

4th Edition

Laboratory Manual of

Simplified Numerical Analysis (C++ version)

Proofreading powered by
various **AI-driven** proprietary software

**Amjad Ali
Muhammad Ishaq
Hamayun Farooq
Muhammad Umar**

Laboratory Manual of Simplified Numerical Analysis (C++ Version)

Accessible through: <https://zenodo.org/record/8264208>

Cite as:

Ali, Amjad, Ishaq, Muhammad, Farooq, Hamayun, & Umar, Muhammad. (2023).
Laboratory Manual of Simplified Numerical Analysis (C++ Version). Zenodo.
<https://doi.org/10.5281/zenodo.8264208>

For availability of the codes, please visit:

[GitHub - DrAmjadAli11/SimplifiedNumericalAnalysis](https://github.com/DrAmjadAli11/SimplifiedNumericalAnalysis)

<https://github.com/DrAmjadAli11/SimplifiedNumericalAnalysis>

Principal Book

Simplified Numerical Analysis

Fourth Edition

www.TimeRenders.com.pk

Companion Books

Laboratory Manual of Simplified Numerical Analysis (C++ Version)

Laboratory Manual of Simplified Numerical Analysis (MATLAB® Version)

Laboratory Manual of Simplified Numerical Analysis (Python Version)

Laboratory Manual of Simplified Numerical Analysis (C++ Version)

Fourth Edition

Amjad Ali, Ph.D.

Bahauddin Zakariya University (BZU), Multan

Muhammad Ishaq, Ph.D.

COMSATS University Islamabad, Vehari Campus

Hamayun Farooq, Ph.D.

Government Degree College, Muzaffar Garh

Muhammad Umar, Ph.D.

University of Heidelberg, Heidelberg

Esteemed Panel of the Supporters

Ms. Aniq Faizan, Bahauddin Zakariya University, Multan

Ms. Amna Waheed, Bahauddin Zakariya University, Multan

Dr. Zainab Bukhari, Times Institute, Multan

Ms. Tayyaba Saeed, Bahauddin Zakariya University, Multan

Ms. Syeda Zahra Kazmi, Bahauddin Zakariya University, Multan

Ms. Zoha Saleem, Bahauddin Zakariya University, Multan

Laboratory Manual of Simplified Numerical Analysis (C++ Version)

A Companion book of the principal book:

Simplified Numerical Analysis (Fourth Edition)

©2023, Amjad Ali, Ph.D. (The Principal Author)

ISBN: 978-969-7821-14-3

Typeset: Mostly personally by **Dr. Amjad Ali** (The Principal Author), also contributed by the supporters.

Title Design: **Mr. Muhammad Rizwan Qadeer** (mrizwanqadeer@gmail.com)

The proofreading is powered by various **AI-driven** proprietary software.



www.timerender.com.pk

Distribution Point:

238-B (PRIDE), Near Girls Comprehensive School, Gulgasht Colony, Multan, Pakistan.

Cell Phone: +923486981925, timerenderpublishers@gmail.com

Table of Contents

Chapter 1: Preliminary Concepts in Numerical Analysis (1)

- 1.1 Introduction
- 1.2 Number Systems and Representations
- 1.3 The Round-off Error
- 1.4 The Truncation Error

Computing Resources	1
Chapter Summary	6
Chapter Exercises	8

Chapter 2: Solution of a Nonlinear Equation in One Variable (9)

Corridor I: BASICS

- 2.1 Introduction
- 2.2 Bracketing Methods
 - 2.2.1 The Bisection Method (or Bolzano Method)
 - 2.2.2 The False-Position Method (or Regula-Falsi Method)
- 2.3 Open Methods
 - 2.3.1 The Fixed-Point Iteration Method
 - 2.3.2 The Newton-Raphson Method
 - 2.3.3 The Secant Method

Corridor II: ANALYSIS

- 2.4 Convergence Analysis
 - The Bisection Method
 - The Regula-Falsi Method
 - The Secant Method
 - The Newton-Raphson Method
 - The Fixed-Point Iteration Method
- 2.5 Further Discussions

Corridor III: PROGRAMMING ARCADE

2.6 Algorithms and Implementations	11
The Newton-Raphson Method.....	11
The Fixed-Point Iteration Method.....	16
The Secant Method	18
The Bisection Method	20
The Regula-Falsi Method	23
Chapter Summary	28
Chapter Exercises	32

Chapter 3: Polynomial Interpolation (37)

Corridor I: BASICS

- 3.1 Introduction
- 3.2 The Newton's Divided Difference Interpolation
- 3.3 The Lagrange Interpolation
- 3.4 Deriving the Lagrange Interpolation Formula from the Newton's Divided-Difference Formula
- 3.5 Interpolation Formulas for Equally Spaced Nodes
- 3.6 Hermite Interpolation
- 3.7 Spline Interpolation
 - 3.7.1 Linear Spline
 - 3.7.2 Quadratic Spline
 - 3.7.3 Cubic Spline

Corridor II: ANALYSIS

- 3.8 Error of Interpolation

Corridor III: PROGRAMMING ARCADE

- 3.9 Algorithms and Implementations 38
 - The Newton's Divided Difference Interpolation Formula 39

- Chapter Summary 42
- Chapter Exercises 45

Chapter 4: Numerical Integration (49)

Corridor I: BASICS

- 4.1 Introduction
- 4.2 The Trapezoidal Rule
- 4.3 The Simpson's 1/3 Rule
- 4.4 Generalized Closed Newton-Cotes Quadrature

Corridor II: ANALYSIS

- 4.5 Truncation Error of the Trapezoidal Rule
- 4.6 Truncation Error of the Simpson's 1/3 Rule
- 4.7 Further Discussions
- 4.8 The Gaussian Quadrature

Corridor III: PROGRAMMING ARCADE

- 4.9 Algorithms and Implementations 53
 - The Composite Trapezoidal Rule 53
 - The Composite Simpson's 1/3 Rule 56
 - The Composite Simpson's 3/8 Rule 58

- Chapter Summary 61
- Chapter Exercises 64

Chapter 5: Numerical Differentiation (69)

- 5.1 Introduction
- 5.2 Finite Difference Approximations of Derivatives using the Taylor Series
 - 5.2.1 First Order Derivatives
 - 5.2.2 Second Order Derivatives
- 5.3 Listing of the Derivative Formulas

Chapter 6: Direct Linear Solvers (71)

Corridor I: BASICS

- 6.1 Introduction to Linear Systems
- 6.2 Solving Linear Systems using the Gaussian Elimination Method
- 6.3 Pivoting Strategies
 - Partial Pivoting
 - Scaled Partial Pivoting
 - Complete Pivoting
- 6.4 The Gauss-Jordan Method
- 6.5 Solving Linear Systems using the LU Factorization Method
 - 6.5.1 The Doolittle's Method
 - 6.5.2 The Crout's Method
 - 6.5.3 The Cholesky's Method

Corridor II: ANALYSIS

- 6.6 Operation Count Analysis
- 6.7 Matrix Inversion

Corridor III: PROGRAMMING ARCADE

6.8 Algorithms and Implementations	73
The Gaussian Elimination Method with Partial Pivoting	74
Solving $AX = B$ using the Doolittle's Method	77
Solving $AX = B$ using the Crout's Method	81
Solving $AX = B$ using the Cholesky's Method	86
Performing Operation Count Analysis	89
Chapter Summary	101
Chapter Exercises	103

Chapter 7: Iterative Linear Solvers (107)

Corridor I: BASICS

- 7.1 Vector Norms and Distances
- 7.2 Convergence Criteria for Linear Solvers
- 7.3 Basic Methods
 - 7.3.1 The Jacobi Method
 - 7.3.2 The Gauss-Seidel Method

7.3.3 The SOR Method

Corridor II: ANALYSIS

- 7.4 Matrix Norms and Conditioning
- 7.5 Iteration Matrix and Matrix Form of a Solver

Corridor III: PROGRAMMING ARCADE

- 7.6 Algorithms and Implementations 108
 - The Jacobi Method 109
 - Modification in the Jacobi Method's algorithm for the Gauss-Seidel Method 110
 - Modification in the Jacobi Method's algorithm for the SOR Method 110

- Chapter Summary 114
- Chapter Exercises 118

Chapter 8: Eigenvalues and Eigenvectors (119)

Corridor I: BASICS

- 8.1 Basic Definitions and Concepts
- 8.2 General Approach of Finding Eigenvalues and Eigenvectors
- 8.3 Some Numerical Methods for Eigenvalues
 - The Power Method
 - The Householder Method
 - The QR Factorization Method
 - The Sturm Method

Corridor II: ANALYSIS

- 8.4 Further Discussions
 - The Power Theorem
 - The Gerschgorin Circle Theorems
 - The Singular Value Decomposition (SVD)

Corridor III: PROGRAMMING ARCADE

- 8.5 Algorithms and Implementations 120
 - The Power Method 121

- Chapter Summary 124
- Chapter Exercises 125

Chapter 9: Numerical Solution of Ordinary Differential Equations (ODEs) (127)

Corridor I: BASICS

- 9.1 Introduction
- 9.2 Solving IVPs using Single Step Methods and Multistep Methods
 - The Euler Method

	The Mid-point Method (an RK2 method of Order 2)	
	The Modified/Improved Euler Method (an RK2 method of Order 2)	
	The RK Method of order 4 (RK4)	
9.3	Solving IVPs using Predictor-Corrector Methods	
	The Adams-Bashforth-Moulton Method of Order 4	
9.4	Solving Systems of ODEs and Higher Order ODEs	
	Using the Classical RK4 Method	
9.5	Solving Linear BVPs using the Finite Difference Method	

Corridor II: ANALYSIS

9.6	Some Theoretical Concepts and Error Analysis	
-----	----------------------------------------------	--

Corridor III: PROGRAMMING ARCADE

9.7	Algorithms and Implementations	129
	Euler method	130
	Mid-point method	135
	Modified/Improved Euler method	138
	RK method of order 4 (RK4)	142
	Adams-Bashforth method of order 4	145
	Adams-Bashforth-Moulton method of order 4	149
	RK4 method for a system of two ODEs	154
	RK4 method for a system of three ODEs	156
	RK4 method for Second Order ODE	160
	RK4 method for Third Order ODE	163
	Linear FDM for BVP	167
	Chapter Summary	172
	Chapter Exercises	173
	Bibliography	177
	The End	178

Numerical
Analysis
is
the
mathematics
of
Scientific
Computing

Preliminary Concepts in Numerical Analysis

- 1.1 Introduction
- 1.2 Number Systems and Representations
- 1.3 The Round-off Error
- 1.4 The Truncation Error

To unleash the topics of this Chapter, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



Computing Resources

The numerical methods are devised just to be used on computers. It makes no sense to study a numerical method without considering its practicality using some computing tools. A variety of numerical computing tools, both freeware and proprietary, are available. The students are advised to understand the algorithmic (step-by-step) style of the numerical methods they learn. This book suggests the following resources for beginners.

- (1) **C++:** The numerical methods can be programmed in any programming language, especially C++, FORTRAN, and Python. The book discusses a wide variety of C++ programs of the numerical methods in this book. One can modify the as per need. Several C++ IDEs (Integrated Development Environments) are available, such as Dev-C++, and Code::Blocks for Windows and GNU-C++ for Linux operating system. One can even find C++ Apps (apps is an acronym for computer application software) for Android or iOS devices. Some online C++ IDEs are also available, which can be used for executing C++ programs without installing them.

- (2) **Python:** There are several free Python IDEs available for the Desktop use (such as Spyder, Jupyter, and PyCharm) or On-line use (such as Google Colab). It is quite a pertinent skill of the day that the students of computational sciences are familiar with programming in Python. The companion website of this book (www.timerender.com.pk) shares a Python Library having a variety of codes for the numerical methods discussed in this book.
- (3) **MATLAB®:** It is a proprietary software, by The MathWorks, Inc., available in both Desktop and Online versions. MATLAB® offers a wide variety of built-in functions and programming capabilities for mathematical computations (both symbolic and numeric, although more suitable and expert for numeric computations), for all modern areas of science and engineering. The book discusses a wide variety of MATLAB® programs and MATLAB® built-in functions for the numerical methods in this book.
- (4) **GNU-Octave:** It is an open-source (and freeware) version of MATLAB®, available in both Desktop and Online versions. Most of the MATLAB® codes and built-in functions discussed in this book can be executed in GNU-Octave and Octave-online.
- (5) **MATHEMATICA®:** It is a proprietary software by Wolfram Research. It is one of the best Computer Algebra Systems (CAS) available. It offers an extensive variety of built-in functions and programming capabilities for mathematical computations (both symbolic and numeric), for all modern areas of science and engineering.
- (6) **MAPLE®:** It is a proprietary software by Maplesoft for mathematical computations (both symbolic and numeric), for all modern areas of science and engineering. It is also one of the best Computer Algebra Systems (CAS).
- (7) **Spread-Sheet:** A spread-sheet software (such as Excel by Microsoft®) can be used for computations involved in simple numerical methods. The companion website of this book (www.timerender.com.pk) may shares a spread-sheet workbook having a variety of sheets for most of the numerical methods discussed in this book.
- (8) **Various Math Solver Tools:** Wolfram|Alpha, Symbolab, and Microsoft® Math Solver are three of the advanced tools for math education to be used as calculators. These are extensive, feature-rich, online tools, accessible both through the web browser and the relevant android/iOS apps. These tools provide automated step by step solutions to algebra and calculus problems covering from middle school through college. The premier versions of these tools are freely available, whereas professional (pro) versions are not free.
- (9) **Various Other Online Tools/Websites:** There are various other online tools and websites that offer basic computing facilities for numerical and symbolic computations. Examples include:
 - **AtoZmath.com** [<https://atozmath.com/>]
 - **CalculatorSoup®** [<https://www.calculatorsoup.com/>].
 - **Keisan - CASIO®** [<https://keisan.casio.com/>]



Question 06: What are the significant figures (or significant digits) of an approximate number?

Significant figures of a number (that approximates a true value) are the digits that are used to express the number meaningfully. The significant figures are counted for a number that approximates some other number to express the degree of precision in the approximate number.

The significant figures begin with the leftmost nonzero digit and end with the rightmost correct digit. The rightmost zeros, which are exact, are also significant. That is,

- All the nonzero digits (i. e., 1, 2, 3, ..., 9) are significant.
- Zeroes appearing anywhere between two nonzero digits are significant (e.g., in 3005.00102 there are nine significant digits).
- Leading zeros (i.e., left to the first nonzero digit) are not significant (e.g., the number 0.000081 has only two significant digits, namely 8 and 1). The leading zeros are used to fix the decimal place.
- Trailing zeroes are significant if they are exact with regard to some true value. Trailing zeros may or may not be significant. It depends on the context; how the number is approximated or obtained by rounding-off some other number.

Remark: The significant figures of a number can easily be identified by using its normalized scientific notation. The digits in the fractional part (or mantissa) are regarded as significant figures. For example, each of the numbers 42.134, 6.0013, and 0.0015784 has five significant figures, which can be identified easily by converting these numbers into their **normalized scientific notation** as:

$$\begin{aligned} 42.134 &= 0.42134 \times 10^2 \\ 6.0013 &= 0.60013 \times 10 \\ 0.0015784 &= 0.15784 \times 10^{-2} \end{aligned}$$

Remarks:

- 6500 has 2 significant figures (i.e., the digits 6 and 5) if it has been obtained by rounding-off a number to the nearest 100 (e.g., by rounding-off the numbers 6497 or 6543.88 to the nearest hundred). In fact, any number in the interval (6450, 6550) gives 6500, when rounded to the nearest 100.
- 6500 has 3 significant figures (i.e., the digits 6, 5, and the following 0) if it has been obtained by rounding-off a number to the nearest 10 (e.g., by rounding-off the numbers 6497 or 6504.99 to the nearest ten). In fact, any number in the interval [6495, 6505] gives 6500, when rounded to the nearest 10.

