, the best Lagrange formula to find the interpolating polynomial to approximate the function is the quadratic polynomial

$$f(x) = p_2(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2),$$

which implies that

$$f(x) = p_2(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} f(x_0) + \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} f(x_1) + \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} f(x_2),$$

or

$$f(x) = p_2(x) = 45.72[x^2 - 0.9x + 0.2] - 200.84[x^2 - 0.7x + 0.1]$$
+ 167.4[x^2 - 0.6x + 0.08].

Thus

$$f(x) = p_2(x) = 12.28x^2 - 1.0x^2 + 2.452. (4.34)$$

Take x = 0.45 in the above polynomial (4.34), we have

$$f(0.45) \approx p_2(0.45) = 4.4887.$$

The exact value of f(0.45) = 4.4790, so, the absolute error is 0.0097. Now to compute an error bound of the approximation, we use the following formula

$$|f(x) - p_2(x)| = \frac{|f^{(3)}(\eta(x))|}{3!} |(x - x_0)(x - x_1)(x - x_2)|.$$
(4.35)

Taking the third derivative of the given function, we get

$$f'(x) = 3e^{3x} - 2\sin 2x, \qquad f''(x) = 9e^{3x} - 4\cos 2x,$$
$$f'''(x) = 27e^{3x} + 8\sin 2x.$$

Thus

$$|f^{(3)}(\eta(x))| = |27e^{3\eta(x)} + 8\sin 2(\eta(x))|, \quad for \quad \eta(x) \in (0.2, 0.5),$$

and it gives

$$|f^{(3)}(0.2)| = 52.3126$$
 and $|f^{(3)}(0.5)| = 127.7374$.

The value of $f^{(3)}(\eta(x))$ can not be computed exactly because $\eta(x)$ is not known. But we can bound the error by computing the largest possible value for $|f^{(3)}(\eta(x))|$. So bound $|f^{(3)}(x)|$ on [0.2, 0.5] can be obtain

$$M = \max_{0.2 \le x \le 0.5} |f^{(3)}(x)| = 127.7374,$$

and so for $|f^{(4)}(\eta(x))| \leq M$, we have (4.35) as follows

$$|f(x) - p_2(x)| \le (127.7374)(0.000625)/6 = 0.0133,$$

which is the required error bound for the approximation.

Example 4.19 Determining spacing h in a table of equally spaced value of the function $f(x) = e^x$ between smallest point a = 1 and largest point b = 2, so that interpolation with a second-degree polynomial in this table will yield a desired accuracy.

Solution. Suppose that the given table contains the function values $f(x_i)$, for the points $x_i = 1 + ih$, i = 0, 1, ..., n, where $n = \frac{(2-1)}{h}$. If $x \in [x_{i-1}, x_{i+1}]$, then we approximate the function f(x) by degree 2 polynomial $p_2(x)$ which interpolates f(x) at x_{i-1}, x_i, x_{i+1} . Then the error formula (4.30) for these data points becomes

$$|f(x) - p_2(x)| = \left| \frac{(x - x_{i-1})(x - x_i)(x - x_{i+1})}{3!} \right| \left| f'''(\eta(x)) \right|$$

where $\eta(x) \in (x_{i-1}, x_{i+1})$. Since the point $\eta(x)$ is unknown and so, we can not estimate $f'''(\eta(x))$, therefore, let

$$|f'''(\eta(x))| \le M = \max_{1 \le x \le 2} |f'''(x)|.$$

Then

$$|f(x) - p_2(x)| \le \frac{M}{6} |(x - x_{i-1})(x - x_i)(x - x_{i+1})|.$$

Since $f(x) = e^x$ and $f'''(x) = e^x$, therefore

$$|f'''(\eta(x))| \le M = e^2 = 7.3891.$$

Now to find the maximum value of $|(x-x_{i-1})(x-x_i)(x-x_{i+1})|$, take $t=(x-x_i)/h$, gives

$$\max_{x \in [x_{i-1}, x_{i+1}]} |(x - x_{i-1})(x - x_i)(x - x_{i+1})| = \max_{t \in [-h, h]} |(t - h)t(t + h)| = \max_{t \in [-h, h]} |t(t^2 - h^2)|,$$

using the linear change of variables $t = x - x_i$. As we see the function $H(t) = t^3 - th^2$ vanishes at t = -h and t = h, so the maximum value of |H(t)| on [-h, h] must occurs at one of the extreme of H(t), which can be found by solving the equation

$$H'(t) = 3t^2 - h^2 = 0$$
, gives, $t = \pm h/\sqrt{3}$.

Hence

$$\max_{x \in [x_{i-1}, x_{i+1}]} |(x - x_{i-1})(x - x_i)(x - x_{i+1})| = \frac{2h^3}{3\sqrt{3}}.$$

Thus, for any $x \in [1, 2]$, we have

$$|f(x) - p_2(x)| \le \frac{(2h^3/3\sqrt{3})e^2}{6} = \frac{h^3e^2}{9\sqrt{3}},$$

if $p_2(x)$ is chosen as the polynomial of degree 2 which interpolates $f(x) = e^x$ at the three tabular points nearest x. If we wish to obtain six decimal place accuracy this way, we would have to choose h so that

 $\frac{h^3 e^2}{9\sqrt{3}} < 5 \times 10^{-7}, \quad which \ implies \ that \quad h^3 < 10.5483 \times 10^{-7},$

and it gives h = 0.01.

If the derivatives of f can be uniformly bounded by a constant, then we can choose n appropriately large in order to force the error term to be as small as we want. So, we can find a very accurate interpolating polynomial. In particular, if we look at the case when we choose to use equally spaced points (but without pre-determining n), then bounding the derivatives allows us to choose n large enough (alternatively h small enough) to give an accurate polynomial interpolant. The form of the error bound Theorem 4.2 is not as useful as the following one for quick calculations.

Theorem 4.3 (Error Bounds for Lagrange Interpolation at Equally Spaced Points)

Assume that f(x) is defined on the interval [a,b], which contains equally spaced points $x_k = x_0 + hk$. Additionally, assume that f(x) and the derivatives of f(x) up to the order (n+1), are continuous and bounded on the special intervals $[x_0, x_1]$, $[x_0, x_2]$ and $[x_0, x_3]$, respectively; that is

$$|f^{(n+1)}(x)| \le M$$
 for $x_0 \le x \le x_n$,

for n = 1, 2, 3. Then error bounds for linear, quadratic and cubic polynomials are:

$$|E_1(x)| \leq \frac{h^2}{8}M \qquad for \qquad x_0 \leq x \leq x_1,$$

$$|E_2(x)| \leq \frac{h^3}{9\sqrt{3}}M \qquad for \qquad x_0 \leq x \leq x_2,$$

$$|E_3(x)| \leq \frac{h^4}{24}M \qquad for \qquad x_0 \leq x \leq x_3.$$

Continue in the similar manner for the interval $[x_0, x_n]$, for n = 1, 2, ..., n, we have

$$|E_n(x)| \le \frac{M}{4(n+1)} \left(\frac{b-a}{n}\right)^{n+1} = \frac{M}{4(n+1)} h^{n+1}, \quad \text{for} \quad x_0 \le x \le x_n,$$
 (4.36)

the general error bound formula.

Example 4.20 Find an error bound if $f(x) = \sin x$ is approximated by an interpolation polynomial with ten equally spaced data points in [0, 1.6875].

Solution. Given n = 9 and a = 0, b = 1.6875, then

$$M = \max_{0 \le x \le 1.6875} |f^{(10)}(x)| = \max_{0 \le x \le 1.6875} |-\sin x| \le 1, \quad \forall \ x \in [0, 1.6875].$$

Hence, the interpolation error (using Theorem 4.36) can be bounded by

$$|E_9(x)| = |\sin x - p_9(x)| \le \frac{1}{40} \left(\frac{1.6875}{9}\right)^{10} \approx 1.34 \times 10^{-9},$$

for all $x \in [0, 1.6875]$.

Note that $f^{(n)}(x) = \pm \sin x$ for even n and $f^{(n)}(x) = \pm \cos x$ for odd n, so we have a uniform bound on $f^{(n)}(x)$ for all n. That is $|f^{(n)}(x)| \le 1$ for all x and for all n. In such instance, we can force interpolation error to 0 by increasing the number of interpolation data points.

On the other hand, if the derivatives of all order for the function f are continuous but we cannot uniformly the derivatives, then increasing the interpolation data points may not result in smaller errors.

Example 4.21 Find an error bound if $f(x) = \frac{1}{x^2 + 1}$ is approximated by an interpolation polynomial with ten equally spaced data points in the interval [-5, 5].

Solution. Given n = 9 and a = -5, b = 5, then

$$M = \max_{-5 \le x \le 5} |f^{(10)}(x)|,$$

One can find the maximum value of $|f^{(10)}(x)|$ is a very large value. Note that the term

$$\frac{1}{40} \left(\frac{10}{9} \right)^{10} \approx 0.0717,$$

so this term will not be small enough to guarantee a reasonable bound on the error. Furthermore, if we allow n to get larger, then the magnitude of the derivatives of f also get very large very fast. So adding more interpolation points can increase the oscillation of the interpolating polynomial. \bullet

Example 4.22 Determine an appropriate step size h to construct a table of $f(x) = x + \ln(x+1)$ on the interval [2,3] when the error for linear Lagrange interpolation is to be bounded by 9×10^{-4} . Then use the constructed table to find the approximation of f(2.6) using linear Lagrange polynomial and compute the absolute error.

Solution. Since we know that error bound formula for the linear Lagrange polynomial is

$$|E| \le \frac{Mh^2}{8}, \quad h = x_1 - x_0,$$
 (4.37)

and it is given that

$$|E_1| \le 9 \times 10^{-4},$$

so

$$\frac{Mh^2}{8} \le 9 \times 10^{-4}$$
, where $M = \max_{2 \le x \le 3} |f''(x)|$.

Since

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

so

$$M = \max_{2 \le x \le 3} \left| \frac{-1}{(x+1)^2} \right| = \frac{1}{9}.$$

Thus

$$h^2 \le 648 \times 10^{-4} = 0.0648$$
, gives $h = 0.25$.

Hence we obtain the required table

for the given function f(x) on [2,3] with step size h = 0.25.

As x = 2.6, so take two points $x_0 = 2.5$ and $x_1 = 2.75$, then linear Lagrange polynomial $p_1(x)$

$$f(x) = p_1(x) = \frac{(x - x_1)}{(x_0 - x_1)} f(x_0) + \frac{(x - x_0)}{(x_1 - x_0)} f(x_1),$$

and interpolating f(2.6) at these points is

$$f(2.6) \approx p_1(2.6) = \frac{(2.6 - 2.75)}{(2.5 - 2.75)}(3.7528) + \frac{(2.6 - 2.5)}{(2.75 - 2.5)}(4.0718) = 3.8804.$$

To compute absolute error in the linear interpolation of f(x), we use the error formula as follows

$$|f(2.6) - p_1(2.6)| = |(2.6 + \ln(2.6)) - 3.8804| = |3.8809 - 3.8804| = 0.0005,$$

the required error.

Example 4.23 (a) Let $f(x) = (x+1) \ln(x+1)$ be the function defined over the interval [1,2]. Find the approximations of $(1.9 \ln 1.9)$ using linear, quadratic and cubic Lagrange interpolating polynomials for equally spaced data points defined over the interval [1,2]. Compute the error bounds for linear, quadratic and cubic Lagrange interpolating polynomials for equally spaced data points. Also, compute absolute error for each case.

(b) Determine the step size h and the number of points to be used in the tabulation of the given function $f(x) = (x+1)\ln(x+1)$ in [1,2] so that linear, quadratic and cubic interpolations will be correct to six decimal places.

Solution. (a) For linear Lagrange polynomial, we have h = 2 - 1 = 1, so using $x_0 = 1, x_1 = 2$ and x = 0.9, in the linear Lagrange formula, we have

$$f(0.9) \approx p_1(0.9) = L_0(0.9)f(1) + L_1(0.9)f(2) = 1.3863L_0(0.9) + 3.2958L_1(0.9).$$
 (4.38)

The Lagrange coefficients can be calculate as follows:

$$L_0(0.9) = \frac{(0.9-2)}{(1-2)} = 1.1$$
 and $L_1(0.9) = \frac{(0.9-1)}{(2-1)} = -0.1$.

Putting these values of the Lagrange coefficients in (4.38), we have

$$f(0.9) \approx p_1(0.9) = 1.3863(1.1) + 3.2958(-0.1) = 1.1954.$$

Now for quadratic Lagrange polynomial, we have $h = \frac{2-1}{2} = 0.5$, so using $x_0 = 1, x_1 = 1.5, x_2 = 2$ and x = 0.9, in the quadratic Lagrange formula, we have

$$f(0.9) \approx p_2(0.9) = L_0(0.9)f(1) + L_1(0.9)f(1.5) + L_2(0.9)f(2)$$

= 1.3863 $L_0(0.9) + 2.2907L_1(0.9) + 3.2958L_2(0.9)$. (4.39)

The Lagrange coefficients can be calculate as follows:

$$L_0(0.9) = \frac{(0.9 - 1.5)(0.9 - 2)}{(1 - 1.5)(1 - 2)} = 1.32,$$

$$L_1(0.9) = \frac{(0.9 - 1)(0.9 - 2)}{(1.5 - 1)(1, 5 - 2)} = -0.44,$$

$$L_2(0.9) = \frac{(0.9 - 1)(0.9 - 1.5)}{(2 - 1)(2 - 1.5)} = 0.12.$$

Putting these values of the Lagrange coefficients in (4.39), we have

$$f(0.9) \approx p_2(0.9) = 1.3863(1.32) + 2.2907(-0.44) + 3.2958(0.12) = 1.2175.$$

For cubic Lagrange polynomial, we have $h = \frac{2-1}{3} = 1/3$, so using $x_0 = 1, x_1 = 4/3, x_2 = 5/3$, $x_3 = 2$ and x = 0.9, in the cubic Lagrange formula, we have

$$f(0.9) \approx p_3(0.9) = L_0(0.9)f(1) + L_1(0.9)f(4/3) + L_2(0.9)f(5/3) + L_3(0.9)f(2)$$

= 1.3863 $L_0(0.9) + 1.9770L_1(0.9) + 2.6156L_2(0.9) + 3.2958L_3(0.9)$. (4.40)

The Lagrange coefficients can be calculate as follows:

$$L_0(0.9) = \frac{(0.9 - 4/3)(0.9 - 5/3)(0.9 - 2)}{(1 - 4/3)(1 - 5/3)(1 - 2)} = 1.6449,$$

$$L_1(0.9) = \frac{(0.9 - 1)(0.9 - 5/3)(0.9 - 2)}{(4/3 - 1)(4/3 - 5/3)(4/3 - 2)} = -1.1385,$$

$$L_2(0.9) = \frac{(0.9 - 1)(0.9 - 4/3)(0.9 - 2)}{(5/3 - 1)(5/3 - 4/3)(5/3 - 2)} = 0.6435,$$

$$L_3(0.9) = \frac{(0.9 - 1)(0.9 - 4/3)(0.9 - 5/3)}{(2 - 1)(2 - 4/3)(2 - 5/3)} = -0.1495.$$

Putting these values of the Lagrange coefficients in (4.40), we have

$$f(0.9) \approx p_3(0.9) = 1.386(1.6449) + 1.977(-1.1385) + 2.616(0.6435) + 3.296(-0.1495) = 1.220.$$

The derivatives of the given function $f(x) = (x+1)\ln(x+1)$ are as follows:

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

Now for error bound of linear Lagrange polynomial, we use the formula

$$|E_1| \le \frac{Mh^2}{8},$$

where h = 2 - 1 = 1 and $M = \max_{1 \le x \le 2} |f''(x)| = \max_{1 \le x \le 2} \left| \frac{1}{x+1} \right| = \frac{1}{2}$. So

$$|E_1| \le \frac{(1/2)(1)^2}{8} = \frac{1}{16} = 0.0625,$$

the error bound for the linear Lagrange polynomial.

Similarly, for error bound of quadratic Lagrange polynomial, we use the formula

$$|E_2| \le \frac{Mh^3}{9\sqrt{3}},$$

where
$$h = (2-1)/2 = 1/2$$
 and $M = \max_{1 \le x \le 2} |f'''(x)| = \max_{1 \le x \le 2} \left| -\frac{1}{(x+1)^2} \right| = \frac{1}{4}$. So

$$|E_2| \le \frac{(1/4)(1/2)^3}{9\sqrt{3}} = \frac{(1/32)}{9\sqrt{3}} = 0.0020,$$

the error bound for the quadratic Lagrange polynomial.

Finally, for error bound of cubic Lagrange polynomial, we use the formula

$$|E_3| \le \frac{Mh^4}{24},$$

where
$$h = (2-1)/3 = 1/3$$
 and $M = \max_{1 \le x \le 2} |f^{(4)}(x)| = \max_{1 \le x \le 2} \left| \frac{2}{(x+1)^3} \right| = \frac{1}{4}$. Thus

$$|E_3| \le \frac{(1/4)(1/3)^4}{24} = \frac{1}{7776} = 0.0001,$$

the error bound for the cubic Lagrange polynomial. Finally,

$$|f(0.9) - p_1(0.9)| = |1.2195 - 1.1954| = 0.0241,$$

$$|f(0.9) - p_2(0.9)| = |1.2195 - 1.2175| = 0.0020,$$

$$|f(0.9) - p_3(0.9)| = |1.2195 - 1.2199| = 0.0004,$$

are respectively, the absolute error for linear, quadratic and cubic polynomials.

(b) Since we know that the upper bound of error in linear polynomial is

$$|E_1| \le \frac{Mh^2}{8}$$
 and $M = \frac{1}{2}$,

therefore,

$$\frac{h^2}{16} \le 5 \times 10^{-6}$$
, or $h^2 \le 80 \times 10^{-6}$.

This gives

$$h \le 0.0089$$
 and $n = 112.36 \approx 113$.

As the upper bound of error in quadratic polynomial is

$$|E_2| \le \frac{Mh^3}{9\sqrt{3}} \quad and \quad M = \frac{1}{4},$$

therefore,

$$\frac{h^3}{36\sqrt{3}} \le 5 \times 10^{-6}, \quad or \quad h^3 \le 311.7691 \times 10^{-6}.$$

This gives

$$h < 0.0678$$
 and $n = 14.7476 \approx 15$.

Finally, as the upper bound of error in cubic polynomial is

$$|E_2| \le \frac{Mh^4}{24} \quad and \quad M = \frac{1}{4},$$

therefore,

$$\frac{h^4}{96} \le 5 \times 10^{-6}$$
, or $h^4 \le 480 \times 10^{-6}$.

This gives

$$h < 0.1480$$
 and $n = 6.7560 \approx 7$.

Thus we need, respectively, 114 points, 16 points and 8 points for the linear, quadratic and cubic interpolations.

Example 4.24 Find the cubic Lagrange interpolating polynomial to find the approximation of f(x) if f(1) = 0.5, f(4/3) = 1, f(5/3) = 2 and f(2) = 3. If $|f^{(4)}| \le \frac{1}{5}$ for $1 \le x \le 2$, then show that the error estimate is $|f(x) - p_3(x)| \le \frac{1}{27}$.

Solution. The given number of points are, $x_0 = 1, x_1 = 4/3, x_2 = 5/3, x_3 = 2$,, therefore the cubic Lagrange interpolating polynomial for the approximation of the given function is:

$$f(x) = p_3(x) = L_0(x)f(x_0) + L_1(x)f(x_1) + L_2(x)f(x_2) + L_3(x)f(x_3),$$

and using the given function values, gives

$$f(x) = p_3(x) = L_0(x)(0.5) + L_1(x)(1) + L_2(x)(2) + L_3(x)(3).$$

The Lagrange coefficients can be calculate as follows:

$$L_0(1.5) = \frac{(x-4/3)(x-5/3)(x-2)}{(1-4/3)(1-5/3)(1-2)} = \frac{(x^3-5x^2+74/9x-40/9)}{-2/9},$$

$$L_1(1.5) = \frac{(x-1)(x-5/3)(x-2)}{(4/3-1)(4/3-5/3)(4/3-2)} = \frac{(x^3-14/3x^2+7x-10/3)}{2/27},$$

$$L_2(1.5) = \frac{(x-1)(x-4/3)(x-2)}{(5/3-1)(5/3-4/3)(5/3-2)} = \frac{(x^3-13/3x^2+6x-8/3)}{-2/27},$$

$$L_3(1.5) = \frac{(x-1)(x-4/3)(x-5/3)}{(2-1)(2-4/3)(2-5/3)} = \frac{(x^3-4x^2+47/9x-20/9)}{2/9}.$$

Putting these values of the Lagrange coefficients in the above equation, we get

$$f(x) = p_3(x) = \frac{1}{4}(-9x^3 + 45x^2 - 62x + 28),$$

which is the required cubic interpolating polynomial for the approximation of the function. Since we know the error of cubic Lagrange polynomial is

$$f(x) - p_3(x) = \frac{f^{(4)}(\eta(x))}{4!}(x - x_0)(x - x_1)(x - x_2)(x - x_3),$$

so

$$|f(x) - p_3(x)| = \frac{|f^{(4)}(\eta(x))|}{4!} |(x - x_0)||(x - x_1)||(x - x_2)||(x - x_3)|,$$

and $|f^{(4)}| \le \frac{1}{5}$ for $1 \le x \le 2$, we obtain

$$|f(x) - p_3(x)| \le \frac{1}{120}|(x - x_0)||(x - x_1)||(x - x_2)||(x - x_3)|.$$

Now for $1 \le x \le 2$, we deduce that

$$|x-1| \le 1$$
, $|x-4/3| \le 2/3$, $|x-5/3| \le 2/3$, $|x-2| \le 1$.

Hence, the possible error in the cubic Lagrange polynomial is

$$|f(x) - p_3(x)| \le \frac{1}{120}(1)(2/3)(2/3)(1) = 0.0037.$$

While the Lagrange interpolation formula is at the heart of polynomial interpolation, it is not, by any stretch of the imagination, the most practical way to use it. Using this interpolation formula, there is no restriction on the spacing and order of the tabulating points x_0, x_1, \ldots, x_n . Also, the value of y (the dependent variable) can be calculated at any point x within minimum and maximum values of x_0, x_1, \ldots, x_n . But just consider for a moment that if we have to add an addition data point in the previous Example 4.32, in order to find cubic polynomial $p_3(x)$, we have to repeat the whole process again because we can not use the solution of the quadratic polynomial $p_2(x)$ in the

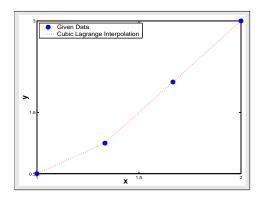


Figure 4.3: Graphical solution of the Example 4.24.

construction of the cubic polynomial $p_3(x)$. Therefore one can note that the Lagrange method is not particular efficient for large values of n, the degree of the polynomial. When n is large and the data for x is ordered, some improvement in efficiency can be obtained by considering only the data pairs in the vicinity of the x values for which f(x) is sought.

One will be quickly convinced that there must be better techniques available. In the following section we discuss some of the more practical approaches to polynomial interpolation. In using the following scheme the construction of the difference table plays an important role. In using the Lagrange interpolation scheme there is no need to construct difference table.

4.2.2 Newton's General Interpolating Formula

Since we noted in the previous section that for a small number of data point one can easily use the Lagrange formula of the interpolating polynomial. However, for a large number of data points there will be many multiplication and more significantly, whenever a new data point is added to an existing set, the interpolating polynomial has to be completely recalculated. Here, we describe an efficient way of organizing the calculations so as to overcome these disadvantages.

Let us consider the *nth*-degree polynomial $p_n(x)$ that agrees with the function f(x) at the distinct numbers x_0, x_1, \ldots, x_n . The divided differences of f(x) with respect to x_0, x_1, \ldots, x_n are derived to express $p_n(x)$ in the form

$$f(x) = p_n(x) = a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1) + \cdots + a_n(x - x_0)(x - x_1) \cdots (x - x_{n-1}),$$

$$(4.41)$$

for appropriate constants a_0, a_1, \ldots, a_n .

Now to determine the constants, firstly, by evaluating $p_n(x)$ at x_0 , we have

$$p_n(x_0) = a_0 = f(x_0) (4.42)$$

Similarly, when $p_n(x)$ is evaluated at x_1 , then

$$p_n(x_1) = a_0 + a_1(x_1 - x_0) = f(x_1),$$

which implies that

$$a_1 = \frac{f(x_1) - f(x_0)}{x_1 - x_0}. (4.43)$$

Divided Differences

Now we express the interpolating polynomial in terms of divided difference.

Firstly, we define the Zeroth divided difference at the point x_i by

$$f[x_i] = f(x_i), \tag{4.44}$$

which is simply the value of the function f(x) at x_i .

The first-order or first divided difference at the points x_i and x_{i+1} can be defined by

$$f[x_i, x_{i+1}] = \frac{f[x_{i+1}] - f[x_i]}{x_{i+1} - x_i} = \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i}.$$
(4.45)

In general, the *nth divided difference* $f[x_i, x_{i+1}, \dots, x_{i+n}]$ is defined by

$$f[x_i, x_{i+1}, \dots, x_{i+n}] = \frac{f[x_{i+1}, x_{i+2}, \dots, x_{i+n}] - f[x_i, x_{i+1}, \dots, x_{i+n-1}]}{x_{i+n} - x_i}.$$
 (4.46)

By using this definition, (4.42) and (4.43) can be written as

$$a_0 = f[x_0];$$
 $a_1 = f[x_0, x_1],$

respectively. Similarly, one can have the values of other constants involving in (4.41) such as

$$a_2 = f[x_0, x_1, x_2],$$

 $a_3 = f[x_0, x_1, x_2, x_3],$
 $\cdots = \cdots$
 $\cdots = \cdots$
 $a_n = f[x_0, x_1, \dots, x_n].$

Putting the values of these constants in (4.41), we get

Table 4.1: Divided difference table for a function y = f(x)

		Zero	First	Second	Third
		Divided	Divided	Divided	Divided
k	x_k	Difference	Difference	Difference	Difference
0	x_0	$f[x_0]$			
1	x_1	$f[x_1]$	$f[x_0, x_1]$		
2	x_2	$f[x_2]$	$f[x_1, x_2]$	$f[x_0, x_1, x_2]$	
3	x_3	$f[x_3]$	$f[x_2, x_3]$	$f[x_1, x_2, x_3]$	$f[x_0, x_1, x_2, x_3]$

$$f(x) = p_n(x) = f[x_0] + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1) + \dots + f[x_0, x_1, \dots, x_n](x - x_0)(x - x_1) \dots (x - x_{n-1}),$$

$$(4.47)$$

which can also be written as

$$f(x) = p_n(x) = f[x_0] + \sum_{k=1}^n f[x_0, x_1, \dots, x_k](x - x_0)(x - x_1) \cdots (x - x_{k-1}).$$
 (4.48)

This type of polynomial is known as the Newton's interpolatory divided difference polynomial. Table 4.1 shows the divided difference for a function f(x).

Example 4.25 Construct the fourth divided differences table for $f(x) = 4x^4 + 3x^3 + 2x^2 + 10$ for the values x = 3, 4, 5, 6, 7, 8.

Solution. The result are listed in Table 4.2.

From the results in Table 4.2, one can note that the nth divided difference for the nth polynomial equation is always constant and the (n+1)th divided difference is always zero for the nth polynomial equation.

		Table 4.2. Divided differences table for $f(x) = e^{-x}$ at given points										
			Zeroth	First	Second	Third	Fourth	Fifth				
			Divided	Divided	Divided	Divided	Divided	Divided				
	k	x_k	Difference	Difference	Difference	Difference	Difference	difference				
Ì	0	3	433									
	1	4	1258	825								
	2	5	2935	1677	426							
	3	6	5914	2979	651	75						
	4	7	10741	4827	924	91	4					
	5	8	18058	7317	1245	107	4	0				

Table 4.2: Divided differences table for $f(x) = e^x$ at given points

Using the following MATLAB command one can construct the Table 4.2 as follows:

>>
$$x = [3\ 4\ 5\ 6\ 7\ 8];$$

>> $y = 4*x.^4+3*x.^3+2*x.^2+10;$
>> $D = divdiff(x,y);$

Divided differences are now can be used to write the Newton's interpolating polynomial. Starting with the constant interpolating polynomial

$$p_0(x) = f[x_0].$$

Linear Newton's Interpolating Polynomial

The linear Newton's interpolating polynomial passing through two points $(x_0, f(x_0))$ and $(x_1, f(x_1))$ can be written as

$$f(x) = p_1(x) = f[x_0] + (x - x_0)f[x_0, x_1].$$

Quadratic Newton's Interpolating Polynomial

The quadratic Newton's interpolating polynomial passing through the points $(x_0, f(x_0))$, $(x_1, f(x_1))$ and $(x_2, f(x_2))$ can be written in terms of divided differences as

$$f(x) = p_2(x) = f[x_0] + (x - x_0)f[x_0, x_1] + (x - x_0)(x - x_1)f[x_0, x_1, x_2].$$

This polynomial can also be written as

$$f(x) = p_2(x) = p_1(x) + (x - x_0)(x - x_1)f[x_0, x_1, x_2],$$

that is, the interpolating polynomial of degree 2 makes full use of the polynomial of degree 1, simply adding one extra term to $p_1(x)$. This is one of the advantages of the Newton's polynomial over Lagrange polynomial.

Cubic Newton's Interpolating Polynomial

Similarly, the cubic Newton's interpolating polynomial passing through the points $(x_0, f(x_0))$, $(x_1, f(x_1)), (x_2, f(x_2))$ and $(x_3, f(x_3))$ can be written in terms of divided differences as

$$f(x) = p_3(x) = f[x_0] + (x - x_0)f[x_0, x_1] + (x - x_0)(x - x_1)f[x_0, x_1, x_2] + (x - x_0)(x - x_1)(x - x_2)f[x_0, x_1, x_2, x_3].$$

This polynomial can also be written as

$$f(x) = p_3(x) = p_2(x) + (x - x_0)(x - x_1)(x - x_2)f[x_0, x_1, x_2, x_3],$$

that is, the interpolating polynomial of degree 3 makes full use of the polynomial of degree 2, simply adding one extra term to $p_2(x)$. Note that using linear polynomial in quadratic polynomial, the starting point x_0 for both polynomials should be same.

Nth Degree Newton's Interpolating Polynomial

Repeating this entire process again, $p_3(x)$, $p_4(x)$ and higher degree interpolating polynomials can be consecutively obtained in the same way. In general, the interpolating polynomial $p_n(x)$ passing through the points $(x_i, f(x_i))(i = 0, 1, ..., n)$, can be written in terms of divided differences as

$$f(x) = p_n(x) = f[x_0] + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1) + \dots + f[x_0, x_1, \dots, x_n](x - x_0)(x - x_1) \dots (x - x_{n-1}),$$

$$(4.49)$$

which can also be written as

$$f(x) = p_n(x) = f[x_0] + \sum_{k=1}^n f[x_0, x_1, \dots, x_k](x - x_0)(x - x_1) \cdots (x - x_{k-1}), \tag{4.50}$$

or

$$f(x) = p_n(x) = f[x_0] + \sum_{k=0}^{n} f[x_0, x_1, \dots, x_k] \prod_{i=0}^{k-1} (x - x_i).$$
 (4.51)

This type of polynomial is known as the Newton's interpolatory divided difference polynomial.

Theorem 4.4 (Newton's Interpolating Polynomial)

Suppose that x_0, x_1, \ldots, x_n are (n+1) distinct points in the interval [a, b]. There exists a unique polynomial $p_n(x)$ of degree at most n with the property that

$$f(x_i) = p_n(x_i), \text{ for } i = 0, 1, \dots, n.$$

The Newton's form of this polynomial is

$$f(x) = p_n(x) = a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1) + \dots + a_n(x - x_0)(x - x_1) + \dots + a_n(x - x_0)(x - x_1),$$

where

$$a_k = f[x_0, x_1, x_2, \dots, x_k], \text{ for } k = 0, 1, 2, \dots, n.$$

One can easily show that (4.50) is simply a rearrangement of the Lagrange form defined by (4.14). For example, the Newton divided difference interpolation polynomial of degree one is

$$f(x) = p_1(x) = f[x_0] + f[x_0, x_1](x - x_0),$$

which implies that

$$f(x) = p_1(x) = f(x_0) + \left(\frac{f(x_1) - f(x_0)}{x_1 - x_0}\right)(x - x_0)$$

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

$$= \left(\frac{x - x_1}{x_0 - x_1}\right)f(x_0) + \left(\frac{x - x_0}{x_1 - x_0}\right)f(x_1),$$

which is the Lagrange interpolating polynomial of degree one. Similarly, one can show the equivalent for the nth-degree polynomial.

Example 4.26 Find the Lagrange and the Newton forms of the interpolating polynomial for the following data

Write both polynomials in the form $a + bx + cx^2$ to verify that they are identical as functions.

Solution. With $x_0 = 0, x_1 = 1$ and $x_2 = 3$, we obtain the quadratic Lagrange interpolating polynomial

$$f(x) = p_2(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} f(x_0) + \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} f(x_1) + \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} f(x_2)$$

$$= \frac{(x - 1)(x - 3)}{(0 - 1)(0 - 3)} (1) + \frac{(x - 0)(x - 3)}{(1 - 0)(1 - 3)} (2) + \frac{(x - 0)(x - 1)}{(3 - 0)(3 - 1)} (3).$$

After simplifying, we get

$$f(x) = p_2(x) = 1 + \frac{7}{6}x - \frac{1}{6}x^2,$$

which is the quadratic Lagrange interpolating polynomial.