- 6500 has 4 significant figures if it has been obtained by rounding-off a number to the nearest whole number (e.g., by rounding-off the numbers 6499.8 or 6500.47 to the nearest whole number). In fact, any number in the interval $[6499.5, 6500.5]$ gives 6500, when rounded to the nearest whole number.
- 70500 has at least 3 significant figures (i.e., the digits 7, 5, and the 0 in between these). Depending upon the context, as just explained, it may have 3 to 5 significant figures.
- 0.00364300 has 4 significant figures (i.e., the digits 3, 6, 4, and 3) if it has been obtained by rounding-off a number to 4 significant figures (e.g., by rounding-off the numbers 0.003642859 or 0.0036432099 to 4 significant figures). Usually, in that case, the approximate number is written as 0.003643, without any non-significant trailing zero. In fact, any number in the interval $[0.0036426, 0.0036435]$ gives 0.003643, when rounded to 4 significant figures.
- 0.00364300 has 5 significant figures (i.e., the digits 3, 6, 4, 3, and the following 0) if it has been obtained by rounding-off a number to 5 significant figures (e.g., by rounding-off the numbers 0.003642978001 or 0.003643049 to 5 significant figures). Usually, in that case, the approximate number is written as 0.0036430, without any non-significant trailing zero. In fact, any number in the interval $[0.00364295, 0.00364306]$ gives 0.0036430, when rounded to 5 s.f.
- 0.00364300 has 6 significant figures (i.e., the digits 3, 6, 4, 3, and the following two 0s) if it has been obtained by rounding-off a number to 6 significant figures (e.g., by rounding-off the numbers 0.003642998001 or 0.003643001 to 6 significant figures). In fact, any number in the interval $[0.003642995, 0.003643006]$ gives 0.00364300, when rounded to 6 significant figures.

■

Remark:

An approximation x^* to a number x is called accurate to t significant figures if there are exactly t digits in the mantissa of x^* that agree with the first t digits of the mantissa of x , where x has the same exponent as x^* . Suppose that the number x is represented in the following form

$$x = \pm 0.d_1d_2d_3 \cdots d_t d_{t+1} \cdots \times 10^e$$

Then, the number x^* is accurate to t significant figures to the number x if it can be written in the following form

$$x^* = \pm 0.d_1d_2d_3 \cdots d_t d'_{t+1} \cdots \times 10^e$$

■

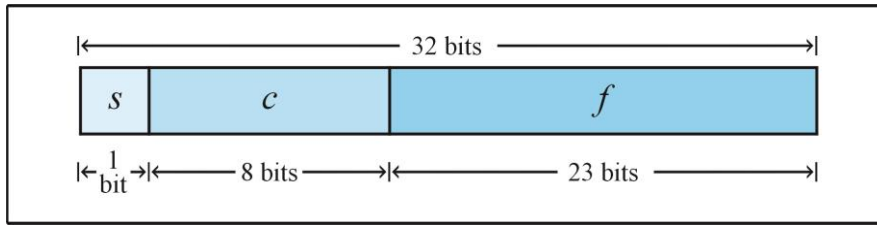


Fig. (1.3): According to the IEEE 754 standard, single-precision floating point representation of a binary real number $x = \pm 1. b_2 b_3 b_4 \dots \times 2^e$ is $(1 - 2s) \times 2^{c-127} \times (1 + f)$.

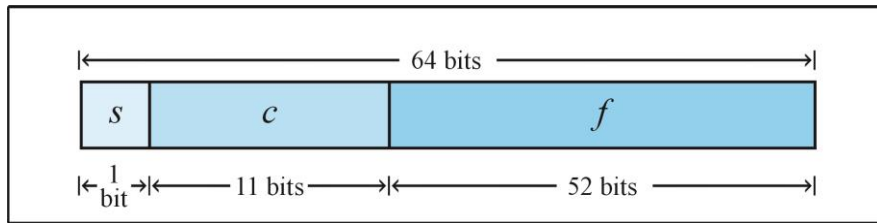


Fig. (1.4): According to the IEEE 754 standard, double-precision floating point representation of a binary real number $x = \pm 1. b_2 b_3 b_4 \dots \times 2^e$ is $(1 - 2s) \times 2^{c-1023} \times (1 + f)$.

Here, s is used for the **sign** of the number (0 means positive, 1 means negative). c in the exponent is called the **biased exponent**. f is the mantissa minus 1 (the hidden bit).

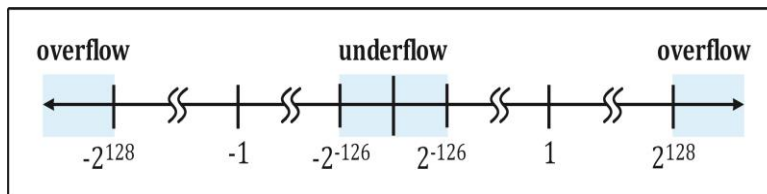


Fig. (1.5): Overflow/Underflow for single-precision floating-point representation

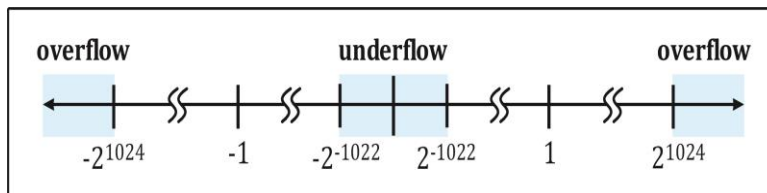


Fig. (1.6): Overflow/Underflow for double-precision floating-point representation



Chapter Summary

- The **numerical methods** obtain some approximate solution of the problems, usually in the numeric form, in contrast to the **analytic** or **exact methods**, which obtain the exact solution of the problem.
- **Numerical Analysis** is the field of deriving, analyzing, and implementing the numerical methods.
- The most common approach followed by the numerical methods is the **iterative** approach. According to this, choose an initial approximation or guess to the solution and apply a set of simple computational steps to obtain a better approximation. Repeatedly apply the same set of steps to the better approximations, ultimately obtaining a sufficiently accurate solution and then stop the repetition. Each course of repetition of the set of computational steps is called **iteration**. Geometrically, a root of an equation $f(x) = 0$ is the point where the graph of $f(x)$ intersects the x -axis.
- For selecting a numerical method from several choices, the characteristics of **accuracy**, **efficiency**, and **robustness** are taken into consideration.
- The numerical analysis may be regarded as the “mathematics of scientific computing”.
- Errors can be quantified as:
 - Absolute Error = $|\text{True value} - \text{Approximate value}|$
 - Relative Error = $\frac{\text{absolute error}}{|\text{True value}|} = \frac{|\text{True value} - \text{Approximate value}|}{|\text{True value}|}$
 - Percentage Relative Error = $\frac{\text{absolute error}}{|\text{True value}|} \times 100 \%$
- The errors can be categorized in three major categories in regard to their sources: **Data Error** or **Inherent Error** (quite unrelated to the numerical methods; occur as blunders, mistakes, model simplification, or data uncertainty), **Round-off Error** (occurs due to number approximation by humans and computers), **Truncation Error** (occurs due to approximation of a mathematical procedure to avoid insignificance), and **Discretization error** (occurs due to approximation of a continuous function by a set of discrete data points).
- **Significant figures** of a real number (which is an approximation of the true value) are the digits that are used to express the number meaningfully. The significant digits begin with the leftmost nonzero digit and end with the rightmost correct digit. The rightmost zeros, which are exact are also significant.
- An approximation x^* to a number x is called accurate to t significant figures if there are exactly t digits in the mantissa of x^* that agreed with the first t digits of the mantissa of x having the same exponent or characteristics.
- **Accuracy** of an approximate value is a measure of how much the approximate value agrees with the true value. **Precision, on the other hand**, has nothing to do with how much the approximate value agrees with the true value. Precision is only concerned about the size of the number.
- The following four are the commonly used number systems, even supported by the computer architectures.

1. Decimal number system (base 10)
 2. Binary number system (base 2)
 3. Octal number system (base 8)
 4. Hexadecimal number system (base 16)
- Any nonzero real decimal number x can be represented in floating-point form: $x = \pm 0.d_1d_2d_3 \dots \times 10^e$. Here $d_i, i = 1, 2, \dots$ are digits from 0 to 9 with $d_1 \neq 0$, called most significant digit and e is an integer that might be positive, negative or zero, called an **exponent** or *characteristic*. The number $0.d_1d_2d_3 \dots$, may be denoted by m , is called the finite normalized **mantissa**. For numbers in the decimal system with base 10, $\frac{1}{10} \leq m < 1$. That is, $m \in \left[\frac{1}{10}, 1\right)$.
 - For numbers in the binary system, the floating-point representation of a number x can be given by, $x = \pm 0.b_1b_2b_3 \dots \times 2^e = \pm m \times 2^e$, where each of b_i is a bit, either 0 or 1, with $b_1 \neq 0$, and $\frac{1}{2} \leq m < 1$.
 - The numbers that are representable precisely in a computer are called **machine numbers**. The real numbers with a non-terminating fractional part (such as $1/3$) cannot be represented, precisely. So many other numbers (for example, 0.01) also has not a precise representation in computer (i.e., a machine number).
 - If the number lies within the allowable range of the possible numbers according to the precision level of the computer, then it is rounded to a nearby machine number (incurring the round-off error) for storing it. The rounding options involve **correct rounding** (round to nearest machine number), *rounding up*, *rounding down* or towards zero, etc.
 - There are commonly two ways to terminate the mantissa of a number to obtain its nearest machine number, namely, correct **chopping** and correct **rounding**. The chopping or rounding of the number to the nearest machine number (representable in a computer) for representation in computers (for storage or for using in computations) causes the error in a number called the **round-off error**.
 - The floating-point form of a number x representable in a computer can be regarded as consisting of the three parts: $x = \pm m \times \beta^e = \mathbf{sign} \times \mathbf{mantissa} \times (\mathbf{base})^{\mathbf{exponent}}$

The sign is either positive (+) or negative (−), the finite normalized mantissa is from the interval $\left[\frac{1}{\beta}, 1\right)$, and the integer exponent either positive, negative, or zero as a power of the base.

- An account on the **IEEE Binary Floating-Point Arithmetic Standard 754-1985** for representing the real numbers in computers can be found under Question 13 in this chapter.
 - If a number x^* is accurate to t significant figures in approximating a number x then the relative error is bounded above by 5×10^{-t} . That is, $\frac{|x-x^*|}{|x|} \leq 5 \times 10^{-t}$
- If an iterative process is to be stopped when the successive approximations become accurate to t significant figures, the relative error bound might be set as 5×10^{-t} . Thus, the relative error is computed after every iteration using the result of the current iteration and that of the previous iteration. If the relative error is smaller than the bound of 5×10^{-t} , then it ensures that the approximation the accurate to t significant digits.
- Whenever two nearly equal numbers are subtracted, some loss of significance might occur. The risk of loss of significance can be eliminated by avoiding the subtraction through some mathematical manipulation.

Chapter Exercises

Exercise 01: Compute the absolute error E_a and relative error E_r in an approximation of x by x^*

- (i) $x = \log_{10} 2, x^* = 0.301$ (ii) $x = 17/6, x^* = 2.8333$
 (iii) $x = \sqrt{\pi}, x^* = 1.77245$ (iv) $x = e^{-1}, x^* = 0.36787$

Exercise 02: Write the following numbers in floating-point form and identify their mantissa and exponent:

- (i) $x = -23.500128$ (ii) $x = 658.000012$ (iii) $x = 0.010023$
 (iv) $x = -0.0000782$ (v) $x = \frac{1}{234.24}$ (vi) $x = 541000$

Exercise 03: Simplify the following expression by performing the computations

- (a) Exactly
 (b) Using four-digit chopping arithmetic
 (c) Using four-digit rounding arithmetic
 (d) Compute the relative errors

- (i) $\frac{7}{4} - \frac{5}{3}$ (ii) $\frac{5}{4} \left(\frac{2}{3} + 4 \right)$ (iii) $\frac{\pi - 1}{\frac{4}{3}}$
 (iv) $10\pi - 2e + 1$ (v) $\left(\frac{432 - 0.0012}{101} \right)$ (vi) $\left(\frac{2}{9} \right) \cdot \left(\frac{9}{7} \right)$

Consider π and e expressed with fifteen significant digits as the exact numbers.

Exercise 04: Calculate the roundoff error if chopping and rounding is used to write the following numbers accurate to four decimal digits:

- (i) $355/113$ (ii) $\sqrt{3/142}$ (iii) $\sqrt[3]{\ln 2}$

Exercise 05: We want to round-off each the following numbers to three decimal places. For which number, the result of “round-off by chopping” and “round-off by rounding-rule” will be the same:

- (A) 5.5555 (B) 3.3575 (C) 5.5565 (D) 4.4555

Exercise 06: Find the absolute and relative errors involved in rounding 4.9997 to 5.000.

Exercise 07: Suppose a real number x is represented approximately by 0.6032 with the relative error is at most 0.1%. What is x ?

Exercise 08: Suppose that a number is accurate to n significant figures and a_1 is the first significant figure than show that the relative error is bounded above by $\frac{1}{a_1} \times 10^{1-n}$.

Exercise 09: Show that if a number is rounded off to n digits than the relative error is bounded by $\frac{1}{2} \times 10^{1-n}$.



Solution of a Nonlinear Equation in One Variable

Corridor I: BASICS

Let's plan it

- 2.1 Introduction
- 2.2 Bracketing Methods
 - 2.2.1 The Bisection Method (or Bolzano Method)
 - 2.2.2 The False-Position Method (or Regula-Falsi Method)
- 2.3 Open Methods
 - 2.3.1 The Fixed-Point Iteration Method
 - 2.3.2 The Newton-Raphson Method
 - 2.3.3 The Secant Method

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



Corridor II: ANALYSIS

Let's think deep

2.4 Convergence Analysis

- The Bisection Method
- The Regula-Falsi Method
- The Secant Method
- The Newton-Raphson Method
- The Fixed-Point Iteration Method

2.5 Further Discussions

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



Corridor III: PROGRAMMING ARCADE

Let's do it

2.6 Algorithms and Implementations

- The Newton-Raphson Method
- The Fixed-Point Iteration Method
- The Secant Method
- The Bisection Method
- The Regula-Falsi Method

To cross-check the results/output of the computer programs you would execute, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



2.6 Algorithms and Implementations

Question 36: Write down the algorithm (pseudo code) of the Newton's method to solve $f(x) = 0$. The algorithm should perform a fixed number of iterations.

Algorithm: To solve $f(x) = 0$ using the following iterative formula (given an initial approximation x_0):

$$x_k = x_{k-1} - \frac{f(x_{k-1})}{f'(x_{k-1})}, \quad \text{for } k = 1, 2, 3, \dots$$

INPUTS: $\{x_0$: a real value as the initial approximation x_0 sufficiently close to the root
 N : an integer as the maximum number of iterations

OUTPUT: $\{xn$: a real value as the approximate solution
 (on completing N iterations)

Step 1 Receive the inputs as stated above

Step 2 Set $xn = x_0$ (initialize xn with the initial approximation)

Step 3 for $k = 1, 2, 3, \dots, N$ perform Steps 4-6

Step 4 Set $xp = xn$ $\{xp$ is to keep a copy of the approximation xn ,
 because xn is going to be updated.

Step 5 Set fxp as the value of $f(xp)$
 Set $dfxp$ as the value of $f'(xp)$

Step 6

$$xn = xp - \frac{fxp}{dfxp} \quad \left\{ \begin{array}{l} \text{Computing a new} \\ \text{approximation to the root} \end{array} \right.$$

end for (Go to Step 4 for the next iteration)

Step 7 Print the output: xn

[Additionally, the initial approximation (x_0), number of iterations (k), and $f(xn)$ can be printed]

STOP.

Remark: In the algorithm, it is assumed that neither any pitfall of the method will occur, nor $f(x)$ will be equal to zero (or the machine-epsilon) in any iteration for the given problem and initial approximation.

Problem 18: Write a C++ program to find a real root of $f(x) = 4x + \sin x - e^x = 0$ using the Newton-Raphson method. Take initial approximation as $x_0 = 0$. Here $f'(x) = 4 + \cos x - e^x$. The program should perform a fixed number of iterations.

```

1  # include <iostream>
2  # include <cmath>
3  using namespace std ;
4  # define N 100                                // maximum number of iterations
5
6  int main ( )
7  {
8      int k ;
9      double xn , xp , fxp , dfxp ;
10
11     cout<< "Enter the initial approximation x0: " ;
12     cin>> xn ;
13
14     //----- Processing Section -----//
15
16     for ( k=1 ; k<=N ; k++ )
17     {
18         xp = xn ;
19         fxp = 4*xp + sin(xp) - exp(xp) ;
20         dfxp = 4 + cos(xp) - exp(xp) ;
21         xn = xp - fxp / dfxp ;
22     }
23
24     //----- Output Section -----//
25
26     cout<< endl << "An approximate root of the given function is " << x << endl ;
27     cout<< endl << N << " iterations completed." ;
28
29 }

```

$$x_k = x_{k-1} - \frac{f(x_{k-1})}{f'(x_{k-1})}$$

Remark: In the program of Problem 18, the code segment of line 26 can be placed just before line 22 to print the latest result on completion of each of the iterations.

Remark: In the program, it is assumed that neither any pitfall of the method will occur, nor $f(x)$ will be equal to zero (or the machine-epsilon) in any iteration for the given problem and initial approximation.

Remark: In the program, it is assumed that neither any pitfall of the method will occur, nor $f(x)$ will be equal to zero (or the machine-epsilon) in any iteration.

Remark: The algorithm in Question 36 (likewise 18) has a shortcoming that on completion of the given fixed number of N iterations the solutions might not have been converged (the desired accuracy might not have been achieved). Moreover, the algorithm has a shortcoming if the convergence has been achieved (or divergence has occurred) in few iterations, even then the iterations would not stop immediately; the algorithm will complete the fixed number of iterations. These shortcomings in the algorithm can be addressed by incorporating the two convergence criteria such that if the convergence is achieved

(i. e., error < tolerance), then no more iterations will be performed, however, the number of iterations would not exceed the maximum limit on the number of iterations. Such an indispensable modification regarding the stopping criteria is adopted throughout the subsequent part of the book.

Question 37: Write down the algorithm (pseudo code) of the Newton's method to solve $f(x) = 0$.

Algorithm: To solve $f(x) = 0$ using the following iterative formula (given an initial approximation x_0):

$$x_k = x_{k-1} - \frac{f(x_{k-1})}{f'(x_{k-1})}, \quad \text{for } k = 1, 2, 3, \dots$$

INPUTS: $\left\{ \begin{array}{l} \mathbf{x0}$: a real value as the initial approximation x_0 sufficiently close to the root
 \mathbf{TOL} : a real value as the tolerance (permissible error)
 \mathbf{N} : an integer as the maximum number of iterations

OUTPUT: $\left\{ \begin{array}{l} \mathbf{xn}$: a real value as the approximate solution
 (either on convergence or on completing \mathbf{N} iterations – whichever happens first)

Step 1 Receive the inputs as stated above

Step 2 Set $\mathbf{xn} = \mathbf{x0}$ (initialize \mathbf{xn} with the initial approximation)

Step 3 for $k = 1, 2, 3, \dots, N$ perform Steps 4-8

Step 4 Set $\mathbf{xp} = \mathbf{xn}$ $\left\{ \begin{array}{l} \mathbf{xp}$ is to keep a copy of the approximation \mathbf{xn} ,
 because \mathbf{xn} is going to be updated.

Step 5 Set \mathbf{fxp} as the value of $f(\mathbf{xp})$
 Set \mathbf{dfxp} as the value of $f'(\mathbf{xp})$

Step 6 $\mathbf{xn} = \mathbf{xp} - \frac{\mathbf{fxp}}{\mathbf{dfxp}}$ $\left\{ \begin{array}{l} \text{Computing a new} \\ \text{approximation to the root} \end{array} \right.$

Step 7 Set $\mathbf{err} = |\mathbf{xn} - \mathbf{xp}|$ (or $\mathbf{err} = |\mathbf{xn} - \mathbf{xp}|/|\mathbf{xp}|$)

Step 8 $\left. \begin{array}{l} \text{if } (\mathbf{err} < \mathbf{TOL}) \text{ then} \\ \text{Exit/Break the loop} \end{array} \right\} \begin{array}{l} \text{This means that the consecutive} \\ \text{approximations are nearly the same.} \\ \text{Therefore, stop iterations.} \end{array}$

end for (Go to Step 4 for the next iteration)

Step 9 Print the output: \mathbf{xn}

[Additionally, the initial approximation ($\mathbf{x0}$), number of iterations (\mathbf{k}), $f(\mathbf{xn})$, and error (\mathbf{err}) can be printed]

if ($\mathbf{err} < \mathbf{TOL}$) OUTPUT ('The desired accuracy achieved; Solution converged.')

else OUTPUT ('The number of iterations exceeded the maximum limit.') because $\mathbf{k} > \mathbf{N}$

STOP.

Remark: In the algorithm, it is assumed that neither any pitfall of the method will occur, nor $f(x)$ will be equal to zero (or the machine-epsilon) in any iteration for the given problem and initial approximation.

Problem 20: Write a C++ program to find a real root of $f(x) = 4x + \sin x - e^x = 0$ using the Newton-Raphson method. Take initial approximation as $x_0 = 0$. Here $f'(x) = 4 + \cos x - e^x$. The iterations of the method should stop when either the approximation is accurate within 10^{-5} , or the number of iterations exceed 100, whichever happens first.

```

1  # include <iostream>
2  # include <cmath>
3  using namespace std ;
4  # define TOL 0.000001          // error tolerance
5  # define N 100                // maximum number of iterations
6
7  int main ( )
8  {
9      int k ;
10     double x0, xn , xp , fxp , dfxp , err ;
11
12     cout<< "Enter the initial approximation x0: " ;
13     cin>> x0 ;
14     xn = x0 ;
15
16     //----- Processing Section -----//
17
18     for ( k=1 ; k<=N ; k++ )
19     {
20         xp = xn ;
21         fxp = 4*xp + sin(xp) - exp(xp) ;
22         dfxp = 4 + cos(xp) - exp(xp) ;
23         xn = xp - fxp / dfxp ;
24         err = abs(xn - xp) ;
25
26         cout<< "After " << k << "iterations, the approximate root = " << xn ;
27         cout<< "    f(x) = " << fxp << "    Error = " << err << endl ;
28
29         if ( err < TOL )    break ;
30     }
31
32     if ( err < TOL )
33         cout<< "The desired accuracy achieved; Solution converged." ;
34     else
35         cout<< "The number of iterations exceeded the maximum limit." ;
36 }

```

Remark: This program is based on the assumption that neither any pitfall of the method will occur, nor $f(x)$ will be equal to zero (or machine-epsilon) in any iteration for the given problem and data.

Problem 22: Write a C++ program to find a real root of the equation $f(x) = 4x + \sin x - e^x = 0$ using the Newton-Raphson method. Take initial approximation as $x_0 = 0$. Here $f'(x) = 4 + \cos x - e^x$. Use #define directive to evaluate $f(x)$ and $f'(x)$ wherever required. The iterations of the method should stop when either the approximation is accurate within 10^{-5} , or the number of iterations exceed 100, whichever happens first.

```

1  # include <iostream>
2  # include <cmath>
3
4  using namespace std ;
5
6  # define TOL 0.000001                // error tolerance
7  # define N 100                      // maximum number of iterations
8  # define fval(x) 4*x + sin(x) - exp(x) // Evaluating f(x)
9  # define dfval(x) 4 + cos(x) - exp(x) // Evaluating f'(x)
10
11 int main ( )
12 {
13     int k ;
14     double x0, xn , xp , fxp , dfxp , err ;
15     cout<< "Enter the initial approximation x0: " ;
16     cin>> x0 ;
17     xn = x0 ;
18
19     //----- Processing Section -----//
20
21     for ( k=1 ; k<=N ; k++ )
22     {
23         xp = xn ;
24         fxp = fval(xp) ;
25         dfxp = dfval(xp) ;
26         xn = xp - fxp / dfxp ;            $x_k = x_{k-1} - \frac{f(x_{k-1})}{f'(x_{k-1})}$ 
27         err = abs(xn - xp) ;           Error =  $|x_k - x_{k-1}|$ 
28
29         cout<< "After " << k << " iterations, the approximate root = " << xn ;
30         cout<< "    f(x) = " << fxp << "    Error = " << err << endl ;
31
32         if ( err < TOL )    break ;
33     }
34
35     if ( err < TOL )    cout<<"The desired accuracy achieved; Solution converged.";
36     else                cout<< "The maximum number of iterations exceeded the limit." ;
37 }

```

Question 38: Write down the algorithm (pseudo code) of the Fixed-Point Iteration method to solve $f(x) = 0$.

Algorithm: To solve $f(x) = 0 \Leftrightarrow x = g(x)$, using the following iterative formula (given an initial approximation x_0)

$$x_k = g(x_{k-1}), \quad \text{for } k = 1, 2, 3, \dots$$

INPUTS: $\left\{ \begin{array}{l} \mathbf{x0}$: a real value as the initial approximation x_0 sufficiently close to the root
 \mathbf{TOL} : a real value as the tolerance (permissible error)
 \mathbf{N} : an integer as the maximum number of iterations

OUTPUT: $\left\{ \begin{array}{l} \mathbf{xn}$: a real value as the approximate solution
 ((either on convergence or on completing \mathbf{N} iterations – whichever happens first)

Step 1 Receive the inputs as stated above

Step 2 Set $\mathbf{xn} = \mathbf{x0}$ (initialize \mathbf{xn} with the initial approximation)

Step 3 for $k = 1, 2, 3, \dots, N$ perform Steps 4-7

Step 4 Set $\mathbf{xp} = \mathbf{xn}$ $\left\{ \begin{array}{l} \mathbf{xp}$ is to keep a copy of the approximation \mathbf{xn} ,
 because \mathbf{xn} is going to be updated.

Step 5

Set \mathbf{xn} as the value of $g(\mathbf{xp})$ $\left\{ \begin{array}{l} \text{Computing a new} \\ \text{approximation to the root} \end{array} \right.$

Step 6 Set $\mathbf{err} = |\mathbf{xn} - \mathbf{xp}|$ (or $\mathbf{err} = |\mathbf{xn} - \mathbf{xp}|/|\mathbf{xp}|$)

Step 7

if ($\mathbf{err} < \mathbf{TOL}$) then
 Exit/Break the loop } This means that the consecutive
 approximations are nearly the same.
 Therefore, stop iterations.

end for (Go to Step 4 for the next iteration)

Step 8 Print the output: \mathbf{xn}

[Additionally, the initial approximation ($\mathbf{x0}$), number of iterations (\mathbf{k}), $f(\mathbf{xn})$, and error (\mathbf{err}) can be printed]

if ($\mathbf{err} < \mathbf{TOL}$) OUTPUT ('The desired accuracy achieved; Solution converged.')

else OUTPUT ('The number of iterations exceeded the maximum limit.')

STOP.

Problem 24: Write a C++ program to find a real root of the equation $f(x) = 4x + \sin x - e^x = 0$ using the Fixed-Point Iteration method. Take $x = g(x) = \frac{1}{4}(e^x - \sin x)$ and $x_0 = 0$ as an initial approximation. The iterations of the method should stop when either the approximation is accurate within 10^{-5} , or the number of iterations exceeds 100, whichever happens first.

```

1  # include <iostream>
2  # include <cmath>
3  using namespace std ;
4
5  # define TOL 0.000001          // error tolerance
6  # define N 100                // maximum number of iterations
7
8  int main ( )
9  {
10     int k ;
11     double x0 , xp , xn , err ;
12
13     cout<< "Enter the initial approximation x0: " ;
14     cin>> x0 ;
15     xn = x0 ;
16
17     //----- Processing Section -----//
18
19     for ( k=1 ; k<=N ; k++)
20     {
21
22         xp = xn ;
23
24         xn = 0.25 * ( exp(xp) - sin(xp) ) ;    // Computing g(x) at the current approx.
25
26         err = abs(xn - xp) ;                  Error = |xk - xk-1|
27
28         if ( err < TOL )
29             break ;
30
31     }
32
33     if ( err < TOL )
34         cout<< "The desired accuracy achieved; Solution converged." ;
35     else
36         cout<< "The number of iterations exceeded the maximum limit." ;
37
38 }

```

Question 39: Write down the algorithm (pseudo code) of the Secant method to solve $f(x) = 0$.

Algorithm: To solve $f(x) = 0$ using the iterative formula (given the root containing interval):

$$x_k = x_{k-1} - \frac{f(x_{k-1})(x_{k-1} - x_{k-2})}{f(x_{k-1}) - f(x_{k-2})}, \quad \text{for } k = 2, 3, 4, \dots$$

INPUTS: $\left\{ \begin{array}{l} \mathbf{a}$ and \mathbf{b} : two real values as the initial approximations sufficiently close to the root
 \mathbf{TOL} : a real value as the tolerance (permissible error)
 \mathbf{N} : an integer as the maximum number of iterations

OUTPUT: $\left\{ \begin{array}{l} \mathbf{xn}$: a real value as the approximate solution
(either on convergence or on completing \mathbf{N} iterations – whichever happens first)

Step 1 Receive the inputs as stated above

Step 2 Set $\mathbf{xn} = \mathbf{b}$ (initialize \mathbf{xn} with any of the two endpoints)

Step 3 Set $\mathbf{x0} = \mathbf{a}$
Set $\mathbf{x1} = \mathbf{b}$
Set $\mathbf{fx0}$ as the value of $f(\mathbf{x0})$
Set $\mathbf{fx1}$ as the value of $f(\mathbf{x1})$

Step 4 for $k = 2, 3, \dots, N + 1$ perform Steps 5-10

Step 5 Set $\mathbf{xp} = \mathbf{xn}$ $\left\{ \begin{array}{l} \mathbf{xp}$ is to keep a copy of the approximation \mathbf{xn} ,
because \mathbf{xn} is going to be updated.

Step 6

$$\mathbf{xn} = \mathbf{x1} - \frac{\mathbf{fx1}(\mathbf{x1} - \mathbf{x0})}{\mathbf{fx1} - \mathbf{fx0}} \quad \left\{ \begin{array}{l} \text{Computing a new} \\ \text{approximation to the root} \end{array} \right.$$

Step 7 Set \mathbf{fxn} as the value of $f(\mathbf{xn})$

Step 8 Set $\mathbf{err} = |\mathbf{xn} - \mathbf{xp}|/|\mathbf{xp}|$ (or $\mathbf{err} = |\mathbf{xn} - \mathbf{xp}|$)

Step 9

if ($\mathbf{err} < \mathbf{TOL}$) then
Exit/Break the loop } This means that the consecutive approximations are nearly the same. Therefore, stop iterations.

else
Set $\mathbf{x0} = \mathbf{x1}$
Set $\mathbf{fx0} = \mathbf{fx1}$
Set $\mathbf{x1} = \mathbf{xn}$
Set $\mathbf{fx1} = \mathbf{fxn}$ } preparing two approximations for the next iteration

end for (Go to Step 5 for the next iteration)

Step 10 Print the output: \mathbf{xn}

[Additionally, the initial approx. ($\mathbf{x0}$ and $\mathbf{x1}$), number of iterations ($k - 1$), $f(\mathbf{xn})$, and error (\mathbf{err}) can be printed]

if ($\mathbf{err} < \mathbf{TOL}$) OUTPUT ('The desired accuracy achieved; Solution converged.')

else OUTPUT ('The number of iterations exceeded the maximum limit.')

STOP.