Now we construct the divided differences table for the given data points. The result of the divided difference is listed in Table 4.3.

Since the Newton's interpolating polynomial of degree 2 is defined as

$$f(x) = p_2(x) = f[x_0] + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1).$$

By using Table 4.3, we have

$$f(x) = p_2(x) = 1 + (1)(x - 0) + \left(-\frac{1}{6}\right)(x - 0)(x - 1).$$

It gives Newton polynomial of the form

$$f(x) = p_2(x) = 1 + \frac{7}{6}x - \frac{1}{6}x^2,$$

which show that both polynomials are identical as functions.

Table 4.3: Divided differences table for the Example 4.26

		Zeroth Divided	First Divided	Second Divided
k	x_k	Difference	Difference	Difference
0	0	1		
1	1	2	1	
2	3	3	$\frac{1}{2}$	$-\frac{1}{6}$

Example 4.27 Show that the Newton's interpolating polynomial $p_2(x)$ of degree 2 satisfies the interpolation conditions

$$p_2(x_i) = f(x_i), \qquad i = 0, 1, 2.$$

Solution. Since the Newton's interpolating polynomial of degree 2 is

$$f(x) = p_2(x) = f[x_0] + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1).$$

First take $x = x_0$, we have

$$p_2(x_0) = f[x_0] + 0 + 0 = f(x_0).$$

Now take $x = x_1$, we have

$$p_2(x_1) = f[x_0] + f[x_0, x_1](x_1 - x_0) + 0 = f(x_0) + \frac{f(x_1) - f(x_0)}{x_1 - x_0}(x_1 - x_0),$$

it gives

$$p_2(x_1) = f(x_0) + f(x_1) - f(x_0) = f(x_1).$$

Finally, take $x = x_2$, we have

$$p_2(x_2) = f[x_0] + f[x_0, x_1](x_2 - x_0) + f[x_0, x_1, x_2](x_2 - x_0)(x_2 - x_1),$$

which can be written as

$$p_2(x_2) = f[x_0] + f[x_0, x_1](x_2 - x_0) + \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0}(x_2 - x_0)(x_2 - x_1).$$

It gives

$$p_2(x_2) = f[x_0] + f[x_0, x_1](x_2 - x_1 + x_1 - x_0) + f[x_1, x_2](x_2 - x_1) - f[x_0, x_1](x_2 - x_1),$$

or

$$p_2(x_2) = f[x_0] + f[x_0, x_1](x_1 - x_0) + f[x_1, x_2](x_2 - x_1).$$

From (4.45), we have

$$p_2(x_2) = f[x_0] + \frac{f(x_1) - f(x_0)}{x_1 - x_0}(x_1 - x_0) + \frac{f(x_2) - f(x_1)}{x_2 - x_1}(x_2 - x_1),$$

which gives

$$p_2(x_2) = f(x_0) + f(x_1) - f(x_0) + f(x_2) - f(x_1) = f(x_2).$$

Example 4.28 The cubic Newton's polynomial $p_3(x) = 2 - (x+1) + x(x+1) - 2x(x+1)(x-1)$ interpolates the first four points in the following table:

By adding one additional term (3,10) to $p_3(x)$, find Newton's polynomial $p_4(x)$ that interpolates the whole table and then use it to find the approximation of f(0.5).

Solution. Since the Newton's polynomial for the whole table data points is the four degree Newton's interpolating polynomial and it can be written as

$$f(x) = p_4(x) = p_3(x) + x(x+1)(x-1)(x-2)f[x_0, x_1, x_2, x_3].$$

Now to find to find fourth divided difference $f[x_0, x_1, x_2, x_3]$, we have to construct the required divided differences table. The result are listed in Table 4.4.

Thus the Newton's interpolating polynomial passing through all the given data points is

$$f(x) = p_4(x) = 2 - (x+1) + x(x+1) - 2x(x+1)(x-1) + 2x(x+1)(x-1)(x-2).$$

Thus at x = 0.5, we get

$$f(0.5) \approx p_4(0.5) = 3.1250$$

the required approximation of the function.

_	Table 1.1. Divided differences table for $f(w) = e^{-av}$ given points										
		Zeroth	First	Second	Third	Fourth					
		Divided	Divided	Divided	Divided	Divided					
k	x_k	Difference	Difference	Difference	Difference	difference					
0	-1	2									
1	0	1	-1								
2	1	2	1	1							
3	2	-7	-9	-5	-2						
4	3	10	17	13	6	2					

Table 4.4: Divided differences table for $f(x) = e^x$ at given points

Example 4.29 Consider the following table of date points

Find the third divided difference f[3,1,5,6] and use it to find the Newton's form of the interpolating polynomial. Find approximation of f(2).

Solution. The third divided differences for the given data points are listed in Table 4.5. The cubic

Table 4.5: Divided difference table for a function y = f(x)

		Zero	First	Second	Third
		Divided	Divided	Divided	Divided
k	x_k	Difference	Difference	Difference	Difference
0	$x_0 = 3$	$f[x_0] = 1$			
1	$x_1 = 1$	$f[x_1] = -3$	$f[x_0, x_1] = 2$		
2	$x_2 = 5$	$f[x_2] = 2$	$f[x_1, x_2] = 5/4$	$f[x_0, x_1, x_2] = -3/8$	
3	$x_3 = 6$	$f[x_3] = 4$	$f[x_2, x_3] = 2$	$f[x_1, x_2, x_3] = 3/20$	$f[x_0, x_1, x_2, x_3] = 7/40$

Newton's interpolating polynomial passing through the given can be written as

$$f(x) = p_3(x) = f[x_0] + (x - x_0)f[x_0, x_1] + (x - x_0)(x - x_1)f[x_0, x_1, x_2] + (x - x_0)(x - x_1)(x - x_2)f[x_0, x_1, x_2, x_3],$$

so using Table 4.5, we have

$$f(x) = p_3(x) = 1 + 2(x - x_0) - \frac{3}{8}(x - x_0)(x - x_1) + \frac{7}{40}(x - x_0)(x - x_1)(x - x_2),$$

or

$$f(x) = p_3(x) = \frac{1}{40} [7x^3 - 78x^2 + 301x - 350].$$

Thus at x = 2, we get

$$f(2) \approx p_3(2) = \frac{1}{40} [7(2)^3 - 78(2)^2 + 301(2) - 350] = -\frac{1}{10},$$

the required approximation of the function at x=2.

Example 4.30 Consider the following table of date points

Is this table form a polynomial? If so, what maximum degree of the polynomial can be obtain to approximate f(0).

Solution. The result are listed in Table 4.6. Since all the second order differences are the same

Table 4.6: Divided difference table for a function y = f(x)

		Zero	First	Second	Third
		Divided	Divided	Divided	Divided
k	x_k	Difference	Difference	Difference	Difference
0	$x_0 = -1$	$f[x_0] = 8$			
1	$x_1 = 3$	$f[x_1] = 0$	$f[x_0, x_1] = -2$		
2	$x_2 = 2$	$f[x_2] = -1$	$f[x_1, x_2] = 1$	$f[x_0, x_1, x_2] = 1$	
3	$x_3 = -2$	$f[x_3] = 15$	$f[x_2, x_3] = -4$	$f[x_1, x_2, x_3] = 1$	$ f[x_0, x_1, x_2, x_3] = 0 $
4	$x_4 = 4$	$f[x_3] = 3$	$f[x_2, x_3] = -2$	$f[x_1, x_2, x_3] = 1$	$f[x_1, x_2, x_3, x_4] = 0$

they equal 1, so that means the maximum degree of the polynomial will be quadratic. By using Table 4.6, it can be written as

$$f(x) = p_2(x) = 8 - 2(x+1) + 1(x+1)(x-3) = x^2 - 4x + 3.$$

Thus at x = 0, we get

$$f(0) \approx p_2(0) = 3$$

the required approximation of the function at x = 0.

Example 4.31 Repeat the Example 4.30 by adding f(0) = 5, that is, the following table of the data points of the form

Solution. The result are listed in Table 4.7. Since we have the value of the fifth order differences, so that means the maximum degree of the polynomial will be five. By using Table 4.7, it can be written as

$$f(x) = p_5(x) = 5 - 3x + 0.3333x(x+1) + 0.3333x(x+1)(x-3) + 0.1667x(x+1)(x-3)(x-2) - 0.0417x(x+1)(x-3)(x-2)(x+2),$$

$$f(x) = p_5(x) = \frac{1}{10000} \left[-417x^5 + 2501x^4 - 416x^3 - 5002x^2 - 15834x - 50000 \right],$$

the required fifth degree polynomial.

lai	Table 4.7. Divided difference table for a function $g = f(x)$									
k	x_k	ODD	1DD	2DD	3DD	4DD	5DD			
0	0	5								
1	-1	8	-3							
2	3	0	-2	0.3333						
3	2	-1	1	1	0.3333					
4	-2	15	-4	1	0	0.1667				
5	4	3	-2	1	0	0	-0.0417			

Table 4.7: Divided difference table for a function y = f(x)

Example 4.32 Construct the divided differences table for $f(x) = x^3 + 7x^2 + 1$ using the values x = 1, 2, 3, 4, 5. If the approximation of f(3.5) by linear Newton's polynomial is 134, then find the best approximation of f(3.5) by using quadratic Newton's polynomial.

Solution. The result are listed in Table 4.8. Since we know the quadratic Newton's polynomial has a form

$$f(x) = p_2(x) = p_1(x) + (x - x_0)(x - x_1)f[x_0, x_1, x_2].$$

After choosing the best points $x_0 = 3, x_1 = 4, x_3 = 5$, and Table 4.11, we get

$$f(3.5) \approx p_2(3.5) = p_1(3.5) + (3.5 - 3)(3.5 - 4)f[3, 4, 5] = 134 - (0.25)(19) = 129.25,$$

and the absolute error

$$|f(3.5) - p_2(3.5)| = |129.625 - 129.25| = 0.375.$$

From the results in Table 4.8, one can note that the nth divided difference for the nth degree polynomial equation is always constant and the (n+1)th divided difference is always zero for the nth degree polynomial equation.

Using the following MATLAB command one can construct the Table 4.8 as follows:

$$>> x = [1\ 2\ 3\ 4\ 5];\ y = x.^3 + 7*x.^2 + 1;\ D = divdiff(x,y);$$

Table 4.8: Divided differences table for $f(x) = x^3 + 7x^2 + 1$

		Zeroth	First	Second	Third	Fourth
		Divided	Divided	Divided	Divided	Divided
k	x_k	Difference	Difference	Difference	Difference	Difference
0	1	9				
1	2	37	28			
2	3	91	54	13		
3	4	177	86	16	1	
4	5	301	124	19	1	0

The main advantage of the Newton divided difference form over the Lagrange form is that polynomial $p_n(x)$ can be calculated from polynomial $p_{n-1}(x)$ by adding just one extra term, since it follows from (4.50) that

$$f(x) = p_n(x) = p_{n-1}(x) + f[x_0, x_1, \dots, x_n](x - x_0)(x - x_1) \cdots (x - x_{n-1}). \tag{4.52}$$

Program 4.2

MATLAB m-file for the Divided Differences of Function

function D=divdiff(x,y)

m = length(x); D = zeros(m, m); D(:, 1) = y(:);

for j=2:m; for i=j:m

D(i,j) = (D(i,j-1) - D(i-1,j-1))/(x(i) - x(i-j+1)); end; end

- **Example 4.33** (a) Construct the divided difference table for the function $f(x) = \ln(x+2)$ in the interval $0 \le x \le 3$ for the stepsize h = 1.
- (b) Use Newton divided difference interpolation formula to construct the interpolating polynomials of degree 2 and degree 3 to approximate ln(3.5).
- (c) Compute error bounds for the approximations in part (b).

Solution. (a) The results of the divided differences are listed in Table 4.9.

(b) Firstly, we construct the second degree polynomial $p_2(x)$ by using the quadratic Newton interpolation formula as follows

$$f(x) = p_2(x) = f[x_0] + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1),$$

then with the help of the divided differences Table 4.9, we get

$$f(x) = p_2(x) = 0.6932 + 0.4055(x - 0) - 0.0589(x - 0)(x - 1),$$

which implies that

$$f(x) = p_2(x) = -0.0568x^2 + 0.4644x + 0.6932$$
 and $p_2(1.5) = 1.2620$,

with possible absolute error

$$|f(1.5) - p_2(1.5)| = |1.2528 - 1.2620| = 0.0072.$$

Now to construct the cubic interpolatory polynomial $p_3(x)$ that fits at all four points. We only have to add one more term to the polynomial $p_2(x)$:

$$f(x) = p_3(x) = p_2(x) + f[x_0, x_1, x_2, x_3](x - x_0)(x - x_1)(x - x_2),$$

and this gives

$$f(x) = p_3(x) = p_2(x) + 0.0089(x^3 - 3x^2 + 2x)$$
 and $f(1.5) \approx p_3(1.5) = 1.2620 - 0.0033 = 1.2587$,

<u> ra</u>	Table 4.9: Divide differences table for the Example 4.55									
		Zeroth	First	Second	Third					
		Divided	Divided	Divided	Divided					
k	x_k	Difference	Difference	Difference	Difference					
0	0	0.6932								
1	1	1.0986	0.4055							
2	2	1.3863	0.2877	- 0.0589						
3	3	1.6094	0.2232	- 0.0323	0.0089					

Table 4.9: Divide differences table for the Example 4.33

with possible absolute error

$$|f(1.5) - p_3(1.5)| = |1.2528 - 1.2587| = 0.0059.$$

We note that the estimated value of f(1.5) by cubic interpolating polynomial is more closer to the exact solution than the quadratic polynomial.

(c) Now to compute the error bounds for the approximations in part (b), we use the error formula (4.30). For the polynomial $p_2(x)$, we have

$$|f(x) - p_2(x)| = \frac{|f'''(\eta(x))|}{3!} |(x - x_0)(x - x_1)(x - x_2)|.$$

The third derivative of the given function is given as

$$f'''(x) = \frac{2}{(x+2)^3}$$
 and $|f'''(\eta(x))| = \left|\frac{2}{(\eta(x)+2)^3}\right|$, for $\eta(x) \in (0,2)$.

Then

$$M = \max_{0 \le x \le 2} \left| \frac{2}{(x+2)^3} \right| = 0.25,$$

and

$$|f(1.5) - p_2(1.5)| \le (0.375)(0.25)/6 = 0.0156,$$

which is the required error bound for the approximation $p_2(1.5)$. Since the error bound for the cubic polynomial $p_3(x)$ is

$$|f(x) - p_3(x)| = \frac{|f^{(4)}(\eta(x))|}{4!} |(x - x_0)(x - x_1)(x - x_2)(x - x_3)|.$$

Taking the fourth derivative of the given function, we obtain

$$f^{(4)}(x) = \frac{-6}{(x+2)^4} \quad and \quad |f^{(4)}(\eta(x))| = \Big|\frac{-6}{(\eta(x)+2)^4}\Big|, \quad for \quad \eta(x) \in (0,3).$$

Since

$$|f^{(4)}(0)| = 0.375$$
 and $|f^{(4)}(3)| = 0.0096$,

so
$$|f^{(4)}(\eta(x))| \le \max_{0 \le x \le 3} \left| \frac{-6}{(x+2)^4} \right| = 0.375$$
 and it gives

$$|f(1.5) - p_3(1.5)| \le (0.5625)(0.375)/24 = 0.0088,$$

which is the required error bound for the approximation $p_3(1.5)$.

Note that in the above Example 4.33, we used the value of the quadratic polynomial $p_2(1.5)$ in calculating the cubic polynomial $p_3(1.5)$. It was possible because the initial value for both polynomials was the same as $x_0 = 0$. But the situation will be quite different if the initial point for both polynomials will be different. For example, if we have to find the approximate value of $\ln(4.5)$, then the suitable data points for the quadratic polynomial will be $x_0 = 1, x_1 = 2, x_2 = 3$ and for the cubic polynomial will be $x_0 = 0, x_1 = 1, x_2 = 2, x_3 = 3$. So for getting the best approximation of $\ln(4.5)$ by the cubic polynomial $p_3(2.5)$, we can not use the value of the quadratic polynomial $p_2(2.5)$ in the cubic polynomial $p_3(2.5)$. The best way is to use the following cubic polynomial form

$$p_3(2.5) = f[0] + (2.5 - 0)f[0, 1] + (2.5 - 0)(2.5 - 1)f[0, 1, 2] + (2.5 - 0)(2.5 - 1)(2.5 - 2)f[0, 1, 2, 3],$$

which gives

$$p_3(2.5) = 0.6932 + 1.0137 - 0.2208 + 0.0166 = 1.5027.$$

Using the following MATLAB command one can reproduce the results of the Example 4.33 as follows:

>>
$$x = [0 \ 1 \ 2 \ 3]; \ y = log(x + 2);$$

>> $x0 = 1.5; \ Y = Ndivf(x, y, x0);$

Program 4.3

MATLAB m-file for Linear Newton's Interpolation Method function Y=Ndivf(x,y,x0)

m = length(x); D = zeros(m, m); D(:, 1) = y(:);

for j=2:m; for i=j:m;

D(i,j) = (D(i,j-1) - D(i-1,j-1))/(x(i) - x(i-j+1)); end; end;

Y = D(m, m) * ones(size(x0));

for i = m - 1 : -1 : 1; Y = D(i, i) + (x0 - x(i)) * Y; end

Example 4.34 Consider the points $x_0 = 0$, $x_1 = 0.4$, $x_2 = 0.7$ and for a function f(x), the divided differences are $f[x_2] = 6$, $f[x_1, x_2] = 10$, $f[x_0, x_1, x_2] = 50/7$. Use linear and quadratic Newton's polynomials to find the approximations of f(0.5).

Solution. First we construct the complete divided differences table for the given data points. Since we know that the second divided difference is defined as

$$f[x_0, x_1, x_2] = \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0}.$$

Using given values, we have

$$\frac{50}{7} = \frac{10 - f[x_0, x_1]}{0.7 - 0}.$$

Solving for $f[x_0, x_1]$, we have, $f[x_0, x_1] = 5$. Now we need to find the values of the zeroth divided differences $f[x_0]$ and $f[x_1]$ which can be obtained by using the first-order divided differences $f[x_0, x_1]$

and $f[x_1, x_2]$. Firstly, we find the value of $f[x_1]$ as follows

$$f[x_1, x_2] = \frac{f[x_2] - f[x_1]}{x_2 - x_1}$$

$$10 = \frac{6 - f[x_1]}{0.7 - 0.4}$$

$$f[x_1] = 6 - 10(0.3) = 3.$$

The other zeroth divided difference $f[x_0]$ can be computed as follows

$$f[x_0, x_1] = \frac{f[x_1] - f[x_0]}{x_1 - x_0}$$

$$5 = \frac{3 - f[x_0]}{0.4 - 0}$$

$$f[x_0] = 3 - 5(0.4) = 1,$$

which completes the divided differences table as shown by Table 4.10. Using the linear Newton's

Table 4.10: Divided differences table for the Example 4.34

		Zeroth Divided	First Divided	Second Divided
k	x_k	Difference	Difference	Difference
0	0	1		
1	0.4	3	5	
2	0.7	6	10	$\frac{50}{7}$

polynomial and the divided differences Table 4.10, we get

$$f(0.5) \approx p_1(0.5) = 3 + (10)(0.5 - 0.4) = 3 + 1 = 4$$

and then the quadratic Newton's polynomial and the divided differences Table 4.10, we have

$$f(0.5) \approx p_2(0.5) = 1 + 5(0.5 - 0) + \frac{50}{7}(0.5 - 0)(0.5 - 0.4) = 3.8571,$$

the approximations of f(0.5) using linear and quadratic Newton's polynomials respectively.

Example 4.35 Let $x_0 = 0.5, x_1 = 0.7, x_2 = 0.9, x_3 = 1.1, x_4 = 1.3, x_5 = 1.5$. Construct the divided difference table for the function $f(x) = e^x$. Use Newton polynomial $p_5(x)$ of degree five to approximate the function $f(x) = e^x$ at x = 0.6 when $p_4(0.6) = 1.9112$. Also, compute an error bound for your approximation.

Solution. Since the fifth-degree Newton polynomial $p_5(x)$ is defined as

$$f(x) = p_5(x) = p_4(x) + (x - x_0)(x - x_1)(x - x_2)(x - x_3)(x - x_4)f[x_0, x_1, x_2, x_3, x_4, x_5],$$

and using the given data points, we have

$$f(0.6) \approx p_5(0.6) = p_4(0.6) + (0.1)(-0.1)(-0.3)(-0.5)(-0.7)f[0.5, 0.7, 0.9, 1.1, 1.3, 1.5].$$

Now we compute the fifth-order divided differences of the function as follows: Thus

	Table 4.11: Divided differences table for $f(x) = e^x$ at the given points										
			Zeroth	First	Second	Third	Fourth	Fifth			
			Divided	Divided	Divided	Divided	Divided	Divided			
	k	x_k	Difference	Difference	Difference	Difference	Difference	difference			
Г	0	0.5	1.6487								
	1	0.7	2.0138	1.8252							
	2	0.9	2.4596	2.2293	1.0102						
	3	1.1	3.0042	2.7228	1.2339	0.3728					
	4	1.3	3.6693	3.3257	1.5071	0.4553	0.1032				
İ	5	1.5	4.4817	4.0620	1.8408	0.4553	0.1260	0.0228			

$$p_5(0.6) = 1.9112 + (0.0010)(0.0228) = 1.9112228.$$

Since the error bound for the fifth-degree polynomial $p_5(x)$ is

$$|f(x) - p_5(x)| = \frac{|f^{(6)}(\eta(x))|}{6!} |(x - x_0)(x - x_1)(x - x_2)(x - x_3)(x - x_4)(x - x_5)|.$$

Taking the sixth derivative of the given function, we have

$$f^{(6)}(x) = e^x,$$

and

$$|f^{(6)}(\eta(x))| = e^{(\eta(x))}, \quad for \quad \eta(x) \in (0.5, 1.5).$$

Since

$$|f^{(6)}(0.5)| = 1.6487$$
 and $|f^{(6)}(1.5)| = 4.4817$,

so

$$|f^{(6)}(\eta(x))| \le \max_{0.5 \le x \le 1.5} |e^x| = 4.4817,$$

therefore, we get

$$|f(0.6) - p_5(0.6)| \le (0.00095)(4.4817)/720 = 0.000006,$$

which is the required error bound for the approximation $p_5(0.6)$.

Example 4.36 If f(x) = p(x)q(x), then show that

$$f[x_0, x_1] = p(x_1)q[x_0, x_1] + q(x_0)p[x_0, x_1].$$

Also, find the values of p[0,1] and q[0,1] when f[0,1] = 4, f(1) = 5 and p(1) = q(0) = 2.

Solution. The first-order divided difference can be written as

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

Now using $f(x_1) = p(x_1)q(x_1)$ and $f(x_0) = p(x_0)q(x_0)$ in the above formula, we have

$$f[x_0, x_1] = \frac{p(x_1)q(x_1) - p(x_0)q(x_0)}{x_1 - x_0}.$$

Adding and subtracting the term $p(x_1)q(x_0)$, we obtain

$$f[x_0, x_1] = \frac{p(x_1)q(x_1) - p(x_1)q(x_0) + p(x_1)q(x_0) - p(x_0)q(x_0)}{x_1 - x_0},$$

which can be written as

$$f[x_0, x_1] = p(x_1) \frac{q(x_1) - q(x_0)}{x_1 - x_0} + q(x_0) \frac{p(x_1) - p(x_0)}{x_1 - x_0}.$$

Thus

$$f[x_0, x_1] = p(x_1)q[x_0, x_1] + q(x_0)p[x_0, x_1].$$

Given $x_0 = 0, x_1 = 1, f(1) = 5$, and f[0, 1] = 4, we obtain

$$f[0,1] = \frac{f(1) - f(0)}{1 - 0} = f(1) - f(0),$$

or

$$f[0,1] = 4 = 5 - f(0), \quad gives \quad f(0) = 1,$$

Also

$$f(1) = 5 = p(1)q(1) = 2q(1)$$
, gives $q(1) = 5/2$,

and

$$f(0) = 1 = p(0)q(0) = 2p(0), qives p(0) = 1/2.$$

Hence

$$p[0,1] = \frac{p(1) - p(0)}{1 - 0} = p(1) - p(0) = 2 - 1/2 = 3/2,$$

and

$$q[0,1] = \frac{q(1) - q(0)}{1 - 0} = q(1) - q(0) = 5/2 - 2 = 1/2.$$

In the case of the Lagrange interpolating polynomial we derived an expression for the truncation error in the form given by (4.30), namely, that

$$R_{n+1}(x) = \frac{f^{(n+1)}(\eta(x))}{(n+1)!} L_n(x),$$

where
$$L_n(x) = (x - x_0)(x - x_1) \cdots (x - x_n)$$
.

For the Newton's divided difference formula, we obtain, following the same reasoning as above

$$f(x) = f[x_0] + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1) + \cdots + f[x_0, x_1, \dots, x_n](x - x_0)(x - x_1) \cdots (x - x_{n-1}) + f[x_0, x_1, \dots, x_n, x](x - x_0)(x - x_1) \cdots (x - x_{n-1})(x - x_n),$$

which can also be written as

$$f(x) = p_n(x) + f[x_0, x_1, \dots, x_n, x](x - x_0)(x - x_1) \cdots (x - x_n), \tag{4.53}$$

or

$$f(x) - p_n(x) = L_n(x)f[x_0, x_1, \dots, x_n, x].$$
(4.54)

Since the interpolation polynomial agreeing with f(x) at x_0, x_1, \ldots, x_n is unique, it follows that these two error expressions must be equal.

Theorem 4.5 Let $p_n(x)$ be the polynomial of degree at most n that interpolates a function f(x) at a set of n+1 distinct points x_0, x_1, \ldots, x_n . If x is a point different from the points x_0, x_1, \ldots, x_n , then

$$f(x) - p_n(x) = f[x_0, x_1, \dots, x_n, x] \prod_{j=0}^{n} (x - x_j).$$
(4.55)

One can easily shows that the relationship between the divided differences and the derivative. From (4.46), we have

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

Now apply the Mean Value Theorem to the above equation implies that when the derivative f' exists, then we have

$$f[x_0, x_1] = f'(\eta(x)),$$

for unknown point $\eta(x)$ lies between x_0 and x_1 . Following theorem generalizes this result.

Theorem 4.6 (Divided Differences and Derivatives)

Suppose that $f \in C^n[a,b]$ and x_0, x_1, \ldots, x_n are distinct number in [a,b]. Then for some point $\eta(x)$ in the interval (a,b) spanned by x_0, \ldots, x_n exists with

$$f[x_0, x_1, \dots, x_n] = \frac{f^{(n)}(\eta(x))}{n!}.$$
(4.56)

Example 4.37 Let $f(x) = x \ln x$, and the points $x_0 = 1.1, x_1 = 1.2, x_2 = 1.3$. Compute the best approximate value for unknown point $\eta(x)$ by using the relation (4.56).

Solution. Given $f(x) = x \ln x$, then

$$f(1.1) = 1.1 \ln(1.1) = 0.1048,$$

 $f(1.2) = 1.2 \ln(1.2) = 0.2188,$
 $f(1.3) = 1.3 \ln(1.3) = 0.3411.$

Since the relation (4.56) for the given data points is

$$f[x_0, x_1, x_2] = \frac{f''(\eta(x))}{2!}. (4.57)$$

To compute the value of the left-hand side of the relation (4.57), we have to find the values of the first-order divided differences

$$f[x_0, x_1] = \frac{f(x_1) - f(x_0)}{x_1 - x_0} = \frac{0.2188 - 0.1048}{1.2 - 1.1} = 1.1400,$$

and

$$f[x_1, x_2] = \frac{f(x_2) - f(x_1)}{x_2 - x_1} = \frac{0.3411 - 0.2188}{1.3 - 1.2} = 1.2230.$$

Using these values, we can compute the second-order divided difference as

$$f[x_0, x_1, x_2] = \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0} = \frac{1.2230 - 1.1400}{1.3 - 1.1} = 0.4150.$$

Now we calculate the right-hand side of the relation (4.57) for the given points and which gives us

$$\frac{f''(x_0)}{2} = \frac{1}{2x_0} = 0.4546, \quad \frac{f''(x_1)}{2} = \frac{1}{2x_1} = 0.4167, \quad \frac{f''(x_2)}{2} = \frac{1}{2x_2} = 0.3846.$$

We note that the left-hand side of (4.57) is nearly equal to the right-hand side when $x_1 = 1.2$. Hence the best approximate value of $\eta(x)$ is 1.2.

Properties of Divided Differences

Now we discuss some of the nice properties of the divided differences as follows:

1. Divided difference of a constant is zero. Let f(x) = a, then

$$f[x_0, x_1] = \frac{f(x_1) - f(x_0)}{x_1 - x_0} = \frac{a - a}{x_1 - x_0} = 0.$$

2. Divided difference of h(x) = af(x), a is constant, is the divided difference of f(x) multiplied by a. Let h(x) = af(x), then

$$h[x_0, x_1] = \frac{h(x_1) - h(x_0)}{x_1 - x_0} = \frac{af(x_1) - af(x_0)}{x_1 - x_0} = a\frac{f(x_1) - f(x_0)}{x_1 - x_0} = af[x_0, x_1].$$

3. Divided difference obeys linear property.

Let
$$F(x) = af_1(x) + bf_2(x)$$
, then

$$F[x_0, x_1] = \frac{F[x_0] - F[x_1]}{x_0 - x_1} = \frac{(af_1(x_0) + bf_2(x_0)) - (af_1(x_1) + bf_2(x_1))}{x_0 - x_1}$$

$$= \frac{(af_1(x_0) - af_1(x_1)) + (bf_1(x_0) - bf_2(x_1))}{x_0 - x_1}$$

$$= a\left(\frac{f_1(x_0) - f_1(x_1)}{x_0 - x_1}\right) + b\left(\frac{f_2(x_0) - bf_2(x_1)}{x_0 - x_1}\right)$$

$$= af_1[x_0, x_1] + bf_2[x_0, x_1].$$

- 4. If $p_n(x)$ is a polynomial of degree n, then the divided differences of order n is always constant and $(n+1), (n+2), \ldots$ are identically zero.
- 5. The divided difference is a symmetric function of its arguments, that is, $f[x_0, x_1] = f[x_1, x_0]$. Thus if (t_0, t_1, \ldots, t_n) is a permutation of (x_0, x_1, \ldots, x_n) , then

$$f[t_0, t_1, \dots, t_n] = f[x_0, x_1, \dots, x_n],$$

This can be verify easily, since the divided differences on the both sides of the above equation are the coefficient of x^n in the polynomial of degree at most n that interpolates f(x) at the n+1 distinct points t_0, t_1, \ldots, t_n and x_0, x_1, \ldots, x_n . These two polynomials are, of course, the same.

6. The interpolating polynomial of degree n can be obtained by adding a single term to the polynomial of degree (n-1) expressed in the Newton form.

$$p_n(x) = p_{n-1}(x) + f[x_0, \dots, x_n] \prod_{j=0}^{n-1} (x - x_j).$$

- 7. The divided difference $f[x_0, \ldots, x_{n-1}]$ is the coefficient of x^{n-1} in the polynomial that interpolates $(x_0, f_0), (x_1, f_1), \ldots, (x_{n-1}, f_{n-1})$.
- 8. A sequence of divided differences may be constructed recursively from the formula

$$f[x_0,\ldots,x_n] = \frac{f[x_1,\ldots,x_n] - f[x_0,\ldots,x_{n-1}]}{x_n - x_0},$$

and the zeroth-order divided difference is defined by

$$f[x_i] = f(x_i), \quad i = 0, 1, \dots, n.$$

9. The another useful property of the divided difference can be obtained by using the definitions of the divided differences (4.46) and (4.49) which can be extended to the case where some or all of the points x_i are coincident, provided that f(x) is sufficiently differentiable. For example, define

$$f[x_0, x_0] = \lim_{\epsilon \to 0} f[x_0, x_0 + \epsilon] = \lim_{\epsilon \to 0} \frac{f(x_0 + \epsilon) - f(x_0)}{\epsilon} = f'(x_0).$$

$$f[x_0, x_0, x_0] = \lim_{\epsilon \to 0} f[x_0, x_0, x_0 + \epsilon] = \lim_{\epsilon \to 0} \frac{f(x_0, x_0 + \epsilon) - f[x_0, x_0]}{\epsilon}$$

$$= \lim_{\epsilon \to 0} \frac{\frac{f(x_0 + \epsilon) - f(x_0)}{\epsilon} - f'(x_0)}{\epsilon}$$

$$= \lim_{\epsilon \to 0} \frac{f(x_0 + \epsilon) - f(x_0) - \epsilon f'(x_0)}{\epsilon^2} \qquad \left(\frac{0}{0} \text{ form}\right)$$

$$= \lim_{\epsilon \to 0} \frac{f'(x_0 + \epsilon) - f'(x_0)}{2\epsilon} \qquad (\text{using } L'H\hat{o}\text{pital's rule})$$

$$= \frac{1}{2} \left[\lim_{\epsilon \to 0} \frac{f'(x_0 + \epsilon) - f'(x_0)}{\epsilon}\right] = \frac{f''(x_0)}{2}.$$

For an arbitrary, $n \ge 1$, let all the points in the Theorem 4.6 approach x_0 . This leads to the definition

$$f[x_0, x_0, \dots, x_0] = \frac{f^{(n)}(x_0)}{n!},$$

where the left hand side denotes the *nth* divided difference, all of whose points are x_0 .

Example 4.38 Let $f(x) = e^{-x}$ and let $x_0 = 0, x_1 = 1$. Using (4.56) and the above divide difference property 9, calculate $f[x_0, x_1, x_0]$, $f[x_0, x_0, x_1, x_1]$ and $f[x_0, x_1, x_1, x_1]$.

Solution. By using (4.56), we have

$$f[x_0, x_0] = \frac{1}{1!}f'(x_0) = f'(x_0).$$

Therefore

$$f[x_0, x_1, x_0] = f[x_0, x_0, x_1] = \frac{f[x_0, x_1] - f[x_0, x_0]}{x_1 - x_0} = \frac{f[x_0, x_1] - f'(x_0)}{x_1 - x_0}.$$

Using definition of the first-order divided difference of f(x) at points x_0 and x_1 , we have

$$f[x_0, x_1] = \frac{f[x_1] - f[x_0]}{x_1 - x_0}, \quad gives \quad f[0, 1] = \frac{0.368 - 1}{1 - 0} = -0.632.$$

Also

$$f'(x_0) = -e^{-x_0}$$
 and $f'(0) = -1$.

Using these values, we obtain the value of the second divided difference as

$$f[0,1,0] = \frac{-0.632 + 1}{1 - 0} = 0.368.$$

Now to find the value of the third divided difference which is defined as

$$f[x_0, x_0, x_1, x_1] = \frac{f[x_0, x_1, x_1] - f[x_0, x_0, x_1]}{x_1 - x_0},$$

and after simplifying, we have

$$f[x_0, x_0, x_1, x_1] = \frac{f'(x_1) - 2f[x_0, x_1] + f'(x_0)}{(x_1 - x_0)^2}.$$

Thus

$$f[0,0,1,1] = \frac{-0.368 - 2(-0.632) - 1}{(1-0)^2} = -0.014.$$

Finally, the other third divided difference is defined as

$$f[x_0, x_1, x_1, x_1] = \frac{f[x_1, x_1, x_1] - f[x_0, x_1, x_1]}{x_1 - x_0},$$

or

$$f[x_0, x_1, x_1, x_1] = \frac{f''(x_1)/2! - (f'(x_1) - f[x_0, x_1])/(x_1 - x_0)}{x_1 - x_0}$$

After simplifying, we have

$$f[x_0, x_0, x_1, x_1] = \frac{(x_1 - x_0)f''(x_1) - 2f'(x_1) + 2f[x_0, x_1]}{2(x_1 - x_0)^2}.$$

By using the value of $f''(1) = e^{-1} = 0.368$, we have

$$f[0, 1, 1, 1] = \frac{(1-0)(0.368) - 2(-1) + 2(-0.632)}{2(1-0)^2} = 1.104,$$

the required value.

Example 4.39 *Let* $f(x) = \ln(x+2)$.

- (a) Compute the divided differences f[0,0,1], f[0,1,1] and f[0,0,1,1].
- (b) Compute the approximation of ln(2.5) by using cubic Newton's interpolating polynomial.

Solution. (a) Using $f(x) = \ln(x+2)$ and $x_0 = 0, x_1 = 1$, we find the third-order divided difference f[0,0,1] as follows:

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

and so using the given values, we get

$$f[0,0,1] = \frac{f[0,1] - f'(0)}{1 - 0} = f(1) - f(0) - f'(0) = 1.0986 - 0.6932 - 0.5 = -0.0946.$$

Similarly, we obtain

$$f[0,1,1] = \frac{f[1,1] - f[0,1]}{1 - 0} = f'(1) - f(1) + f(0) = 0.3333 - 1.0986 + 0.6932 = -0.0721,$$

and

$$f[0,0,1,1] = \frac{f[0,1,1] - f[0,0,1]}{1-0} = -0.0721 + 0.0946 = 0.0225.$$

(b) The cubic Newton's interpolating polynomial has the following form

$$p_3(x) = f[x_0] + (x - x_0)f[x_0, x_0] + (x - x_0)(x - x_0)f[x_0, x_0, x_1] + (x - x_0)(x - x_0)(x - x_1)f[x_0, x_0, x_1, x_1],$$

so using values of part (a) and $x = 0.5$, we get

$$\ln(2.5) \approx p_3(0.5) = f(0) + (0.5 - 0)f'(0) + (0.5 - 0)(0.5 - 0)f[0, 0, 1] + (0.5 - 0)(0.5 - 0)(0.5 - 1)f[0, 0, 1, 1],$$
or

$$\ln(2.5) \approx p_3(0.5) = 0.6932 + 0.25 - 0.0237 + 0.0090 = 0.9286,$$

the required approximation of $\ln(3.5)$ and

$$|\ln(2.5) - p_3(0.5)| = |0.9163 - 0.9285| = 0.0122,$$

the possible absolute error in the approximation.

Example 4.40 Let $f(x) = x^3$ and $\alpha \neq 1$, then find the value of α such that $f[\alpha, 1, 1] = 1$.

Solution. Using $f(x) = x^3$ and $x_0 = \alpha, x_1 = 1$, we find the third-order divided difference $f[\alpha, 1, 1]$ as follows:

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

and so using the given values, we get

$$1 = f[\alpha, 1, 1] = \frac{f'(1) - f[\alpha, 1]}{1 - \alpha}, \quad gives \quad 1 - \alpha = 3 - f[\alpha, 1].$$

From this we get

$$2 + \alpha = \frac{f[1] - f[\alpha]}{1 - \alpha} = \frac{1 - \alpha^3}{1 - \alpha} = (\alpha^2 + \alpha + 1),$$

and after simplifying, we get $\alpha^2 = 1$, which means $\alpha = -1$ (because $\alpha \neq 1$) is the required value. \bullet

Example 4.41 If $f(x) = \frac{2}{x}$, find the third-order divided difference f[a, b, c, d] of the function f(x) and then show that the third divided difference f[1, 1, 1, 2] = -1.

Solution. Since we know that

$$f[a,b] = \frac{f(b) - f(a)}{b - a} = \frac{\frac{2}{b} - \frac{2}{a}}{b - a} = \frac{2(a - b)}{ab(b - a)} = -\frac{2}{ab},$$

and

$$f[a,b,c] = \frac{f[b,c] - f[a,b]}{c-a} = \frac{-\frac{2}{bc} + \frac{2}{ab}}{c-a} = \frac{2(c-a)}{abc(c-a)} = -\frac{2}{abc}.$$

Thus

$$f[a,b,c,d] = \frac{f[b,c,d] - f[a,b,c]}{d-a} = \frac{\frac{2}{bcd} - \frac{2}{abc}}{d-a} = -\frac{2}{abcd},$$

is the required third-order divided difference f[a, b, c, d] of the given function. using this result, can find the value of f[1, 1, 1, 2] as follows:

$$f[1, 1, 1, 2] = -\frac{2}{(1)(1)(1)(2)} = -1,$$

is the value of the third-order divided difference f[1,1,1,2] of the given function. We can also reproduced this value of f[1,1,1,2] by using divide difference property 6 as follows:

$$f[1,1,1,2] = \frac{f[1,1,2] - f[1,1,1]}{2-1} = f[1,1,2] - f[1,1,1]$$
$$= \frac{f[1,2] - f[1,1]}{2-1} - \frac{f''(1)}{2!}$$

$$= \frac{f(2) - f(1)}{2 - 1} - \frac{f'(1)}{1!} - \frac{f''(1)}{2!}$$

$$= f(2) - f(1) - f'(1) - \frac{f''(1)}{2}.$$

$$ave, f'(x) = -\frac{2}{x^2} \text{ and } f''(x) = \frac{4}{x^3}. \text{ Thus}$$

Since
$$f(x) = \frac{2}{x}$$
, so we have, $f'(x) = -\frac{2}{x^2}$ and $f''(x) = \frac{4}{x^3}$. Thus

$$f[1, 1, 1, 2] = f(2) - f(1) - f'(1) - \frac{f''(1)}{2} = 1 - 2 + 2 - 2 = -1,$$

is the required value.

4.3 Interpolation with Spline Functions

In the previous sections we studied the use of interpolation polynomials for approximating the values of the functions on closed intervals. An alternative approach is divide the interval into a collection of subintervals and construct a different approximating polynomial on each subinterval. Approximation by polynomial of this type is called *piecewise polynomial approximation*. Here, we will discuss some of the examples of a piecewise curve fitting techniques; the use of the piecewise linear interpolation.

Definition 4.1 (Spline Function)

Let $a = x_0 < x_1 < x_2 \cdots < x_n = b$. A function $s : [a, b] \to \mathbf{R}$ is a spline or spline function of degree m with points x_0, x_1, \ldots, x_n if:

1. A function s is a piecewise polynomial such that, on each subinterval $[x_k, x_{k+1}]$, s has degree at most m.

2. A function s is m-1 times differentiable everywhere.

A spline is a flexible drafting device that can be constrained to pass smoothly through a set of plotted data points. Spline functions are a mathematical tool which is an adaptation of this idea.

4.3.1 Piecewise Linear Interpolation

It is the one of the simplest piecewise polynomial interpolation for the approximation of the function, called *linear spline*. The linear spline is continuous function and the basic of it is simply connect consecutive points with straight lines. Consider the set of seven data points (x_0, y_0) , (x_1, y_1) , (x_2, y_2) , $(x_3, y_3), (x_4, y_4), (x_5, y_5)$ and (x_6, y_6) which define six subintervals. These intervals are denoted as $[x_0, x_1], [x_1, x_2], [x_2, x_3], [x_3, x_4], [x_4, x_5]$ and $[x_5, x_6],$ where $x_0, x_1, x_2, x_3, x_4, x_5,$ and x_6 are distinct x-values. If we use a straight line on each subinterval (see Figure 4.4) then we can interpolate the data with a piecewise linear function, where

$$s_k(x) = p_k(x) = \frac{(x - x_{k+1})}{(x_k - x_{k+1})} y_k + \frac{(x - x_k)}{(x_{k+1} - x_k)} y_{k+1},$$

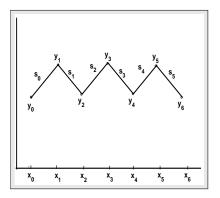


Figure 4.4: Linear spline.

or

$$s_k(x) = y_k + \frac{(y_{k+1} - y_k)}{(x_{k+1} - x_k)}(x - x_k).$$

It gives us

$$s_k(x) = A_k + B_k(x - x_k), (4.59)$$

where the values of the coefficients A_k and B_k are given as

$$A_k = y_k$$
 and $B_k = \frac{(y_{k+1} - y_k)}{(x_{k+1} - x_k)}$. (4.60)

Note that the linear spline must be continuous at given points x_0, x_1, \ldots, x_n and

$$s(x_k) = f(x_k) = y_k,$$
 for $k = 0, 1, ..., n$.

Example 4.42 Find the values of unknown coefficients a and b so that the following function is a linear spline.

$$s(x) = \begin{cases} a - x, & 0 \le x \le 1, \\ 3x - b, & 1 \le x \le 2, \\ 2x + 1, & 2 \le x \le 3. \end{cases}$$

Solution. Since the given function is a linear spline, so s must be continuous at the internal points 1 and 2. Continuity at x = 1 implies that

$$\lim_{x \to 1^{-}} s(x) = \lim_{x \to 1^{+}} s(x),$$

$$\lim_{x \to 1^{-}} a - x = \lim_{x \to 1^{+}} 3x - b,$$

$$a - 1 = 3 - b,$$

and it gives an equation of the from

$$a + b = 4$$
.

Now continuity at x = 2 implies that

$$\lim_{x \to 2^{-}} s(x) = \lim_{x \to 2^{+}} s(x),$$

$$\lim_{x \to 2^{-}} 3x - b = \lim_{x \to 2^{+}} 2x + 1,$$

$$6 - b = 5$$

and it gives b = 1. Using this value of b, we get a = 3, and so

$$s(x) = \begin{cases} 3 - x, & 0 \le x \le 1, \\ 3x - 1, & 1 \le x \le 2, \\ 2x + 1, & 2 \le x \le 3, \end{cases}$$

is the linear spline function.

Example 4.43 Find the linear splines which interpolates the following data

Find the approximation of the function $y(x) = \frac{2}{x+1}$ at x = 2.9. Compute absolute error.

Solution. Given $x_0 = 1.0, x_1 = 2.0, x_2 = 3.0, x_3 = 4.0$, then using (4.60), we have

$$A_0 = y_0 = 1.0, \ A_1 = y_1 = 0.67, \ A_2 = y_2 = 0.50, \ A_3 = y_3 = 0.4,$$

and

$$B_0 = \frac{(y_1 - y_0)}{(x_1 - x_0)} = \frac{(0.67 - 1.0)}{(2.0 - 1.0)} = -0.33,$$

$$B_1 = \frac{(y_2 - y_1)}{(x_2 - x_1)} = \frac{(0.50 - 0.67)}{(3.0 - 2.0)} = -0.17,$$

$$B_2 = \frac{(y_3 - y_2)}{(x_3 - x_2)} = \frac{(0.40 - 0.50)}{(4.0 - 3.0)} = -0.10.$$

Now using (4.59), the linear splines for three subintervals are define as

$$s(x) = \begin{cases} s_0(x) &= 1.0 - 0.33(x - 1.0) = 1.33 - 0.33x, & 1 \le x \le 2, \\ s_1(x) &= 0.67 - 0.17(x - 2.0) = 1.01 - 0.17x, & 2 \le x \le 3, \\ s_2(x) &= 0.50 - 0.10(x - 3.0) = 0.80 - 0.10x, & 3 \le x \le 4. \end{cases}$$

The value x = 2.9 lies in the interval [2, 3], so

$$f(2.9) \approx s_1(2.9) = 1.01 - 0.17(2.9) = 0.517.$$

Also,

$$|f(2.9) - s_1(2.9)| = |0.513 - 0.517| = 0.004,$$

the required absolute error.

Using MATLAB command window, we can reproduce above results as follows:

$$>> X = [1\ 2\ 3\ 4];\ Y = [1\ 0.67\ 0.50\ 0.40];$$

 $>> x = 2.9; s = LSpline(X, Y, x);$

236 4.7 Exercises

Program 4.4

MATLAB m-file for the Linear Spline Functions

function LS=LSpline(X,Y,x)

n = length(X); for i = n-1:-1:1

D = x - X(i); if (D >= 0); break; end; end

D = x - X(i); if (D < 0); i = 0; D = x - X(1); end

M = (Y(i+1) - Y(i))/(X(i+1) - X(i)); LS = Y(i) + M * D; end

4.4 Exercises

- 1. Use the Lagrange interpolation formula based on the points $x_0 = 0, x_1 = 1, x_2 = 2.5$ to find the equation of the quadratic polynomial to approximate $f(x) = \frac{2}{x+2}$ at x = 2.3.
- 2. Let $f(x) = \cos(x\pi/4)$, where x is in radian. Use the quadratic Lagrange interpolation formula based on the points $x_0 = 0, x_1 = 1, x_2 = 2$ and $x_3 = 4$ to find the polynomial $p_2(x)$ to approximate the function f(x) at x = 0.5 and x = 3.5.
- 3. Let $f(x) = x + 2\ln(x+2)$. Use the quadratic Lagrange interpolation formula based on the points $x_0 = 0, x_1 = 1, x_2 = 2$ and $x_3 = 3$ to approximate f(0.5) and f(2.8). Also, compute the error bounds for your approximations.
- 4. Consider the function $f(x) = e^{x^2}$ and x = 0, 0.25, 0.5, 1. Then use the suitable Lagrange interpolating polynomial to approximate f(0.75). Also, compute an error bound for your approximation.
- 5. Let $f(x) = x^4 2x + 1$. Use cubic Lagrange interpolation formula based on the points $x_0 = -1, x_1 = 0, x_2 = 2$ and $x_3 = 3$ to find the polynomial $p_3(x)$ to approximate the function f(x) at x = 1.1. Also, compute an error bound for your approximation.
- 6. Construct the Lagrange interpolation polynomials for the following functions and compute the error bounds for the approximations:

(a)
$$f(x) = x + 2^{x+1}$$
, $x_0 = 0, x_1 = 1, x_2 = 2.5, x_3 = 3.5$

(b)
$$f(x) = 3x^3 + 2x^2 + 1$$
, $x_0 = 1, x_1 = 2, x_2 = 3$

$$\begin{array}{ll} (a) & f(x)=x+2^{x+1}, & x_0=0, x_1=1, x_2=2.5, x_3=3. \\ (b) & f(x)=3x^3+2x^2+1, & x_0=1, x_1=2, x_2=3. \\ (c) & f(x)=\cos x-\sin x, & x_0=0, x_1=0.25, x_2=0.5, x_3=1. \end{array}$$

7. Consider the following table:

- (a) Construct divided difference table for the tabulated function.
- (b) Compute the Newton interpolating polynomials $p_2(x)$ and $p_3(x)$ at x=2.2.

8. Consider the following table:

- (a) Construct divided difference table for the tabulated function.
- (b) Compute the Newton interpolating polynomials $p_3(x)$ and $p_4(x)$ at x = 2.5, 3.5.
- 9. Consider the following table of the $f(x) = \sqrt{x}$:

- (a) Construct the divided difference table for the tabulated function.
- (b) Find the Newton interpolating polynomials $p_3(x)$ and $p_4(x)$ at x=5.9.
- (c) Compute error bounds for your approximations in part (b).
- 10. Let $f(x) = e^x \sin x$, with $x_0 = 0, x_1 = 2, x_2 = 2.5, x_3 = 4, x_4 = 4.5$. Then
 - (a) Construct the divided-difference table for the given data points.
 - (b) Find the Newton divided difference polynomials $p_2(x)$, $p_3(x)$ and $p_4(x)$ at x=2.4.
 - (c) Compute error bounds for your approximations in part (b).
 - (d) Compute the actual error.
- 11. Show that if x_0, x_1 and x_2 are distinct then

$$f[x_0, x_1, x_2] = f[x_1, x_2, x_0] = f[x_2, x_0, x_1]$$

12. The divided difference form of the interpolating polynomial $p_3(x)$ is

$$p_3(x) = f[x_0] + (x - x_0)f[x_0, x_1] + (x - x_0)(x - x_1)f[x_0, x_1, x_2, x_0] + (x - x_0)(x - x_1)(x - x_2)f[x_0, x_1, x_2, x_3]$$

By expressing these divided differences in terms of the function values $f(x_i)(i = 0, 1, 2, 3)$, verify that $p_3(x)$ does pass through the points $(x_i, f(x_i))(i = 0, 1, 2, 3)$.

- 13. Let $f(x) = x^2 + e^x$ and $x_0 = 0, x_1 = 1$. Use the divided differences to find the value of the second divided difference $f[x_0, x_1, x_0]$.
- 14. Which of the following functions are linear splines?

(a)
$$s(x) = \begin{cases} x, & 0 \le x \le 1\\ 2x - 1, & 1 \le x \le 2\\ x + 2, & 2 \le x \le 4 \end{cases}$$

(b)
$$s(x) = \begin{cases} 2-x, & 0 \le x \le 1\\ 2x-1, & 1 \le x \le 2\\ x+1, & 2 \le x \le 4 \end{cases}$$

238 4.7 Exercises

15. Find the linear spline which interpolates the data:

What are its values at x = 0.55, 1.15 and 2.5 ?

16. Find the linear spline which interpolates the data:

$$(0,0), (0.2,0.18), (0.3,0.26), (0.5,0.41)$$

What are its values at x = 0.15, 0.25, and 0.45?

17. Find the linear splines which interpolate the following data:

Compare interpolated values at $x = 0.5, 11.5, \text{ and } 30.5 \text{ to } f(x) = \sqrt[4]{x}.$

18. Find the linear splines which interpolate the following data:

$$(0,1), (2,0.9976), (3,0.9945), (4,0.9903)$$

Compare interpolated values at $x = 1.5, 2.5, \text{ and } 3.5 \text{ to } f(x) = \cos(2x).$

19. Find the linear splines which interpolate the following data:

Compare interpolated values at $x = 2.5, 5.5, \text{ and } 10.5 \text{ to } f(x) = \sqrt{x+1}.$

Chapter 5

Numerical Differentiation and Integration

5.1 Introduction

In this chapter we deal with techniques for approximating numerically the two fundamental operations of the calculus, differentiation and integration. Both of these problems may be approached in the same way. Although both numerical differentiation and numerical integration formulas will be discussed, it should be noted that numerical differentiation is inherently much less accurate than numerical integration, and its application is generally avoided whenever possible. Nevertheless, it has been used successfully in certain applications.

Engineers are frequently confronted with the problem of differentiating functions which are defined in tabular or graphical form rather than as explicit functions. The interpretation of experimentally obtained data is a good example of this. A similar situation involves the integration of functions which have explicit forms that are difficult or impossible to integrate in terms of elementary functions. Graphical techniques, employing the construction of tangents to curves and the estimation of areas under curves, are commonly used in solving such problems, when great accuracy is not a prerequisite for the results. However, there are occasions when a higher degree of accuracy is desired, and, for these, various numerical methods are available.

Important Points of the Chapter 5

- I. In this chapter we shall find the approximate solutions of derivative (first- and second-order) and antiderivative (definite integral only).
- II. Given data points should be equally spaced only (length of each subinterval should be same). Smaller the length of the interval better the approximation.
- III. Numerical methods for differentiation and integration can be derived using Lagrange interpolating polynomial at equally-spaced data points.
- IV. Error term for each numerical method will be discuss which helps to look for the maximum 239

240 5.2 Numerical Differentiation

error in the approximation.

V. Two-point formula (for first derivative) and three-point formulas (for first and second derivatives) for numerical differentiation and Trapezoidal and Simpson's rules for numerical integration will be discuss here.

5.2 Numerical Differentiation

Firstly, we discuss the numerical process for approximating the derivative of the function f(x) at the given point. A function f(x), known either explicitly or as a set of data points, is replaced by a simpler function. A polynomial p(x) is the obvious choice of approximating function, since the operation of differentiation is then easily performed. The polynomial p(x) is differentiated to obtain p'(x), which is taken as an approximation to f'(x) for any numerical value of x. Geometrically, this is equivalent to replacing the slope of f(x), at x, by that of p(x). Here, numerical differentiation are derived by differentiating interpolating polynomials.

We now turn our attention to the numerical process for approximating the derivative of a function f(x) at x, that is

$$f'(x) = \lim_{h \to 0} \frac{f(x+h) - f(x)}{h}$$
, provided the limit exits. (5.1)

In principle, it is always possible to determine an analytic form (5.1) of a derivative for a given function. In some cases, however, the analytic form is very complicated, and a numerical approximation of the derivative may be sufficient for our purpose.

The formula (5.1) provides an obvious way to get an approximation to f'(x); simply compute

$$D_h f(x) = \frac{f(x+h) - f(x)}{h},$$
 (5.2)

for small values of stepsize h, called numerical differentiation formula for (5.1).

Numerical differentiation is a much less satisfactory process because the seemingly obvious approximations are not always a good as they seem. Therefore, this process should, for this reason, be avoided if at all possible. We study it mainly as a means to end of solving differential equations by various numerical methods, based on the approximations we shall obtain for the derivatives of a function. We study it also, because it often happens that the thing we want to differentiate is not known function. We may, for instance, be given a table of speeds of a body observed at certain times, and wish to estimate its acceleration at these time.

Numerical differentiation is useful in estimating the derivative of a function when either function f(x) is difficult to differentiate easily, or it is not known as explicit expression in x but the values of the function are described only in terms of tabulated data. Generally, it is considered that numerical differentiation is basically an *unstable process* which means that small errors made in the initial computations may cause greatly magnified errors in the final result. In fact, we may not always expect reasonable results even when the original data are known to be more accurate.

Here, we shall derive some formulas for estimating derivatives but we should avoid as far as possible, numerically calculating derivatives higher than the first, as the error in their evaluation increases

with their orders. In spite of some inherent shortcomings, numerical differentiation is important to derive formulas for solving integrals and the numerical solution of both ordinary and partial differential equations.

There are three different approaches for deriving the numerical differentiation formulas. The first approach is based on the Taylor expansion of a function about a point, the second is to use difference operators, and the third approach to numerical differentiation is to fit a curve with a simple form to a function, and then to differentiate the curve-fit function. For example, the polynomial interpolation or spline methods of the Chapter 4 can be used to fit a curve to tabulated data for a function and the resulting polynomial or spline can then be differentiated. When a function is represented by a table of values, the most obvious approach is to differentiate the Lagrange interpolation formula

$$f(x) = p_n(x) + \frac{f^{(n+1)}(\eta(x))}{(n+1)!} \prod_{i=0}^{n} (x - x_i),$$
 (5.3)

where the first term $p_n(x)$ of the right hand side is the Lagrange interpolating polynomial of degree n and the second term is its error term.

It is interesting to note that the process of numerical differentiation may be less satisfactory than interpolation the closeness of the ordinates of f(x) and $p_n(x)$ on the interval of interest does not guarantee the closeness of their respective derivatives. Note that the derivation and analysis of formulas for numerical differentiation is considerably simplifies when the data is equally spaced. It will be assumed, therefore, that the points x_i are given by $x_i = x_0 + ih$, (i = 0, 1, ..., n) for some fixed tabular interval h.

5.3 Numerical Differentiation Formulas

Here, we will find the approximation of first and second derivative of a function at a given arbitrary point x. For the approximation of the first derivative of a function we will use two-point formula, three-point formula, and Richardson's extrapolation formula. While for second derivative approximation we will discuss three-point formula only.

5.3.1 First Derivative Numerical Formulas

To obtain general formula for approximation of the first derivative of a function f(x), we consider that $\{x_0, x_1, \ldots, x_n\}$ are (n+1) distinct equally spaced points in some interval I and function f(x) is continuous and its (n+1)th derivatives exist in the given interval, that is, $f \in C^{n+1}(I)$. Then by differentiating (5.3) with respect to x and at $x = x_k$, we have

$$f'(x_k) = \sum_{i=0}^{n} f(x_i) L_i'(x_k) + \frac{f^{(n+1)}(\eta(x_k))}{(n+1)!} \prod_{\substack{i=0\\i\neq k}}^{n} (x_k - x_i).$$
 (5.4)

The formula (5.4) is called the (n+1)-point formula to approximate $f'(x_k)$. From this formula we can obtain many numerical differentiation formulas but here we shall discuss only three formulas to approximate (5.1) at given point $x = x_k$. First one is called the two-point formula which we can get from (5.4) by taking n = 1 and k = 0. The second numerical differentiation formula is called the three-point formula which can be obtained from (5.4) when n = 2 and k = 0, 1, 2. Finally, we will discuss the five-point formula to approximate (5.1) by using (5.4) when n = 4 and k = 0, 1, 2, 3, 4.

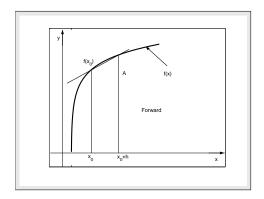


Figure 5.1: Forward-difference approximations.

Two-point Formula

Consider two distinct points x_0 and x_1 , then, to find the approximation of (5.1), the first derivative of a function at given point, take $x_0 \in (a, b)$, where $f \in C^2[a, b]$ and that $x_1 = x_0 + h$ for some $h \neq 0$ that is sufficiently small to ensure that $x_1 \in [a, b]$. Consider the linear Lagrange interpolating polynomial $p_1(x)$ which interpolate f(x) at the given points is

$$f(x) = p_1(x) = \left(\frac{x - x_1}{x_0 - x_1}\right) f(x_0) + \left(\frac{x - x_0}{x_1 - x_0}\right) f(x_1). \tag{5.5}$$

By taking derivative of (5.5) with respect to x and at $x = x_0$, we obtain

$$f'(x)|_{x=x_0} \approx p'_1(x)|_{x=x_0} = -\frac{f(x_0)}{x_0 - x_1} + \frac{f(x_1)}{x_1 - x_0}.$$

Simplifying the above expression, we have

$$f'(x_0) \approx -\frac{f(x_0)}{h} + \frac{f(x_0 + h)}{h},$$

which can be written as

$$f'(x_0) \approx \frac{f(x_0 + h) - f(x_0)}{h} = D_h f(x_0).$$
 (5.6)

It is called the two-point formula for smaller values of h. For h > 0, sometime the formula (5.6) is also called the two-point forward-difference formula because it involves only differences of a function values forward from $f(x_0)$. The two-point forward-difference formula has a simple geometric interpretation as the slope of the forward secant line, as shown in Figure 5.1.

If h < 0, then the formula (5.6) is also called the two-point backward-difference formula, which can be written as

$$f'(x_0) \approx \frac{f(x_0) - f(x_0 - h)}{h}.$$
 (5.7)

In this case, a value of x behind the point of interest is used. The formula (5.7) is useful in cases where the independent variable represents time. If x_0 denotes the present time, the backward-difference formula uses only present and past samples, it does not rely on future data samples that may not yet be available in a real time application.

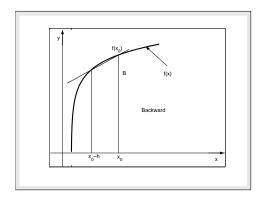


Figure 5.2: Backward-difference approximations.

The geometric interpretation of the two-point backward-difference formula, as the slope of the backward secant line, is shown in Figure 5.2.

Example 5.1 Let $f(x) = e^x$ and h = 0.1, h = 0.01. Use two-point forward difference formula to approximate f'(2). For which value of h we have better approximation and why?

Solution. Using the formula (5.6), with $x_0 = 2$, we have

$$f'(2) \approx \frac{f(2+h) - f(2)}{h}.$$

Then for h = 0.1, we get

$$f'(2) \approx \frac{f(2.1) - f(2)}{0.1} \approx \frac{e^{2.1} - e^2}{0.1} = 7.7712.$$

Similarly, by using h = 0.01, we obtain

$$f'(2) \approx \frac{(e^{2.01} - e^2)}{0.01} = 7.4262.$$

Since the exact solution of $f'(2) = e^2$ is, 7.3891, so the corresponding actual errors with h = 0.1 and h = 0.01 are, -0.3821 and -0.0371 respectively. This shows that the approximation obtained with h = 0.01 is better than the approximation with h = 0.1.

The above results can be easily achieved with MATLAB commands as follows:

$$>> x0 = 2.0; h = 10.^-(1:2); df = (exp(x0+h) - exp(x))./h;$$

Similarly, by using the formula (5.7), with $x_0 = 2$, we have

$$f'(2) \approx \frac{f(2) - f(2-h)}{h},$$

then for h = 0.1, we have

$$f'(2) \approx \frac{f(2) - f(1.9)}{0.1} = \frac{e^2 - e^{1.9}}{0.1} = 7.0316.$$

For h = 0.01, we have

$$f'(2) \approx \frac{e^2 - e^{1.99}}{0.01} = 7.3522.$$

The corresponding actual errors with h=0.1 and h=0.01 are, 0.3575 and 0.0369 respectively, which shows that the approximation with h=0.01 is better than the approximation with h=0.1. Note that the both errors for h=0.1 and h=0.01 by using the backward-difference formula is better than the forward-difference formula for the same values of h. The MATLAB symbolic toolbox contains commands for obtaining the symbolic derivative of symbolically written function which are as follows:

>>
$$syms \ x;$$

>> $f = x.^2 \cdot cos(x); \ df = diff(f, 1); \ subs(ddf, 1);$

Error Term of Two-point Formula

The formula (5.6) is not very useful, therefore, let us attempt to find the error involves in our first numerical differentiation formula (5.6). Consider the error term for the linear Lagrange polynomial which can be written as

$$f(x) - p_1(x) = \frac{f''(\eta(x))}{2!} \prod_{i=0}^{1} (x - x_i),$$

for some unknown point $\eta(x) \in (x_0, x_1)$. By taking derivative of above equation with respect to x and at $x = x_0$, we have

$$f'(x_0) - p'_1(x_0) = \left(\frac{d}{dx}f''(\eta(x))\Big|_{x=x_0}\right) \frac{(x-x_0)(x-x_1)}{2} + \frac{f''(\eta(x_0))}{2} \left(\frac{d}{dx}(x^2 - x(x_0+h) - xx_0 + x_0(x_0+h))\Big|_{x=x_0}\right).$$

Since $\frac{d}{dx}f''(\eta(x)) = 0$ only if $x = x_0$, so error in the forward-difference formula (5.6) is

$$E_F(f,h) = f'(x_0) - D_h f(x_0) = -\frac{h}{2} f''(\eta(x)), \text{ where } \eta(x) \in (x_0, x_1),$$
 (5.8)

which is called the *error formula* of the two-point formula (5.6). Hence the formula (5.6) can be written as

$$f'(x_0) = \frac{f(x_0 + h) - f(x_0)}{h} - \frac{h}{2}f''(\eta(x)), \quad \text{where} \quad \eta \in (x_0, x_1).$$
 (5.9)

The formula (5.9) is more useful than the formula (5.6) because now on a large class of function, an error term is available along with the basic numerical formula. Note that the error term in (5.9) has two parts; a power of h and a factor involving some higher-order derivative of f(x) which

gives us an indication of the class of function to which the error estimate is applicable. The h term in the error makes the entire expression converge to zero as h approaches zero. The rapidity of this convergence will depend on the power of h. These remarks apply to many error estimates in numerical analysis. There will usually be a power of h and a factor telling us to what smoothness class of the function must belong so that the estimate is valid.

Note that the formula (5.9) may also be derived from the Taylor's theorem. Expansion of function $f(x_1)$ about x_0 as far as term involving h^2 gives

$$f(x_1) = f(x_0) + hf'(x_0) + \frac{h^2}{2!}f''(\eta(x)). \tag{5.10}$$

From this the result follows by subtracting $f(x_0)$ both sides and dividing both sides by h and put $x_1 = x_0 + h$.