Problem 26: Write a C++ program to find a real root of the equation $f(x) = 4x + \sin x - e^x = 0$ using the Secant method. Take initial approximation as $x_0 = 0$ and $x_1 = 1$. The iterations of the method should stop when either the approximation is accurate within 10^{-5} , or the number of iterations exceeds 100, whichever happens first.

```

1  # include <iostream>
2  #include<cmath>
3  using namespace std ;
4  # define TOL 0.000001          // error tolerance
5  # define N 100                // maximum number of iterations
6
7  int main ( )
8  {
9      int k ;
10     double x0 , x1 , xn , xp , fx0 , fx1 , fxn , fxp , err ;
11     cout<< "Enter the first initial approximation: " ;      cin>> a ;
12     cout<< "Enter the second initial approximation: " ;    cin>> b ;
13
14     //----- Processing Section -----//
15     xn = b ;
16     x0 = a ;
17     x1 = b ;
18     fx0 = 4*x0 + sin(x0) - exp(x0) ;    // Evaluating f(x) at x0
19     fx1 = 4*x1 + sin(x1) - exp(x1) ;    // Evaluating f(x) at x1
20
21     for ( k=2 ; k<=N+1 ; k++)
22     {
23         xp = xn ;
24         xn = x1 - (fx1 * (x1 - x0)) / (fx1 - fx0) ;
25         fxn = 4*xn + sin(xn) - exp(xn) ;
26         err = abs(xn - xp)/abs(xn) ;
27
28         
$$x_k = x_{k-1} - \frac{f(x_{k-1})(x_{k-1} - x_{k-2})}{f(x_{k-1}) - f(x_{k-2})}$$

29         
$$\text{Error} = \frac{|x_k - x_{k-1}|}{|x_k|}$$

30
31     cout<< "After " << k-1 << " iterations, the approximate root = " << xn ;
32     cout<< "   f(x) = " << fxn << "   Error = " << err << endl ;
33
34     if ( err < TOL )    break ;
35     else
36     {
37         x0 = x1 ;
38         fx0 = fx1 ;
39         x1 = xn ;
40         fx1 = fxn ;
41     }
42 }

```

```

43   if ( err < TOL )
44       cout << "The desired accuracy achieved; Solution converged." ;
45   else
46       cout << "The number of iterations exceeded the maximum limit." ;
47   }

```

Question 40: Write down the algorithm (pseudo code) of the Bisection method to solve $f(x) = 0$.

Algorithm: To solve $f(x) = 0$ using the iterative formula (given the root containing interval):

$$x_k = x_{k-2} + \frac{x_{k-1} - x_{k-2}}{2}, \quad \text{for } k = 2, 3, 4, \dots$$

INPUTS: $\left\{ \begin{array}{l} \mathbf{a} \text{ and } \mathbf{b}: \text{ two real values as the initial approximations bracketing the root} \\ \mathbf{TOL}: \text{ a real value as the tolerance (permissible error)} \\ \mathbf{N}: \text{ an integer as the maximum number of iterations} \end{array} \right.$

OUTPUT: $\left\{ \begin{array}{l} \mathbf{xn}: \text{ a real value as the approximate solution} \\ \text{(either on convergence or on completing } \mathbf{N} \text{ iterations – whichever happens first)} \end{array} \right.$

Step 1 Receive the inputs as stated above

Step 2 Set $\mathbf{xn} = \mathbf{b}$ (initialize \mathbf{xn} with any of the two endpoints)

Step 3 Set $\mathbf{x0} = \mathbf{a}$
Set $\mathbf{x1} = \mathbf{b}$
Set $\mathbf{fx0}$ as the value of $f(\mathbf{x0})$
Set $\mathbf{fx1}$ as the value of $f(\mathbf{x1})$

Step 4 for $\mathbf{k} = 2, 3, \dots, \mathbf{N} + 1$ perform Steps 5-10

Step 5 Set $\mathbf{xp} = \mathbf{xn}$ $\left\{ \begin{array}{l} \mathbf{xp} \text{ is to keep a copy of the approximation } \mathbf{xn}, \\ \text{because } \mathbf{xn} \text{ is going to be updated.} \end{array} \right.$

Step 6

$$\mathbf{xn} = \mathbf{x0} + \frac{\mathbf{x1} - \mathbf{x0}}{2} \quad \left\{ \begin{array}{l} \text{Computing a new} \\ \text{approximation to the root} \end{array} \right.$$

Step 7 Set \mathbf{fxn} as the value of $f(\mathbf{xn})$

Step 8 Set $\mathbf{err1} = |\mathbf{xn} - \mathbf{xp}|/|\mathbf{xn}|$ (or $\mathbf{err} = |\mathbf{xn} - \mathbf{xp}|$)
Set $\mathbf{err2} = |\mathbf{fxn}|$
Set $\mathbf{err} = \min(\mathbf{err1}, \mathbf{err2})$

Step 9

$\left. \begin{array}{l} \text{if } (\mathbf{err} < \mathbf{TOL}) \text{ then} \\ \text{Exit/Break the loop} \end{array} \right\}$ This means that either $f(\mathbf{xn})$ is close to zero, or the consecutive approximations are nearly the same. Therefore, stop iterations.

```

else if ( $f(x_0)f(x_n) < 0$ ) then
    Set  $x_1 = x_n$ 
    Set  $f_{x_1} = f_{x_n}$ 
else
    Set  $x_0 = x_n$ 
    Set  $f_{x_0} = f_{x_n}$ 
end for      (Go to Step 5 for the next iteration)

```

} Adjusting one endpoint
of the interval such that
half of the interval will be
used in the next iteration

Step 10 Print the output: x_n

[Additionally, the starting interval $[a, b]$, number of iterations $(k - 1)$, $f(x_n)$, and error (err) can be printed]

if ($err < TOL$) OUTPUT ('The desired accuracy achieved; Solution converged.')

else OUTPUT ('The number of iterations exceeded the maximum limit.')

STOP.

Remark: While using a bracketing method, there might arise a situation in which the two consecutive approximations to the roots are not sufficiently close to each other (i.e., the sequence of successive approximations has not converged), but the function values at the approximations are sufficiently close to zero (i.e., $|f(x_k)| < \text{tolerance}$). Therefore, there is no point to proceed the iterations further. The iterations should be stopped. Therefore, the algorithm of a bracketing method (the Bisection, or Regula-Falsi method) should include both of the convergence criteria of testing the convergence of the roots, and closeness of the function values to zero. The iterations should be terminated on whichever criterion is met first, ensuring the convergence. To accommodate this in the algorithm, the two kinds of errors are computed and the minimum of the two errors is found to compare with the tolerance:

Set $err1 = |x_n - x_p|/|x_n|$ (or $err = |x_n - x_p|$)

Set $err2 = |f_{x_n}|$

Set $err = \min(err1, err2)$

Problem 28: Write a C++ program to find a real root of the equation $f(x) = 4x + \sin x - e^x = 0$ in $[0, 1]$ using the Bisection method. The two function values at the endpoints of the interval have opposite signs. The iterations of the method should stop when either the approximation is accurate within 10^{-5} , or the number of iterations exceeds 100, whichever happens first.

```

1 # include <iostream>
2 # include <cmath>
3 using namespace std ;
4

```

```

5 # define TOL 0.000001 // error tolerance
6 # define N 100 // maximum number of iterations
7
8 int main (
9 {
10 int k ;
11 double x0 , x1 , xn , xp , fx0 , fx1 , fxn , fxp , err , err1 , err2 ;
12 cout<< "Enter the left endpoint of the interval: " ; cin>> a ;
13 cout<< "Enter the right endpoint of the interval: " ; cin>> b ;
14
15 //----- Processing Section -----//
16
17 xn = b ;
18 x0 = a ;
19 x1 = b ;
20 fx0 = 4*x0 + sin(x0) - exp(x0) ; // Evaluating f(x) at x0
21 fx1 = 4*x1 + sin(x1) - exp(x1) ; // Evaluating f(x) at x1
22
23 for ( k=2 ; k<=N+1 ; k++)
24 {
25 xp = xn ;
26
27 xn = x0 + (x1 - x0) / 2 ;  $x_k = x_{k-2} + \frac{x_{k-1} - x_{k-2}}{2}$ 
28 fxn = 4*xn + sin(xn) - exp(xn) ;
29
30 err1 = abs(xn - xp)/abs(xn) ; Error 1 =  $|x_k - x_{k-1}|/|x_k|$ 
31 err2 = abs(fxn) ; Error 2 =  $|f(x_k)|$ 
32 err = min(err1, err2) ;
33
34 cout<< "After " << k-1 << " iterations, the approximate root = " << xn ;
35 cout<< " f(x) = " << fxn << " Error = " << err1 << endl ;
36
37 if ( err < TOL ) break ;
38 else if ( fx0*fxn < 0 )
39 {
40 x1 = xn ;
41 fx1 = fxn ;
42 }
43 else
44 {
45 x0 = xn ;
46 fx0 = fxn ;
47 }
48
49 }

```



```

50
51   if ( err < TOL )
52       cout << "The desired accuracy achieved; Solution converged." ;
53   else
54       cout << "The number of iterations exceeded the maximum limit." ;
55 }

```

Question 41: Write down the algorithm (pseudo code) of the Regula-Falsi method to solve $f(x) = 0$.

Algorithm: To solve $f(x) = 0$ using the iterative formula (given the root containing interval):

$$x_k = x_{k-1} - \frac{f(x_{k-1})(x_{k-1} - x_{k-2})}{f(x_{k-1}) - f(x_{k-2})}, \quad \text{for } k = 2, 3, 4, \dots$$

INPUTS: $\left\{ \begin{array}{l} \mathbf{a} \text{ and } \mathbf{b}: \text{ two real values as the initial approximations bracketing the root} \\ \mathbf{TOL}: \text{ a real value as the absolute error tolerance} \\ \mathbf{N}: \text{ an integer as the maximum number of iterations} \end{array} \right.$

OUTPUT: $\left\{ \begin{array}{l} \mathbf{xn}: \text{ a real value as the approximate solution} \\ \text{(either on convergence or on completing } \mathbf{N} \text{ iterations – whichever happens first)} \end{array} \right.$

Step 1 Receive the inputs as stated above

Step 2 Set $\mathbf{xn} = \mathbf{b}$ (initialize \mathbf{xn} with any of the two endpoints)

Step 3 Set $\mathbf{x0} = \mathbf{a}$
Set $\mathbf{x1} = \mathbf{b}$
Set $\mathbf{fx0}$ as the value of $f(\mathbf{x0})$
Set $\mathbf{fx1}$ as the value of $f(\mathbf{x1})$

Step 4 for $k = 2, 3, \dots, N + 1$ perform Steps 5-10

Step 5 Set $\mathbf{xp} = \mathbf{xn}$ $\left\{ \begin{array}{l} \mathbf{xp} \text{ is to keep a copy of the approximation } \mathbf{xn}, \\ \text{because } \mathbf{xn} \text{ is going to be updated.} \end{array} \right.$

Step 6

$$\mathbf{xn} = \mathbf{x1} - \frac{\mathbf{fx1}(\mathbf{x1} - \mathbf{x0})}{\mathbf{fx1} - \mathbf{fx0}} \quad \left\{ \begin{array}{l} \text{Computing a new} \\ \text{approximation to the root} \end{array} \right.$$

Step 7 Set \mathbf{fxn} as the value of $f(\mathbf{xn})$

Step 8 Set $\mathbf{err1} = |\mathbf{xn} - \mathbf{xp}|/|\mathbf{xn}|$ (or $\mathbf{err} = |\mathbf{xn} - \mathbf{xp}|$)

Set $\mathbf{err2} = |\mathbf{fxn}|$

Set $\mathbf{err} = \min(\mathbf{err1}, \mathbf{err2})$

Step 9

```

if (err < TOL) then
    Exit/Break the loop } This means that either  $f(x_n)$  is close to
                           zero, or the consecutive approximations are
                           nearly the same. Therefore, stop iterations.
else if ( $f(x_0)f(x_n) < 0$ ) then
    Set  $x_1 = x_n$ 
    Set  $fx_1 = fx_n$ 
else
    Set  $x_0 = x_n$ 
    Set  $fx_0 = fx_n$ 
} Adjusting one endpoint
  of the interval such that
  a shorter interval will be
  used in the next iteration
end for      (Go to Step 5 for the next iteration)

```

Step 10 Print the output: x_n

[Additionally, the starting interval $[a, b]$, number of iterations $(k - 1)$, $f(x_n)$, and error (*err*) can be printed]

if (*err* < *TOL*) OUTPUT ('The desired accuracy achieved; Solution converged.')

else OUTPUT ('The number of iterations exceeded the maximum limit.')

STOP.

Problem 30: Write a C++ program to find a real root of the equation $f(x) = 4x + \sin x - e^x = 0$ in $[0, 1]$ using the Regula-Falsi method. The two function values at the endpoints of the interval have opposite signs. The iterations of the method should stop when either the approximation is accurate within 10^{-5} , or the number of iterations exceeds 100, whichever happens first.

```

1  # include <iostream>
2  # include <cmath>
3  using namespace std ;
4
5  # define TOL 0.000001          // error tolerance
6  # define N 100                // maximum number of iterations
7
8  int main ( )
9  {
10     int k ;
11     double x0 , x1 , xn , xp , fx0 , fx1 , fxn , fxp , err , err1 , err2 ;
12     cout << "Enter the left endpoint of the interval: " ;      cin >> a ;
13     cout << "Enter the right endpoint of the interval: " ;    cin >> b ;
14
15     //----- Processing Section -----//
16
17     xn = b ;
18     x0 = a ;
19     x1 = b ;
20     fx0 = 4*x0 + sin(x0) - exp(x0) ;    // Evaluating f(x) at x0
21     fx1 = 4*x1 + sin(x1) - exp(x1) ;    // Evaluating f(x) at x1
22

```

```

23   for ( k=2 ; k<=N+1 ; k++)
24   {
25       xp = xn ;
26
27       xn = x1 - (fx1 * (x1 - x0)) / (fx1 - fx0) ;
28       fxn = 4*xn + sin(xn) - exp(xn) ;
29
30       err1 = abs(xn - xp)/abs(xn) ;
31       err2 = abs(fxn) ;
32       err = min(err1, err2) ;
33
34       cout<< "After " << k-1 << " iterations, the approximate root = " << xn ;
35       cout<< "   f(x) = " << fxn << "   Error = " << err1 << endl ;
36
37       if ( err < TOL )   break ;
38       else if ( fx0*fxn < 0 )
39       {
40           x1 = xn ;
41           fx1 = fxn ;
42       }
43       else
44       {
45           x0 = xn ;
46           fx0 = fxn ;
47       }
48   }
49
50
51   if ( err < TOL )
52       cout<< "The desired accuracy achieved; Solution converged." ;
53   else
54       cout<< "The number of iterations exceeded the maximum limit." ;
55 }

```

Problem 32: Write a C++ program to find a real root of the equation $f(x) = 4x + \sin x - e^x = 0$ in $[0, 1]$ using the Regula-Falsi method. The two function values at the endpoints of the interval have opposite signs. Use #define directive to evaluate $f(x)$ wherever required. The iterations of the method should stop when either the approximation is accurate within 10^{-5} , or the number of iterations exceeds 100, whichever happens first.

```

1  # include <iostream>
2  # include <cmath>
3  using namespace std ;
4
5  # define TOL 0.000001

```

// error tolerance

```

6 # define N 100 // maximum number of iterations
7 # define fval(x) 4 * x + sin(x) - exp(x) // Evaluating f(x)
8
9 int main ( )
10 {
11     int k ;
12     double x0 , x1 , xn , xp , fx0 , fx1 , fxn , fxp , err , err1 , err2 ;
13     cout<< "Enter the left endpoint of the interval: " ; cin>> a ;
14     cout<< "Enter the right endpoint of the interval: " ; cin>> b ;
15
16     //----- Processing Section -----//
17
18     xn = b ;
19     x0 = a ;
20     x1 = b ;
21     fx0 = fval(x0) ; // Evaluating f(x) at x0
22     fx1 = fval(x1) ; // Evaluating f(x) at x1
23
24     for ( k=2 ; k<=N+1 ; k++)
25     {
26         xp = xn ;
27
28         xn = x1 - (fx1 * (x1 - x0)) / (fx1 - fx0) ;  $x_k = x_{k-1} - \frac{f(x_{k-1})(x_{k-1} - x_{k-2})}{f(x_{k-1}) - f(x_{k-2})}$ 
29         fxn = fval(xn) ;
30
31         err1 = abs(xn - xp)/abs(xn) ; Error 1 =  $|x_k - x_{k-1}|/|x_k|$ 
32         err2 = abs(fxn) ; Error 2 =  $|f(x_k)|$ 
33         err = min(err1, err2) ;
34
35         cout<< "After " << k-1 << " iterations, the approximate root = " << xn ;
36         cout<< " f(x) = " << fxn << " Error = " << err1 << endl ;
37
38         if ( err < TOL ) break ;
39         else if ( fx0*fxn < 0 )
40         {
41             x1 = xn ;
42             fx1 = fxn ;
43         }
44         else
45         {
46             x0 = xn ;
47             fx0 = fxn ;
48         }
49     }
50 }

```

```
51
52   if ( err < TOL )
53       cout<< "The desired accuracy achieved; Solution converged." ;
54   else
55       cout<< "The number of iterations exceeded the maximum limit." ;
56 }
```

■

Remark: An interesting online calculator by CASIO® at <https://keisan.casio.com> has the following webpage to approximate the root of a non-linear equation using different methods.