Note that for a linear function, f(x) = ax + b, the approximate formula (5.6) is exact; that is, it yields the correct value of first derivative of the function f(x) for any nonzero value of h.

Example 5.2 Let $f(x) = x^3$ be defined in the interval [0.2, 0.3]. Use the error formula (5.8) of two-point formula for the approximation of f'(0.2) to compute a value of η .

Solution. Since the exact value of the first derivative of the function at $x_0 = 0.2$ is

$$f'(x) = 3x^2$$
 and $f'(0.2) = 3(0.2)^2 = 0.12$,

and the approximate value of f'(0.2) using two point formula is

$$f'(0.2) \approx \frac{f(0.3) - f(0.2)}{0.1} = \frac{(0.3)^3 - (0.2)^3}{0.1} = 0.19,$$

so error E can be calculated as

$$E = 0.12 - 0.19 = -0.07.$$

Using the formula (5.6) and $f''(\eta) = 6\eta$, we have

$$-0.07 = -\frac{0.1}{2}6\eta,$$

and solving for η , we get $\eta = 0.233$.

Example 5.3 Let $f(x) = x^2 \cos x$ and h = 0.1. Then

- (a) Compute the approximate value of f'(1) using forward difference two-point formula (5.6).
- (b) Compute the error bound for your approximation using the formula (5.8).
- (c) Compute the absolute error.
- (d) What best maximum value of stepsize h required to obtain the approximate value of f'(1) correct to 10^{-2} .

Solution. (a) Given $x_0 = 1, h = 0.1$, then by using the formula (5.6), we have

$$f'(1) \approx \frac{f(1+0.1) - f(1)}{0.1} = \frac{f(1.1) - f(1)}{0.1} = D_h f(1).$$

Thus

$$f'(1) \approx \frac{(1.1)^2 \cos(1.1) - (1)^2 \cos(1)}{0.1} \approx \frac{0.5489 - 0.5403}{0.1} = 0.0860,$$

which is the required approximation of f'(x) at x = 1.

(b) To find the error bound, we use the formula (5.8), which gives

$$E_F(f,h) = -\frac{0.1}{2}f''(\eta(x)), \text{ where } \eta(x) \in (1,1.1),$$

or

$$|E_F(f,h)| = \left| -\frac{0.1}{2} \right| |f''(\eta(x))|, \text{ for } \eta \in (1,1.1).$$

The second derivative f''(x) of the function can be found as

$$f(x) = x^2 \cos x$$
, gives $f''(x) = (2 - x^2) \cos x - 4x \sin x$.

The value of the second derivative $f''(\eta(x))$ cannot be computed exactly because $\eta(x)$ is not known. But one can bound the error by computing the largest possible value for $|f''(\eta(x))|$. So bound |f''| on [1,1.1] can be obtain

$$M = \max_{1 \le x \le 1.1} |(2 - x^2)\cos x - 4x\sin x| = 3.5630,$$

at x = 1.1. Since $|f''(\eta(x))| \leq M$, therefore, for h = 0.1, we have

$$|E_F(f,h)| \le \frac{0.1}{2}M = 0.05(3.5630) = 0.1782,$$

which is the possible maximum error in our approximation.

The MATLAB symbolic toolbox contains commands for obtaining the symbolic derivative of symbolically written function which are as follows:

>>
$$syms\ x;$$

>> $f = x.^2 \cdot scos(x);\ df = diff(f, 1);\ ddf = diff(f, 2);\ subs(df, 1);$

(c) Since the exact value of the derivative f'(1) is 0.2392, therefore the absolute error |E| can be computed as follows:

$$|E| = |f'(1) - D_h f(1)| = |0.2391 - 0.0860| = 0.1531.$$

(d) Since the given accuracy required is 10^{-2} , so

$$|E_F(f,h)| = |-\frac{h}{2}f''(\eta(x))| \le 10^{-2},$$

for $\eta(x) \in (1, 1.1)$. This gives

$$\frac{h}{2}M \leq 10^{-2}, \quad or \quad h \leq \frac{(2\times 10^{-2})}{M}.$$

Using M = 3.5630, we obtain

$$h \le \frac{2}{356.3000} = 0.0056,$$

which is the best maximum value of h to get the required accuracy.

The truncation error in the approximation of (5.9) is roughly proportional to stepsize h used in its computation. The situation is made worse by the fact that the round-off error in computing the approximate derivative (5.6) is roughly proportion to $\frac{1}{h}$. The overall error therefore is of the form

$$E = ch + \frac{\delta}{h},$$

where c and δ are constants. This places serve restriction on the accuracy that can be achieved with this formula.

Now we discuss little more about the role of the round-off error in the numerical differentiation. Consider the formula (5.6) which is

$$f'(x_0) \approx \frac{f(x_0 + h) - f(x_0)}{h} = D_h f(x_0).$$

If h is small, then we can reasonably assume that $f(x_0)$ and $f(x_0 + h)$ have similar magnitude and, therefore, similar round-off errors. Let the actual function values used in the computation be denoted by \tilde{f}_0 and \tilde{f}_1 with

$$f(x_i) - \tilde{f}_i = \epsilon_i,$$
 for $i = 0, 1,$

the errors in the function values. Thus the actual quantity calculated is

$$\tilde{D}_h f(x_0) = \frac{\tilde{f}_1 - \tilde{f}_0}{h}.$$

For the error in this quantity, replace \tilde{f}_i by $f(x_i) - \epsilon_i$, for i = 0, 1, we obtain

$$f'(x_0) - \tilde{D}_h f(x_0) = f'(x_0) - \frac{(f(x_1) - \epsilon_1) - (f(x_0) - \epsilon_0)}{h}$$
$$= f'(x_0) - \frac{f(x_1) - f(x_0)}{h} + \frac{\epsilon_1 - \epsilon_0}{h}.$$

Then the overall error is given by

$$|f'(x_0) - \tilde{D}_h f(x_0)| \le |f'(x_0) - \frac{f(x_1) - f(x_0)}{h}| + \left|\frac{\epsilon_1 - \epsilon_0}{h}\right| \le \left|\frac{h}{2} f''(\eta(x_0))\right| + \frac{2\delta}{h},$$

the errors ϵ_0, ϵ_1 are generally random in some interval $[-\delta, \delta]$.

It is the second term on the right side of this error bound which leads to the growth of error as $h \to 0$. If

$$|f(x_i) - \tilde{f}_i| \le \frac{1}{2} \times 10^{-t} = \delta,$$

and t is the required decimal digits of accuracy, then the maximum rounding error in the two-point formula is $\frac{10^{-t}}{h}$. While the truncation error $\left|\frac{h}{2}f''(\eta(x_0))\right|$ decreases with h, the rounding error increases. The total error, E(h), therefore has a minimum with respect to h. If

$$E(h) = E_{trunc} + E_{round} = \frac{h}{2}M + \frac{10^{-t}}{h},$$

where $M = \max_{x_0 \le x \le x_1} |f''(\eta(x_0))|$, then

$$\frac{dE}{dh} = \frac{M}{2} - \frac{10^{-t}}{h^2}.$$

A minimum of E(h) satisfies the equation $\frac{dE}{dh} = 0$, that is

$$\frac{dE}{dh} = \frac{M}{2} - \frac{10^{-t}}{h^2} = 0.$$

solving for h, we obtain

$$h = h_{opt} = \sqrt{\frac{2}{M} \times 10^{-t}},$$

which gives the optimal value for h. Thus the minimum error is

$$E(h_{opt}) = \frac{M}{2} \sqrt{\frac{2}{M} \times 10^{-t}} + \frac{10^{-t}}{\sqrt{\frac{2}{M} \times 10^{-t}}} = \sqrt{2M \times 10^{-t}}.$$

Example 5.4 Consider $f(x) = x^2 \cos x$ and $x_0 = 1$. To show the effect of rounding error, the values \tilde{f}_i are obtained by rounding $f(x_i)$ to seven significant digits, compute the total error for h = 0.1 and also, find the optimum h.

Solution. Given $|\epsilon_i| \leq \frac{1}{2} \times 10^{-7} = \delta$ and h = 0.1. Now to calculate the total error, we use

$$E(h) = \frac{h}{2}M + \frac{10^{-t}}{h},$$

where

$$M = \max_{1 \le x \le 1.1} |(2 - x^2)\cos x - 4x\sin x| = 3.5630.$$

Then

$$E(h) = \frac{0.1}{2}(3.5630) + \frac{10^{-7}}{0.1} = 0.17815 + 0.000001 = 0.178151.$$

Now to find the optimum h, we use

$$h = h_{opt} = \sqrt{\frac{2}{M} \times 10^{-t}} = \sqrt{\frac{2}{3.5630} \times 10^{-7}} = 0.00024,$$

which is the smallest value of h, below which the total error will begin to increase.

Note that for

$$h = 0.00024,$$
 $E(h) = 0.000844,$
 $h = 0.00015,$ $E(h) = 0.000934,$
 $h = 0.00001,$ $E(h) = 0.010018.$

A similar effect is present for all numerical differentiation formulas.

The numerical differentiation formulas are often judge by the power of h in the error term. Since h is always small, so the higher power of h involve will give better approximation. In this assessment, the formula (5.9) acts poorly, as the error term involved h of power one. A superior formulas can be obtained by deriving some useful three-point formulas together with error terms which involved h^2 using the formula (5.4) by taking n=2.

Three-point Central Difference Formula

Consider the quadratic Lagrange interpolating polynomial $p_2(x)$ to the three distinct equally spaced points x_0, x_1 , and x_2 , with $x_1 = x_0 + h$ and $x_2 = x_0 + 2h$, for smaller value h, we have

$$f(x) = p_2(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} f(x_0) + \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} f(x_1) + \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} f(x_2).$$

Now taking the derivative of the above expression with respect to x and then take $x = x_k$, for k = 0, 1, 2, we have

$$f'(x_k) \approx \frac{(2x_k - x_1 - x_2)}{(x_0 - x_1)(x_0 - x_2)} f(x_0) + \frac{(2x_k - x_0 - x_2)}{(x_1 - x_0)(x_1 - x_2)} f(x_1) + \frac{(2x_k - x_0 - x_1)}{(x_2 - x_0)(x_2 - x_1)} f(x_2).$$
(5.11)

Three different numerical differentiation formulas can be obtained from (5.11) by putting $x_k = x_0$, or $x_k = x_1$ or $x_k = x_2$, which are use to find the approximation of the first derivative of a function defined by the formula (5.1) at the given point. Firstly, we take $x_k = x_1$, then the formula (5.11) becomes

$$f'(x_1) \approx \frac{(2x_1 - x_1 - x_2)}{(x_0 - x_1)(x_0 - x_2)} f(x_0) + \frac{(2x_1 - x_0 - x_2)}{(x_1 - x_0)(x_1 - x_2)} f(x_1) + \frac{(2x_1 - x_0 - x_1)}{(x_2 - x_0)(x_2 - x_1)} f(x_2).$$

After, simplifying, and replacing $x_0 = x_1 - h, x_2 = x_1 + h$, we obtain

$$f'(x_1) \approx \frac{f(x_1+h) - f(x_1-h)}{2h} = D_h f(x_1).$$
 (5.12)

It is called the *three-point central-difference formula* for finding the approximation of the first derivative of a function at the given point x_1 .

Note that the formulation of the formula (5.12) uses data points that are centered about the point of interest x_1 even though it does not appear in the right side of (5.12).

The geometric interpretation of the central-difference formula is shown in Figure 5.3.

Error Formula of Central Difference Formula

The formula (5.12) is not very useful, therefore, let us attempt to find the error involve in the formula (5.12) for numerical differentiation. Consider the error term for the quadratic Lagrange polynomial which can be written as

$$f(x) - p_2(x) = \frac{f'''(\eta(x))}{3!} \prod_{i=0}^{2} (x - x_i),$$

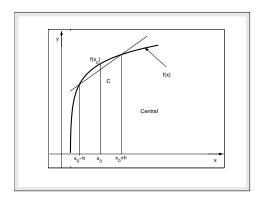


Figure 5.3: Central-difference approximations.

for some unknown point $\eta(x) \in (x_0, x_2)$. By taking derivative of the above equation with respect to x and then taking $x = x_1$, we have

$$f'(x_1) - p_2'(x_1) = \left(\frac{d}{dx}f'''(\eta(x))\Big|_{x=x_1}\right) \frac{(x-x_0)(x-x_1)(x-x_2)}{6}$$

$$+ \frac{f'''(\eta(x_1))}{6} \left((x-x_1)(x-x_2) + (x-x_0)(x-x_2) + (x-x_0)(x-x_1)\Big|_{x=x_1}\right).$$

Since $\frac{d}{dx}f'''(\eta(x)) = 0$ only if $x = x_1$, therefore the error formula of the central-difference formula (5.12) can be written as

$$E_C(f,h) = f'(x_1) - D_h f(x_1) = -\frac{h^2}{6} f'''(\eta(x_1)),$$
(5.13)

where $\eta(x_1) \in (x_1 - h, x_1 + h)$. Hence the formula (5.12) can be written as

$$f'(x_1) = \frac{f(x_1 + h) - f(x_1 - h)}{2h} - \frac{h^2}{6}f'''(\eta(x_1)), \tag{5.14}$$

where $\eta(x_1) \in (x_1 - h, x_1 + h)$. The formula (5.14) is more useful than the formula (5.12) because now on a large class of function, an error term is available along with the basic numerical formula. Note that for a quadratic function, $f(x) = ax^2 + bx + c$, the approximate formula (5.12) is exact; that is, it yields the correct value of first derivative of a function f(x) for any nonzero value of h.

Example 5.5 Let $f(x) = x^2 + \cos x$ and h = 0.1. Then

- (a) Compute the approximate value of f'(1) by using three-point central difference formula (5.12).
- (b) Compute the error bound for your approximation using (5.13).
- (c) Compute the absolute error.
- (d) What is the best maximum value of stepsize h required to obtain the approximate value of f'(1) correct to 10^{-2} .

Solution. (a) Given $x_1 = 1, h = 0.1$, then using the formula (5.12), we have

$$f'(1) \approx \frac{f(1+0.1) - f(1-0.1)}{2(0.1)} = \frac{f(1.1) - f(0.9)}{0.2} = D_h f(1).$$

Then

$$f'(1) \approx \frac{[(1.1)^2 + \cos(1.1)] - [(0.9)^2 + \cos(0.9)]}{0.2} \approx \frac{1.6636 - 1.4316}{0.2} = 1.1600.$$

(b) By using the error formula (5.13), we have

$$E_C(f,h) = -\frac{(0.1)^2}{6}f'''(\eta(x_1)), \quad for \quad \eta(x_1) \in (0.9, 1.1),$$

or

$$|E_C(f,h)| = \left| -\frac{(0.1)^2}{6} \right| |f'''(\eta(x_1))|, \quad for \quad \eta(x_1) \in (0.9, 1.1).$$

Since

$$f'''(\eta(x_1)) = \sin \eta(x_1).$$

This formula cannot be computed exactly because $\eta(x_1)$ is not known. But one can bound the error by computing the largest possible value for $|f'''(\eta(x_1))|$. So bound |f'''| on [0.9, 1.1] is

$$M = \max_{0.9 \le x \le 1.1} |\sin x| = 0.8912,$$

at x = 1.1. Thus, for $|f'''(\eta(x_1))| \leq M$ and h = 0.1, gives

$$|E_C(f,h)| \le \frac{0.01}{6}M = \frac{0.01}{6}(0.8912) = 0.0015,$$

which is the possible maximum error in our approximation.

(c) Since the exact value of the derivative f'(1) is, 0.2391, therefore, the absolute error |E| can be computed as follows

$$|E| = |f'(1) - D_h f(1)| = |(2 - \sin 1) - 1.1600| = |1.1585 - 1.1600| = 0.0015.$$

(d) Since the given accuracy required is 10^{-2} , so

$$|E_C(f,h)| = \left| -\frac{h^2}{6}f'''(\eta(x_1)) \right| \le 10^{-2},$$

for $\eta(x_1) \in (0.9, 1.1)$. Then

$$\frac{h^2}{6}M \le 10^{-2}.$$

Solving for h and taking M = 0.8912, we obtain

$$h^2 \le \frac{6}{0.8912} = 0.0673$$
, and $h \le 0.2594$.

So the best value of h is 0.25.

To get above results using the MATLAB commands, we do the following:

>>
$$x0 = 1.0$$
; $h = 0.1$;
>> $df = (x0 + h).^2 + \cos(x0 + h) - (x0 - h).^2 + \cos(x0 - h)./(2.*h)$;

Note that the formula (5.14) may also be derived from the Taylor's theorem. The second degree Taylor's expansion f(x) about x_1 , for $f(x_1 + h)$ and $f(x_1 - h)$, gives

$$f(x_1 + h) = f(x_1) + hf'(x_1) + \frac{h^2}{2!}f''(x_1) + \frac{h^3}{3!}f'''(\eta_1(x)),$$

and

$$f(x_1 - h) = f(x_1) - hf'(x_1) + \frac{h^2}{2!}f''(x_1) - \frac{h^3}{3!}f'''(\eta_2(x)).$$

Subtracting above two equations, the results is

$$f(x_1 + h) - f(x_1 - h) = 2hf'(x_1) + \frac{h^3}{3!} \left[f'''(\eta_1(x)) + f'''(\eta_2(x)) \right].$$

Since f'''(x) is continuous, the intermediate value theorem can be used to find a value of $\eta(x)$ so that

$$\frac{f'''(\eta_1(x)) + f'''(\eta_2(x))}{2} = f'''(\eta(x)),$$

which is the required formula (5.14).

Three-point Forward and Backward Difference Formulas

Similarly, the two other three-point formulas can be obtained by taking $x_k = x_0$ and $x_k = x_2$ in the formula (5.11). Firstly, by taking $x_k = x_0$ in the formula (5.11) and then after simplifying, we have

$$f'(x_0) \approx \frac{-3f(x_0) + 4f(x_0 + h) - f(x_0 + 2h)}{2h} = D_h f(x_0), \tag{5.15}$$

which is called the three-point forward-difference formula which is use to approximate the formula (5.1) at given point $x = x_0$. The error term of this approximation formula can be obtain in the similar way as we obtained for the central-difference formula and it is

$$E_F(f,h) = \frac{h^2}{3}f'''(\eta(x_0)), \tag{5.16}$$

where $\eta(x_0) \in (x_0, x_0 + 2h)$. Similarly, taking $x_k = x_2$ in the formula (5.11), and after simplifying, we obtain

$$f'(x_2) \approx \frac{f(x_2 - 2h) - 4f(x_2 - h) + 3f(x_2)}{2h} = D_h f(x_2), \tag{5.17}$$

which is called the three-point backward-difference formula which is use to approximate the formula (5.1) at given point $x = x_2$. It has the error term of the form

$$E_B(f,h) = \frac{h^2}{3}f'''(\eta(x_2)), \tag{5.18}$$

where $\eta(x_2) \in (x_2 - 2h, x_2)$.

Note that the backward-difference formula (5.17) can be obtained from the forward-difference formula by replacing h with -h. Also, note that the error in (5.12) is approximately half the error in (5.15) and (5.17). This is reasonable since in using the central-difference formula (5.12) data is being examined on both sides of point x_1 , and for others in (5.15) and (5.17) only on one side. Note that in using the central-difference formula, a function f(x) needs to be evaluated at only two points, whereas in using the other two formulas, we need the values of a function at three points. The approximations in using the formulas (5.15) and (5.17) are useful near the ends of the required interval, since the information about a function outside the interval may not be available. Thus the central-difference formula (5.12) is superior to both the forward-difference formula (5.15) and the backward-difference formula (5.17). The central-difference represents the average of the forward-difference and the backward-difference.

Example 5.6 Consider the following table for set of data points

- (a) Use three-point formula for smaller value of h to find approximation of f'(3).
- (b) The function tabulated is $\ln x$, find error bound and absolute error for the approximation of f'(3).
- (c) What is the best maximum value of stepsize h required to obtain the approximate value of f'(3) within the accuracy 10^{-4} .

Solution. (a) For the given table of data points, we can use all three-points formulas as for the central difference we can take

$$x_0 = x_1 - h = 2$$
, $x_1 = 3$, $x_2 = x_1 + h = 4$, gives $h = 1$,

for the forward difference formula we can take

$$x_0 = 3$$
, $x_1 = x_0 + h = 3.9$, $x_2 = x_0 + 2h = 4.8$, gives $h = 0.9$,

and for the backward difference formula we can take

$$x_0 = x_2 - 2h = 1.6$$
, $x_1 = x_2 - h = 2.3$, $x_2 = 3$, gives $h = 0.7$.

Since we know that smaller the vale of h better the approximation of the derivative of the function, therefore, for the given problem, backward difference is the best formula to find approximation of f'(3) as

$$f'(3) \approx \frac{f(1.6) - 4f(2.3) + 3f(3)}{2(0.7)} \approx \frac{[0.47 - 4(0.83) + 3(1.10)]}{1.4} = 0.3214.$$

(b) Using error term of backward difference formula, we have

$$E_B(f,h) = \frac{h^2}{3}f'''(\eta), \quad or \quad |E_B(f,h)| \le \frac{h^2}{3}|f'''(\eta)|.$$

Taking $|f'''(\eta(x_2))| \le M = \max_{1.6 \le x \le 3} |f'''(x)| = \max_{1.6 \le x \le 3} |2/x^3| = 0.4883$. Thus using h = 0.7, we obtain

$$|E_B(f,h)| \le \frac{(0.7)^2}{3}(0.4883) = 0.0798,$$

the required error bounds for the approximations. To compute the absolute error we do as

$$|E| = |f'(3) - 0.3214| = |0.3333 - 0.3214| = 0.0119.$$

(c) Since the given accuracy required is 10^{-4} , so

$$|E_B(f,h)| = \left|\frac{h^2}{3}f'''(\eta)\right| \le 10^{-4},$$

for $\eta \in (1.6,3)$. Then

$$\frac{h^2}{3}M \le 10^{-4}.$$

Solving for h by taking M = 0.4883, we obtain

$$h^2 \le \frac{3 \times 10^{-4}}{0.4883} = 0.0248,$$

and so h = 0.025 the best maximum value of h.

Example 5.7 Use the best three-point formula to find approximation of f'(1.5) using the following table for set of data points

Solution. The best three-point formula for this problem is the central difference because this formula does not need the value of f(1.5) which is not given in the table while the other two need this value. Since $x_1 = 1.5$, so by taking $x_1 + h = 2$ and $x_1 - h = 1$, gives h = 0.5, we have

$$f'(1.5) \approx \frac{(f(2) - f(1))}{2(0.5)} = \frac{(0.69 - 0.0)}{1} = 0.69,$$

the required approximation of f'(1.5).

Example 5.8 Use the three-point formulas (5.12), (5.15) and (5.17) to approximate the first derivative of the function $f(x) = e^x$ at x = 2, take h = 0.1. Also, compute the error bound for each approximation.

Solution. Given $f(x) = e^x$ and h = 0.1, then

Central-difference formula:

$$f'(2) \approx \frac{(f(2.1) - f(1.9))}{2h} = \frac{(e^{2.1} - e^{1.9})}{0.2} = 7.4014.$$

Forward-difference formula:

$$f'(2) \approx \frac{-3f(2) + 4f(2.1) - f(2.2)}{2h} \approx \frac{-3e^2 + 4e^{2.1} - e^{2.2}}{0.2} = 7.3625.$$

Backward difference formula:

$$f'(2) \approx \frac{f(1.8) - 4f(1.9) + 3f(2)}{2h} \approx \frac{e^{1.8} - 4e^{1.9} + 3e^2}{0.2} = 7.3662.$$

Since the exact solution of the first derivative of the given function at x=2 is 7.3891, so the corresponding actual errors are, -0.0123, 0.0266 and 0.0229 respectively. This shows that the approximate solution got by using the central-difference formula is closer to exact solution as compared with the other two difference formulas.

The error bounds for the approximations got by (5.12), (5.15), and (5.17) are as follows:

Central-difference formula:

$$E_C(f,h) = -\frac{h^2}{6}f'''(\eta(x_1)), \quad or \quad |E_C(f,h)| \le \frac{h^2}{6}|f'''(\eta(x_1))|.$$

Taking $|f'''(\eta(x_1))| \le M = \max_{1.9 \le x \le 2.1} |e^x| = e^{2.1}$ and h = 0.1, we obtain

$$|E_C(f,h)| \le \frac{(0.1)^2}{6}e^{2.1} = 0.0136.$$

Forward-difference formula:

$$E_F(f,h) = \frac{h^2}{3}f'''(\eta(x_0)), \quad or \quad |E_F(f,h)| \le \frac{h^2}{3}|f'''(\eta(x_0))|.$$

Taking $|f'''(\eta(x_0))| \le M = \max_{2 \le x \le 2} |e^x| = e^{2.2}$ and h = 0.1, we obtain

$$|E_F(f,h)| \le \frac{(0.1)^2}{3}e^{2.2} = 0.0301.$$

Backward difference formula:

$$E_B(f,h) = \frac{h^2}{3}f'''(\eta(x_2)), \quad or \quad |E_B(f,h)| \le \frac{h^2}{3}|f'''(\eta(x_2))|.$$

Taking $|f'''(\eta(x_2))| \le M = \max_{1.8 \le x \le 2} |e^x| = e^2$ and h = 0.1, we obtain

$$|E_B(f,h)| \le \frac{(0.1)^2}{3}e^2 = 0.0246.$$

Thus we got the required error bounds for the approximations.

Example 5.9 Let $f(x) = (x \ln x + x)$ and x = 0.9, 1.3, 1.6, 2.1, 2.5, 3.1. Then

- (a) Find the approximate value of f'(1.9) using three-point formula for smaller value of h.
- (b) Compute the error bound for your approximation.
- (c) Compute the absolute error.
- (d) What is the best maximum value of stepsize h required to obtain the approximate value of f'(1.9) within the accuracy 10^{-2} .

Solution. (a) For the given data points we can use all three-points difference formulas with central difference at $x_0 = 1.3, x_1 = 1.9, x_2 = 2.5$, forward difference at $x_0 = 1.9, x_1 = 2.5, x_2 = 3.1$, and backward difference at $x_0 = 1.3, x_1 = 1.6, x_2 = 1.9$. But the value of h = 0.3 for the backward difference formula is smaller than both the other formulas with h = 0.6. So the best three-point formula for this case is the following backward difference formula

$$f'(x_2) \approx \frac{f(x_2 - 2h) - 4f(x_2 - h) + 3f(x_2)}{2h} = D_h f(x_2).$$

Thus using $x_2 = 1.9, x_2 - h = 1.6$, and $x_2 - 2h = 1.3$, we have

$$f'(1.9) \approx \frac{f(1.3) - 4f(1.6) + 3f(1.9)}{2(0.3)},$$

and using $f(x) = x \ln x + x$, we obtain

$$f'(1.9) \approx \frac{(1.3 \ln 1.3 + 1.3) - 4(1.6 \ln 1.6 + 1.6) + 3(1.9 \ln 1.9 + 1.9)}{0.6} \approx 2.6527,$$

the required approximation of f'(1.9).

(b) By using the backward difference error

$$E_B(f,h) = \frac{(0.3)^2}{3} f'''(\eta(x_2)), \quad for \quad \eta(x_1) \in (1.3, 1.9),$$

we get

$$|E_C(f,h)| = \left|\frac{(0.3)^2}{3}\right| |f'''(\eta(x_1))|, \quad for \quad \eta(x_2) \in (1.3, 1.9).$$

Since

$$f'''(\eta(x_2)) = -1/(\eta(x_2))^2.$$

This formula cannot be computed exactly because $\eta(x_2)$ is not known. But one can bound the error by computing the largest possible value for $|f'''(\eta(x_2))|$. So bound |f'''| on [1.3, 1.9] is

$$M = \max_{\substack{1.3 \le x \le 1.9}} |-1/x^2| = 1/(1.3)^2 = 0.5917.$$

Thus, for $|f'''(\eta(x_2))| \leq M$ and h = 0.1, gives

$$|E_B(f,h)| \le \frac{0.09}{3}M = \frac{0.09}{3}(0.5917) = 0.0178,$$

which is the possible maximum error in our approximation.

(c) Since the exact value of the derivative f'(1.9) is, 2.6419, therefore, the absolute error |E| can be computed as follows

$$|E| = |f'(1) - D_h f(1.9)| = |2.6419 - 2.6527| = 0.0108.$$

(d) Since the given accuracy required is 10^{-2} , so

$$|E_B(f,h)| = \left|\frac{h^2}{3}f'''(\eta(x_2))\right| \le 10^{-2}$$

for $\eta(x_2) \in (1.3, 1.9)$. Then

$$\frac{h^2}{3}M \le 10^{-2}$$
.

Solving for h and taking M = 0.5917, we obtain

$$h^2 \le \frac{3}{59.17} = 0.0507$$
, or $h \le 0.2252$,

and so h = 0.2 the best maximum value of h.

5.3.2 Second Derivative Numerical Formula

It is also possible to estimate second and higher order derivatives numerically. Formulas for higher derivatives can be found by differentiating the interpolating polynomial repeatedly or using the Taylor's theorem. Since the *two-point* and *three-point* formulas for the approximation of the first derivative of a function were derived by differentiating the Lagrange interpolation polynomials for f(x) but the derivation of the higher-order can be tedious. Therefore, we shall use here the Taylor's theorem for finding the *three-point* central-difference formulas for finding approximation of the second derivative f''(x) of a function f(x) at the given point $x = x_1$. The process used to obtain numerical formulas for first and second derivatives of a function can be readily extended to third- and higher-order derivatives of a function.

Three-point Central Difference Formula

To find the three-point central-difference formula for the approximation of the second derivative of a function at given point, we use the third-order Taylor's theorem by expanding a function f(x) about a point x_1 and evaluate at $x_1 + h$ and $x_1 - h$. Then

$$f(x_1 + h) = f(x_1) + hf'(x_1) + \frac{1}{2}h^2f''(x_1) + \frac{1}{6}h^3f'''(x_1) + \frac{1}{24}h^4f^{(4)}(\eta_1(x)),$$

and

$$f(x_1 - h) = f(x_1) - hf'(x_1) + \frac{1}{2}h^2f''(x_1) - \frac{1}{6}h^3f'''(x_1) + \frac{1}{24}h^4f^{(4)}(\eta_2(x)),$$

where $(x_1 - h) < \eta_2(x) < x_1 < \eta_1(x) < (x_1 + h)$.

By adding these equations and simplifies, we have

$$f(x_1 + h) + f(x_1 - h) = 2f(x_1) + h^2 f''(x_1) + \frac{(f^{(4)}(\eta_1(x)) + f^{(4)}(\eta_2(x)))}{24} h^4.$$

Solving this equation for $f''(x_1)$, we obtain

$$f''(x_1) = \frac{f(x_1 - h) - 2f(x_1) + f(x_1 + h)}{h^2} - \frac{h^4}{24} \left[f^{(4)}(\eta_1(x)) + f^{(4)}(\eta_2(x)) \right].$$

If $f^{(4)}$ is continuous on $[x_1 - h, x_1 + h]$, then by using the Intermediate Value Theorem, the above equation can be written as

$$f''(x_1) = \frac{f(x_1 - h) - 2f(x_1) + f(x_1 + h)}{h^2} - \frac{h^4}{12}f^{(4)}(\eta(x_1)).$$

Then the following formula

$$f''(x_1) \approx \frac{f(x_1 - h) - 2f(x_1) + f(x_1 + h)}{h^2} = D_h^2 f(x_1), \tag{5.19}$$

is called the three-point central-difference formula for the approximation of the second derivative of a function f(x) at the given point $x = x_1$.

Example 5.10 Let $f(x) = x \ln x + x$ and x = 0.9, 1.3, 2.1, 2.5, 3.2. Then find the approximate value of $f''(x) = \frac{1}{x}$ at x = 1.9. Also, compute the absolute error.

Solution. Given $f(x) = x \ln x + x$, then one can easily find second derivative of the function as

$$f'(x) = \ln x + 2$$
 and $f''(x) = \frac{1}{x}$.

To find the approximation of $f''(x) = \frac{1}{x}$ at the given point $x_1 = 1.9$, we use the three-point formula (5.19)

$$f''(x_1) \approx \frac{f(x_1+h) - 2f(x_1) + f(x_1-h)}{h^2} = D_h^2 f(x_1).$$

Taking the three points 1.3, 1.9 and 2.5 (equally spaced), giving h = 0.6, we have

$$f''(1.9) \approx \frac{f(2.5) - 2f(1.9) + f(1.3)}{0.36}$$

$$\approx \frac{((2.5 \ln 2.5 + 2.5) - 2(1.9 \ln 1.9 + 1.9) + (1.3 \ln 1.3 + 1.3))}{0.36}$$

$$\approx \frac{4.7907 - 6.2391 + 1.6411}{0.36} = 0.5353 = D_h^2 f(1.9).$$

Since the exact value of f''(1.9) is $\frac{1}{1.9} = 0.5263$, therefore, the absolute error |E| can be computed as follows:

$$|E| = |f''(1.9) - D_h^2 f(1.9)| = |0.5263 - 0.5353| = 0.009.$$

Note that the *error term* of the three-point central-difference formula (5.19) for the approximation of the second derivative of a function f(x) at the given point $x = x_1$ is of the form

$$E_C(f,h) = -\frac{h^2}{12}f^{(4)}(\eta(x_1)), \tag{5.20}$$

for some unknown point $\eta(x_1) \in (x_1 - h, x_1 + h)$.

Note that for a cubic function, $f(x) = ax^3 + bx^2 + cx + d$, the difference formula (5.19) is exact; that is, it yields the correct value of f''(x) for any nonzero value of stepsize h.

Example 5.11 Let $f(x) = x^2 + \cos x$. Then

- (a) Compute the approximate value of f''(x) at x = 1, taking h = 0.1 using (5.19).
- (b) Compute the error bound for your approximation using (5.20).
- (c) Compute the absolute error.
- (d) What is the best maximum value of stepsize h required to obtain the approximate value of f''(1) within the accuracy 10^{-2} .

Solution. (a) Given $x_1 = 1, h = 0.1$, then the formula (5.19) becomes

$$f''(1) \approx \frac{f(1+0.1) - 2f(1) + f(1-0.1)}{(0.1)^2} = D_h^2 f(1),$$

or

$$f''(1) \approx \frac{f(1.1) - 2f(1) + f(0.9)}{0.01}$$

$$\approx \frac{[(1.1)^2 + \cos(1.1)] - 2[1^2 + \cos(1)] + [(0.9)^2 + \cos(0.9)]}{0.01}$$

$$\approx \frac{1.6636 - 3.0806 + 1.4316}{0.01} \approx 1.4600 = D_h^2 f(1).$$

(b) To compute the error bound for our approximation in part (a), we use the formula (5.20) and have

$$E_C(f,h) = -\frac{h^2}{12}f^{(4)}(\eta(x_1)), \quad for \quad \eta(x_1) \in (0.9, 1.1),$$

or

$$|E_C(f,h)| = \left| -\frac{h^2}{12} \right| |f^{(4)}(\eta(x_1))|, \quad for \quad \eta(x_1) \in (0.9, 1.1).$$

The fourth derivative of the given function at $\eta(x_1)$ is

$$f^{(4)}(\eta(x_1)) = \cos \eta(x_1),$$

and it cannot be computed exactly because $\eta(x_1)$ is not known. But one can bound the error by computing the largest possible value for $|f^{(4)}(\eta(x_1))|$. So bound $|f^{(4)}|$ on the interval (0.9, 1.1) is

$$M = \max_{0.9 \le x \le 1.1} |\cos \eta(x_1)| = 0.4536,$$

at x = 1.1, Thus, for $|f^{(4)}(\eta(x))| \leq M$, we have

$$|E_C(f,h)| \le \frac{h^2}{12}M.$$

Taking M = 0.4536 and h = 0.1, we obtain

$$|E_C(f,h)| \le \frac{0.01}{12}(0.4536) = 0.0004,$$

which is the possible maximum error in our approximation.

By using the MATLAB symbolic commands:

$$>> syms \ x;$$

 $>> f = x.^2 + cos(x); \ ddf = diff(f, 2); \ ddddf = diff(f, 4); \ subs(ddf, 1);$

(c) Since the exact value of f''(1) is

$$f''(1) = (2 - 1^2)\cos 1 - 4(1)\sin 1 = -2.8256,$$

therefore, the absolute error |E| can be computed as follows:

$$|E| = |f''(1) - D_h^2 f(1)| = |1.4597 - 1.4600| = 0.0003.$$

(d) Since the given accuracy required is 10^{-2} , so

$$|E_C(f,h)| = \left| -\frac{h^2}{12} f^{(4)}(\eta(x_1)) \right| \le 10^{-2},$$

for $\eta(x_1) \in (0.9, 1.1)$. Then for $|f^{(4)}(\eta(x_1))| \leq M$, we have

$$\frac{h^2}{12}M \le 10^{-2}$$
.

Solving for h^2 , we obtain

$$h^2 \le \frac{(12 \times 10^{-2})}{M} = \frac{(12 \times 10^{-2})}{0.4536} = 0.2646,$$

and it gives the value of h as

$$h \le 0.5144$$
.

Thus the best maximum value of h is 0.5.

Using the following MATLAB commands, we can easily achieved the above results:

>>
$$x0 = 1.0$$
; $h = 0.1$;
>> $ddf = ((x0 + h).^2 + cos(x0 + h) - 2. *x0.^2 + cos(x0) + ...$
 $(x0 - h).^2 + cos(x0 - h))./(h.^2)$

The central-difference formula is probably the most used approximation for derivatives. Many real problems are modeled by second-order differential equations, involving either ordinary or partial derivatives. These equations cannot be solved analytically. To solve the equations numerically requires the replacement of the second-order derivatives by the difference formula (5.19).

Example 5.12 The function f(x) satisfies a given equation $f''(x) = x^2 f(x)$ and the conditions f(0) = 1, f(0.2) = 3. Use the central-difference formula (5.19) for f''(x) and the stepsize h = 0.1 to estimate the value of f(0.1).

Solution. Given $x_1 = 0.1$ and h = 0.1, then using the central-difference formula for the second derivative of a function

$$f''(x_1) \approx \frac{f(x_1 - h) - 2f(x_1) + f(x_1 + h)}{h^2},$$

we obtain

$$f''(0.1) = (0.1)^2 f(0.1) \approx \frac{f(0) - 2f(0.1) + f(0.2)}{0.01},$$

which is equal to

$$(0.01)f(0.1) \qquad \approx \quad \frac{1 - 2f(0.1) + 3}{0.01}$$

$$(0.0001)f(0.1) \approx (4-2f(0.1)).$$

Solving for f(0.1), we obtain

$$f(0.1) \approx \frac{4}{2 + 0.0001} \approx 1.9999,$$

the required solution.

Example 5.13 Consider following set of data points

Use the table, find the best approximation of f'(0.75) and the worst approximations of f'(0.1) and f''(0.6) by using three-point formulas.

Solution. For the best approximation of f'(0.75), we have to take small value of h = 0.15, so using the central difference three-point formula (5.12), we get

$$f'(0.75) \approx \frac{f(0.9) - f(0.6)}{2(0.15)} \approx \frac{1.52 - 1.43}{0.3} = 0.3,$$

while the exact value of f'(0.75) is 0.3184. For the worst approximation of f'(0.1), we have to take big value of h = 0.5, so using the forward difference three-point formula (5.15), we get

$$f'(0.1) \approx \frac{-3f(0.1) + 4f(0.6) - f(1.1)}{2(0.5)} \approx \frac{-3(1.1) + 4(1.43) - 1.55}{1} \approx 0.87.$$

Similarly, for the worst approximation of f''(0.6), we have to take big value of h = 0.6, so using the central difference three-point formula (5.19), we get

$$f''(0.6) \approx \frac{f(1.2) - 2f(0.6) + f(0.0)}{0.36}$$

$$\approx \frac{1.56 - 2(1.43) + 1.0}{0.36} \approx -0.8333,$$

the required worst approximation.

5.4 Formulas for Computing Derivatives

For convenience, we collect following some useful central-difference, forward-difference and backward-difference formulas for computing different orders derivatives.

262 5.5 Numerical Integration

5.4.1 Central Difference Formulas

The central-difference formula (5.12) for first derivative $f'(x_1)$ of a function required that a function can be computed at points that lies on both sides of x_1 . The Taylor series can be used to obtain central-difference formulas for higher derivatives. The most usable are those of order $O(h^2)$ and $O(h^4)$ and are given as follows:

$$f'(x_0) = \frac{f_1 - f_{-1}}{2h} + O(h^2)$$

$$f'(x_0) = \frac{-f_2 + 8f_1 - 8f_{-1} + f_{-2}}{12h} + O(h^4)$$

$$f''(x_0) = \frac{f_1 - 2f_0 + f_{-1}}{h^2} + O(h^2)$$

$$f''(x_0) = \frac{-f_2 + 16f_1 - 30f_0 + 16f_{-1} - f_{-2}}{12h^2} + O(h^4)$$

$$f'''(x_0) = \frac{f_2 - 2f_1 + 2f_{-1} - f_{-2}}{2h^3} + O(h^2)$$

$$f'''(x_0) = \frac{-f_3 + 8f_2 - 13f_1 + 13f_{-1} - 8f_{-2} + f_{-3}}{8h^3} + O(h^4)$$

$$f^{(4)}(x_0) = \frac{f_2 - 4f_1 + 6f_0 - 4f_{-1} + f_{-2}}{h^4} + O(h^2)$$

$$f^{(4)}(x_0) = \frac{-f_3 + 12f_2 - 39f_1 + 56f_0 - 39f_{-1} + 12f_{-2} - f_{-3}}{6h^4} + O(h^4)$$

5.4.2 Forward and Backward Difference Formulas

If a function cannot be evaluated at points that lie on both sides of x_0 , then the central-difference formula cannot be used to approximate the derivatives of a function. When a function can be evaluated at equally spaced points that lie to the right (or left) of point x_0 , then forward-(or backward-)difference formula can be used. These formulas can be derived by using the Taylor series, the Lagrange interpolating polynomials, or the Newton interpolating polynomials. Some of them are mostly usable to find derivatives of a function are as follows:

$$f'(x_0) = \frac{-3f_0 + 4f_1 - f_2}{2h} + O(h^2)$$

$$f''(x_0) = \frac{3f_0 - 4f_{-1} + f_{-2}}{2h} + O(h^2)$$

$$f'''(x_0) = \frac{2f_0 - 5f_1 + 4f_2 - f_3}{h^2} + O(h^2)$$

$$f'''(x_0) = \frac{2f_0 - 5f_{-1} + 4f_{-2} - f_{-3}}{h^2} + O(h^2)$$

$$f'''(x_0) = \frac{-5f_0 + 18f_1 - 24f_2 + 14f_3 - 3f_4}{2h^3} + O(h^2)$$

$$f'''(x_0) = \frac{5f_0 - 18f_{-1} + 24f_{-2} - 14f_{-3} + 3f_{-4}}{2h^3} + O(h^2)$$

$$f^{(4)}(x_0) = \frac{3f_0 - 14f_1 + 26f_2 - 24f_3 + 11f_4 - 2f_5}{h^4} + O(h^2)$$

$$f^{(4)}(x_0) = \frac{3f_0 - 14f_{-1} + 26f_{-2} - 24f_{-3} + 11f_{-4} - 2f_{-5}}{h^4} + O(h^2)$$

5.5 Numerical Integration

Numerical integration has a history extending back to the invention of calculus and before. It is used to integrate tabulated functions or to integrate functions whose integrals are either impossible or very difficult to obtain analytically. Even when analytical integration is easy, numerical integration may save time and effort if only the numerical value of the integral is desired. Consequently, numerical methods of integration represent a natural alternative whenever conventional methods fail to yield a solution.

Now for numerical integration, we wish to find an approximation to the definite integral

$$I(f) = \int_{a}^{b} f(x)dx,$$
(5.21)

assuming that f(x) is integrable. If $f(x) \ge 0$ on the given interval [a, b], then geometrically, the integral (5.21) is equivalent to replacing the *area* under the graph of f(x), the x-axis and between the ordinates x = a and x = b.

The fundamental theorem of calculus shows that integration is the inverse process to differentiation. If we can find a function F(x), called the antiderivative of f(x), that is, F'(x) = f(x), then we can evaluate the integral (5.21) using the relation

$$I(f) = \int_{a}^{b} f(x)dx = F(b) - F(a).$$
 (5.22)

Sometimes considerable skill is required to obtain F(x), perhaps by making a change of variable or integrating by parts. But in many cases F(x) cannot be found by elementary methods. In such cases the computation of the integral (5.21) by means of formula (5.22) may be either too difficult or practically impossible. Even if F(x) can be found, it may still more convenient to use a numerical method to estimate the integral (5.21) if the evaluation of F(x) required a great deal of computation. Moreover, in practical applications, a function f(x) often given in tabular form and then the entire concept of antiderivative is meaningless. An obvious approach is to replace a function f(x) in the integral (5.21) by an approximating polynomial p(x), that is

$$I(f) = \int_{a}^{b} f(x)dx \approx \int_{a}^{b} p(x)dx.$$

Numerical integration formulas are derived by integrating interpolation polynomials. Therefore, different interpolation formulas will leads to different numerical integration methods.

The definite integral (5.21) may be interpreted as the area under the curve of y = f(x) from a to b as shown by Figure 5.4.

It should be noted that any areas beneath the x-axis are counted as negative. Many numerical methods for integration are based on using this interpretation of the integral to derive approximations to it by dividing the interval [a,b] into a number of smaller subintervals. By making simple approximations to the curve y = f(x) in the small subinterval its area may be obtained and on summing all the contributions we obtain an approximation to a integral in the interval [a,b]. Variations of this technique are derived by taking groups of subintervals and fitting different degree polynomials as approximations for each of these groups. The lead of accuracy obtained is dependent on the number of intervals used and the nature the approximation function.

There are several methods available in the literature for numerical integration but the most commonly methods may be classified into two groups.

264 5.5 Numerical Integration

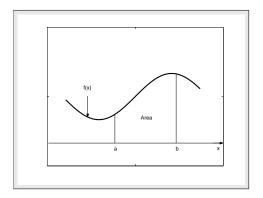


Figure 5.4: Definite integral for f(x)

- (a) The Newton-Cotes formulas that employ functional values at equally spaced data points.
- (b) The Gaussian quadrature formulas that employ unequally spaced data points determined by certain properties of orthogonal polynomials.

Firstly, we shall discuss the Newton-Cotes formulas which has two different types, called, the *closed Newton-Cotes* formulas and the *open Newton-Cotes* formulas. In the first type, we shall discuss in some details the two mostly usable formulas, called the *Trapezoidal rule* and the *Simpson's rule* which can be derived by integrating the Lagrange interpolating polynomials of degree 1 and 2 respectively. In the second type we shall consider some good formulas. The use of the closed Newton-Cotes and other integration formulas of order higher than the Simpson's rule is seldom necessary in most engineering applications and can be use for those cases where extremely high accuracy is required.

Example 5.14 Find the numerical integration rule using the points $x_0 = 0, x_1 = 1$ and $x_2 = 2$ for approximating the integral

$$\int_0^3 f(x)dx,$$

for $f(x) = 1, x, x^2$. Use the resulting rule to estimate the integral

$$\int_0^3 (x^2+1)dx.$$

Solution. Consider

$$I = \int_0^3 f(x)dx \approx a_0 f(x_0) + a_1 f(x_1) + a_2 f(x_2),$$

then for f(x) = 1, we have

$$I = \int_0^3 1 dx \approx a_0 f(0) + a_1 f(1) + a_2 f(2),$$

and it gives

$$3 \approx a_0 + a_1 + a_2$$
.

Similarly, for f(x) = x and $f(x) = x^2$, we get

$$\frac{9}{2} \approx a_1 + 3a_2 \quad and \quad 9 \approx a_1 + 9a_2,$$

respectively. Now solving the following linear system

we obtain $a_0 = 3/4$, $a_1 = 0$, $a_2 = 9/4$. Thus the required integration rule is

$$\int_0^3 f(x)dx \approx \frac{3f(0) + 9f(2)}{4},$$

and with $f(x) = x^2 + 1$, we get

$$\int_0^3 (x^2 + 1)dx \approx \frac{3(0+1) + 9(4+1)}{4} = 12.$$

One can easily compute the true integral as

$$\int_0^3 (x^2 + 1)dx = (x^3/3 + x) \Big|_0^3 = 12.$$

Thus the above integration rule gave the exact solution for the given function.

Using MATLAB symbolic toolbox we can easily obtained symbolic integration of function as

>>
$$syms \ x;$$

>> $f = x.^2 + 1; \ I = int(f);$
>> $a = 0; b = 3; \ A = int(f, a, b); \ subs(A);$

Example 5.15 Find α such that the integration formula $\int_0^1 \frac{f(x)}{\sqrt{x}} \approx Af(0) + Bf(\alpha) + Cf(1)$ may be exact for polynomials up to degree 3.

Solution. Consider

$$I = \int_0^1 \frac{f(x)}{\sqrt{x}} = Af(0) + Bf(\alpha) + Cf(1),$$

then by taking $f(x) = 1, x, x^2$ and x^3 , we get

Subtracting second equation from third equation, gives, $B\alpha(\alpha - 1) = -4/15$, and then subtracting third equation from fourth equation, gives, $B\alpha^2(\alpha - 1) = -4/35$. Thus

$$-4/35 = B\alpha^{2}(\alpha - 1) = \alpha[B\alpha(\alpha - 1)] = \alpha(-4/15), \quad or \quad -4/35 = \alpha(-4/15),$$

and solving for α , it gives, $\alpha = 3/7$.

5.6 Newton-Cotes Formulas

The usual strategy in developing formulas for numerical integration is similar to that for numerical differentiation. We pass a polynomial through points of a function and then integrate this polynomial approximation to a function. This allows us to integrate a function known only as a table of values. Some common formulas based on polynomial interpolation are referred to as the Newton-Cotes formulas.

An (n+1)-point Newton-Cotes formula for approximating the definite integral (5.21) is obtained by replacing the integrand f(x) by the *nth*-degree Lagrange polynomial that interpolates the values of f(x) at equally spaced data points

$$a = x_0 < x_1 < \ldots < x_n = b.$$

Note that if the end-points a and b of the given interval [a, b] are in the set of interpolating points; then the Newton-Cotes formulas are called *closed*; otherwise, it is said to be *open*.

5.6.1 Closed Newton-Cotes Formulas

An (n+1)-point closed Newton-Cotes formula used points $x_i = x_0 + ih$, for, i = 0, 1, 2, ..., n, where $x_0 = a$, $x_n = b$ and $h = \frac{b-a}{n}$, has the form (see Figure 5.5)

$$\int_{a}^{b} f(x)dx = \int_{x_0}^{x_n} f(x)dx \approx \sum_{i=0}^{n} a_i f(x_i),$$
 (5.23)

where

$$a_{i} = \int_{x_{0}}^{x_{n}} L_{i}(x)dx = \int_{x_{0}}^{x_{n}} \prod_{\substack{j=0\\j\neq i}}^{n} \frac{(x-x_{j})}{(x_{i}-x_{j})} dx.$$
 (5.24)

The following theorem describes the error analysis associated with the above closed Newton-Cotes formulas.

Theorem 5.1 (Close Newton-Cotes Formulas)

Suppose that $\sum_{i=0}^{n} a_i f(x_i)$ denotes the (n+1)-point closed Newton-Cotes formula with $x_0 = a, x_n = b$, and h = (b-a)/n. There exists $\eta(x) \in (a,b)$ for which

$$\int_{a}^{b} f(x)dx = \sum_{i=0}^{n} a_{i}f(x_{i}) + \frac{h^{n+3}f^{(n+2)}(\eta(x))}{(n+2)!} \int_{0}^{n} t^{2}(t-1)\cdots(t-n)dt,$$
 (5.25)

if n is even and $f \in C^{n+2}[a,b]$. For $f \in C^{n+1}[a,b]$, and n is odd, then

$$\int_{a}^{b} f(x)dx = \sum_{i=0}^{n} a_{i}f(x_{i}) + \frac{h^{n+2}f^{(n+1)}(\eta(x))}{(n+1)!} \int_{0}^{n} t(t-1)\cdots(t-n)dt.$$
 (5.26)

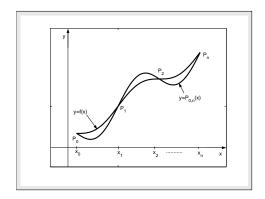


Figure 5.5: Close Newton-Cotes approximation

Different numerical integration formulas can be obtained by using the formulas (5.25) and (5.26) to approximate the definite integral (5.21). By using the formula (5.26) for n = 1, we have well-known numerical integration formula, called, the *Trapezoidal rule*. Similarly, by using the formula (5.25) for n = 2, we have one of the best integration rule called, the *Simpson's rule*. We shall discuss the formulation of both these rules and also discuss about their error terms. Later we shall also consider some more closed Newton-Cotes formulas.

Simple Trapezoidal Rule

It is one of the oldest and good numerical method for approximating the definite integral (5.21). It is based on approximating a function in each subinterval by a straight line.

To derive the Trapezoidal rule for one-strip (one interval), let us consider the first degree Lagrange interpolating polynomial with equally spaced data points, that is, $x_0 = a, x_1 = b$ and $h = x_1 - x_0$, then

$$f(x) = p_1(x) = \left(\frac{x - x_1}{x_0 - x_1}\right) f(x_0) + \left(\frac{x - x_0}{x_1 - x_0}\right) f(x_1). \tag{5.27}$$

Taking integral on both sides of (5.27) with respect to x between the limits x_0 and x_1 , we have

$$\int_{x_0}^{x_1} f(x)dx \approx \frac{f(x_0)}{x_0 - x_1} \int_{x_0}^{x_1} (x - x_1)dx + \frac{f(x_1)}{x_1 - x_0} \int_{x_0}^{x_1} (x - x_0)dx,$$

which implies that

$$\int_{x_0}^{x_1} f(x)dx \approx \frac{f(x_0)}{x_0 - x_1} \left[\frac{(x - x_1)^2}{2} \Big|_{x_0}^{x_1} \right] + \frac{f(x_1)}{x_1 - x_0} \left[\frac{(x - x_0)^2}{2} \Big|_{x_0}^{x_1} \right] \approx \frac{(x_1 - x_0)}{2} [f(x_0) + f(x_1)],$$

and by taking $h = x_1 - x_0$, we get

$$\int_{a=x_0}^{b=x_1} f(x)dx \approx T_1(f) = \frac{h}{2} [f(x_0) + f(x_1)]. \tag{5.28}$$

Then $T_1(f)$ is called the *simple Trapezoidal rule* or the Trapezoidal rule for one trapezoid or one strip and can be use for the approximation of the definite integral (5.21). The reason for calling this formula the Trapezoidal rule is that when f(x) is a function with positive values, the integral (5.21) is approximated by the area in the trapezoid, see Figure 5.6.

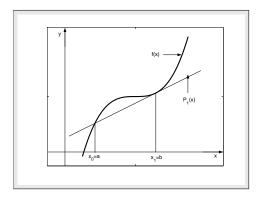


Figure 5.6: Simple Trapezoidal rule.

Example 5.16 Approximate the following integral

$$\int_{1}^{2} \frac{1}{x+1} dx,$$

using the simple Trapezoidal rule and compute the absolute error.

Solution. Given $f(x) = \frac{1}{x+1}$ and h = 1, so using the simple Trapezoidal rule (5.28), gives

$$T_1(f) = \frac{1}{2}[f(1) + f(2)] = 0.4167.$$

The exact solution of the given integral is

$$I(f) = \ln(3/2) = 0.4055$$
, so $|E_{T_1}(f)| = |I(f) - T_1(f)| = |0.4055 - 0.4167| = 0.0112$,

is the required absolute error.

Composite Trapezoidal Rule

It is evident that the Newton-Cotes formulas produce accurate approximations to the definite integral (5.21) only when the limits a and b are close together, that is, the integration interval is not large. Formulas based on low-degree interpolating polynomials are clearly unsuitable since it is then necessary to use large values of h. Also, note that higher-order Newton-Cotes formulas will not necessarily produce more accurate approximations to the given integral. This difficulty can be avoided by using a piecewise approach; the integration interval is divided into subintervals and low-order formulas are applied on each of these. The corresponding integration rules are said to be in *composite form*, and the most suitable formula of this type make use of the Trapezoidal rule. The interval [a, b] is partitioned into n subintervals (x_{i-1}, x_i) , $i = 1, 2, \ldots, n$ with $a = x_0$ and $b = x_n$ of equal width h = (b - a)/n and the rule for a single interval (the simple rule (5.28)) is applied to each subinterval or a grouping of subintervals (see Figure 5.7). Since the Trapezoidal rule requires only one interval for application, there is no restriction on the integer n. We define the composite Trapezoidal rule in the form of the following theorem.

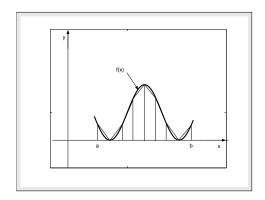


Figure 5.7: Composite Trapezoidal rule.

Theorem 5.2 (Composite Trapezoidal Rule)

Let $f \in C^2[a,b]$, n may be odd or even, h = (b-a)/n, and $x_i = a+ih$ for each $i = 0,1,2,\ldots,n$. Then the composite Trapezoidal rule for n subintervals can be written as

$$\int_{a=x_0}^{b=x_n} f(x)dx \approx T_n(f) = \frac{h}{2} \left[f(a) + 2 \sum_{i=1}^{n-1} f(x_i) + f(b) \right].$$
 (5.29)

Proof. Since for the composite form of the Trapezoidal rule, the interval is divided into n equal subintervals of width h so that $h = \frac{b-a}{n}$, and (n+1) distinct points $a = x_0 < x_1 < x_2 \ldots < x_n = b$, then we have

$$\int_{a}^{b} f(x)dx = \int_{x_{0}}^{x_{1}} f(x)dx + \int_{x_{1}}^{x_{2}} f(x)dx + \dots + \int_{x_{n-1}}^{x_{n}} f(x)dx.$$

Applying the Trapezoidal rule (5.28) for one strip to each of these integral, we have

$$\int_{a=x_0}^{b=x_n} f(x)dx \approx \frac{h}{2}[f(x_0) + f(x_1)] + \frac{h}{2}[f(x_1) + f(x_2)] + \dots + \frac{h}{2}[f(x_{n-1}) + f(x_n)].$$

Note that each of the interior point is counted twice and therefore has a coefficient of two whereas the endpoints are counted once and therefore has a coefficient one.

Example 5.17 Evaluate the integral $\int_0^1 e^{4x} dx$ by using the Trapezoidal rule with n = 1, 2, 4, 8. Also compute the corresponding absolute errors.

Solution. For n = 1, we use the formula (5.28) for h = 1, as follows

$$T_1(f) = \frac{1}{2} [f(0) + f(1)] = 27.7991.$$

For n = 2, using the formula (5.29) and h = 0.5, we have

$$T_2(f) = \frac{0.5}{2} [f(0) + 2f(0.5) + f(1)] = 17.5941.$$

For n = 4, using the formula (5.29) and h = 0.25, we have

$$T_4(f) = \frac{0.25}{2} \left[f(0) + 2[f(0.25) + f(0.5) + f(0.75)] + f(1) \right] = 14.4980.$$

Finally, for n = 8, using (5.29) and h = 0.125, we have

$$T_8(f) = \frac{0.125}{2} \Big[f(0) + 2[f(0.125) + f(0.25) + f(0.375) + f(0.5) + f(0.625) + f(0.75) + f(0.875)] + f(1) \Big] = 13.6776.$$

Since the exact value of the given integral is

$$I(f) = \frac{1}{4}[e^4 - 1] = 13.4000.$$

So the corresponding absolute errors are, 14.3991, 4.1941, 1.0980 and 0.2776, respectively, which decrease by a factor of about *four* at each stage.

Example 5.18 Suppose that $f(0.25) = f(0.75) = \alpha$. Find α if the composite Trapezoidal rule with n = 2 gives the value 2 for $\int_0^1 f(x) dx$ and with n = 4 gives the value of 1.75.

Solution. For n = 2, using the formula (5.29) and h = 0.5, we have

$$\int_0^1 f(x) \ dx = 2 \approx T_2(f) = \frac{0.5}{2} \Big[f(0) + 2f(0.5) + f(1) \Big],$$

which is equivalent to

$$8 \approx f(0) + 2f(0.5) + f(1)$$
.

For n = 4, using the formula (5.29) and h = 0.25, we have

$$\int_{0}^{1} f(x) dx = 1.75 \approx T_4(f) = \frac{0.25}{2} \left[f(0) + 2(2\alpha) + 2f(0.5) + f(1) \right],$$

which is equals to

$$8(1.75) \approx f(0) + 2f(0.5) + f(1) + 4\alpha$$
, or $8(1.75) \approx 8 + 4\alpha$.

(using $8 \approx f(0) + 2f(0.5) + f(1)$). Solving for α , we get $\alpha \approx 1.5$, the required value.

Note that the Trapezoidal rule of integration involves no restriction relative to number of data points involved. This is not the case with elaborate methods yet to discussed. This is one reason why it is one of the favorite numerical integration method used in mathematics and engineering.

Error Terms for Trapezoidal Rule

Now we discuss the error formulas for the Trapezoidal rule. This formula will lead to a better understanding of the method, showing both their weakness and strengths, and it will allow improvements of the method. We discuss the error for the simple Trapezoidal rule (5.28) in the from of the following theorem and then we use it to define the error for the composite Trapezoidal rule (5.29).

Theorem 5.3 (Error term for Simple Trapezoidal Rule)

Let $f \in C^2[a,b]$, and h = (b-a). The local error that the simple Trapezoidal rule (5.28) makes in estimating the definite integral (5.21) is

$$E_{T_1}(f) = -\frac{h^3}{12}f''(\eta(x)), \tag{5.30}$$

where $\eta(x) \in (a,b)$.

Proof. Consider two points $a = x_0 < x_1 = b$ with $h = x_1 - x_0$. From the linear Lagrange interpolation formula with error term, we have

$$f(x) = p_1(x) + \frac{f''(\eta(x)(x))}{2!} \prod_{i=0}^{1} (x - x_i).$$
 (5.31)

By integrating (5.31) with respect to x and between x_0 and x_1 , we have

$$\int_{x_0}^{x_1} f(x)dx = \int_{x_0}^{x_1} p_1(x)dx + \frac{1}{2} \int_{x_0}^{x_1} f''(\eta(x)) \prod_{i=0}^{1} (x - x_i)dx.$$
 (5.32)

The error term for the Trapezoidal rule of one strip can be obtained as follows:

$$E_{T_1}(f) = \frac{1}{2} \int_{x_0}^{x_1} f''(\eta(x))(x - x_0)(x - x_1) dx.$$

Note that $f''(\eta(x))$ is a continuous function of x and the term $(x-x_0)(x-x_1)$ is negative on (a,b), therefore, by using the Mean Value Theorem for integrals, we have

$$E_{T_1}(f) = \frac{f''(\eta(x))}{2} \int_{x_0}^{x_1} (x - x_0)(x - x_1) dx.$$

Now to solve the integral on the right side of above equation, we use the change of variable

$$x - x_0 = uh$$
, $x - x_1 = (u - 1)h$, and $dx = hdu$

then we have

$$E_{T_1}(f) = \frac{f''(\eta(x))}{2} \int_0^1 huh(u-1)hdu = \frac{f''(\eta(x))h^3}{2} \int_0^1 (u^2 - u)du.$$

Thus

$$E_{T_1}(f) = -\frac{h^3 f''(\eta(x))}{12},\tag{5.33}$$

for some $\eta(x) \in (a,b)$. Then the formula (5.32) can be written as

$$\int_{a=r_0}^{b=x_1} f(x)dx = \frac{h}{2}[f(x_0) + f(x_1)] - \frac{h^3}{12}f''(\eta(x)), \tag{5.34}$$

for $\eta(x) \in (a,b)$, which is the simple Trapezoidal rule with its error term.

The formula (5.33) indicates that the local error of the Trapezoidal rule is proportional to second derivative f''. So, if the Trapezoidal rule is used to integrate each of $f(x) = 1, x, x^2, x^3, \ldots$, then results have no error for f(x) = 1 and f(x) = x but there are error for x^2 and higher power of x.

Example 5.19 Compute the local error for the Trapezoidal rule (5.28) using the integral

$$\int_1^2 \frac{1}{x+1} dx.$$

Solution. Given $f(x) = \frac{1}{x+1}$, and [a,b] = [1,2], then the second derivative of the function is

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

Since the error formula for the simple Trapezoidal rule is

$$E_{T_1}(f) = -\frac{h^3}{12}f''(\eta(x)), \quad where \quad \eta(x) \in (1,2).$$

This formula cannot be computed exactly because $\eta(x)$ is not known. But one can bound the error by computing the largest possible value for $|f''(\eta(x))|$. Bound $|f''(\eta(x))|$ on [1, 2] is

$$M = \max_{1 \le x \le 2} \left| \frac{2}{(x+1)^3} \right| = 0.25.$$

Then, for $|f''(\eta(x))| \leq M$, we have

$$|E_{T_1}(f)| \le \frac{h^3}{12}M.$$

Using M = 0.25 and h = 1, we get

$$|E_{T_1}(f)| \le \frac{0.25}{12} = 0.0208.$$

Comparing this with the absolute error 0.0112, this bound is about 2 times the actual error.

Example 5.20 Let $f(x) = x^3 + 1$ defined on the interval [0.1, 0.2]. Find the value of unknown point $\eta(x)$ by using the local error for the simple Trapezoidal rule (5.28).

Solution. Given $f(x) = x^3 + 1$, and [a, b] = [0.1, 0.2], we use the formula (5.28) for h = 0.1, as follows

$$ApproxValue = \frac{0.1}{2} \left[f(0.1) + f(0.2) \right] = \frac{0.1}{2} \left[\left[(0.1)^3 + 1 \right] + \left[(0.2)^3 + 1 \right] \right] = 0.10045.$$

We know that

$$ExactValue = \int_{0.1}^{0.2} (x^3 + 1) \ dx = (x^4/4 + x) \Big|_{0.1}^{0.2} = [(0.2)^4/4 + 0.2] - [(0.1)^4/4 + 0.1] = 0.100375,$$

so we have the error

$$E = (ExactValue) - (ApproxValue) = 0.100375 - 0.10045 = -0.000075.$$

since the second derivative of the given function is f''(x) = 6x, so by using the local error for the Trapezoidal rule (5.28), we have

$$-0.000075 = -\frac{(0.1)^3}{12}(6\eta(x)),$$

gives the value of $\eta(x) = 0.15$.