<https://keisan.casio.com/menu/system/000000001000>

■■■

Chapter Summary

- The root-finding problem refers to find some appropriate value $x = \alpha$ in the domain of a function f such that $f(\alpha) = 0$. Every such possible value α is called a **root of the equation** $f(x) = 0$.
- Geometrically, a root of an equation $f(x) = 0$ is the point where the graph of f intersects the x -axis.
- An iterative numerical method to approximate the root starts with some appropriate or reasonable estimation (also called **initial approximation** or guess) of the exact root and attempts to refine the approximation, iteratively. The iterations are repeated until a desired level of accuracy is achieved.
- Let x_0 denotes the initial approximation and x_1, x_2, x_3, \dots denote the successive iterative solutions to an exact root α of the equation $f(x) = 0$. The sequence $\{x_k\}_{k=0}^{\infty}$ of the successive approximations is said to **converge** to the exact root α , if the successive approximations approach α . In such a case, the iterative method is also said to converge. In other words, the iterative method is said to be **convergent** for a given initial approximation if the corresponding sequence of successive approximations is convergent to the exact solution. Under certain conditions, it is possible for an iterative method that the sequence of successive approximations might **diverge** from a desired exact root α .
- **Stopping Criteria:** The most common convergence criterion to stop the iterative process is based on the comparison of the estimated error with the error tolerance. For this purpose, the current approximation is considered as the true solution and the previous approximation is considered as the approximate solution for estimating the error and any appropriate one of the following criteria is used,

$$(1) \quad |x_k - x_{k-1}| \leq \tau$$

$$(2) \quad \left| \frac{x_k - x_{k-1}}{x_k} \right| \leq \tau$$

$$(3) \quad \left| \frac{x_k - x_{k-1}}{x_k} \right| \times 100 \leq \tau$$

Here x_k and x_{k-1} denote the current and previous approximations, respectively, and τ denotes the tolerance.

- **Another Stopping Criterion:** Note that the values of the function f tend to zero with the progress of the iterative process. Thus, falling of the difference between the function values and zero beyond a certain level might also indicate convergence.
- The numerical methods of finding a root of $f(x) = 0$ can be categorized as bracketing methods and open methods.
- **Bracketing methods** start with an interval containing a root and squeeze down the interval, iteratively. Two well-known root bracketing methods are the Bisection method and the Regula-Falsi (False-Position) method.

- **Open Methods** are those who obtain successive single approximations irrespective of their location at any side of the root. Some of the well-known open methods are the Fixed-Point Iteration method, the Newton-Raphson method (Newton's method), and the Secant method.

- A bracketing method for finding a root/zero of a continuous function f starts with an interval $[a, b]$ containing a root. The opposite signs of $f(a)$ and $f(b)$ ensure (due to the Intermediate value theorem) that there exists a root α of $f(x) = 0$ in (a, b) . To get closer to the root α , first a point $c \in (a + b)$ is chosen. If $f(c) = 0$, then c is the exact root. Otherwise, either of the intervals $[a, c]$ or $[c, b]$ is chosen as the squeezed interval containing the root. The root lies in $[a, c]$ if $f(a)f(c) < 0$, or in $[c, b]$ if $f(c)f(b) < 0$. The selected interval is relabeled as $[a, b]$ and the process is repeated. This way, a sequence of points c_1, c_2, c_3, \dots , is formed. The iterations are performed until the approximations of the root of $f(x)$ in two consecutive iterations are sufficiently close to each other.

- **The Bisection method** selects $c \in (a + b)$, as the midpoint of the interval $[a, b]$, using the formula

$$c = \frac{(a + b)}{2}$$

- **The Regula-Falsi method** selects $c \in (a + b)$, as the point where the line segment joining $f(a)$ and $f(b)$ intersects the x -axis, using the formula

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

- For the Bisection method, the error-bound is given by,

$$|\alpha - c_k| \leq \frac{b - a}{2^k}, \quad \text{for } k = 1, 2, 3, \dots,$$

Here α is the exact root of the equation $f(x) = 0$ in (a, b) and $c_k = \frac{a_{k-1} + b_{k-1}}{2}$ is the midpoint of the interval in k th iteration.

- The formula to determine the maximum number of iterations N of the Bisection method after which the error associated with any point in the squeezed interval is not greater than a given permissible absolute error τ_a is as below:

$$N \geq \frac{\log(b - a) - \log(\tau_a)}{\log(2)}$$

This formula tells that, for an interval of unit length, it is sure that after 10, 14, 17, and 20 iterations the length of the squeezed interval (or the absolute error) is not greater than 10^{-3} , 10^{-4} , 10^{-5} , and 10^{-6} , respectively.

- **The Fixed-Point Iteration method** is an open method that approximates a root of the equation $f(x) = 0$ by rearranging the equation $f(x) = 0$ to get an appropriate form $x = g(x)$ and generating a sequence of successive approximations $\{x_k\}_{k=1}^{\infty}$ by the iterative formula $x_k = g(x_{k-1})$, for $k = 1, 2, 3, \dots$. The said sequence may

- converge but could be different for different forms of $x = g(x)$,

- converge but could be different for different choices of the initial approximation x_0 for a particular form of $x = g(x)$, or
 - diverge for some unsuitable form of $x = g(x)$ or an initial approximation x_0 .
- Suppose that f is a continuous function and the equation $f(x) = 0$ has a real root α . Suppose that the equation $f(x) = 0$ can be rearranged in the form $x = g(x)$ such that α is a fixed-point of the function g , and g and g' are continuous in some neighbourhood I around α . If

$$|g'(x)| \leq K < 1, \quad \text{for all } x \in I,$$

then for any initial approximation $x_0 \in I$, the sequence $\{x_k\}_{k=1}^{\infty}$ of successive approximations, generated by the iterative formula $x_k = g(x_{k-1})$, for $k = 1, 2, 3, \dots$, converges to the solution α .

- To find a root of a non-linear equation $f(x) = 0$ **the Newton-Raphson method** requires an initial solution x_0 and considers the x -intercept of the tangent line to the function $f(x)$ at $x = x_0$ as the new approximation. Then, the x -intercept of the tangent line to the function at the new approximation is considered as the next approximation. This way, the process is repeated with the successive approximations until sufficient convergence is achieved. The formula to generate the sequence of successive approximations based on the said approach is given by

$$x_k = x_{k-1} - \frac{f(x_{k-1})}{f'(x_{k-1})}, \quad \text{for } k = 1, 2, 3, \dots$$

- A sufficient condition of convergence for the Newton-Raphson method: Suppose that α is a root of the equation $f(x) = 0$. Suppose that I is a neighbourhood of α such that $f(x)$, $f'(x)$ and $f''(x)$ are continuous on I . If $|f(x)f''(x)| \leq |f'(x)|^2$, for all $x \in I$, then for an initial approximation $x_0 \in I$, the sequence $\{x_k\}_{k=1}^{\infty}$ of successive approximations, generated by the Newton's formula, converges to the solution α .
- The iterative formula of **the Secant method** for solving $f(x) = 0$ (with $x = x_0$ and $x = x_1$ as the initial approximations) is given by

$$x_k = x_{k-1} - \frac{f(x_{k-1})(x_{k-1} - x_{k-2})}{f(x_{k-1}) - f(x_{k-2})}, \quad \text{for } k = 2, 3, 4, \dots$$

- Comparison of the False-Position method and the Secant method:
 - The False-Position method is a bracketing method, whereas the Secant method is an open method.
 - The False-Position method keeps the root bracketed by selects out the root bracketing subintervals out the two subintervals obtains in each of the iterations. On the other hand, the Secant method selects the two most recent approximations out of the three available approximations in any iteration to proceed to the next iteration.
 - The False-Position method always converges, whereas the Secant method may not converge for certain situations.

- If the Secant method is convergent, it converges faster than the False-Position method. That is, it has a higher convergence rate than that of the False-Position method.
- The order/rate of convergence of the Bisection method is 1 (i.e., linear) and the asymptotic error constant is $(1/2)$
- The order/rate of convergence of the False-Position or Regula-Falsi method is 1 (i.e., linear) and the asymptotic error constant is $-\frac{1}{2} \frac{f''(\alpha)}{f'(\alpha)} \varepsilon_0$
- The order/rate of convergence of the Fixed-Point Iteration method is 1 (i.e., linear) and the asymptotic error constant is the maximum value of the function $g'(x)$ in some neighbourhood around the solution α .
- The order/rate of convergence of the Newton-Raphson method is 2 (i.e., quadratic) and the asymptotic error constant is $-\frac{1}{2} \frac{f''(\alpha)}{f'(\alpha)}$
- The order/rate of convergence of the Secant method is 1.62 (i.e., superlinear).
- The Newton-Raphson method may fail to converge to a root in different situations including where $f'(x)$ or $f''(x)$ becomes zero at any approximation.
- The Newton-Raphson method converges to a multiple root very slowly (instead of exhibiting quadratic convergence).
- **The Aitken's Δ^2 method** offers a technique for accelerating the convergence of any sequence that is linearly convergent. From the given sequence $\{x_k\}_{k=1}^{\infty}$ that linearly converges to α , another sequence $\{\bar{x}_k\}_{k=1}^{\infty}$ that also converges to α with possibly improved convergence rate is constructed by using the Aitken's acceleration formula given as

$$\bar{x}_k \cong x_k - \frac{(\Delta x_k)^2}{\Delta^2 x_k}$$



Chapter Exercises

Exercise 01: Find a real root of the following equations using the Bisection method accurate to four decimal places.

- (i) $\log(x) - \cos x = 0$
- (ii) $e^{-x} - 10x = 0$
- (iii) $x^3 + x^2 - 1 = 0$

Exercise 02: Find a real root of the following equations using the Bisection method accurate to three decimal places.

- (i) $x^6 - x^4 - x^3 - 1 = 0$
- (ii) $x^3 - \sin x + 1 = 0$
- (iii) $x \log_{10} x = 4.77$

Exercise 03: Approximate the solution of the following equations using the Regula-Falsi method accurate to three decimal places.

- (i) $3x + \sin x - e^x = 0$
- (ii) $4x^3 - 1 - e^{(x^2/2)} = 0$
- (iii) $x^2 = (e^{-2x} - 2)/x$

Exercise 04: Find the approximation to a real root of the equation $2 \sin x - \frac{e^x}{4} - 1 = 0$ starting with $[-5, -3]$ using the Regula-Falsi method.

Exercise: Find a real root of each of the following equations using (a) the Bisection method, (b) the Regula-Falsi method, (c) the Newton's method, (d) the Secant method. Choose the initial approximation/s in the given interval. Assume that the tolerance for the approximate root is 0.001. The numeric values should not be rounded to less than 5 decimal places. (x is in radians, wherever applicable).

- (i) $\cos x - xe^x = 0$, in $[0, 1]$
- (ii) $\cos x - x + 2 = 0$, in $[1, 2]$
- (iii) $e^x - x - 3 = 0$, in $[1, 2]$
- (iv) $\ln(x) + x - 4 = 0$, in $[2, 3]$
- (v) $4x + \sin x - e^x = 0$, in $[0, 1]$.

Exercise 06: Find a real root of the Chebyshev polynomial of degree four, $T_4(x) = 8x^4 - 8x^2 + 1$ using the Newton's method accurate to four decimal places.

Exercise 07: Find a root of the *Laquerre* polynomial of degree four, $L_4(x) = x^4 - 16x^3 + 72x^2 - 96x + 24$ using the Newton's method accurate to four decimal places.

Exercise 08: Find a root of the following equations using the Newton's method accurate to 4 decimal places.

- (i) $2x + 3 \cos x - e^x = 0,$
- (ii) $x^2 - 4x + 4 - \ln x = 0$
- (iii) $\tan x - 6 = 0$

Exercise 09: Find the roots accurate to within 10^{-3} of the Legendre polynomial $P_4(x) = x^4 - \frac{6}{7}x^2 + \frac{3}{35}$ on each interval, using the Secant method.

- (i) $[-1, -0.5]$
- (ii) $[-0.5, 0]$
- (iii) $[0, 0.5]$
- (iv) $[0.5, 1]$

Exercise 10: Approximate the value of $\sqrt[3]{4}$ using the Secant method accurate to 10^{-4} .

Exercise 11: Find a real root of the following equations using the Secant method accurate to 10^{-3} .

- (i) $x^3 - 2x + 2 = 0$
- (ii) $10 - 2x + \sin x = 0$
- (iii) $2e^{-3x} + 1 - 3e^{-3x} = 0$

Exercise 12: Use the Fixed-Point method to find a root of the following, accurate to 3 decimal places.

- (i) $e^x - 2x^2$ for $0 \leq x \leq 2$
- (ii) $xe^x = 0$ for $1 \leq x \leq 2$
- (iii) $x^2 - \sin x - x = 0$

Exercise 13: Find the solutions of the following equations using the fixed-point method accurate to 10^{-3} .

- (i) $x = \tan x$
- (ii) $x = \cos x$
- (iii) $x = \sin(x + 2)$

Exercise 14: Find the solution of the equation (relevant to the vibrating beam),

$$\cos x \cosh x = 1$$

near $x = -\frac{3}{2}\pi$ using the Newton-Raphson method.

Exercise 15: The velocity V , in meters per second (m/s), of a free falling sky diver is expressed as:

$$V = \frac{gm}{D_c} \left(1 - \exp\left(\frac{-D_c t}{m}\right) \right)$$

Here m is the mass of the falling body in kilograms (kg), D_c is the drag coefficient in kilogram per second (kg/s), t is the time in seconds (s), and $g = 9.8m/s^2$ is the gravitation acceleration. If the velocity of a body of mass $85kg$ is $40m/s$ after 5 seconds of free fall, then calculate the drag coefficient.

Hint for the Solution:

Given $m = 80kg$, $V = 40m/s$, $g = 9.8m/s^2$, and $t = 5s$, the equation takes the form:

$$40 = \frac{(9.8)(85)}{D_c} \left(1 - \exp\left(\frac{-5D_c}{85}\right) \right)$$

or

$$f(D_c) = D_c + 17 \ln(1 - 0.04802D_c) = 0$$

Solve this equation for D_c , using any appropriate iterative method. To obtain an initial guess of D_c , a trick is to calculate V for different assumed values of D_c . The values of the D_c , which produce values of V close to 40, can offer reasonable initial guess of D_c . While using an iterative method, approximate error should be calculated at each iteration. [$\exp(x) = e^x$]

Exercise 16: The velocity V , in meters per second (m/s), of a free falling sky diver is expressed as:

$$V = \frac{gm}{D_c} \left(1 - \exp\left(\frac{-D_c t}{m}\right) \right)$$

Here m is the mass of the falling body in kilograms (kg), D_c is the drag coefficient in kilogram per second (kg/s), t is the time in seconds (s), and $g = 9.8m/s^2$ is the gravitation acceleration. If the velocity of a falling body with drag coefficient of $18 kg/s$ is $50m/s$ after 7 seconds of free fall, then calculate the mass m of the body, accurate to 0.0001. [$\exp(x) = e^x$]

Hint for the Solution:

Given $D_c = 18kg/s$, $V = 50m/s$, $g = 9.8m/s^2$, and $t = 7s$, the equation takes the form:

$$50 = \frac{(9.8)m}{18} \left(1 - \exp\left(\frac{-126}{m}\right) \right)$$

or

$$f(m) = m \ln \left(1 - \frac{91.83673}{m} \right) + 126 = 0$$

Solve this equation for m , using any appropriate iterative method. To obtain an initial guess of m , a trick is to calculate V for different assumed values of m . The values of the m , which produce values of V close to 50, can offer reasonable initial guess of m . While using an iterative method, approximate the error at each iteration.