Error Term for Composite Trapezoidal Rule

The global error of the Trapezoidal rule (5.29) equals the sum of n local errors of the Trapezoidal rule (5.28), that is

$$E_{T_n}(f) = -\frac{h^3}{12}f''(\eta_1(x)) - \frac{h^3}{12}f''(\eta_2(x)) - \dots - \frac{h^3}{12}f''(\eta_n(x)),$$

which can also written as

$$E_{T_n}(f) = -\frac{h^3}{12} \sum_{i=1}^n f''(\eta_i(x)), \text{ for } \eta_i(x) \in (x_{i-1}, x_i),$$

or

$$E_{T_n}(f) = -\frac{h^3}{12} n f''(\eta(x)),$$

where $f''(\eta(x))$ is the average of the *n* individual values of the second derivative. Since $n = \frac{b-a}{h}$, thus the global error in the composite Trapezoidal rule (5.29) is

$$E_{T_n}(f) = -\frac{h^2}{12}(b-a)f''(\eta(x)), \quad \eta(x) \in (a,b).$$
(5.35)

Hence

$$\int_{a}^{b} f(x)dx = \frac{h}{2} \left[f(a) + 2 \sum_{i=1}^{n-1} f(x_i) + f(b) \right] - \frac{h^2}{12} (b - a) f''(\eta(x)), \tag{5.36}$$

for $\eta(x) \in (a,b)$, is the composite Trapezoidal rule with its error term.

Note that whereas the simple Trapezoidal rule (5.28) has a truncation error of order h^3 , the composite Trapezoidal rule (5.29) has an error of order h^2 . This means that when h is halved and the number of subintervals is doubled the error decreases by a factor of approximately four (assuming that $f''(\eta(x))$ remains fairly constant throughout [a,b]). Of course, it is also possible to express the truncation error in terms of n rather than n. Since $n = \frac{b-a}{n}$, it follows that the global truncation error (5.35) is of order $O(n^2)$.

Example 5.21 (a) Find approximation of $\int_1^2 f(x) dx$, taking h = 0.2 by using following set of data points

The function tabulated is xe^{-x} , compute error bound and the absolute error for the approximation using Trapezoidal rule.

(b) How many subintervals approximate the given integral to an accuracy of at least 10^{-6} ?

Solution. (a) Given h = 0.2, so we have the select following set of data points for Trapezoidal rule as

so the composite Trapezoidal rule (5.29) for six points can be written as

$$\int_{1}^{2} f(x) dx \approx T_{5}(f) = \frac{h}{2} \Big[f(x_{0}) + 2 \Big(f(x_{1}) + f(x_{2}) + f(x_{3}) + f(x_{4}) \Big) + f(x_{5}) \Big],$$

and by using the given values, we get

$$\int_{1}^{2} f(x) dx \approx 0.1 \Big[0.368 + 2(0.361 + 0.355 + 0.323 + 0.298) + 0.271 \Big] = 0.3313.$$

The second derivative of the function $f(x) = xe^{-x}$ can be obtain as

$$f'(x) = (1-x)e^{-x}$$
 and $f''(x) = (x-2)e^{-x}$.

Since $\eta(x)$ is unknown point in (1,2), therefore, the bound |f''| on [1,2] is

$$M = \max_{1 \le x \le 2} |f''(x)| = \max_{1 \le x \le 2} |(x-2)e^{-x}| = 0.3679,$$

at x = 1. Thus the error formula (5.35) becomes

$$|E_{T_5}(f)| \le \frac{(0.2)^2(1)}{12}(0.3679) = 0.0012,$$

which is the possible maximum error in our approximation. We can easily computed the exact value of the given integral as

$$\int_{1}^{2} xe^{-x} dx = (-xe^{-x} - e^{-x})\Big|_{1}^{2} = 0.3298.$$

Thus the absolute error |E| in our approximation is given as

$$|E| = |0.3298 - T_5(f)| = |0.3298 - 0.3313| = 0.0015.$$

(b) To find the minimum subintervals for the given accuracy, we use the formula (5.35) such that

$$|E_{T_n}(f)| \le \frac{|-(b-a)^3|}{12n^2} M \le 10^{-6},$$

where h = (b - a)/n. Since M = 0.3679, then solving for n^2 , we obtain

$$n^2 \ge 30658.3333$$
, gives $n \ge 175.0952$.

Hence to get the required accuracy, we need 176 subintervals or 177 points.

Example 5.22 Consider the integral $I(f) = \int_1^2 \ln(x+1)dx$; n=6

- (a) Compute the approximation of the integral using the composite Trapezoidal rule.
- (b) Compute the error bound for your approximation using the formula (5.35).
- (c) Compute the absolute error.
- (d) How many subintervals approximate the given integral to an accuracy of at least 10^{-4} using the

composite Trapezoidal rule ?

Solution. (a) Given $f(x) = \ln(x+1)$, n = 6, and so $h = \frac{2-1}{6} = \frac{1}{6}$, then the composite Trapezoidal rule (5.29) for n = 6, can be written as

$$T_6(f) = \frac{1/6}{2} \left[\ln(1+1) + 2\left(\ln\left(\frac{7}{6}+1\right) + \ln\left(\frac{8}{6}+1\right) + \ln\left(\frac{9}{6}+1\right) + \ln\left(\frac{10}{6}+1\right) + \ln\left(\frac{11}{6}+1\right) \right) + \ln(2+1) \right].$$

Hence

$$\int_{1}^{2} \ln(x+1)dx \approx T_{6}(f) = \frac{1}{12}[0.6932 + 2(4.5591) + 1.0986] = 0.9092.$$

(b) The second derivative of the function can be obtain as

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

Since $\eta(x)$ is unknown point in (1,2), therefore, the bound |f''| on [1,2] is

$$M = \max_{1 \le x \le 2} |f''(x)| = \left| \frac{-1}{(x+1)^2} \right| = 0.25.$$

Thus the error formula (5.35) becomes

$$|E_{T_6}(f)| \le \frac{(1/6)^2}{12}(0.25) = 0.0006,$$

which is the possible maximum error in our approximation in part (a).

(c) The absolute error |E| in our approximation is given as

$$|E| = |(3 \ln 3 - 2 \ln 2 - 1) - T_6(f)| = |0.9095 - 0.9092| = 0.0003.$$

(d) To find the minimum subintervals for the given accuracy, we use the formula (5.35) such that

$$|E_{T_n}(f)| \le \frac{|-(b-a)^3|}{12n^2} M \le 10^{-4},$$

where h = (b - a)/n. Since M = 0.25, then solving for n^2 , we obtain

$$n^2 \ge 208.3333$$
, gives $n \ge 14.4338$.

Hence to get the required accuracy, we need 15 subintervals.

To use MATLAB command for the composite Trapezoidal rule, first we define a function m-file as fn.m for the function as follows:

function
$$y = fn(x)$$

 $y = \log(x+1);$
 $>> T6 = TrapezoidalR('fn', 1, 2, 6)$

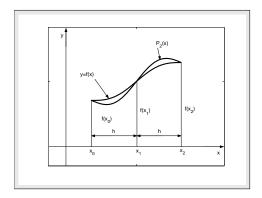


Figure 5.8: Simple Simpson's rule.

Program 5.1

MATLAB m-file for the Composite Trapezoidal Rule

function TN=TrapezoidR(fn,a,b,n);

h=(b-a)/n; s=(feval(fn,a)+feval(fn,b))/2;

for k=1:n-1; x = a + h * k s=s+feval(fn,x); end TN = s * h;

Suppose that, due to rounding, the f_i are in error by at most $\frac{1}{2} \times 10^{-t}$. Then we see from (5.29) that the error in the composite Trapezoidal rule due to rounding is not greater than

$$\frac{1}{2}h[1+2+2+\cdots+2+1]\times 10^{-t} = nh(\frac{1}{2}\times 10^{-t}) = \frac{1}{2}(b-a)\times 10^{-t}.$$

Thus, rounding errors do not seriously affect the accuracy of this quadrature rule. This is generally true of numerical integration unlike numerical differentiation, as we saw earlier in the chapter.

Simple Simpson's Rule

The Trapezoidal rule approximates the area under a curve by the area of trapezoid formed by connecting two points on the curve by straight line. The Simpson's rule gives a more accurate approximation since it consists of connecting three points on the curve by second-degree parabola and the area under the parabola to obtain the approximate area under the curve, see Figure 5.8. Let us consider the second-degree Lagrange interpolating polynomial, with equally spaced base points, that is, $x_0 = a, x_1 = a + h$ and $x_2 = a + 2h$, with h = (b - a)/2, then

$$f(x) = p_2(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} f(x_0) + \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} f(x_1) + \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} f(x_2).$$

Taking integral on both sides of the above equation with respect to x between the limits x_0 and

 x_2 , we have

$$\int_{x_0}^{x_2} f(x)dx \approx \frac{f(x_0)}{(x_0 - x_1)(x_0 - x_2)} \int_{x_0}^{x_2} (x - x_1)(x - x_2)dx
+ \frac{f(x_1)}{(x_1 - x_0)(x_1 - x_2)} \int_{x_0}^{x_2} (x - x_0)(x - x_2)dx
+ \frac{f(x_2)}{(x_2 - x_0)(x_2 - x_1)} \int_{x_0}^{x_2} (x - x_0)(x - x_1)dx,$$

which implies that

$$\int_{a}^{b} f(x)dx \approx \frac{f(x_0)}{2h^2} I_1 + \frac{f(x_1)}{-h^2} I_2 + \frac{f(x_2)}{2h^2} I_3,$$

where

$$I_1 = \int_{x_0}^{x_2} (x - x_1)(x - x_2) dx; \ I_2 = \int_{x_0}^{x_2} (x - x_0)(x - x_2) dx; \ I_3 = \int_{x_0}^{x_2} (x - x_0)(x - x_1) dx.$$

Solving above three integrals by using integration by parts, we obtain the values of I_1 , I_2 and I_3 as follows

$$I_1 = \frac{2h^3}{3}$$
, $I_2 = -\frac{4h^3}{3}$, $I_3 = \frac{2h^3}{3}$.

By using these values, we have

$$\int_{a}^{b} f(x)dx \approx \frac{f(x_0)}{2h^2} \left(\frac{2h^3}{3}\right) + \frac{f(x_1)}{-h^2} \left(\frac{-4h^3}{3}\right) + \frac{f(x_2)}{2h^2} \left(\frac{2h^3}{3}\right).$$

Simplifying, gives

$$\int_{a}^{b} f(x)dx \approx S_{2}(f) = \frac{h}{3}[f(x_{0}) + 4f(x_{1}) + f(x_{2})]. \tag{5.37}$$

which is called the *simple Simpson's rule* or Simpson's rule for two strips (or 3 points).

Example 5.23 Approximate the following integral

$$\int_{1}^{2} \frac{1}{x+1} dx,$$

using simple Simpson's rule. Compute the actual error.

Solution. Since $f(x) = \frac{1}{x+1}$ and h = (2-1)/2 = 0.5, then by using Simpson's rule (5.37), we have

$$S_2(f) = \frac{0.5}{3} [f(1) + 4f(1.5) + f(2)] = (0.1667)[0.5 + 1.6 + 0.3333] = 0.4056.$$

Hence

$$\int_{1}^{2} \frac{1}{x+1} dx \approx S_2(f) = 0.4056.$$

Since the exact solution of the given integral is, 0.4055, therefore, the actual error is

$$E_{S_2} = I(f) - S_2(f) = -0.0001.$$

To compare this error with the error got by using the simple Trapezoidal rule, the error in Simpson's rule is much smaller than for the Trapezoidal rule by a factor of about 123, a significant increase in accuracy.

Example 5.24 Use simple Simpson's rule to show that

$$\int_{1}^{1.6} \frac{2}{x} dx < 1 < \int_{1}^{1.7} \frac{2}{x} dx.$$

Solution. Given $f(x) = \frac{2}{x}$ and take h = (1.6 - 1)/2 = 0.3, then by using Simpson's rule (5.37), we have

$$S_2(f) = \frac{0.3}{3} [f(1) + 4f(1.3) + f(1.6)] = (0.1)[2 + 6.1538 + 1.25] = 0.9404.$$

Now taking h = (1.7 - 1)/2 = 0.35, then by using Simpson's rule (5.37), we have

$$S_2(f) = \frac{0.35}{3} [f(1) + 4f(1.35) + f(1.7)] = (0.1167)[2 + 5.9260 + 1.1764] = 1.0623.$$

Hence

the required result.

Example 5.25 Let f be defined by

$$f(x) = \begin{cases} x^2 - x + 1, & \text{if } 0 \le x \le 1, \\ 2x - 1, & \text{if } 1 \le x \le 2. \end{cases}$$

Approximate the integral $\int_0^2 f(x)dx$ by using Simpson's rule with n=2.

Solution. Since one can know that

$$\int_{0}^{2} f(x)dx = \int_{0}^{1} f(x)dx + \int_{1}^{2} f(x)dx,$$

and we are given

$$\int_0^2 f(x)dx = \int_0^1 (x^2 - x + 1)dx + \int_1^2 (2x - 1)dx.$$

First we find the approximation of the first integral on the right hand side of above equation for n = 2, using the formula (5.37) and h = 0.5, we have

$$I_1(f) = \frac{0.5}{3} [f(0) + 4f(0.5) + f(1)] = \frac{0.5}{3} [1 + 3 + 1] = 0.8333.$$

Now we find the approximation of the second integral on the right hand side of above equation for n = 2, using the formula (5.37) and h = 0.5, we have

$$I_2(f) = \frac{0.5}{3} [f(1) + 4f(1.5) + f(2)] = \frac{0.5}{3} [1 + 8 + 3] = 2.0000.$$

Hence

$$\int_0^2 f(x)dx = I_1(f) + I_2(f) \approx 0.8333 + 2.000 = 2.83333,$$

the required approximation of the given integral.

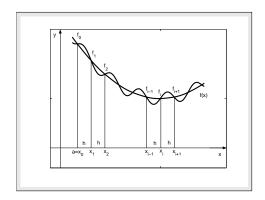


Figure 5.9: Composite Simpson's Rule.

Composite Simpson's Rule

Just as with the simple Trapezoidal rule (5.28), the simple Simpson's rule (5.37) can be improved by dividing the integration interval [a, b] into a number of subintervals of equal width h; where $h = \frac{b-a}{n}$. Since the simple Simpson's rule (5.37) requires a interval consisting of three points (pair of strips). In practice, we usually take more than three points and add the separate results for the different pairs of strips (see Figure 5.9). Since the simple Simpson's rule requires a pair of strips for application, so there is restriction on the integer n, which must be even. We define the composite Simpson's rule in the form of the following theorem.

Theorem 5.4 (Composite Simpson's Rule)

Let $f \in C^4[a,b]$, n be even, h = (b-a)/n, and $x_i = a+ih$ for each $i = 0,1,2,\ldots,n$. Then the composite Simpson's rule for n subintervals can be written as

$$\int_{a}^{b} f(x)dx \approx S_{n}(f) = \frac{h}{3} \left[f(a) + 2 \sum_{i=1}^{n/2-1} f(x_{2i}) + 4 \sum_{i=1}^{n/2} f(x_{2i-1}) + f(b) \right].$$
 (5.38)

Proof. Since for the composite form of the Simpson's rule, the interval is divided into n equal subintervals of width h so that $h = \frac{b-a}{n}$. For this rule to work, n must be even number and the total number of (n+1) distinct points $a = x_0 < x_1 < x_2 ... < x_n = b$ should be odd. The total integral can be represented as

$$\int_{x_0}^{x_n} f(x)dx = \int_{x_0}^{x_2} f(x)dx + \int_{x_2}^{x_4} f(x)dx + \dots + \int_{x_{n-2}}^{x_n} f(x)dx.$$

Substitute the simple Simpson's rule (5.37) for the individual integral yields

$$\int_{x_0}^{x_n} f(x)dx \approx \frac{h}{3} [f(x_0) + 4f(x_1) + f(x_2)] + \frac{h}{3} [f(x_2) + 4f(x_3) + f(x_4)]$$

$$+ \dots + \frac{h}{3} [f(x_{n-2}) + 4f(x_{n-1}) + f(x_n)].$$

To avoid the repetition of terms, we summed them. Note that each of the odd interior point is counted four and therefore has a coefficient of four whereas each of the even interior point is counted two and therefore has a coefficient of two. The endpoints are counted once and therefore has a coefficient one.

Example 5.26 Let f be defined by

$$f(x) = \begin{cases} x^2 - x + 1, & \text{if } 0 \le x \le 1, \\ 2x - 1, & \text{if } 1 \le x \le 2. \end{cases}$$

Approximate the integral $\int_0^2 f(x)dx$ by using Simpson's rule with n=4.

Solution. Since one can know that

$$\int_0^2 f(x)dx = \int_0^1 f(x)dx + \int_1^2 f(x)dx,$$

and we are given

$$\int_0^2 f(x)dx = \int_0^1 (x^2 - x + 1)dx + \int_1^2 (2x - 1)dx = I_1(f) + I_2(f).$$

First we find the approximation of the first integral on the right hand side of above equation for n = 4, using the formula (5.38) and h = 0.25, we have

$$I_1(f) \approx \frac{0.5}{3} \Big[f(0) + 4(f(0.25) + f(0.75)) + 2f(0.5) + f(1) \Big] \approx \frac{0.25}{3} \Big[8 \Big] \approx 0.8333.$$

Now we find the approximation of the second integral on the right hand side of above equation for n = 4, using the formula (5.38) and h = 0.25, we have

$$I_2(f) \approx \frac{0.5}{3} \Big[f(1) + 4(f(1.25) + f(1.75)) + 2f(1.5) + f(2) \Big] \approx \frac{0.25}{3} \Big[24 \Big] \approx 2.0000.$$

Hence

$$\int_{0}^{2} f(x)dx = I_{1}(f) + I_{2}(f) \approx 0.8333 + 2.000 = 2.83333,$$

the required approximation of the given integral.

Example 5.27 Suppose that f(1) = 0.5, f(1.2) = 0.9, $[f(1.25) + f(1.75)] = \alpha$, f(1.5) = 1.5, f(1.6) = 1.65, f(1.95) = 1.95 and f(2) = 2. Find the approximate value of α if the best composite Simpson's rule gives the value, 1.35, for the integral $\int_{1}^{2} f(x) \ dx$.

Solution. Since we need equally spaced data points, so take f(1) = 0.5, $[f(1.25) + f(1.75)] = \alpha$, f(1.5) = 1.5 and f(2) = 2, gives n = 4, so $h = \frac{2-1}{4} = 0.25$. By using the composite formula (5.38) for n = 4, we have

$$\int_{1}^{2} f(x) dx \approx \frac{0.25}{3} \Big[f(1) + 4[f(1.25) + f(1.75)] + 2f(1.5) + f(2) \Big].$$

Now using the given values, we obtain

$$1.35 \approx \frac{1}{12}[0.5 + 4(\alpha) + 2(1.5) + 2], \quad or \quad 12(1.35) - 5.5 \approx 4\alpha,$$

gives $\alpha \approx 2.675$.

Example 5.28 Evaluate the integral $\int_0^1 e^{4x} dx$ by using the Simpson's rule with n = 2, 4, 8. Also, compute the corresponding actual errors.

Solution. For n = 2, using the formula (5.37) and h = 0.5, we have

$$S_2(f) = \frac{0.5}{3} [f(0) + 4f(0.5) + f(1)] = 14.1924.$$

For n = 4, using the formula (5.38) and h = 0.25, we have

$$S_4(f) = \frac{0.25}{3} \left[f(0) + 4[f(0.25) + f(0.75)] + 2f(0.5) + f(1) \right] = 13.4659.$$

For n = 8, using the formula (5.38) and h = 0.125, we have

$$S_8(f) = \frac{0.125}{3} \Big[f(0) + 4[f(0.125) + f(0.375) + f(0.625) + f(0.875)] + 2[f(0.25) + f(0.5) + f(0.75)] + f(1) \Big] = 13.4041.$$

Note that the exact value of the given integral is 13.39995, and so the corresponding errors are, 0.79245, 0.06595 and 0.00411 respectively, which decrease by a factor of about 16 at each stage. \bullet

To use the MATLAB command for the composite Simpson's rule, first we define a function m-file as fn.m for the function as follows:

```
function y = fn(x)

y = exp(4 * x);

>> S2 = SimpsonR('fn', 0, 1, 2)
```

Program 5.2

MATLAB m-file for the Composite Simpson's Rule function SN=SimpsonR(fn,a,b,n);

h=(b-a)/n; s=feval(fn,a)+feval(fn,b);

for k=1:2:n-1; s = s + 4 * feval(fn, a + k * h); end

for k=2:2:n-2; s = s + 2*feval(fn, a + k*h); end; SN = (s*h)/3;

One can note that these results are more accurate than those obtained in the previous Example 5.17 using the Trapezoidal rule (for same number of function evaluations). The main disadvantage of the Simpson's rule is that it can only be used when the given interval [a, b] is divided in to an even number of subintervals.

Error Terms for Simpson's Rule

Now we discuss the local error and the global error formulas for Simpson's rule.

Theorem 5.5 (Error Term for Simple Simpson's Rule)

Let $f \in C^4[a,b]$, and h = (b-a)/2. The local error that the Simpson's rule makes in estimating the definite integral (5.21) is

$$E_{S_2}(f) = -\frac{h^5}{90} f^{(4)}(\eta(x)), \tag{5.39}$$

where $\eta(x) \in (a,b)$.

Proof. Consider three equally spaced points $a = x_0 < x_1 < x_2 = b$. From the quadratic Lagrange interpolation formula with error term, we have

$$f(x) = p_2(x) + \frac{f'''(\eta(x))}{3!} \prod_{i=0}^{2} (x - x_i),$$
 (5.40)

and by integrating (5.40) with respect to x, we have

$$\int_{a}^{b} f(x)dx = \int_{a}^{b} p_{2}(x)dx + \frac{1}{6} \int_{a}^{b} f'''(\eta(x)(x)) \prod_{i=0}^{2} (x - x_{i})dx.$$
 (5.41)

The second term on the right hand side of the (5.41)

$$E_{S_2}(f) = \frac{1}{6} \int_a^b f'''(\eta(x)(x))(x - x_0)(x - x_1)(x - x_2)dx, \tag{5.42}$$

is called the error term of the Simpson's rule for n=2.

In this way it provides only an $O(h^4)$ error term, involving $f'''(\eta(x))$. By approaching the problem in another way, a higher-order term involving $f^{(4)}(\eta(x))$ can be obtained.

Consider the two intervals from x_{i-1} to x_i and x_i to x_{i+1} , with $h = (x_i - x_{i-1})$, and also, assume that F(x) is the indefinite integral of a function f(x) which we are trying to integrate, then the exact value of the integral from x_{i-1} to x_{i+1} is

$$I_i(f) = \int_{x_{i-1}}^{x_{i+1}} f(x)dx = F(x_{i+1}) - F(x_{i-1}). \tag{5.43}$$

The approximated value calculated by using the Simpson's rule is

$$S_{2i}(f) = \frac{h}{3} [f(x_{i-1}) + 4f(x_i) + f(x_{i+1})]. \tag{5.44}$$

Then the error define in using the Simpson's rule on this two intervals is

$$E_i(f) = I_i(f) - S_{2i}(f) (5.45)$$

Now by expanding f(x) about $x = x_i$ to get $f(x_{i-1})$ in terms of a function and derivatives at $x = x_i$ by using the Taylor's series, we have

$$f(x_{i-1}) = f(x_i) + (x_{i-1} - x_i)f'(x_i) + \frac{(x_{i-1} - x_i)^2}{2!}f''(x_i) + \frac{(x_{i-1} - x_i)^3}{3!}f'''(x_i) + \frac{(x_{i-1} - x_i)^4}{4!}f^{(4)}(x_i) + \cdots$$
(5.46)

Since we know that $h = (x_i - x_{i-1})$, or $-h = (x_{i-1} - x_i)$, therefore

$$f(x_{i-1}) = f(x_i) - hf'(x_i) + \frac{h^2}{2!}f''(x_i) - \frac{h^3}{3!}f'''(x_i) + \frac{h^4}{4!}f^{(4)}(x_i) - \cdots$$
 (5.47)

Similar way we expand $f(x_{i+1})$ about x_i , we get

$$f(x_{i+1}) = f(x_i) + (x_{i+1} - x_i)f'(x_i) + \frac{(x_{i+1} - x_i)^2}{2!}f''(x_i) + \frac{(x_{i+1} - x_i)^3}{3!}f'''(x_i) + \frac{(x_{i+1} - x_i)^4}{4!}f^{(4)}(x_i) + \cdots,$$
 (5.48)

or

$$f(x_{i+1}) = f(x_i) + hf'(x_i) + \frac{h^2}{2!}f''(x_i) + \frac{h^3}{3!}f'''(x_i) + \frac{h^4}{4!}f^{(4)}(x_i) + \cdots$$
 (5.49)

So by using (5.47) and (5.49), we get (5.44) of the form

$$S_{2i}(f) = \frac{h}{3} \Big[\{ f(x_i) - hf'(x_i) + \frac{h^2}{2!} f''(x_i) - \frac{h^3}{3!} f'''(x_i) + \frac{h^4}{4!} f^{(4)}(x_i) - \dots \} + \{ 4f(x_i) \}$$

$$+ \{ f(x_i) + hf'(x_i) + \frac{h^2}{2!} f''(x_i) + \frac{h^3}{3!} f'''(x_i) + \frac{h^4}{4!} f^{(4)}(x_i) + \dots \} \Big],$$

which implies that

$$S_{2i}(f) = \frac{h}{3} \left[6f(x_i) + \frac{2h^2}{2!} f''(x_i) + \frac{2h^4}{4!} f^{(4)}(x_i) + \cdots \right]$$

$$= 2hf(x_i) + 2f''(x_i) \frac{h^3}{3!} + 2f^{(4)}(x_i) \frac{h^5}{3(4!)} + \cdots$$
(5.50)

Now we use the same procedure for F(x) as we used for f(x) to get

$$F(x_{i-1}) = F(x_i) - hF'(x_i) + \frac{h^2}{2!}F''(x_i) - \frac{h^3}{3!}F'''(x_i) + \frac{h^4}{4!}F^{(4)}(x_i) - \frac{h^5}{5!}F^{(5)}(x_i) + \cdots,$$
 (5.51)

and

$$F(x_{i+1}) = F(x_i) + hF'(x_i) + \frac{h^2}{2!}F''(x_i) + \frac{h^3}{3!}F'''(x_i) + \frac{h^4}{4!}F^{(4)}(x_i) + \frac{h^5}{5!}F^{(5)}(x_i) + \cdots$$
 (5.52)

Then by using (5.51) and (5.52), we got (5.43) of the form

$$I_i(f) = 2hF'(x_i) + 2F'''(x_i)\frac{h^3}{3!} + 2F^{(5)}(x_i)\frac{h^5}{5!} + \cdots$$
 (5.53)

But we know that

$$F'(x) = f(x), \ F''(x) = f'(x), \ F'''(x) = f''(x), \ F^{(4)}(x) = f'''(x), \ F^{(5)}(x) = f^{(4)}(x),$$

therefore, (5.53) can be written as

$$I_i(f) = 2hf(x_i) + 2f''(x_i)\frac{h^3}{3!} + 2f^{(4)}(x_i)\frac{h^5}{5!} + \cdots$$
 (5.54)

So by using (5.50) and (5.54), we get (5.45) of the form

$$E_{i}(f) = -\left[2hf(x_{i}) + 2f''(x_{i})\frac{h^{3}}{3!} + 2f^{(4)}(x_{i})\frac{h^{5}}{3(4!)} + \cdots\right]$$

$$+ \left[2hf(x_{i}) + 2f''(x_{i})\frac{h^{3}}{3!} + 2f^{(4)}(x_{i})\frac{h^{5}}{5!} + \cdots\right]$$

$$= -2f^{(4)}(x_{i})\frac{h^{5}}{3(4!)} + 2f^{(4)}(x_{i})\frac{h^{5}}{5!} + \text{higher terms in } h^{7} + \cdots$$

$$(5.55)$$

Assuming that h is small, we may neglect the terms in h^7 and above, and get the approximate error as follows:

$$E_i(f) \approx f^{(4)}(x_i)h^5\left[\frac{2}{5!} - \frac{2}{3(4!)}\right] \approx -\frac{1}{90}f^{(4)}(x_i)h^5,$$
 (5.56)

which is the desired local error for Simpson's rule.

Since (5.56) indicates that the error of the Simpson's rule is proportional to fourth derivative $f^{(4)}$. If the Simpson's rule is used to integrate $f(x) = 1, x, x^2$ and x^3 , then the results have no error. In more general terms, the Newton-Cotes closed formula of odd order n is exact if the integrand is a polynomial of order n or less, whereas that of an even n is exact when the integrand is a polynomial of order n + 1 or less.

Example 5.29 Compute the local error for the Simpson's rule using the following integral

$$\int_{1}^{2} \frac{1}{x+1} dx.$$

Solution. Given $f(x) = \frac{1}{x+1}$, and [a,b] = [1,2], then the fourth derivative of the function can be obtain as

$$f' = \frac{-1}{(x+1)^2}$$
, $f'' = \frac{2}{(x+1)^3}$, $f''' = \frac{-6}{(x+1)^4}$, $f^{(4)} = \frac{24}{(x+1)^5}$.

Since the error formula for the Simpson's rule is

$$E_{S_2}(f) = -\frac{h^5}{90} f^{(4)}(\eta(x)), \quad where \quad \eta(x) \in (1, 2),$$

or

$$|E_{S_2}(f)| = \left| -\frac{h^5}{90} \right| \left| f^{(4)}(\eta(x)) \right|, \quad for \quad \eta(x) \in (1, 2).$$

This formula cannot be computed exactly because $\eta(x)$ is not known. But one can bound the error by computing the largest possible value for $|f^{(4)}|$. Bound $|f^{(4)}|$ on [1,2] is

$$M = \max_{1 \le x \le 2} = \left| \frac{24}{(x+1)^5} \right| = 0.75.$$

Then for $|f^{(4)}(\eta(x))| \leq M$, we have

$$|E_{S_2}(f)| \le \frac{h^5}{90}M.$$

Taking M = 0.75 and h = 0.5, we get

$$|E_{S_2}(f)| \le \frac{(0.03125)}{90}(0.75) = 0.0003.$$

Comparing this with the actual error -0.0001, this bound is about 3 times the actual error.

Error Term for Composite Simpson's Rule

Since the composite Simpson's rule (5.38) requires that the given interval [a, b] is divided into even number of subintervals and each application of the simple Simpson's rule requires two subintervals, therefore, the global error of the composite Simpson's rule (5.38) is the sum of $\frac{n}{2}$ local truncation error of the simple Simpson's rule with $n = \frac{b-a}{2}$ that is

error of the simple Simpson's rule with
$$n = \frac{b-a}{h}$$
, that is,

$$E_{S_n}(f) = -\frac{h^5}{90} f^{(4)}(\eta_1(x)) - \frac{h^5}{90} f^{(4)}(\eta_2(x)) - \dots - \frac{h^5}{90} f^{(4)}(\eta_{n/2}(x)),$$

which implies that

$$E_{S_n}(f) = -\frac{h^5}{90} \left(\frac{n}{2}\right) \left[\frac{\sum_{i=1}^{n/2} f^{(4)}(\eta_i(x))}{n/2} \right].$$

Thus by using the Intermediate Value Theorem, we have

$$E_{S_n}(f) = -\frac{(b-a)}{180} h^4 f^{(4)}(\eta(x)), \tag{5.57}$$

for $\eta(x) \in (a,b)$ and nh = b - a. Then the formula (5.57) is known as the *global error* of the Simpson's rule.

Note that the truncation error in the composite Simpson's rule is of order h^4 . This means that when h is halved and the number of subintervals is doubled the error decreases by a factor of approximately 16, considerably better than the composite Trapezoidal rule.

Example 5.30 (a) Find the approximation of $\int_0^{1.2} f(x) dx$ taking h = 0.3 by using the following set of data points:

The function tabulated is $x + \cos x$, compute error bound and absolute error for the approximation by Simpson's rule.

(b) How many subintervals approximate the given integral to within accuracy of 10^{-6} ?

Solution. (a) Given h = 0.3, so to select the following set of data points for Simpson's rule as

so the composite Simpson's rule (5.38) for five points can be written as

$$\int_0^{1.2} f(x) dx \approx S_4(f) = \frac{h}{2} \Big[f(x_0) + 4(f(x_1) + f(x_3)) + 2f(x_2) + f(x_4) \Big],$$

and by using the given values, we get

$$\int_0^{1.2} f(x) dx \approx 0.1 \left[1.0 + 4(1.26 + 1.52) + 2(1.43) + 1.56 \right] = 1.6521.$$

The fourth derivative of the function $f(x) = x + \cos x$ can be obtain as

$$f'(x) = 1 - \sin x$$
, $f''(x) = -\cos x$, $f'''(x) = \sin x$, $f^{(4)}(x) = \cos x$.

Since $\eta(x)$ is unknown point in (0,1.2), therefore, the bound $|f^{(4)}|$ on [0,1.2] is

$$M = \max_{0 \le x \le 1.2} |f^{(4)}| = \max_{0 \le x \le 1.2} |\cos x| = 1.0,$$

at x = 1. Thus the error formula (5.57) becomes

$$|E_{S_4}(f)| \le \frac{(0.2)^4(1.2)}{180}(1.0) = 0.0001,$$

which is the possible maximum error in our approximation.

We can easily computed the exact value of the given integral as

$$\int_0^{1.2} (x + \cos x) \ dx = (x^2/2 + \sin x) \Big|_0^{1.2} = 1.6520.$$

Thus the absolute error |E| in our approximation is given as

$$|E| = |0.3298 - T_5(f)| = |1.6520 - 1.6521| = 0.0001.$$

(b) To find the minimum subintervals for the given accuracy, we use the formula (5.57) such that

$$|E_{S_n}(f)| \le \frac{|-(b-a)^5|}{180n^4} M \le 10^{-6},$$

where h = (b - a)/n. Since M = 1., then solving for n^4 , we obtain

$$n^4 \ge 16589$$
, gives $n \ge 11.3489$.