Exercise 17: The volume V of spherical water-tank in cubic meters can be calculated as:

$$V = \frac{\pi H^2(3R - H)}{3}$$

where H denotes the height of water level in meters from the base of the tank, and R denotes the radius of the spherical tank in meters. If the radius R of a tank is 2.5 meters, then how much water level must be raised in the tank to hold 27 cubic meters of water.

Hint for the Solution:

Given $R = 2.5$ and $V = 27$, and taking $\pi = 3.14159$ the equation takes the form

$$27 = \frac{\pi H^2(7.5 - H)}{3}$$

or

$$f(H) = 3.14159H^3 - 23.56193H^2 + 81 = 0$$

Solve this equation for H , using any appropriate iterative method. Intuitively, appropriate initial guesses for H can be taken from $[0, 2R]$. While using an iterative method, approximate error should be calculated at each iteration.

Exercise 18: Numerically, compare the convergence of the method:

$$x_k = x_{k-1} - 2 \frac{f(x_{k-1})}{f'(x_{k-1})}, \quad \text{for } k = 1, 2, 3, \dots$$

with the Newton-Raphson method on a function with a known double root.

Exercise 19: The ideal gas equation relates the volume (V in L), temperature (T in K), pressure (P in atm), and the amount of gas (number of moles n) by:

$$P = \frac{nRT}{V}$$

where $R = 0.08206$ ($L atm$)/(mol K) is the gas constant.

The van der Waals equation gives the relationship between these quantities for a real gas by

$$\left(P + \frac{n^2 a}{V^2} \right) (V - nb) = nRT$$

where a and b are constants that are specific for each gas.

Calculate the volume of 2 mol CO_2 at temperature of $50^\circ C$, and pressure of 6 atm. For CO_2 , $a = 3.59 (L^2 atm)/mol^2$, and $b = 0.0427 L/mol$. Because CO_2 is a real gas, so we need to use the second equation for the solution. But for solving the second equation for the volume, obtain an appropriate guess of the volume from the first equation: ideal gas equation.

Exercise 20: Golden-ratio corresponds to the order of which method:

- (A) Secant (B) Regula-Falsi (C) Fixed-Point Iteration (D) Newton-Raphson

Exercise 21: Which of the following methods, has an explicit formula that can be used to determine the required number of iterations in advance for achieving a given accuracy:

- (A) Bisection (B) Regula-Falsi (C) Fixed-Point Iteration (D) Newton-Raphson (E) Secant

Exercise 22: The convergence rate of which of the following methods is highest:

- (A) Bisection (B) Regula-Falsi (C) Fixed-Point Iteration (D) Newton-Raphson (E) Secant

■■■

Polynomial Interpolation

Corridor I: BASICS

Let's plan it

- 3.1 Introduction
- 3.2 The Newton's Divided Difference Interpolation
- 3.3 The Lagrange Interpolation
- 3.4 Deriving the Lagrange Interpolation Formula from the Newton's Divided-Difference
- 3.5 Interpolation Formulas for Equally Spaced Nodes
- 3.6 Hermite Interpolation
- 3.7 Spline Interpolation
 - 3.7.1 Linear Spline
 - 3.7.2 Quadratic Spline
 - 3.7.3 Cubic Spline

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



Corridor II: ANALYSIS

Let's think deep

3.8 Error of Interpolation

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)

■ ■ ■

Corridor III: PROGRAMMING ARCADE

Let's do it

3.9 Algorithms and Implementations

The Newton's Divided Difference Interpolation Formula

To cross-check the results/output of the computer programs you would execute, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)

■ ■ ■

3.9 Algorithms and Implementations

Question 21: Write down an algorithm (pseudo code) to interpolate or extrapolate the function at a point using the n th-degree Newton's Divided difference interpolating polynomial.

Algorithm: Given $n + 1$ data points, approximate $f(x)$ at $x = xp$ with $P_n(xp)$.

INPUTS: $\left\{ \begin{array}{l} n: \text{an integer as the degree of interpolating polynomial} \\ x_i, 0 \leq i \leq n: \text{real values as the arbitrary nodes} \\ f_i, 0 \leq i \leq n: \text{real values as the function values corresponding to } x_i \text{ nodes} \\ xp: \text{real values as the entries} \end{array} \right.$

OUTPUT: fxp : a real number as an interpolated value at $x = xp$

Step 1 Receive the inputs as stated above

Step 2 for $i = 0, 1, \dots, n$
 $ddf_{i,0} = f_i$ (Computing zeroth divided differences, $f[x_i] = f_i$)

Step 3 (Computing the divided differences of order 1 to n)

for $j = 1, 2, \dots, n$
 for $i = 0, 1, \dots, n - j$
 $ddf_{i,j} = \frac{ddf_{i+1,j-1} - ddf_{i,j-1}}{x_{i+j} - x_i}$ $\left\{ \begin{array}{l} f[x_i, \dots, x_{i+j}] = \\ \frac{f[x_{i+1}, \dots, x_{i+j}] - f[x_i, \dots, x_{i+j-1}]}{x_{i+j} - x_i} \end{array} \right.$

Step 4 (Evaluating the interpolation polynomial at xp)

Set $pro = 1$
 Set $fxp = ddf_{0,0}$
 for $k = 1, 2, \dots, n$
 $pro = pro \times (xp - x_{k-1})$
 $fxp = fxp + pro \times ddf_{0,k}$ $\left\{ \begin{array}{l} \left(P_n = f[x_0] + \sum_{k=1}^n \left[f[x_0, \dots, x_k] \prod_{t=0}^{k-1} (xp - x_t) \right] \right) \end{array} \right.$

Step 5 Print the output: fxp

STOP.

Problem 14: Write a C++ program for the second order Newton's Divided Difference Interpolation.

```

1  # include <iostream>
2  using namespace std ;
3  # define n 2                                // degree of interpolating polynomial
4  int main ( )
5  {
6    char quit ;
7    int i , j , k ;
8    double xp , fxp , pro ;
9
10 //----- Input Section -----//
11 double ddf[n+1][n+1] , f[n+1] = {5, 10, 12} , x[n+1] = {0, 1, 3} ;
12 cout<< "The Divided Difference Interpolation. " << endl ;
13 cout<< "Enter a real value it which the interpolation is to be obtained: \n" ;
14 cin>> xp ;
15
16 //----- Processing Section -----//
17 // Computing zeroth divided differences, f[x_i ] = f_i
18 for ( i=0 ; i<=n ; i++ )
19     ddf[i][0] = f[i] ;
20
21 // Computing the divided differences of order 1 to n
22 for ( j=1 ; j<=n ; j++ )
23     for ( i=0 ; i<=n-j ; i++ )
24         ddf[i][j] = ( ddf[i+1][j-1] - ddf[i][j-1] ) / ( x[i+j] - x[i] ) ;
25
26 // Evaluating the interpolation polynomial at xp
27 pro = 1 ;    fxp = ddf[0][0] ;
28 for ( k=1 ; k<=n ; k++ )
29     {
30     pro = pro * ( xp - x[k-1] ) ;
31     fxp = fxp + pro * ddf[0][k] ;
32     }
33
34 //----- Output Section -----//
35 cout<< "The interpolate or extrapolate value of function at x = xp is " << endl ;
36 cout<< fxp ;
37 cin>> quit ;
38 return 0 ;    }

```

Problem 16: Write a C++ program for the Newton's Divided Difference Interpolation.

```

1  # include <iostream>
2  using namespace std ;
3  # define n 10                                // degree of interpolating polynomial
4
5  int main ( )
6  {
7      char quit ;
8      int i , j , k ;
9      double ddf[n+1][n+1] , f[n+1] , x[n+1] ;
10     double xp , fxp , pro ;
11
12     //----- Input Section -----//
13
14     cout<< "The Divided Difference Interpolation. " << endl ;
15     cout<< "Enter real values as the arbitrary nodes " << endl ;
16     for ( i=0 ; i<=n ; i++ )
17         cin>> x[i] ;
18     cout<< endl ;
19     cout<< " Enter real values as the function values corresponding to x_i
20     nodes " << endl ;
21     for ( i=0 ; i<=n ; i++ )
22         cin>> f[i] ;
23     cout<< "Enter a real value it which the interpolation is to be obtained: \n" ;
24     cin>> xp ;
25     //----- Processing Section -----//
26     // Computing zeroth divided differences, f[x_i ] = f_i
27
28     for ( i=0 ; i<=n ; i++ )
29         ddf[i][0] = f[i] ;
30     // Computing the divided differences of order 1 to n
31
32     for ( j=1 ; j<=n ; j++ )
33         for ( i=0 ; i<=n-j ; i++ )
34             ddf[i][j] = ( ddf[i+1][j-1] - ddf[i][j-1] ) / ( x[i+j] - x[i] ) ;
35
36     // Evaluating the interpolation polynomial at xp
37
38     pro = 1 ;    fxp = ddf[0][0] ;
39     for ( k=1 ; k<=n ; k++ )

```

```

40     {
41         pro = pro * ( xp - x[k-1] ) ;
42         fxp = fxp + pro * ddf[0][k] ;
43     }
44
45 //----- Output Section -----//
46
47     cout<< "The interpolate or extrapolate value of function at x = xp is " << endl ;
48     cout<< fxp ;
49
50     cin>> quit ;
51     return 0 ;
52
53 }

```

■ ■

Chapter Summary

- **Curve fitting** refers to the process of constructing a curve (a mathematical function) that reasonably fits the given discrete data points along a continuum. The obtained curve offers a simpler alternative to the original function (whose values at discrete points were given) that might be used to estimate the data values at points between the given points (and sometimes beyond the given data points, as well).
- **Regression** and **Interpolation** are the two basic approaches for curve fitting. Regression is the process of deriving a single curve that provides for the general trend of the data (and that curve is not required to pass through any of the data points). Interpolation is the process of fitting a curve (a single function or a piecewise function) that interpolates (passes through) each of the given data points.
- Suppose that the values of a function f at different points $x_0, x_1, x_2, \dots, x_n$ are given. The points x_i are referred to as **nodes** or **arguments** and the $n + 1$ ordered pairs $(x_i, f(x_i)), i = 0, 1, 2, \dots, n$, are referred to as **data points** of f . **Interpolation** (or, more precisely, *polynomial interpolation*) refers to the process of approximating the value of f at any intermediate point to the given data points.
- The interpolation process consists of determining the *unique* polynomial $P_n(x)$ of degree at most n that interpolates (passes through) the given data points, i.e.,

$$P_n(x_i) = f(x_i)$$

And then, the polynomial $P_n(x)$ serves as the formula to approximate the function values at intermediate points to the given data points and, thus, is referred to as **interpolating polynomial**. If the polynomial $P_n(x)$ is used approximate the function values at beyond the given data points, then the process is called **extrapolation**.

- **Newton's Divided Difference Interpolation:** For $n + 1$ arbitrarily spaced data points, $(x_0, f_0), (x_1, f_1), \dots, (x_n, f_n)$, of a function f , the Newton's Divided Difference interpolation formula for the interpolating polynomial $P_n(x)$ of degree at most n is given by

$$P_n(x) = f[x_0] + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1) + \dots + f[x_0, x_1, x_2, \dots, x_n](x - x_0)(x - x_1) \dots (x - x_{n-1})$$

or

$$P_n(x) = f[x_0] + \sum_{k=1}^n f[x_0, x_1, x_2, \dots, x_k](x - x_0)(x - x_1) \dots (x - x_{k-1})$$

Here the **k th divided difference** of the function f with respect to the nodes $x_i, x_{i+1}, \dots, x_{i+k}$ is denoted by $f[x_i, x_{i+1}, \dots, x_{i+k}]$ and is recursively defined by

$$f[x_i, x_{i+1}, \dots, x_{i+k}] = \frac{f[x_{i+1}, x_{i+2}, \dots, x_{i+k}] - f[x_i, x_{i+1}, \dots, x_{i+(k-1)}]}{x_{i+k} - x_i}$$

with $f[x_i] = f(x_i) = f_i$ as the zeroth divided difference.

- **Lagrange Interpolation:** For $n + 1$ arbitrarily spaced data points, $(x_0, f_0), (x_1, f_1), \dots, (x_n, f_n)$, of a function f , the Lagrange interpolation formula for the interpolating polynomial $P_n(x)$ of degree at most n is given by

$$\begin{aligned} P_n(x) &= L_0(x)f(x_0) + L_1(x)f(x_1) + \dots + L_n(x)f(x_n) \\ &= \sum_{k=0}^n L_k(x)f(x_k) \end{aligned}$$

Here $L_k(x)$ denotes the **k th Lagrange coefficient** (also called **cardinal polynomial**) and is defined by

$$L_k(x) = \frac{(x - x_0)(x - x_1) \dots (x - x_{k-1})(x - x_{k+1}) \dots (x - x_n)}{(x_k - x_0)(x_k - x_1) \dots (x_k - x_{k-1})(x_k - x_{k+1}) \dots (x_k - x_n)} = \prod_{\substack{j=0 \\ j \neq k}}^n \frac{x - x_j}{x_k - x_j}$$

and satisfies the **Kronecker delta** equation:

$$L_k(x) = \begin{cases} 1 & \text{for } x = x_k \\ 0 & \text{for all } x, \text{ except } x = x_k \end{cases}$$

- **First Theorem on Interpolation Error:** If $P_n(x)$ is the polynomial of degree at most n that interpolates a function f at $n + 1$ arbitrary nodes x_0, x_1, \dots, x_n in an interval $[a, b]$ and if $f \in C^{(n+1)}[a, b]$, then for each x in $[a, b]$, there exists an ξ in (a, b) for which

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

Here $E(x)$ is the truncation error of the polynomial interpolation.

- A Lagrange interpolation formula can be obtained from the relevant Newton's Divided Difference interpolation formula, after some rearrangements.
- Suppose that $n + 1$ data points, $(x_0, f_0), (x_1, f_1), \dots, (x_n, f_n)$, of a function f are given on the interval $[a, b]$ for consecutively arranged and equispaced nodes $x_0, x_1, x_2, \dots, x_n$, such that

$$a = x_0 < x_1 < x_2 < \dots < x_{n-1} < x_n = b$$

$$\text{with astep size of length } h = x_i - x_{i-1}, \text{ for } i = 1, 2, 3, \dots, n$$

$$\text{and } f(x_i) = f_i$$

The **Newton Forward-Difference Interpolation formula** for the interpolating polynomial $P_n(x)$ of degree at most n is given by

$$P_n(x) = f_0 + \alpha \Delta f_0 + \frac{\alpha(\alpha-1)}{2!} \Delta^2 f_0 + \dots + \frac{\alpha(\alpha-1)(\alpha-2) \dots (\alpha-(n-1))}{n!} \Delta^n f_0$$

where

$$\alpha = \frac{x - x_0}{h}$$

Here the **k th forward-difference** of f at x_i is denoted by $\Delta^k f_i$ and is recursively defined by

$$\Delta^k f_i = \Delta(\Delta^{k-1} f_i) = \Delta^{k-1} f_{i+1} - \Delta^{k-1} f_i \quad \text{for } k = 2, 3, \dots, n$$

with $\Delta f_i = f_{i+1} - f_i$

The **Newton Backward-Difference Interpolation formula** for the interpolating polynomial $P_n(x)$ of degree at most n is given by

$$P_n(x) = f_n + \beta \nabla f_n + \frac{\beta(\beta+1)}{2!} \nabla^2 f_n + \dots + \frac{\beta(\beta+1)(\beta+2) \dots (\beta+(n-1))}{n!} \nabla^n f_n$$

where

$$\beta = \frac{x - x_n}{h}$$

Here the **k th backward-difference** of f at x_i is denoted by $\nabla^k f_i$ and is recursively defined by

$$\nabla^k f_i = \nabla(\nabla^{k-1} f_i) = \nabla^{k-1} f_i - \nabla^{k-1} f_{i-1} \quad \text{for } k = 2, 3, \dots, n$$

with $\nabla f_i = f_i - f_{i-1}$

- There are central difference interpolation formulas also available in the literature, which are more suited for approximation of a function value around mid of the interval of interpolation. Following are the examples of some well-known central difference interpolation formulas:
 - Gauss Forward Difference Interpolation Formula
 - Gauss Backward Difference Interpolation Formula
 - Stirling's Central Difference Interpolation Formula
 - Bessel's Central Difference Interpolation Formula
 - Everett's Central Difference Interpolation Formula



Chapter Exercises

Exercise 01: Find the linear interpolating polynomial passing through the following set of pairs of the points.