Hence to get the required accuracy, we need 12 subintervals or 13 points.

Example 5.31 Consider the integral $I(f) = \int_1^2 \ln(x+1)dx;$ n=6.

- (a) Find the approximation of the give integral using the composite Simpson's rule.
- (b) Compute the error bound for the approximation using the formula (5.57).
- (c) Compute the absolute error.
- (d) How many subintervals approximate the given integral to an accuracy of at least 10^{-4} using the composite Simpson's rule?

Solution. (a) Given $f(x) = \ln(x+1)$, n = 6, and so $h = \frac{2-1}{6} = \frac{1}{6}$, then the composite Simpson's rule (5.38) for n = 6, can be written as

$$S_6(f) = \frac{1/6}{3} \left[\ln(1+1) + 4(\ln\left(\frac{7}{6}+1\right) + \ln\left(\frac{9}{6}+1\right) + \ln\left(\frac{11}{6}+1\right) \right] + \left[2\left(\ln\left(\frac{8}{6}+1\right) + \ln\left(\frac{10}{6}+1\right)\right) + \ln(2+1) \right].$$

Hence

$$\int_{1}^{2} \ln(x+1)dx \approx S_{6}(f) = \frac{1}{18} \Big[0.6932 + 4(2.7309) + 2(1.8281) + 1.0986 \Big] = 0.9095.$$

(b) Since the fourth derivative of the function is

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

Since $\eta(x)$ is unknown point in (1,2), therefore, the bound $|f^{(4)}|$ on [1,2] is

$$M = \max_{1 \le x \le 2} |f^{(4)}(x)| = \left| \frac{-6}{(x+1)^4} \right| = 6/16 = 0.375.$$

Thus the error formula (5.57) becomes

$$|E_{T_6}(f)| \le \frac{(1/6)^4}{180}(0.375) = 0.000002,$$

which is the possible maximum error in our approximation in part (a).

(c) The absolute error |E| in our approximation is given as

$$|E| = |3 \ln 3 - 2 \ln 2 - 1 - S_6(f)| = 0.0000003.$$

(d) To find the minimum subintervals for the given accuracy, we use the error formula (5.57) which is

$$|E_{S_n}(f)| \le \frac{(b-a)^5}{180n^4} M \le 10^{-4}.$$

Since we know M = 0.375, then we have

$$n^4 \ge 20.83333$$
, gives $n \ge 2.136435032$.

Hence to get the required accuracy, we need 4 subintervals (because n should be even) that ensures the stipulated accuracy.

288 5.11 Exercises

To find the approximate value of the given integral within the given accuracy, we use MATLAB command which gives us the approximate solution and the number of subintervals n. First we define m-file as fn.m for the function, so after finding the value of M, it simply compute n using (5.57) and then calls the previously defined SimpsonR function, we have the results:

```
function y = fn(x)

y = \log(x+1);

>> k = ErrorSR('fn', 0, 1, 0.375, 1e-4)
```

Program 5.3

MATLAB m-file for computing Error term of the Composite Simpson's Rule

$$\begin{split} L &= abs(b-a); n = ceil(L*sqrt(sqrt(L*M/180/eps))); \\ &\text{if } mod(n,2) == 1; n = n+1; \text{ end}; \text{ k=SimpsonR(fn,a,b,n)}; \end{split}$$

Example 5.32 Determine the number of subintervals n required to approximate

$$I(f) = \int_0^2 \frac{1}{x+4} dx,$$

with an error less than 10^{-4} using Simpson's rule.

Solution. we have to use the error formula (5.57) which is

$$|E_{S_n}(f)| \le \frac{(b-a)}{180} h^4 M \le 10^{-4}.$$

Given the integrand is $f(x) = \frac{1}{x+4}$, and we have $f^{(4)}(x) = \frac{24}{(x+4)^5}$. The maximum value of $|f^{(4)}(x)|$ on the interval [0,2] is 3/128, and thus $M = \frac{3}{128}$. Using the above error formula, we get

$$\frac{3}{(90 \times 128)} h^4 \le 10^{-4}, \quad or \quad h \le \frac{2}{5} \sqrt[4]{15} = 0.7872.$$

Since $n = \frac{2}{h} = \frac{2}{0.7872} = 2.5407$, so the number of subintervals n required is $n \ge 4$.

5.7 Exercises

- 1. Let $f(x) = (x-1)e^x$ and take h = 0.01. Calculate approximation to f'(2.3) using the two-point forward-difference formula and compute the actual error and an error bound for you approximation.
- 2. Solve the Problem 1 for the $f(x) = (x^2 + x + 1)e^{2x}$ with h = 0.05.

- 3. By using the following data: (1.2, 11.6), (1.29, 13.8), (1.3, 14), (1.31, 14.3), (1.4, 16.8). Compute the best approximations of f'(1.3) using the two-point forward-difference formula.
- 4. Let $f(x) = \sin(x+1)$. Compute the approximation of $f'(\frac{\pi}{4})$ using the two-point formula. Compute actual error and error bound.
- 5. Use the three-point central-difference formula to compute the approximate value for f'(5) with $f(x) = (x^2 + 1) \ln x$ and h = 0.05. Compute the actual error and the error bound for you approximation.
- 6. Use the three-point central-difference formula to compute the approximate value for f'(2) with $f(x) = e^{x/2} + 2\cos x$, and h = 0.01. Compute actual error and error bound.
- 7. Solve the Problem 3 to find the best approximation of f'(1.3) using the three-point forward-difference and backward-difference formulas.
- 8. By using the following data: (1.0, 2.0), (1.5, 1.94), (2.0, 2.25), (3.0, 3.11). Find the best approximate values for f'(1.5), f'(1.0), and f'(3.0) using suitable three-point formulas.
- 9. Use all three-point formulas to compute the approximate value for f'(2) for the derivative of $f(x) = e^{x/2} + x^3$, taking h = 0.1. Also, compute the actual errors and error bounds for your approximation.
- 10. Use all three-point formulas to compute the approximate value for f'(2.2) for the derivative of $f(x) = x^2 e^x x + 1$, taking h = 0.2. Compute actual errors and error bounds.
- 11. Use the most accurate three-point formula to determine approximations that will complete the following table.

X	f(x)	f'(x)
8.1	16.94410	
8.3	17.56492	
8.5	18.19056	
8.7	18.82091	

- 12. The data in the Problem 11 were taken from the function $f(x) = x \ln x$. Compute the actual errors and error bounds.
- 13. Let $f(x) = x + \ln(x+2)$, with h = 0.1. Use the three-point formula to approximate f'(2). Find error bound for your approximation and compare the actual error to the bound.
- 14. Let $f(x) = e^{-2x}$, with points x = 0.25, 0.5, 0.75, 1.25, 1.50. Use the three-point central-difference formula to approximate f'(1.0). Compute error bound for your approximation.
- 15. Approximate the integral $\int_0^2 x^2 e^{-x^2} dx$, using suitable composite rules with n=4, and n=6.

290 5.11 Exercises

16. The following values of a function $f(x) = \tan x/x$ are given

X	f(x)	X	f(x)
1.00	1.5574	1.40	4.1342
1.10	1.7862	1.50	9.4009
1.20	2.1435	1.60	-21.3953
1.30	2.7709		

Find $\int_{1.0}^{1.6} f(x)dx$, using the Trapezoidal rule with h = 0.1.

- 17. Use a suitable composite integration formula to approximate $\int_0^1 \frac{dx}{2e^x 1}$, with n = 5.
- 18. Use a suitable composite integration formula for the approximation of the integral $\int_1^2 \frac{dx}{3-x}$, with n=5. Compute an upper bound for your approximation.
- 19. Use the composite Trapezoidal rule for the approximation of the integral $\int_1^3 \frac{dx}{7-2x}$ with h = 0.5. Also, compute an error term.
- 20. Find the stepsize h so that the absolute value of the error for the composite Trapezoidal rule is less than 5×10^{-4} when it is used to approximate the integral $\int_2^7 \frac{dx}{x}$.
- 21. Estimate the integral $\int_{-1}^{1} \frac{dx}{1+x^2}$ using the Simpson's rules with n=8.
- 22. Use the composite Trapezoidal and the Simpson's rules to approximate the integral $\int_1^2 \frac{dx}{4x+1}$ such that the error does not exceed 10^{-2} .
- 23. Evaluate $\int_0^1 e^{x^2} dx$ by the Simpson's rule choosing h small enough to guarantee five decimal accuracy. How large can h be ?
- 24. Use a suitable composite integration rule to find the best approximate value of the integral $\int_{1}^{2} \sqrt{1 + \sin x} \ dx$, with h = 0.1. Estimate the error bound.

Chapter 6

Numerical Solution of Ordinary Differential Equations

6.1 Introduction

The differential equations are of fundamental importance in engineering mathematics because many physical laws of biology, chemistry, ecology, economics, business, etc., and relations appear mathematically in the form of such equations. We know that many differential equations can be solved explicitly in terms of elementary functions of calculus. For example, the explicit solution of the differential equation

$$\frac{dy}{dx} = e^{x-y},$$

can be obtained easily as

$$y(x) = \ln(e^x + C),$$

and using initial condition y(0) = 1, we get

$$y(x) = \ln(e^x + e - 1).$$

We can use the MATLAB command dsolve (MATLAB's symbolic differential equation solver) which produces the general solution to the differential equation, or the specific solution to an associated initial-value problem.

>>
$$syms \ x \ y;$$

>> $y = dsolve('Dx = exp(x - y)', 'x');$
>> $y = dsolve('Dx = exp(x - y)', 'y(0) = 1', 'x'); \ pretty(y)$

But there are many differential equations which cannot be solved explicitly in terms of the functions of calculus. For example, the solutions of the differential equation of the form

$$\frac{dy}{dx} = e^{-x^2},$$

are the integral, or antiderivatives of e^{-x^2} ,

$$y(x) = \int \frac{e^{-x^2}}{291} \, dx + C,$$

but it is known that these integrals cannot be expressed in terms of the functions of calculus. We will refer to solutions that can be explicitly written in terms of elementary functions or special functions as formula or symbolic solutions-whether they were obtained via hand-calculation or via dsolve. When we cannot solve a differential equation in this way, or if the formula we find to complicated, we turn to numerical methods to solve as initial-value problem. This is similar to a situation in calculus: if we cannot find an antiderivative in terms of elementary functions, we

turn to a numerical method such as Trapezoidal rule or Simpson's rule to evaluate a definite integral.

Important Points of the Chapter 6

I. In this chapter we shall find the approximate solutions of the differential equations.

II. Differential equation may be ordinary differential equation (only one independent variable) or partial differential equation (more than one independent variable).

III. Given data points should be equally spaced only (length of each subinterval should be same). Smaller the length of the interval better the approximation.

IV. Given differential equations may be linear or nonlinear and first degree and first-order.

V. We shall discuss the first-order ordinary differential equations and sets of simultaneous first-order differential equations, since, one can easily find that *nth*-order differential equation may be solved by transforming it to a set of n-simultaneous first-order differential equations. All the specified conditions are on the same endpoints. These are *initial-value problems*. Many numerical methods are discussed for the approximate solutions of such initial value problem.

VI. We shall use single step numerical methods for the approximate solution of the ordinary differential equations.

6.2 Ordinary Differential Equations

Here, we will discuss about the ordinary differential equations and their numerical solutions.

Definition 6.1 (Differential Equation)

An equation which involving functions and their derivatives. For example, the following equations

(a)
$$\frac{dy}{dx} = 3x$$
, (b) $\frac{d^2y}{dx^2} + 4\frac{dy}{dx} + y = 0$,

(c)
$$\frac{dy}{dx} = x^2 + y^2$$
, (d) $\left(\frac{d^3y}{dx^3}\right)^2 - 5\frac{d^2y}{dx^2} + 2y = 5$,

are differential equations.

For the sake of completeness, we shall define some of the standard terms for differential equations.

Definition 6.2 (Dependent Variable)

It is the variable that has being differentiated. For example, in each of above differential equations (a)-(d), y is the dependent variable.

Definition 6.3 (Independent Variable)

It is the variable with respect to which the differentiation is performed. For example, in each of above differential equations (a)-(d), x is the independent variable.

Definition 6.4 (Order of Differential Equation)

The order of the differential equation is the order of the highest derivative involved. For example, the differential equations (a) and (c) are of first-order since the highest derivatives that appear is of first-order, whereas the differential equations (b) and (d) are respectively, the second-order and the third-order.

Definition 6.5 (Degree of Differential Equation)

The degree of the differential equation is the power to which the highest-order derivative is raised. For example, the differential equations (a)-(c) are of degree 1 while the differential equation (d) is of degree 2.

Definition 6.6 (Linear Differential Equation)

An differential equation is linear if

- (1) The dependent variable y and all its derivatives are of the first degree, that is, the power of each term involving y or its derivatives is one.
- (2) Each coefficient depends on only independent variable x or constant.

For example, the above differential equations (a) and (b) are the linear differential equations while the differential equations (c) and (d) are the nonlinear differential equations.

Definition 6.7 (Initial Conditions)

When all of the conditions are given at starting value of independent variable x to solve a given differential equation, is called a initial condition. When the conditions are given at the endpoints of x-values, then the conditions are called the boundary conditions.

6.2.1 Classification of Differential Equations

There are two major types of differential equations, called, ordinary differential equations (ODE) and partial differential equations (PDE). If an equation contains only ordinary derivatives of one or more dependent variables, with respect to a single independent variable, it is then said to be an ordinary differential equation. For example, all the differential equations (a)-(d) are ordinary differential equations because there is only one independent variable, called x.

An equation involving the partial derivatives of one or more dependent variables of two or more independent variables is called it partial differential equation. For example, the following differential equation

$$\frac{\partial^2 y}{\partial x^2} = c \frac{\partial^2 y}{\partial t^2},$$

is the partial differential because it involves two independent variables, x and t. Although partial differential equations are very useful and important, their study demands a good foundation in the theory of ordinary differential equations. Consequently, in this chapter the discussion that follows we shall confine our attention to ordinary differential equations.

As a mathematical form, the ordinary differential equation is a very useful tool. It is used in modeling of a wide variety of physical phenomena, that is, chemical reactions, satellite orbits, vibrating or oscillating systems, electrical networks, and so on. In many cases, the independent variable represents time so that a differential equation describes change, with respect to time, in the system being modeled. The solution of a differential equation will be representation of the state of the system at any point in time and one can use it to study the behavior of the system. Consequently, the problem of finding the solution of a differential equation play an important role in scientific research.

The solution of a differential equation is the function which satisfies the differential equation. In solving a differential equation analytically, one usually compute a general solution containing arbitrary constant. The simplest form of the differential equation is

$$y' = f(x), (6.1)$$

with f(x) a given function. The general solution of this equation is

$$y(x) = \int f(x)dx + C, \tag{6.2}$$

where C is an arbitrary constant. For example, the differential equation of the form

$$y' = \cos x,\tag{6.3}$$

has general solution of the form

$$y(x) = \sin x + C. \tag{6.4}$$

The more general equation is

$$y' = f(x, y(x)). \tag{6.5}$$

Since the general solution of differential equation is depends on an arbitrary constant C, so this constant can be calculated by specifying the value of function y(x) at a particular point x_0

$$y(x_0) = y_0$$
.

The point x_0 is called initial point, and the number y_0 is called the initial value. We call the problem of solving

$$y' = \frac{dy}{dx} = f(x, y); \quad x_0 \le x \le x_n, \quad y(x_0) = y_0,$$
 (6.6)

the initial-value problem (IVP). For example, for finding the solution of the differential equation (6.3) satisfying $y(\pi) = 1$, we have the value of the constant C = 1, so (6.4) becomes

$$y(x) = \sin x + 1,$$

and it is called the particular solution of the differential equation (6.3), or called the solution of the initial-value problem

$$y' = \cos x, \qquad y(\pi) = 1.$$

The main concern of this chapter is approximating the solution to the problem (6.6). The initial-value problems are problems in which the value of the dependent variable y is known at a point x_0 . Such a large number of methods are available to handle problems of this type that one may have difficulty in deciding which to use. Solving initial-value problem numerically we will assume that the solution is being sought on a given finite interval $x_0 \le x \le x_n$ with h = (b - a)/n, where $x_0 = a, x_n = b$ and n be the number of subintervals. In this chapter the most widely used numerical methods are discussed in some details to find the solution of the initial-value problem. If the analytical process of finding a exact solution y(x) is not feasible, it is still useful to know whether a solution exists and unique using numerical methods. To make precise preceding discussion, we give the following theorem which gives a sufficient condition for the existence and uniqueness of the initial-value problem (6.6).

Theorem 6.1 (Existence and Uniqueness Theorem)

Let f(x,y) and $\frac{\partial f}{\partial y}$ be continuous functions of x and y at all points (x,y) in some neighborhood of the initial point (x_0,y_0) . Then there is a unique function y(x) defined on some interval $[x_0-\epsilon,x_0+\epsilon]$ and satisfying

$$y'(x) = f(x, y(x)), \quad y(x_0) = y_0, \quad x \in [x_0 - \epsilon, x_0 + \epsilon], \quad \epsilon > 0$$
 (6.7)

For example, the initial-value problem

$$y'(x) = 2xy^2, y(0) = 1,$$

has a unique solution

$$y(x) = \frac{1}{1 - x^2}, \quad -1 < x < 1,$$

because the both functions

$$f(x,y) = 2xy^2, \qquad \frac{\partial f}{\partial y} = 4xy,$$

are continuous for all (x, y). Note that this example also showed that the continuity of the function f(x, y) and $\frac{\partial f}{\partial y}$ for all (x, y) does not imply the existence of a y(x) that is continuous for all x.

6.3 Numerical Methods for Solving IVP

By a numerical method for solving the initial-value problem (6.6) is meant a procedure for finding approximate values y_0, y_1, \ldots, y_n of the exact solution y(x) at the given points $x_0 < x_1 < \ldots < x_n$. We will let y_i denote the numerical value obtained as approximation to the exact solution $y(x_i)$, with $x_i = x_0 + ih$ for $i = 0, 1, \ldots, n$, where h (constant) is the size of the interval. Numerical methods for differential equations are of great importance to the engineer and physicist because practical problems often lead to differential equations that cannot be solved by any analytical method or to equations for which the solutions in terms of formulas are so complicated that are often prefers to calculate a table of values by applying a numerical method to such an equation.

Two different types of numerical methods are available to solve initial-value problem (6.6). These are called the *single-step* and the *multi-steps* methods. The methods discussed will vary in complexity, since in general, the greater the accuracy of a method, the greater is its complexity. We shall

discuss many numerical methods for solving the approximate solution of the initial-value problems (6.6) and the error analysis of each of the methods is explained in detail. Firstly, we discuss the single-step methods for solving the problem (6.6).

6.4 Single-Step Methods for IVP

This type of method called self-starting, refers to estimate y'(x), from a initial condition $y(x_0) = y_0$ and $y'_0 = f(x_0, y_0)$ from (6.6) and proceed step-wise. In the first-step we compute an approximate value y_1 of the solution y(x) at $x = x_1 = x_0 + h$. In the second-step we compute an approximate value y_2 of that solution at $x = x_2 = x_0 + 2h$ and so on. Although these methods generally use functional evaluation information for x_i and x_{i+1} , they do not retain that information for direct use in future approximations. All the information used by these methods is consequently obtained within the interval over which the solution is being approximated. Among of them we will discuss here, the Euler's method, the Taylor's method of higher-orders, and the Runge-Kutta method of different orders.

6.4.1 Euler's Method

One of the simplest and most straight forward numerical method for solving first-order ordinary differential equation of the form (6.6) is called method of *Euler*. This method is not an efficient numerical method and so seldom used, but it is relatively easy to analysis and many of the ideas involved in the numerical solution of differential equations are introduced most simply with it. In principle, the Euler's method uses the forward difference formula approximation of y'(x) which we discussed in the previous Chapter 5. That is

$$y' = \frac{dy}{dx} \approx \frac{y(x_{i+1}) - y(x_i)}{h},\tag{6.8}$$

where h is the stepsize and it is equal to $x_{i+1} - x_i$. Given that $\frac{dy}{dx} = f(x,y)$ and the initial conditions $x = x_0$, $y(x) = y(x_0)$, we have

$$\frac{y(x_1) - y(x_0)}{h} \approx f(x_0, y(x_0)), \quad \text{or} \quad y(x_1) \approx y(x_0) + hf(x_0, y(x_0)),$$

which shows that $y(x_1)$ is approximately given by $y(x_0) + hf(x_0, y(x_0))$. We can now use this approximation for $y(x_1)$ to estimate $y(x_2)$, that is

$$y(x_2) \approx y(x_1) + h f(x_1, y(x_1)),$$

and so on. In general,

$$y(x_{i+1}) \approx y(x_i) + hf(x_i, y(x_i)), \quad i = 0, 1, \dots, n-1.$$

Taking $y_i \approx y(x_i)$, for each i = 1, 2, ..., n, we have

$$y_{i+1} = yx_i + hf(x_i, y_i), \quad i = 0, 1, \dots, n-1.$$
 (6.9)

This simple integration strategy is known as the Euler's method, or the Euler-Cauchy method. It is called an explicit method because the value of y(x) at the next step is calculated only from the

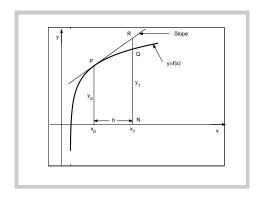


Figure 6.1: Geometrically interpretation of the Euler's method.

value of y(x) at the previous step. Given the approximate formula, one can solve for y_{i+1} in terms of x_i , y_i and $f(x_i, y_i)$, all of which are known. Note that the above formula (6.9) can be derive by using the Taylor series expansion of the unknown solution y(x) to the problem (6.6) about the point $x = x_i$, for each i = 0, 1, ..., n-1

$$y(x_{i+1}) = y(x_i) + (x_{i+1} - x_i)y'(x_i) + \frac{(x_{i+1} - x_i)^2}{2!}y''(\eta_i) = y(x_i) + hy'(x_i) + \frac{h^2}{2!}y''(\eta_i), \quad (6.10)$$

where unknown point η_i lies in the interval (x_i, x_{i+1}) . For the smaller value of stepsize h, the higher power h^2 will be very small and may be neglected. Using $f(x_i, y_i)$ to evaluate $y'(x_i)$ and $y_i \approx y(x_i)$, we have the formula (6.9).

Geometrically interpretation of the method is shown by Figure 6.1.

Example 6.1 Use the Euler's method to find the approximate value of y(1) for the given initial-value problem

$$y' = xy + x$$
, $0 \le x \le 1$, $y(0) = 0$, with $h = 0.1$, 0.2.

Compare your approximate solutions with the exact solution $y(x) = -1 + e^{x^2/2}$.

Solution. Since f(x,y) = xy + x, and $x_0 = 0$, $y_0 = 0$, then

$$y_{i+1} = y_i + h f(x_i, y_i), \quad for \quad i = 0, 1, \dots, 9$$

Then for h = 0.1 and taking i = 0, we have

$$y_1 = y_0 + hf(x_0, y_0) = y_0 + h(x_0y_0 + x_0) = 0 + (0.1)[(0)(0) + (0)] = 0.0000.$$

Similar way, we have other approximations by taking $x_i = x_{i-1} + h, i = 1, 2, ..., 9$, as follows

$$y_2 = 0.0100, \quad y_3 = 0.0302, \quad y_4 = 0.0611, \quad y_5 = 0.1036, \quad y_6 = 0.1587,$$

$$y_7 = 0.2283$$
, $y_8 = 0.3142$, $y_9 = 0.4194$, $y_{10} = 0.5471$,

with possible absolute error

$$|y(1) - y_{10}| = |0.6487 - 0.5471| = 0.1016.$$

Similarly, the approximations for h = 0.2, give

$$y_1 = 0.0000, \quad y_2 = 0.0400, \quad y_3 = 0.1232, \quad y_4 = 0.2580, \quad y_5 = 0.4592,$$

with possible absolute error

$$|y(1) - y_5| = |0.6487 - 0.4592| = 0.1895.$$

It showed that the result for h=0.1 is better than the h=0.2 and for both cases the approximation is not even correct to 1 decimal place. Clearly, the results using this method are inferior to those we will obtain in the coming methods. However, the accuracy of the Euler's method could be considerable improved by using smaller value of h than 0.1.

Note that before calling MATLAB function Euler1 which defined below, we must define MATLAB function fun1 as follows:

$$function f = fun1(x, y)$$

$$f = x * y + y;$$

Given Euler1.m and fun1.m the results obtained manually in the preceding example are reproduced with following MATLAB command:

$$>> [x', y'] = Euler1('fun1', 0, 1, 0, 10); [x', y'] = Euler1('fun1', 0, 1, 0, 5)$$

The same results are obtained with the following statements that define MATLAB command as an inline function object:

$$>> sol = inline('x * y + x', 'x', 'y'); [x', y'] = Euler1(sol, 0, 1, 0, 10); disp([x', y'])$$

Example 6.2 Use the Euler's method to find the approximate value of y(1.4) for the given initial-value problem

$$\frac{1}{x}y' - y^2 = 0$$
, $y(1) = 1$, with $n = 2$

Compare your approximate solutions with the exact solution $y(x) = 2/(3-x^2)$.

Solution. Since $f(x,y) = xy^2$, and $x_0 = 1$, $y_0 = 1$, h = 0.2, then using the Euler's method

$$y_{i+1} = y_i + hf(x_i, y_i), \quad for \quad i = 0, 1, \dots, n-1,$$

for i = 0, we have

$$y_1 = y_0 + hf(x_0, y_0) = y_0 + h(x_0y_0^2) = 1 + (0.2)[(1)(1)] = 1.2.$$

Similar way, we have other approximations by taking i = 1, as follows

$$y_2 = y_1 + hf(x_1, y_1) = y_1 + h(x_1y_1^2) = 1.2 + (0.2)[(1.2)(1.44)] = 1.5456,$$

the required approximation of y(1.4) and

$$|y(1.4) - y_2| = |1.9231 - 1.5456| = 0.3775,$$

is the possible absolute error.

6.4.2 Analysis of the Euler's Method

The preceding Example 6.1 demonstrates that the error in applying the Euler's method is reduced when h is reduced. The question of how well the Euler's method for solving the initial-value problem (6.6) works is closely related to the truncation error of the method. There are two types of such error, local and global truncation error.

In case of local truncation error one consider the size of the error made during one step and for global truncation error, one can consider the errors in the entire interval $x_0 \le x \le x_n$ over which the solution is sought.

We turn to the Taylor series to find an expression that represents the error, we have

$$y(x+h) = y(x) + hy'(x) + \frac{h^2}{2!}y''(\eta(x)),$$

for unknown $\eta(x) \in [x, x+h]$.

If $y(x_{i+1})$ is the true value of y(x), then the Taylor series expression at x_i is

$$y(x_{i+1}) = y(x_i) + hf(x_i, y_i) + \frac{h^2}{2!}y''(\eta(x_i)), \quad \eta(x_i) \in (x_i, x_{i+1})$$

as $y' = f(x_i, y_i)$. The Euler's formula uses the recurrence relation

$$y_{i+1} = y(x_i) + hf(x_i, y_i),$$

to estimate y_{i+1} assuming that $y(x_i)$ is the true solution. The error in y_{i+1} is given by $y_{i+1} - y(x_{i+1})$ which can be written as

$$y_{i+1} - y(x_{i+1}) = (y(x_i) + hf(x_i, y_i)) - (y(x_i) + hf(x_i, y_i) + \frac{h^2}{2!}y''(\eta(x_i))) = -\frac{h^2}{2!}y''(\eta(x_i)), (6.11)$$

for i = 0, 1, ..., n - 1. We call the term $-\frac{h^2}{2}y''(\eta(x_i))$, the *local* truncation error for the Euler's method. It is of order h^2 .

Note that this error term only applies in the region (x_i, x_{i+1}) , hence it is only the error in estimating y_{i+1} from $y(x_i)$. It does not take into account the compounded error from previous estimates. If we assume that the error is increasing linearly with n, then the error will be proportional to nh^2 , but n is dependent on h as $h = \frac{x_n - x_0}{n}$, so the error will be proportional to

$$\frac{x_n - x_0}{h}h^2 = (x_n - x_0)h,$$

which is order h. This error is called the *global* truncation error.

The analysis above leads to important theorem in the analysis of numerical methods.

Theorem 6.2 For the differential equations $\frac{dy}{dx} = f(x, y)$, if the leading term in the local truncation error involves h^{p+1} , for some integer p, then the global error, for small h, is of order h^p , that is

$$y_{i+1} - y(x_{i+1}) \approx ch^p$$
,

where c does not depend on stepsize h.

Note that the Euler's method is called the first-order method because of its local truncation error given by the formula (6.11), since this arises on each application of the method. Thus, in generating the solution point (x_k, y_k) the truncation error appears k times, once for each application of the method.

```
Program 6.1

MATLAB m-file for Euler Method function sol=Euler1(fun1,a,b,y0,n) h=(b-a)/n; x=a+(0:n)*h; y(1)=y0; for k=1:n; y(k+1)=y(k)+h*feval(fun1,x(k),y(k)); end; sol=[x',y'];
```

Since the Euler's method is an iterative method so, it may be converge or diverge. If divergence occurs, then the procedure should be terminated because there may be no solution.

6.4.3 Higher-Order Taylor Methods

The basis for many numerical techniques finding the approximate solution of the initial-value problem can be depend to the Taylor's series, as we used this series in the previous section in finding the Euler's method which also called the Taylor's method of order one. One can, of course, develop the Taylor's method for higher-order to obtain better accuracy, and in general, one expect that higher the order of the method, greater the accuracy for a given stepsize. The Taylor's method is relatively easy to use, however, the necessity of calculating the higher derivatives makes the Taylor's method completely unsuitable. Nevertheless, it is of great theoretical interest because the most of the practical methods attempt to achieve the same accuracy as the Taylor's method of a given order without the disadvantage of having to calculate the higher derivatives. Assuming that the solution y(x) of the initial-value problem (6.6) has (n + 1) continuous derivatives and expanding y(x) in terms of its nth degree Taylor polynomial about x_i , we get

$$y(x_{i+1}) = y(x_i) + hy'(x_i) + \frac{h^2}{2!}y''(x_i) + \cdots + \frac{h^n}{n!}y^{(n)}(x_i) + \frac{h^{n+1}}{(n+1)!}y^{(n+1)}(\eta(x_i)),$$
(6.12)

for some $\eta(x_i) \in (x_i, x_{i+1})$. The derivatives in this expansion are not known explicitly since the solution is not known. However, if f is sufficiently differentiable, they can be obtained by taking the total derivative of (6.6) with respect to x, keeping in mind that f is an implicit function of y. Thus

$$y' = f(x,y) = f$$

$$y'' = f' = f_x + f_y f$$

$$y''' = f'' = f_{xx} + 2f_{xy}f + f_{yy}f^2 + f_x f_y + f_y^2 f$$

$$\vdots$$
(6.13)

Continuing in this manner, we can express any derivative of y in terms of f(x, y) and its partial derivatives. It is already clear, however, that unless f(x, y) is a very simple function, the higher total derivatives become increasingly complex. Now substituting these results into (6.12), gives

$$y(x_{i+1}) = y(x_i) + hf(x_i, y(x_i)) + \frac{h^2}{2!}f'(x_i, y(x_i)) + \cdots$$

+
$$\frac{h^n}{n!}f^{(n-1)}(x_i,y(x_i)) + \frac{h^{n+1}}{(n+1)!}f^{(n)}(\eta(x_i),y(\eta(x_i))).$$
 (6.14)

By taking $y_i \approx y(x_i)$, that the approximation to the exact solution at x_i , for each $i = 0, 1, \dots, n-1$, we have

$$y_{i+1} = y_i + hf(x_i, y_i) + \frac{h^2}{2!}f'(x_i, y_i) + \dots + \frac{h^n}{n!}f^{(n-1)}(x_i, y_i).$$
(6.15)

Then this formula is called the Taylor's method of order n. The last term of (6.14), called remainder, shows that the local error of Taylor's method of order n is

$$E = \frac{h^{n+1}}{(n+1)!} f^{(n)}(\eta_i, y(\eta(x_i))) = \frac{h^{n+1}}{(n+1)!} y^{(n+1)}(\eta(x_i)), \tag{6.16}$$

for some $x_i < \eta(x_i) < x_{i+1}$.

Example 6.3 Use the Taylor's method of order 2 to find the approximate value of y(1) for the given initial-value problem.

$$y' = xy + x$$
, $0 \le x \le 1$, $y(0) = 0$, with $h = 0.2$

Compare your approximate solution with the exact solution $y(x) = -1 + e^{x^2/2}$.

Solution. Since f(x,y) = xy + x, and $x_0 = 0$, $y_0 = 0$, then

$$y_{i+1} = y_i + hf(x_i, y_i) + \frac{h^2}{2}f'(x_i, y_i), \quad for \quad i = 0, 1, 2, 3, 4$$

where $f'(x_i, y_i) = y_i + x_i^2 y_i + x_i^2 + 1$. Then for i = 0, we have

$$y_1 = y_0 + h(x_0y_0 + x_0) + \frac{h^2}{2}(y_0 + x_0^2y_0 + x_0^2 + 1) = 0 + (0.2)(0) + (0.02)(1) = 0.0200,$$

and similar way, we have for i = 1, 2, 3, 4, as follows

$$y_2 = 0.0820, \quad y_3 = 0.1937, \quad y_4 = 0.3694, \quad y_5 = 0.6334,$$

with absolute possible error

$$|y(1) - y_5| = |0.6487 - 0.6334| = 0.0153.$$

It showed that the result is entirely correct to 1 decimal place. Clearly, the result using this method is better than the Euler's method and it could be considerable improved by using smaller value of h than 0.2.