(i) $\{(0.1, \sin(0.1)), (0.2, \sin(0.2))\}$

(ii) $\left\{\left(1.2, \frac{1}{(1.2)^2}\right), \left(1.4, \frac{1}{(1.4)^2}\right)\right\}$

(iii) $\{(1, 7), (2, 4)\}$

(iv) $\left\{(1, e^{-1}), \left(1.5, e^{-\frac{1}{1.5}}\right)\right\}$

Exercise 02: Construct the interpolating polynomial to approximate the following functions at $x = 0.25$. Use the arguments $x_0 = -0.3$, $x_1 = 0$, $x_2 = 0.4$.

(i) $f(x) = \ln(1 + x)$

(ii) $f(x) = e^{-x^2}$

(iii) $f(x) = \tan x^2$

(iv) $f(x) = \frac{1}{\sqrt{x^2-1}}$

Exercise 03: Use the Lagrange Interpolating Polynomial and the Newton's Divided Difference Interpolating polynomial of the appropriate degree to interpolate the following:

(i) Compute $f(1.5)$, given that $f(0.5) = 0.479$, $f(1.0) = 0.841$, $f(2.0) = 0.909$

(ii) Compute $f(3.6)$, given that $f(3.0) = 0.1506$, $f(4.0) = 0.3001$, $f(4.5) = 0.2663$, $f(4.7) = 0.2346$

(iii) Compute $f(2/3)$, given that,

$$f(1.1) = -0.071812, \quad f(1.3) = -0.024750, \quad f(1.7) = 0.334937, \quad f(2.0) = 1.101000$$

Exercise 04: Find the missing value in the following table using the Newton's Divided Difference Interpolating polynomial.

x	-1	1	2	3
$f(x)$	-21	15	?	3

Exercise 05: Find the missing value in the following table using Lagrange Interpolating Polynomial

x	-2	0	2	4	6
$f(x)$	33	5	9	?	113

Exercise 06: Find, for what values of x , y attained extreme values using the data given below

x	3	4	5	6	7	8
y	0.205	0.240	0.259	0.262	0.250	0.224

Exercise 07: Use Lagrange Interpolating Polynomial of the appropriate degree to complete the record of the export of a certain commodity during six years

Year: x	1981	1982	1983	1984	1985	1986
Export: y	43	84	93	?	105	157

Exercise 08: Use the Newton's Divided Difference Interpolating Polynomial to obtain an interpolation that passing through the following points

x	0	0.1	0.3	0.4	0.7	0.8
y	-1.5	-1.27	-0.98	-0.63	-0.22	0.25

Exercise 09: Find a bound for the error associated with linear polynomial interpolation for the following function. Use the arguments $x_0 = 0$, $x_1 = 0.4$.

- (i) $f(x) = \ln(1 + x)$
- (ii) $f(x) = e^{-x^2}$
- (iii) $f(x) = \tan x^2$
- (iv) $f(x) = \frac{1}{\sqrt{x^2-1}}$

Exercise 10: Find a bound for the error associated with quadratic polynomial interpolation for the following function. Use the arguments $x_0 = 0$, $x_1 = 0.1$, $x_2 = 0.4$.

- (i) $f(x) = \sin x + \cos x$
- (ii) $f(x) = x \ln x$
- (iii) $f(x) = x \sin x - x^3 + 2x - 1$
- (iv) $f(x) = \sqrt{x - x^2}$

Exercise 11: Find a bound for the error associated with cubic polynomial interpolation for the following function. Use the arguments $x_0 = 1$, $x_1 = 1.3$, $x_2 = 1.6$, $x_3 = 2.0$

- (i) $f(x) = \sin(e^{-x} - 1)$
- (ii) $f(x) = \ln x - x^4 + x^2 - 1$
- (iii) $f(x) = x^2 e^{-x^2}$
- (iv) $f(x) = \frac{1}{\sqrt{1+x}}$

Exercise 12: Construct the Newton's Forward and Backward Difference Interpolating polynomials passes through the points (0.2, 0.9980), (0.4, 0.9686), (0.6, 0.8443), and (0.8, 0.5358).

Exercise 13: Construct the Newton's Forward and Backward Difference Interpolating polynomials to approximate the following functions at $x = 1.2$ and 2.0 . Use the arguments $x_0 = 1.1$, $x_1 = 1.3$, $x_2 = 1.5$, $x_3 = 1.7$, $x_4 = 1.9$

(i) $f(x) = \ln(1 + x)$

(ii) $f(x) = e^{-x^2}$

(iii) $f(x) = \tan x^2$

(iv) $f(x) = \frac{1}{\sqrt{x^2-1}}$

Exercise 14: Some data of the speed (V) versus drag coefficient (C_d) of a cricket ball is given in the following table: Estimate C_d at $V = 150$ km/h.

V in km/h	C_d
0	0.5
80	0.48
120	0.39
160	0.32

Exercise 15: The mileages covered by a car per liter of fuel at different speeds are shown in the table below:

Speed in km/h	Mileage covered in km/l
60	14.2
75	16.1
90	14.8
105	13.7
120	11.5

Using interpolation, approximate the fuel efficiency of the car at the speed of 100 km/h.

Hint for the Solution: Use any interpolation formula, preferable the Newton's Backward Difference Interpolation formula.

Exercise 16: Some recorded data of number of deaths due to Novel Coronavirus (2019-nCoV) is given in the table below. Use interpolation to determine number of deaths on January 29 and 31, 2020.

Date	Number of Deaths
Jan. 24	16
Jan. 26	24
Jan. 28	26
Jan. 30	43
Feb. 1	45

Hint for the Solution: The given data spans over 9 days. The function values are given for $x = 1, 3, 5, 7, 9$. Find an interpolating polynomial and use it to calculate value at $x = 6$ and $x = 8$ for the desired solutions.

Exercise 17: The census data of Pakistan is given in the following table (source: Pakistan Bureau of Statistics):

Census Year	Population in thousands
1951	33740
1961	42880
1972	65309
1981	84254
1998	132352
2017	207774

Use interpolation to determine the population for the year 2010.

Hint for the Solution: The given data spans over 67 years. The function values are given for $x = 1, 11, 22, 31, 48, 67$. Find an interpolating polynomial and use it to calculate value at $x = 60$ for the desired solution.

Exercise 18: Suppose that a table lists the values of the tangent function for the angles ranging from 0° to 45° in increments of 5° . What is the largest error that we would introduce by performing linear interpolation between successive values in this table?



Numerical Integration

Corridor I: BASICS

Let's plan it

- 4.1 Introduction
- 4.2 The Trapezoidal Rule
- 4.3 The Simpson's 1/3 Rule
- 4.4 Generalized Closed Newton-Cotes Quadrature

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



Question 12: Tabulate Closed Newton-Cotes Integration formulas with relevant features, for both the basic and the composite forms, separately.

Suppose that n data points, (x_j, f_j) , where $f(x_j) = f_j$, of the integrand $f(x)$ are given on the interval $[a, b] = [x_0, x_1]$ for consecutively arranged and equispaced nodes x_j such that $h = (b - a)/n$. The Closed Newton Cotes quadrature formulas for the definite integral $= \int_{x_0}^{x_n} f(x) dx$ are tabulated below.

Numerical Integration Method	Formula	Required number of function values at equidistant points	Interpolating polynomial used for integral evaluation (to derive the formula)
Rectangular Rule	$I = h(f_0)$ (starting-point rule) or $I = h(f_1)$ (end-point rule) or $I = h(f^*)$ (mid-point rule) where $f^* = f\left(\frac{x_0+x_1}{2}\right)$	one	Polynomial of degree 0 (constant function)
Trapezoidal Rule	$I = \frac{h}{2} [f_0 + f_1]$	two	Polynomial of degree 1 (linear polynomial)
Simpson's 1/3 Rule	$I = \frac{h}{3} [f_0 + 4f_1 + f_2]$	three	Polynomial of degree 2 (quadratic polynomial)
Simpson's 3/8 Rule	$I = \frac{3h}{8} [f_0 + 3(f_1 + f_2) + f_3]$	four	Polynomial of degree 3 (cubic polynomial)
Boole's Rule (Milne's Rule)	$I = \frac{2h}{45} [7f_0 + 32f_1 + 12f_2 + 32f_3 + 7f_4]$	five	Polynomial of degree 4
Six-Point Rule	$I = \frac{5h}{288} [19f_0 + 75f_1 + 50f_2 + 50f_3 + 75f_4 + 19f_5]$	six	Polynomial of degree 5
Weddle's Rule	$I = \frac{h}{140} [41f_0 + 216f_1 + 27f_2 + 272f_3 + 27f_4 + 216f_5 + 41f_6]$	seven	Polynomial of degree 6

Numerical Integration Method	Formula (for $n + 1$ data points, $(x_j, f_j), j = 0, 1, 2, \dots, n$, and n subintervals of equal length $h = (x_n - x_0)/n$)	Possible values of n (K represents the number of multiple applications of the formula)	Interpolating polynomial used for integral evaluation (to derive the formula)
Composite Rectangular Rule	$I = h[f_0 + f_1 + f_2 + \dots + f_{n-1}]$ (starting-point rule) or $I = h[f_1 + f_2 + \dots + f_{n-1} + f_n]$ (end-point rule) or $I = h[f_1^* + f_2^* + \dots + f_{n-1}^* + f_n^*]$ (mid-point rule) where $f_j^* = f\left(\frac{x_{j-1} + x_j}{2}\right)$, for $j = 1, 2, 3, \dots, n$	$n = 1, 2, 3, \dots$ (i.e., $n = K$ could be any positive integer)	Piecewise polynomial of degree 0 (piecewise-constant function)
Composite Trapezoidal Rule	$I = \frac{h}{2}[f_0 + 2(f_1 + f_2 + \dots + f_{n-1}) + f_n]$	$n = 1, 2, 3, \dots$ (i.e., $n = K$ could be any positive integer)	Piecewise polynomial of degree 1 (piecewise-linear)
Composite Simpson's 1/3 Rule	$I = \frac{h}{3}[f_0 + 4(f_1 + f_3 + \dots + f_{n-1}) + 2(f_2 + f_4 + \dots + f_{n-2}) + f_n]$	$n = 2, 4, 6, \dots$ (i.e., $n = 2K$, where $K = 1, 2, 3, \dots$)	Piecewise polynomial of degree 2 (piecewise-quadratic)
Composite Simpson's 3/8 Rule	$I = \frac{3h}{8}[f_0 + 3(f_1 + f_2) + 2(f_3) + 3(f_4 + f_5) + 2(f_6) + \dots + 3(f_{n-2} + f_{n-1}) + f_n]$	$n = 3, 6, 9, \dots$ (i.e., $n = 3K$, where $K = 1, 2, 3, \dots$)	Piecewise polynomial of degree 3 (piecewise-cubic)
Composite Boole's Rule (Composite Milne's Rule)	$I = \frac{2h}{45}[7f_0 + 32(f_1 + f_5 + f_9 + \dots + f_{n-3}) + 12(f_2 + f_6 + f_{10} + \dots + f_{n-2}) + 32(f_3 + f_7 + f_{11} + \dots + f_{n-1}) + 14(f_4 + f_8 + f_{12} + \dots + f_{n-4}) + 7f_n]$	$n = 4, 8, 12, \dots$ (i.e., $n = 4K$, where $K = 1, 2, 3, \dots$)	Piecewise polynomial of degree 4
Composite Six-Point Rule	$I = \frac{5h}{288}[19f_0 + 75(f_1 + f_6 + f_{11} + \dots + f_{n-4}) + 50(f_2 + f_7 + f_{12} + \dots + f_{n-3}) + 50(f_3 + f_8 + f_{13} + \dots + f_{n-2}) + 75(f_4 + f_9 + f_{14} + \dots + f_{n-1}) + 38(f_5 + f_{10} + f_{15} + \dots + f_{n-5}) + 19f_n]$	$n = 5, 10, 15, \dots$ (i.e., $n = 5K$, where $K = 1, 2, 3, \dots$)	Piecewise polynomial of degree 5
Composite Weddle's Rule	$I = \frac{h}{140}[41f_0 + 216(f_1 + f_7 + f_{13} + \dots + f_{n-5}) + 27(f_2 + f_8 + f_{14} + \dots + f_{n-4}) + 272(f_3 + f_9 + f_{15} + \dots + f_{n-3}) + 27(f_4 + f_{10} + f_{16} + \dots + f_{n-2}) + 216(f_5 + f_{11} + f_{17} + \dots + f_{n-1}) + 82(f_6 + f_{12} + f_{18} + \dots + f_{n-6}) + 41f_n]$	$n = 6, 12, 18, \dots$ (i.e., $n = 6K$, where $K = 1, 2, 3, \dots$)	Piecewise polynomial of degree 6



Corridor II: ANALYSIS

Let's think deep

- 4.5 Truncation Error of the Trapezoidal Rule
- 4.6 Truncation Error of the Simpson's 1/3 Rule
- 4.7 Further Discussions
- 4.8 The Gaussian Quadrature

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



Corridor III: PROGRAMMING ARCADE

Let's do it

- 4.9 Algorithms and Implementations
 - The Composite Trapezoidal Rule
 - The Composite Simpson's 1/3 Rule
 - The Composite Simpson's 3/8 Rule

To cross-check the results/output of the computer programs you would execute, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



4.9 Algorithms and Implementations

Question 27: Write down the algorithm (pseudo-code) of the Composite Trapezoidal rule for numerical integration of definite integrals.

Algorithm: To approximate the definite integral $I = \int_a^b f(x) dx$ using the formula:

$$I = \frac{h}{2} [f_0 + 2(f_1 + f_2 + f_3 + \cdots + f_{n-1}) + f_n]$$

INPUTS: $\{a$ and b : two real values as the endpoints of the interval of integration
 n : a positive integer as the number of subintervals

OUTPUT: I : a real number as an approximation to the integral

Step 1 Receive the inputs as stated above

Step 2 Set real number $x_0 = a$
 Set real number $x_n = b$
 Set real number $h = (x_n - x_0)/n$
 Set real number f_{x_0} as the value $f(a)$
 Set real number f_{x_n} as the value $f(b)$

Step 3 Set $I = f_{x_0} + f_{x_n}$

Step 4 Set real number $xc = x_0$
 Set real number $sum = 0$
 for $j = 1, 2, \dots, n - 1$
 Set $xc = xc + h$
 Set f_{xc} as the value $f(xc)$
 Set $sum = sum + f_{xc}$ (Forming $f_1 + f_2 + \cdots + f_{n-1}$)
 end for

Step 5 Set $I = (h/2) \times (I + 2 \times sum)$

Step 6 Print the output: I

STOP.