Note that before calling MATLAB functions tayl1 and fun1, we must define MATLAB function dfun1 as follows:

function
$$f = dfun1(x, y)$$

 $f = y + x$. $2 * y + x$. $2 + 1$;

Given tayl1.m, fun1.m and dfun1.m, the results obtained manually in the preceding example are reproduced with following MATLAB commands:

$$>> [x', y'] = tayl1('fun1', 'dfun1', 0, 1, 0, 5); disp([x', y'])$$

Example 6.4 Show that Taylor's method of order 2 for the initial-value problem

$$e^{y}y' - e^{x} = 0$$
, $0 \le x \le 1$, $y(0) = 1$, with $h = 0.5$,

is

$$y_{i+1} = y_i + he^{(x_i - y_i)} \left[1 + \frac{h}{2} \left(1 - e^{(x_i - y_i)} \right) \right], \quad i \ge 0.$$

What are the values of y_0, y_1, y_2 . Compare your approximate solution with the exact solution $y(x) = \ln(e^x + e - 1)$.

Solution. Since the Taylor's method of order 2 is

$$y_{i+1} = y_i + hf(x_i, y_i) + \frac{h^2}{2}f'(x_i, y_i), \text{ for } i \ge 0,$$

and the given function is $f(x,y) = e^{x-y}$ with its first derivative $f'(x,y) = e^{x-y}[1 - e^{x-y}]$. So using these values, we have

$$y_{i+1} = y_i + he^{x_i - y_i} + \frac{h^2}{2} \left[e^{x - y} \left(1 - e^{x - y} \right) \right], \quad \text{for} \quad i \ge 0,$$

or

$$y_{i+1} = y_i + he^{(x_i - y_i)} \left[1 + \frac{h}{2} \left(1 - e^{(x_i - y_i)} \right) \right], \quad i \ge 0.$$

Now for i = 0, we have

$$y_1 = y_0 + he^{(x_0 - y_0)} \left[1 + \frac{h}{2} \left(1 - e^{(x_0 - y_0)} \right) \right],$$

and using $x_0 = 0, y_0$ and h = 0.5, we get y_1 as follows

$$y_1 = 1 + (0.5)e^{(0-1)} \left[1 + \frac{0.5}{2} \left(1 - e^{(0-1)} \right) \right] = 1.2130.$$

Similar way, we have the value of y_2 for taking i = 1, as follows

$$y_2 = y_1 + he^{(x_1 - y_1)} \left[1 + \frac{h}{2} \left(1 - e^{(x_1 - y_1)} \right) \right]$$
$$= 1.2130 + (0.5)e^{(0.5 - 1.2130)} \left[1 + \frac{0.5}{2} \left(1 - e^{(0.5 - 1.2130)} \right) \right] = 1.4893.$$

the required approximation of y(x) at x = 1. Thus

$$|y(1) - y_2| = |1.4899 - 1.4893| = 0.0006,$$

is the possible absolute error in our approximation.

Program 6.2

MATLAB m-file for Taylor's Method of order 2

function sol=tayl1(fun1,dfun1,a,b,y0,n)

h=(b-a)/n; x = a + (0:n) * h; y(1)=y0; for k=1:n

 $y(k+1) = y(k) + h * feval(fun1, x(k), y(k)) + (h.^2 * feval(dfun1, x(k), y(k)))/2;$ end; sol = [x', y'];

Example 6.5 Use the Taylor's method of order 3 to find the approximate value of y(1) for the given initial-value problem

$$4y' - y = 0$$
, $0 \le x \le 1$, $y(0) = 1$, with $n = 2$.

Compare your approximate solution with the exact solution $y(x) = e^{x/4}$.

Solution. Since the Taylor's method of order 3 is

$$y_{i+1} = y_i + hf(x_i, y_i) + \frac{h^2}{2!}f'(x_i, y_i) + \frac{h^3}{3!}f''(x_i, y_i),$$

for i = 0, 1, ..., n - 1, and using the given values $x_0 = 0$, $y_0 = 1$ and f(x, y) = 1/4y, we get, f'(x, y) = 1/16y and f''(x, y) = 1/64y. So using these values we obtain Taylor's method of order 3 of the form

$$y_{i+1} = y_i + h(1/4y_i) + \frac{h^2}{2}(1/16y_i) + \frac{h^3}{6}(1/64y_i).$$

Then for i = 0, we have

$$y_1 = y_0 \left[1 + \frac{h}{4} + \frac{h^2}{32} + \frac{h^3}{384} \right],$$

and by taking $y_0 = 1$, h = 0.5, we get

$$y(0.5) \approx y_1 = 1(1 + 0.125 + 0.0078 + 0.0003) = 1.1331,$$

and similar way, we have other approximation for taking i = 1, as follows

$$y(1) \approx y_2 = y_1(1 + 0.125 + 0.0078 + 0.0003) = 1.1331(1.1331) = 1.2839,$$

the required approximation of y(1) and

$$|y(1) - y_2| = |1.2840 - 1.2839| = 0.0001,$$

is the possible absolute error.

Example 6.6 Show that third order Taylor's method for the given initial-value problem

$$y' - 2x + y = 0, \qquad y(0) = -1,$$

is

$$y_{i+1} = y_i + (2x_-y_i)h + \left(\frac{h^2}{2} - \frac{h^3}{6}\right)[2(1-x_i) + y_i], \quad i = 0, 1, \dots, n-1.$$

Use it to find approximation of y(0.1).

Solution. Using f(x,y) = 2x - y, f'(x,y) = 2(1-x) + y, f''(x,y) = 2(x-1) - y, then the Taylor's method of order 3 gets the form

$$y_{i+1} = y_i + h[2x_i - y_i] + \frac{h^2}{2}[2(1 - x_i) + y_i] + \frac{h^3}{6}[2(x_i - 1) - y_i],$$

and after simplifying, we get

$$y_{i+1} = y_i + (2x_i - y_i)h + \left(\frac{h^2}{2} - \frac{h^3}{6}\right)[2(1 - x_i) + y_i].$$

Now by taking i = 0 in the above formula, we obtain

$$y(x_1) \approx y_1 = y_0 + (2x_0 - y_0)h + \left(\frac{h^2}{2} - \frac{h^3}{6}\right)[2(1 - x_0) + y_0],$$

by using $x_0 = 0$, $y_0 = 1$ and h = 0.1, we obtain, $y(0.1) \approx y_1 = -1 + 0.1 + (0.005 - 0.0002)(2 - 1) = -0.8952$, the required approximation.

In using the Taylor's method, we replace the infinite Taylor series for f(x + h) by a partial sum. The local truncation error is inherent in any algorithm that we might choose.

If we retain term up to and including h^n in the series, then the local truncation error is the sum of all the remaining terms that we do not include by Taylor's method. These terms can be compressed into a single term of the form

$$\frac{h^{n+1}}{(n+1)!} f^{(n+1)}(\eta(x), y(\eta(x))),$$

for some unknown point $\eta(x)$. We say that the local truncation error is order h^{n+1} . An error of this sort is present in each step of the numerical solution. The accumulation of all the many local truncation error gives rise to the *global truncation* error which must be of order h^n because the number of steps necessary to reach on arbitrary point x, having started at x_0 , is $\frac{x-x_0}{h}$. Choosing n large so that this error is small.

6.4.4 Runge-Kutta Methods

Since we studied that the Euler's method is not very useful in practical problems because it requires a very small stepsize for reasonable accuracy. the Taylor's method of higher-order is difficult to use because it needs to obtain higher total derivatives of y(x). An important group of methods which allow us to obtain greater accuracy at each step and yet require only initial value of y(x) to be given with the differential equation are called the Runge-Kutta methods. The Runge-Kutta methods attempt to obtain greater accuracy, and at the same time avoid the need of higher derivatives by evaluating the function f(x, y) at selected points on each subintervals. These methods can be used to generate not only starting values but, in fact, in whole solution. They are self-starting and easy

to program for a digital computer. We shall begin by showing how to derive the simplest formulas in this class. These are of the form

$$y_{i+1} = y_i + (w_1 k_1 + w_2 k_2), (6.17)$$

where

$$k_1 = h f(x_i, y_i)$$
 and $k_2 = h f(x_i + ah, y_i + bk_1)$.

The parameters w_1 , w_2 , a, and b are chosen in order to make the formula (6.17) as accurate as possible, that is, to make the order of accuracy as large as possible. To this end, we substitute the exact value $y(x), y(x_{i+1})$ by the local solution into the formula (6.17) and expand about the point x_i . The parameters are then chosen to make the resulting expansion agree as much as possible with the Taylor series for $y(x_{i+1})$ about x_i . Upon substituting into (6.17), we first expanding $y(x_{i+1})$ in the Taylor series through terms of order h^3 , we obtain

$$y(x_{i+1}) = y(x_i) + hy'(x_i) + \frac{h^2}{2!}y''(x_i) + \frac{h^3}{3!}y'''(x_i) + \cdots$$
(6.18)

Since

$$y' = f(x,y)$$

$$y'' = f'(x_i, y_i) = (f_x + f_y f)_i$$

$$y''' = f''(x_i, y_i) = (f_{xx} + 2f_{xy}f + f_{yy}f^2 + f_x f_y + f_y^2 f)_i + O(h^4).$$
(6.19)

So

$$y(x_{i+1}) = y(x_i) + hf(x_i, y_i) + \frac{h^2}{2!} (f_x + ff_y)_i$$

$$+ \frac{h^3}{3!} (f_{xx} + 2f_{xy}f + f_{yy}f^2 + f_xf_y + f_y^2 f)_i + O(h^4),$$
(6.20)

where the subscripts on f denote partial derivatives with respect to the indicated variables, and the subscript i means that all functions involved are to be evaluated at (x_i, y_i) . Now using the Taylor's expansion for functions of two variables, we find that

$$k_2 = hf(x_i + ah, y_i + bk_1) = h[f + h(af_x + bff_y)]$$

$$+ \frac{h^2}{2}(a^2f_{xx} + 2abff_{xy} + b^2f^2f_{yy}) + O(h^4)]_i.$$
(6.21)

Now we substitute this expression for k_2 into (6.17), gives

$$y_{i+1} = y_i + h[w_1 f(x_i, y_i) + w_2 f(x_i + ah, y_i + bk_1)]$$

$$= y_i + h[(w_1 + w_2)f]_i + h^2 w_2 [(af_x + bff_y)]_i$$

$$+ \frac{h^3}{2} w_2 [a^2 f_{xx} + 2abf f_{xy} + b^2 f^2 f_{yy}]_i + O(h^4).$$
(6.22)

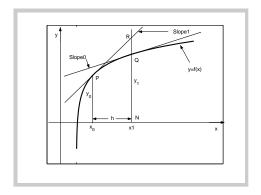


Figure 6.2: Geometrically interpretation of the Modified Euler method.

On comparing (6.20) and (6.22), we see that to make the corresponding powers of h and h^2 agree, we must have

$$w_1 + w_2 = 1$$
 and $a = \frac{1}{2w_2} = b$.

This is a system of two nonlinear equations in the four unknowns a, b, w_1 , and w_2 and its solution can be written in the form

$$b = a = \frac{1}{2w_2}, \quad w_1 = 1 - w_2. \tag{6.23}$$

There are many solutions to (6.23) depending on the choices of w_2 . These choices leads to the numerical method which has order 2 and some of them do correspond to any of the standard numerical integration formulas. Taking the first choice when $w_2 = 1/2$, we have

$$y_{i+1} = y_i + \frac{h}{2} [f(x_i, y_i) + f(x_{i+1}, y_i + hf(x_i, y_i))].$$
(6.24)

Runge-Kutta Method of Order Two (Modified Euler's Method)

The equation (6.24) can be written in a standard form as

$$y_{i+1} = y_i + \frac{h}{2}[k_1 + k_2], (6.25)$$

where

$$k_1 = f(x_i, y_i)$$
 and $k_2 = f(x_{i+1}, y_i + hk_1),$

for each i = 0, 1, ..., n-1. Then the relation (6.25) is called the Runge-Kutta method of order 2 which is also known as the Modified Euler's method. This method corresponds to using the Trapezoidal rule to estimate the integral where a preliminary (full) Euler step is taken to obtain the (approximate) value at x_{i+1} . Geometrically interpretation of the method is shown by Figure 6.2. The local error of this formula is, however, of order h^3 , whereas that of the Euler's method is h^2 . We can therefore expect to be able to use a large stepsize with this formula.

Example 6.7 Use Runge-Kutta method of order two (Modified Euler's method) to find the approximate value of y(1) for the given initial-value problem

$$y' = xy + x$$
, $0 \le x \le 1$, $y(0) = 0$, with $h = 0.2$.

Compare your approximate solution with the exact solution $y(x) = -1 + e^{x^2/2}$.

Solution. Since f(x,y) = xy + x, and $x_0 = 0$, $y_0 = 0$, then for i = 0, we have

$$k_1 = f(x_0, y_0) = (x_0y_0 + x_0) = 0.0000$$

 $k_2 = f(x_1, y_0 + hk_1) = (x_1(y_0 + hk_1) + x_1) = (0 + 0.2) = 0.2000,$

and using these values, we have

$$y_1 = y_0 + \frac{h}{2} [k_1 + k_2] = 0 + 0.1(0 + 0.2000) = 0.0200.$$

Continuing in this manner, we have

$$k_1 = 0.204,$$
 $k_2 = 0.4243,$ then $y_2 = 0.0828,$ $k_1 = 0.4331,$ $k_2 = 0.7017,$ then $y_3 = 0.1963,$ $k_1 = 0.7178,$ $k_2 = 1.0719,$ then $y_4 = 0.3753,$ $k_1 = 1.1002,$ $k_2 = 1.5953,$ then $y_5 = 0.6449,$

with possible error

$$|y(1) - y_5| = |0.6487 - 0.6449| = 0.0039$$

The results of the Example 6.7 can be obtained by using the following MATLAB command as follows:

$$>> sol = mod1('fun1', 0, 1, 0, 5);$$

So there is a significant improvement in accuracy of this method as compared with the Euler's method but the problem of accuracy still remains however since error will accurate from step to step. In particular, since the function f(x,y) calculated repeatedly from values of y(x) which include the accumulated error, these errors may grow in an unpredictable way.

Example 6.8 Use the Runge-Kutta method of order two (the Modified Euler's method) to find the approximate value of y(1.4) for the given initial-value problem

$$xy' + y' - 2y = 0$$
, $y(1) = 4$, with $n = 2$.

and compare your approximate solution with the exact solution $y(x) = (x+1)^2$.

Solution. Since $f(x,y) = \frac{2y}{x+1}$, and $x_0 = 1$, $y_0 = 4$, h = (1.4-1)/2 = 0.2, then for i = 0, we have

$$k_1 = f(x_0, y_0) = \frac{2y_0}{x_0 + 1} = 4,$$

$$k_2 = f(x_1, y_0 + hk_1) = \frac{2(y_0 + hk_1)}{x_1 + 1} = 4.3636,$$

and by using these values, we have

$$y(1.2) \approx y_1 = y_0 + \frac{h}{2}[k_1 + k_2] = 4 + 0.1(4 + 4.3636) = 4.8364.$$

Similar manner, we have the other approximation for taking i = 1, as follows

$$k_1 = 4.3967$$
 and $k_2 = 4.7631$,

and by using these values, we have

$$y(1.4) \approx y_2 = y_1 + \frac{h}{2}[k_1 + k_2] = 4.8364 + 0.1(4.3967 + 4.7631) = 5.7523,$$

the required approximation of y(1.4) and

$$|y(1.4) - y_2| = |5.7600 - 5.7523| = 0.0077,$$

is the possible absolute error.

Example 6.9 Use the Runge-Kutta method of order two (the Modified Euler's method) to find the approximate value of y(1.2) for the given initial-value problem

$$x^2y' - y = 0$$
, $y(1) = 2$, with $n = 2$.

and compare your approximate solution with the exact solution $y(x) = e^{(x-1)/x}$.

Solution. Since $f(x,y) = x^{-2}y$ and $x_0 = 1$, $y_0 = 2$, h = (1.2 - 1)/2 = 0.1, then for i = 0, we have

$$k_1 = f(x_0, y_0) = f(1, 2) = (1)^{-2}(2) = 2,$$

$$k_2 = f(x_1, y_0 + hk_1) = f(1.1, 2.2) = (1.1)^{-2}(2.2) = 1.8182,$$

and by using these values, we have

$$y(1.1) \approx y_1 = y_0 + \frac{h}{2} [k_1 + k_2] = 2 + 0.05(2 + 1.8182) = 2.1909.$$

Similar manner, we have the other approximation for taking i = 1, as follows

$$k_1 = 1.8107$$
 and $k_2 = 1.6035$,

and by using these values, we have

$$y(1.2) \approx y_2 = y_1 + \frac{h}{2} [k_1 + k_2] = 2.1909 + 0.05(1.8107 + 1.6035) = 2.3616,$$

the required approximation of y(1.2) and

$$|y(1.2) - y_2| = |2.3627 - 2.3616| = 0.0011,$$

is the possible absolute error.

Example 6.10 Use the Runge-Kutta method of order two (the Modified Euler's method) to find the approximate value of y(0.2) for the given initial-value problem

$$2\sin xy' - y\sin 2x = 0$$
, $y(0) = 1$, with $n = 2$.

and compare your approximate solution with the exact solution $y(x) = e^{\sin x}$.

Solution. Since $f(x,y) = y \cos x$ and $x_0 = 0$, $y_0 = 1$, h = (0.2 - 0)/2 = 0.1, then for i = 0, we have

$$k_1 = f(x_0, y_0) = f(0, 1) = 1\cos 0 = 1,$$

$$k_2 = f(x_1, y_0 + hk_1) = f(0.1, 1.1) = (1.1)\cos 0.1 = 1.0945,$$

and by using these values, we have

$$y(0.1) \approx y_1 = y_0 + \frac{h}{2} [k_1 + k_2] = 1 + 0.05(1 + 1.0945) = 1.1047.$$

Similar manner, we have the other approximation for taking i = 1, as follows

$$k_1 = 1.0992$$
 and $k_2 = 1.1904$,

and by using these values, we have

$$y(0.2) \approx y_2 = y_1 + \frac{h}{2} [k_1 + k_2] = 1.1047 + 0.05(1.0992 + 1.1904) = 1.2192,$$

the required approximation of y(1.2) and

$$|y(0.2) - y_2| = |1.2198 - 1.2192| = 0.0006,$$

is the possible absolute error.

The modified Euler's method is classified as a predictor-corrector method. This means that in the case of the modified Euler's method the initial-value problem is given by the formula

$$y_{i+1}^{(k)} = y_i + h f(x_i, y_i), (6.26)$$

which is called the predictor and this is corrected by the repeated application of the formula

$$y_{i+1}^{(k+1)} = y_i + \frac{h}{2} \left[f(x_i, y_i) + f(x_{i+1}, y_{i+1}^{(k)}) \right], \quad k = 0, 1, 2, \dots$$
 (6.27)

for each i = 0, 1, ..., n-1. This is called the corrector. There are many predictor-corrector formulas and some provide much greater accuracy than the relatively the modified Euler's method. These methods however, require accurate estimates for a number of initial values of y(x) before they can be used. We shall discuss some of those predictor-corrector formulas latter in the chapter.

Program 6.3

MATLAB m-file for the Modified Euler's Method function sol=mod1(fun1,a,b,y0,n)

$$h = (b-a)/n; x = a + (0:n)*h; y(1) = y0; for k = 1:n k1 = feval(fun1, x(k), y(k)); k2 = feval(fun1, x(k) + h, y(k) + h * k1); y(k+1) = y(k) + h * (k1 + k2)/2; end; sol = [x', y'];$$

310 6.13 Exercises

6.5 Exercises

- 1. Find the general solution of the differential equation $y' = \frac{x}{y}$. In Problems 2 – 20 solve each of the following initial-value problem using indicated method.
- 2. Put following differential equations into a form for numerical solution by Euler's method. (a) y + 2yy' y' = 0; (b) $\ln y' = x^2 y^2$; (c) $y' x^2y' = y$.
- 3. Solve the following initial-value problems using the Euler's method.
- (a) $y' = y + x^2$, x = 0(0.2)1, y(0) = 1.
- (b) $y' = (x-1)(x+y+1), \quad x = 0(0.2)1, \quad y(0) = 1.$
- 4. Solve the following initial-value problems and compare the numerical solutions obtained with the Euler's method using the values of h = 0.1 and h = 0.2. Compare the results to the actual values.
- (a) $y' = 1 + x^2$, $0 \le x \le 1$, y(0) = 0, $y(x) = \tan x$.
- (b) $y' = 2(y+1), \quad 0 \le x \le 1, \quad y(0) = 0, \quad y(x) = e^{2x} 1.$
- (c) $y' = 2(y-1)^2$, $1 \le x \le 2$, y(1) = 0.5, y(x) = (2x-1)/2x.
- 5. Solve the following initial-value problems using the Taylor's method of order two.
- (a) $y' = 2x^2 y$, x = 0(0.2)1, y(0) = -1.
- (b) $y' = 3x^2y$, x = 0(0.2)1, y(0) = 1.
- (c) y' = x/y x, x = 0(0.2)1, y(0) = 2.
- 6. Solve the initial-value problems by using Taylor's method of order three of the Problem 3.
- 7. Solve the following initial-value problems using the Modified Euler's method.
- (a) $y' = y^2x^2$, x = 1(0.2)2, y(1) = -1.
- (b) y' = x y/2x, x = 1(0.02)1.10, y(1) = 0.25.
- (c) $y' = 1/y^2 yx$, x = 1(0.2)2, y(1) = 1.
- 8. Solve the following initial-value problems and compare the numerical solutions obtained with the Modified Euler's method using the values of h = 0.05 and h = 0.1 and compare the results with the actual values.
- (a) $y' = x + \frac{3y}{x}$, $1 \le x \le 2$, y(1) = 0, $y(x) = x^3 x^2$.
- (b) $y' = \sqrt{y}$, $0 \le x \le 1$, y(0) = 1, $y(x) = (x+2)^2/4$.
- (c) y' = 4 3y, $0 \le x \le 1$, y(0) = 5, $y(x) = 4/3 + 11/3e^{-3x}$.

9. Solve the following initial-value problems using the Modified Euler's method.

(a)
$$y' = (x+1)y$$
, $x = 0.5(0.2)1.5$, $y(0.5) = 1$.

(b)
$$y' = -xy^2$$
, $x = 0(0.2)1$, $y(0) = 2$.

(c)
$$y' = t^2 + x^2$$
, $x = 1(0.2)2$, $y(1) = -1$.

Answers to Selected Exercises Chapter 2

1.
$$x_{10} = -1.84141$$
 3. (a) $x_{16} = 0.35173$ (b). $x_{15} = 0.97300$

5. (a)
$$x_8 = 1.53906$$
 (b) $|\alpha - x_8| \le \frac{3-1}{2^8} = 0.00781$

7. **(b)**
$$x_3 = 0.58705$$
 (c) $|\alpha - x_3| \le \frac{(0.41218)^3}{1 - 0.41218} |0.5 - 0| = 0.05957$

9.
$$k=1, \ g'(1)=-\frac{1}{2}\neq 0$$
, a linear convergence.

11. (b)
$$k = \max_{3 \le 4} |g'(x)| = 4/9, x_0 = 3, x_1 = 10/3, n = 8$$

13.
$$x_3 = 0.0667$$
 15. $x_3 = 1.532$, quadratic convergence, **17.** $g'(\alpha) = \frac{1}{2}$

19.
$$x_4 = 0.567$$
 21. $f'''(1) = 6 \neq 0$, **23.** (a) $x_{23} = 1.00$, (b) $x_4 = 0.99$, (c). $x_4 = 1.00$

25.
$$g'(\sqrt{2}) = -0.414214 \neq 0$$
, **27.** $x^{(3)} = 1.08828$, $y^{(3)} = 0.84434$

Chapter 3

3.
$$|AB| = 0$$
 5. $x = -3$, $y = 2$ **7.**(a) **B 9.** $det(A) = 86$

21. (a)
$$\mathbf{x} = [1.16, 0.88, -1.32]^T$$
 23. (c) $\mathbf{x} = [0.6429, -2.2798, -1.9345]^T$

25. (a)
$$\mathbf{x} = [-1.33, 3, 2.33]^T$$
 31. (c) $\mathbf{x} = [-1, 2, 3]^T$

33. (a) WP:
$$\mathbf{x} = [-500.0, 333.7]^T PP : \mathbf{x} = [-499.98, 333.7]^T$$

35. (a)
$$\mathbf{x} = [0.74, -0.84, 1.05, -1.96]^T$$
 37. (c) -5 **39.** (a) -3

41.
$$||A|| = 6, 7, 5.53$$
 43. (c) $K(A) = 39.64$ **45.** (a) 0.5, 11.80

47.
$$\delta x = [-100, 101]^T$$
, $K(A) = 404$ **49.** $\delta x = [-17, 20]^T$, $K(A) = 8004$

51. (a)
$$||T_J|| = 70$$
 (b) $||T_J|| = 0.5$ (c) $||T_J|| = 0.5$

53. (a)
$$\mathbf{x}^{(13)} = [2, 4, 3]^T$$
 (b) $\mathbf{x}^{(22)} = [1, 1, 1]^T$ (c) $\mathbf{x}^{(12)} = [-0.1016, 0.5391, 0.3281]^T$

55.
$$||T_J|| = 0.75$$
, $||X - X^{(20)}|| = 0.0072$

57.
$$A = \begin{bmatrix} 5 & -1 & 1; 0 & 3 & -1; 1 & 2 & 4 \end{bmatrix}, \quad X^{(1)} = \begin{bmatrix} 0.2, -0.333, 0.6167 \end{bmatrix}^T, \quad \|X - X^{(10)}\| = 0.00011$$

312 6.13 Exercises

59. $||T_G|| = 0.6$, $X^{(2)} = [-0.9490, -2.9854, -1.9836]^T$, $||X - X^{(2)}|| = 0.3780$

Chapter 4

- **1.** $p_2(2.3) = 0.4548$ **3.** $p_2(0.5) = 2.3212$, $E_B = 0.0313$, $p_2(2.8) = 5.9399$, $E_B = 0.007111$
- **5.** $p_3(1.1) = -3.686$, $E_B = 3.73308$ **7. (b)** $p_2(2.2) = 10.284$, $p_3(2.2) = 10.1224$
- **9.** (b) $p_3(5.9) = 2.429$, $p_4(5.9) = 2.429$, (c) $E_3 = 0.0005742$, $E_4 = 0.00000004$
- 11. All three divided differences can be expanded as

$$\frac{(x_2 - x_1)f_0 - (x_2 - x_0)f_1 + (x_1 - x_0)f_2}{(x_2 - x_1)(x_2 - x_0)(x_1 - x_0)}$$

- **13.** f[0,1,0] = 0.7183, **15.** $s_0(0.55) = 3.7200$, $s_1(1.15) = 4.0210$, $s_2(2.5) = 5.2500$
- **17.** $s_0(0.15) = 0.1350$, $s_1(0.25) = 0.2200$, $s_2(0.45) = 0.3725$
- **19.** $s_0(2.5) = 1.833$, $s_1(5.5) = 2.5$, $s_2(10.5) = 3.357$, $E_0 = 0.0375$, $E_1 = 0.0495$, $E_2 = 0.034$

Chapter 5

- **1.** (a). 23.1059, -0.16529, 0.16673 (b). 22.77676, 0.16386, 0.16457
- **3.** FD = 30.0, BD = 20.00 **5.** 21.2946, 0.0002, 0.0002
- 7. FD = 21.294, 0.0004, 0.0004 BD = 21.294, 0.0004, 0.0004
- **9.** CD = 13.3697, 0.0106, 0.01, FD = 13.34, 0.02, 0.02, BD = 13.34, 0.02, 0.02
- **11.** FD = f'(8.1) = 3.09, CD = f'(8.3) = 3.12, CD = f'(8.5) = 3.14, BD = f'(2.4) = 3.16
- **13.** Approx = 1.24999992, EB = 0.00000008, AE = 0.00000010097
- **17.** $T_5(f) = 0.49557207$ **19.** $T_4(f) = 0.84166667$, EB = 0.3333333
- **21.** $S_8(f) = 1.57078431$ **23.** h = 0.0543 and n = 18

Chapter 6

- **1.** $y^2 x^2 = c$ **3.** (a) $y(0.25) \approx 1.5876$
- **3.** (a) $y(0.25) \approx 1.5876$ **5.** $y(1) \approx 0.1463$ **7.** $y(2) \approx -0.3075$
- 9. Heun's Method (a) $y(1.5) \approx 7.0236$ Midpoint Method (a) $y(1.5) \approx 6.9937$

Index

absolute error, 4
absolute value, 113
adjacent points, 182
algebraic form, 106
algorithm, 1
approximate area, 276
approximate number, 4
approximating function, 240
approximating functions, 182
approximation polynomials, 181
approximation theory, 181
area, 263
augmented matrix, 78

nth divided difference, 212

backward substitution, 101, 103 backward-difference formula, 242, 253 band matrix, 87 bisection method, 9 Bolzano's method, 9 boundary value problems, 137

central-difference formula, 249 chord, 38 coefficient matrix, 78 cofactor, 90, 91 column matrix, 78 composite form, 268 condition number, 162 consistent system, 77 continuous function, 8, 12, 182 Crout's method, 125 cubic function, 258

definite integral, 263 determinant, 88, 95 diagonal matrix, 85 direct method, 116 discretization error, 1 divided difference, 212 divided differences, 211 Doolittle's method, 118

elementary functions, 181 elimination methods, 137 equivalent system, 102 error bound, 5, 202, 245 error formula, 244 error term, 244 Euclidean, 136 exact number, 4 exponential functions, 181 extrapolation, 181

factorization method, 116 first divided difference, 212 five-point formula, 241 fixed-point, 16 fixed-point method, 17 forward elimination, 103 forward-difference formula, 242, 253 Frobenius norm, 136 full rank, 100

Gauss factorization, 118 Gauss-Jordan method, 115 Gauss-Seidel iterative method, 141 Gaussian elimination method, 101 Gaussian quadrature, 264 geometric interpretation, 242 global error, 273 Graphical techniques, 239

higher derivatives, 257 homogeneous system, 79 313 314 Index

identity matrix, 82 ill-conditioned systems, 161 inconsistent system, 77 integration by parts, 277 interpolating point, 189 interpolating polynomial, 186, 242 interpolation, 181 interpolation conditions, 185 interval-halving method, 9 inverse matrix, 84 invertible matrix, 83, 84 iterative methods, 137

Jacobi method, 138 Jacobian matrix, 65

Lagrange coefficient polynomials, 184
Lagrange coefficients, 189, 191, 195, 201, 210
Lagrange interpolation, 182, 185
Lagrange interpolatory polynomial, 181
Laplace Expansion Theorem, 92
linear combination, 76
linear convergent method, 8
linear equation, 74
linear equations, 73
linear function, 245
linear independent, 76
linear polynomial, 183
linear spline, 233
local error, 272
lower-triangular matrix, 86

Maclaurin's series expansion, 6
matrix inversion method, 99
matrix norm, 135
matrix of cofactor, 92
maximum error, 251
method of elimination, 100
method of tangents, 29
minor, 90
minors, 91
modified Newton's method, 48
multiple root, 8
multiple roots, 43
multiples, 102

LU decomposition, 116

multiplicity, 45

Newton divided difference, 221
Newton divided difference interpolation, 215
Newton interpolation, 221
Newton's method, 30, 38
Newton-Cotes formulas, 264
nonhomogeneous system, 100
nonlinear algebraic equations, 7
nonsingular matrix, 83
numerical differentiation, 239
numerical formula, 250
Numerical integration, 263
numerical integration, 239

order of multiplicity, 50 overdetermined system, 75

partial derivatives, 70
partial pivoting, 114
percentage error, 5
piecewise curve fitting, 233
piecewise linear interpolation, 233
piecewise polynomial, 233
piecewise polynomial approximation, 233
pivot element, 102, 107
pivotal equation, 102
pivoting strategy, 113
pole values, 195
polynomial functions, 181
polynomial interpolation, 210
product matrix, 81

quadratic convergent method, 8 quadratic function, 250 quadratic polynomial, 189 quadrature rule, 276

rank, 100
rank deficient, 100
rate of convergence, 21, 55
rational functions, 181
rectangular array, 79
rectangular matrix, 83
relative error, 5
round-off errors, 6
rounding error, 248

Index 315

scalar matrix, 85, 96 secant line, 242 secant method, 39, 41 significant digits, 248 simple root, 8 Simpson's rule, 276 simultaneous equations, 73 simultaneous linear systems, 75 skew matrix, 87 skew symmetric matrix, 87 slope, 38 sparse matrix, 88 spline, 233 square matrix, 82 strictly diagonally dominant matrix, 132 strictly lower-triangular matrix, 86 strictly upper-triangular matrix, 85 subdiagonal, 88 superdiagonal, 88 symmetric matrix, 86 system of linear equations, 74 system of nonlinear equations, 67 system of two equations, 66

tabulated data, 241 Taylor's series, 283 The fixed point, 26 three-point formula, 241 three-point formulas, 248 total error, 248 transcendental equation, 7 transpose matrix, 83 Trapezoidal rule, 267 triangular form, 101 triangular system, 101 tridiagonal matrix, 88 trigonometric functions, 181 trivial solution, 100 truncation error, 6 two-point formula, 241

underdetermined system, 75 unique solution, 75 upper-triangular matrix, 85

vector norm, 134

Weierstrass approximation theorem, 182

zero matrix, 82 Zeroth divided difference, 212