Problem 16: Write a C++ program to evaluate the integral of $f(x) = \sqrt{x^2 + 1}$ over $[0, 2]$ with 12 subintervals using the Composite Trapezoidal rule.

```
1  #include<iostream>
2  #include<cmath>
3  using namespace std ;
4
5  int main ( )
6  {
7
8      int n , j ;
9      double h , sum , I , x0 , xn , xc , fx0 , fxn , fxc ;
10
11     //----- Input Section -----//
12     cout<< "The Composite Trapezoidal Rule." ;
13     cout<< "\nEnter the lower limit of the integral: " ;
14     cin>> x0 ;
15     cout<< "\nEnter the upper limit of the integral: " ;
16     cin>> xn ;
17     cout<< "\nEnter the number of subintervals n: " ;
18     cin>> n ;
19
20     //----- Processing Section -----//
21
22     h = (xn - x0)/n ;           // finding step size
23     fx0 = sqrt(x0*x0 + 1) ;
24     fxn = sqrt(xn*xn + 1) ;
25     I = fx0 + fxn ;
26     xc = x0 ;
27     sum = 0.0 ;
28
29     for ( j=1 ; j<=n-1 ; j++ )
30     {
31         xc = xc + h ;
32         fxc = sqrt(xc*xc + 1) ;
33         sum = sum + fxc ;
34     }
35
36     I = (h / 2.0) * (I + 2.0 * sum) ;
37
38     //----- Output Section -----//
39     cout<< "The approximate integral = " << I << endl ;
40
41 }
```


Problem 18: Write a C++ program to evaluate the integral of $f(x) = \sqrt{x^2 + 1}$ over $[0, 2]$ with 12 subintervals using the Composite Trapezoidal rule. Use #define directive for evaluating $f(x)$ at the different nodes (i.e., for finding the values of f at the different nodes).

```
1  #include<iostream>
2  #include<cmath>
3  using namespace std ;
4
5  # define fval(x) sqrt(x*x + 1)           // evaluating f(x)
6
7  int main ( )
8  {
9      int n , j ;
10     double h , sum , I , x0 , xn , xc , fx0 , fxn , fxc ;
11
12     //----- Input Section -----//
13     cout<< "The Composite Trapezoidal Rule." ;
14     cout<< "\nEnter the lower limit of the integral: " ;
15     cin>> x0 ;
16     cout<< "\nEnter the upper limit of the integral: " ;
17     cin>> xn ;
18     cout<< "\nEnter the number of subintervals n: " ;
19     cin>> n ;
20
21     //----- Processing Section -----//
22
23     h = (xn - x0)/n ;           // finding step size
24     fx0 = fval(x0) ;
25     fxn = fval(xn) ;
26     I = fx0 + fxn ;
27     xc = x0 ;
28     sum = 0.0 ;
29
30     for ( j=1 ; j<=n-1 ; j++ )
31     {
32         xc = xc + h ;
33         fxc = fval(xc) ;
34         sum = sum + fxc ;
35     }
36
37     I = (h / 2.0) * (I + 2.0 * sum) ;
```

```

38
39 //----- Output Section -----//
40 cout<< "The approximate integral = " << I << endl ;
41 }

```

Question 28: Write down the algorithm (pseudo-code) of the Composite Simpson's 1/3 rule for numerical integration of definite integrals.

Algorithm: To approximate the definite integral $I = \int_a^b f(x) dx$ using the formula:

$$I = \frac{h}{3} [f_0 + 4(f_1 + f_3 + \dots + f_{n-1}) + 2(f_2 + f_4 + \dots + f_{n-2}) + f_n]$$

INPUTS: {**a** and **b**: two real values as the endpoints of the interval of integration
n: a positive even integer as the number of subintervals

OUTPUT: **I**: a real number as an approximation to the integral

Step 1 Receive the inputs as stated above

Step 2 Set real number **x0** = **a**
Set real number **xn** = **b**
Set real number **h** = (**xn** - **x0**)/**n**
Set real number **fx0** as the value $f(a)$
Set real number **fxn** as the value $f(b)$

Step 3 Set $I = fx0 + fxn$

Step 4 Set real number **xc** = **x0**
Set real number **sum1** = **0**
Set real number **sum2** = **0**

for $j = 1, 2, \dots, n - 1$

Set **xc** = **xc** + **h**

Set **fxc** as the value $f(xc)$

if j is odd

Set **sum1** = **sum1** + **fxc** (Forming $f_1 + f_3 + \dots + f_{n-1}$)

else

Set **sum2** = **sum2** + **fxc** (Forming $f_2 + f_4 + \dots + f_{n-2}$)

end for

Step 5 Set $I = (h/3) \times (I + 4 \times sum1 + 2 \times sum2)$

Step 6 Print the output: **I**; **STOP**.

Problem 20: Write a C++ program to evaluate the integral of $f(x) = \sqrt{x^2 + 1}$ over $[0, 2]$ with 12 subintervals using the Composite Simpson's 1/3 rule.

```
1  #include<iostream>
2  #include<cmath>
3  using namespace std ;
4
5  int main ( )
6  {
7      int n , j ;
8      double h , sum1 , sum2 , I , x0 , xn , xc , fx0 , fxn , fxc ;
9
10     //----- Input Section -----//
11
12     cout<< "The Composite Simpson's 1/3 Rule." ;
13     cout<< "\nEnter the lower limit of the integral: " ;
14     cin>> x0 ;
15     cout<< "\nEnter the upper limit of the integral: " ;
16     cin>> xn ;
17     cout<< "\nEnter the number of subintervals n: " ;
18     cin>> n ;
19
20     //----- Processing Section -----//
21
22     h = (xn - x0)/n ;           // finding step size
23     fx0 = sqrt(x0*x0 + 1) ;
24     fxn = sqrt(xn*xn + 1) ;
25     I = fx0 + fxn ;
26     xc = x0 ;
27     sum1 = 0.0 ; sum2 = 0.0 ;
28
29     for ( j=1 ; j<=n-1 ; j++ )
30     {
31         xc = xc + h ;
32         fxc = sqrt(xc*xc + 1) ;
33         if ( j%2 != 0 )
34             sum1 = sum1 + fxc ;
35         else
36             sum2 = sum2 + fxc ;
37     }
38     I = (h / 3.0) * (I + 4 * sum1 + 2 * sum2) ;
```

```

39
40 //----- Output Section -----//
41 cout<< "The approximate integral = " << I << endl ;
42
43 }

```

Question 29: Write down the algorithm (pseudo-code) of the Composite Simpson's 3/8 rule for numerical integration of definite integrals.

Algorithm: To approximate the definite integral $I = \int_a^b f(x) dx$ using the formula:

$$I = \frac{3h}{8} [f_0 + 3(f_1 + f_2 + f_4 + f_5 + \dots + f_{n-2} + f_{n-1}) + 2(f_3 + f_6 + \dots + f_{n-3}) + f_n]$$

INPUTS: $\{a$ and b : two real values as the endpoints of the interval of integration
 n : a positive integer (multiple of 3) as the number of subintervals

OUTPUT: I : a real number as an approximation to the integral

Step 1 Receive the inputs as stated above

Step 2 Set real number $x_0 = a$
Set real number $x_n = b$
Set real number $h = (x_n - x_0)/n$
Set real number f_{x_0} as the value $f(a)$
Set real number f_{x_n} as the value $f(b)$

Step 3 Set $I = f_{x_0} + f_{x_n}$

Step 4 Set real number $x_c = x_0$
Set real number $sum1 = 0$
Set real number $sum2 = 0$

for $j = 1, 2, \dots, n - 1$

Set $x_c = x_c + h$

Set f_{x_c} as the value $f(x_c)$

if j is divisible by 3

Set $sum2 = sum2 + f_{x_c}$ (Forming $f_3 + f_6 + \dots + f_{n-3}$)

else

Set $sum1 = sum1 + f_{x_c}$ (Forming $f_1 + f_2 + f_4 + f_5$
 $+ \dots + f_{n-2} + f_{n-1}$)

end for

Step 5 Set $I = (3 \times h/8) \times (I + 3 \times sum1 + 2 \times sum2)$

Step 6 Print the output: I ; **STOP.**

Problem 22: Write a C++ program to evaluate the integral of $f(x) = \sqrt{x^2 + 1}$ over $[0, 2]$ with 12 subintervals using the Composite Simpson's 3/8 rule.

```
1  #include<iostream>
2  #include<cmath>
3  using namespace std ;
4
5  int main ( )
6  {
7      int n , j ;
8      double h , sum1 , sum2 , I , x0 , xn , xc , fx0 , fxn , fxc ;
9
10     //----- Input Section -----//
11
12     cout<< "The Composite Simpson's 3/8 Rule." ;
13     cout<< "\nEnter the lower limit of the integral: " ;
14     cin>> x0 ;
15     cout<< "\nEnter the upper limit of the integral: " ;
16     cin>> xn ;
17     cout<< "\nEnter the number of subintervals n: " ;
18     cin>> n ;
19
20     //----- Processing Section -----//
21
22     h = (xn - x0)/n ;           // finding step size
23     fx0 = sqrt(x0*x0 + 1) ;
24     fxn = sqrt(xn*xn + 1) ;
25     I = fx0 + fxn ;
26     xc = x0 ;
27     sum1 = 0.0 ; sum2 = 0.0 ;
28
29     for ( j=1 ; j<=n-1 ; j++ )
30     {
31         xc = xc + h ;
32         fxc = sqrt(xc*xc + 1) ;
33         if ( j%3 == 0 )
34             sum2 = sum2 + fxc ;
35         else
36             sum1 = sum1 + fxc ;
37     }
38     I = (3.0 * h / 8.0) * (I + 2 * sum2 + 3 * sum1) ;
```

```
39
40     //----- Output Section -----//
41     cout<< "The approximate integral = " << I << endl ;
42
43 }
```

Remark: Likewise the programs in the solutions of Problem 19, the programmer can modify the programs in the solutions of Problems 20 and 22 to evaluate the function values at the desired nodes through the use of user-defined function and inline function (using #define).



Chapter Summary

- **Numerical integration** or **quadrature** refers to the process of numerically approximating the value of the integral $I = \int_a^b f(x) dx$, by using the values of f at a finite number of sample points. The limits of integration could be finite, semi-finite, or infinite.
- The integral is approximated by a numerical **integration rule** or **quadrature formula**, Q_f , which is a linear combination of certain function values:

$$I = \int_a^b f(x) dx \cong Q_f = \sum_{j=0}^n \omega_j \cdot f(x_j)$$

Here x_i are the ordered points, called the **quadrature nodes** (or simply **nodes**), taken usually within the limits of integration at which the function values $f(x_j)$ are known and ω_j are called the **weights** of the quadrature formula.

- The quadrature formula satisfies the property that

$$I = \int_a^b f(x) dx = Q_f + E_f,$$

where E_f is the **truncation error** (also called the **error term**) associated with the quadrature formula.

- The Newton-Cotes integration formulas are based on the approach that $n + 1$ number of equispaced and ordered nodes are chosen within the limits of integration and the integrand function is replaced by an interpolating polynomial of degree at most n by using the nodes, and then the analytic integration of the polynomial is performed to obtain the formula. A Composite Newton-Cotes integration formula is obtained by applying the relevant Newton-Cotes formula in each of the different consecutive segments of the interval of integration and then summing the integrals over all the segments.
- The examples of Newton-Cotes integration formulas include Trapezoidal rule, Simpson's 1/3 rule, Simpson's 3/8 rule, Boole's rule, Six-Point rule, and Weddle's rule.
- The Trapezoidal rule to numerically integrate the function f over the interval $[a, b]$ is

$$I = \int_a^b f(x) dx \cong \frac{(b-a)}{2} [f(a) + f(b)]$$

- The Composite Trapezoidal rule to integrate a function f over the interval $[a, b]$ is given by,

$$I = \int_a^b f(x) dx \cong \frac{h}{2} [f_0 + 2(f_1 + f_2 + \cdots + f_{n-1}) + f_n]$$

$$\text{where } h = \frac{b-a}{n} = \frac{x_n - x_0}{n}, \quad f_j = f(x_j) \quad \text{and} \quad x_j = x_0 + jh \quad \text{for } j = 0, 1, 2, \dots, n$$

- The Simpson's 1/3 rule to integrate a function f over the interval $[a, b]$ is given by,

$$I = \int_a^b f(x) dx \cong \frac{h}{3} [f_0 + 4f_1 + f_2]$$

$$\text{where } h = \frac{b-a}{2} = \frac{x_2 - x_0}{2}, \quad f_j = f(x_j) \quad \text{and} \quad x_j = x_0 + jh \quad \text{for } j = 0, 1, 2$$

- The Composite Simpson's 1/3 rule to integrate a function f over the interval $[a, b]$ is given by,

$$I = \int_a^b f(x) dx \cong \frac{h}{3} [f_0 + 4(f_1 + f_3 + \cdots + f_{n-1}) + 2(f_2 + f_4 + \cdots + f_{n-2}) + f_n]$$

$$\text{where } h = \frac{b-a}{n} = \frac{x_n - x_0}{n}, \quad f_j = f(x_j) \quad \text{and} \quad x_j = x_0 + jh \quad \text{for } j = 0, 1, 2, \dots, n$$

- A comprehensive summary of the Newton-Cotes formulas and the Composite Newton-Cotes formulas can be found under Question 12 (page 252).

- The error term E_T of order $\mathcal{O}(h^3)$ associated with the Trapezoidal rule in approximating $I = \int_a^b f(x) dx$ is given by,

$$E_T = -\frac{1}{12} h^3 f''(\xi),$$

for some appropriate point ξ in (a, b) and $h = b - a$.

- The error term E_{CT} of order $\mathcal{O}(h^2)$ associated with the Composite Trapezoidal rule in approximating $I = \int_a^b f(x) dx$ is given by,

$$E_{CT} = -\frac{b-a}{12} h^2 f''(\eta),$$

for some appropriate point η in (a, b) and $h = (b-a)/n$, where n is the number of subintervals of $[a, b]$.

- The error term E_S of order $\mathcal{O}(h^5)$ associated with the Simpson's 1/3 rule in approximating $I = \int_a^b f(x) dx$ is given by,

$$E_S = -\frac{1}{90} h^5 f^{(4)}(\xi),$$

for some appropriate point ξ in (a, b) and $h = (b-a)/2$.

- The error term E_{CS} of order $\mathcal{O}(h^4)$ associated with the Composite Simpson's 1/3 rule in approximating $I = \int_a^b f(x) dx$ is given by,

$$E_{CS} = -\frac{b-a}{180} h^4 f^{(4)}(\eta)$$

for some appropriate point η in (a, b) and $h = (b-a)/n$, where n is the number of subintervals of $[a, b]$.

- Suppose I_h denotes the approximate integral using a quadrature formula with step size h , and E_h denotes the associated error. Then, the exact integral $= I_h + E_h$

Similarly, suppose $I_{h/2}$ denotes the approximate integral using the same quadrature formula with a step size $h/2$, and $E_{h/2}$ denotes the associated error. Then, the exact integral $= I_{h/2} + E_{h/2}$

According to the interval halving method, for a Newton-Cotes integration formula with an error of order $\mathcal{O}(h^N)$ an estimate of the error $E_{h/2}$ is given by,

$$E_{h/2} \cong \frac{1}{2^N - 1}(I_{h/2} - I_h)$$

This leads to a better approximation of the integral as below:

$$I \cong I_{h/2} + \frac{1}{2^N - 1}(I_{h/2} - I_h)$$

This corresponds to a special process called **Richardson Extrapolation**, in which two estimates of the solution are used to obtain a third approximation, which is a more accurate one. This approach for numerical integration forms an initial stage of a relatively broader way of numerical integration, called **Romberg Integration**. Recall that for the Composite Trapezoidal rule $N = 2$, and for the Composite Simpson's 1/3 rule $N = 4$.

- There could be several approaches for improving the estimates of the integrals:
 - Using smaller step size (or larger number of subintervals)
 - Using higher-order formula (e.g., using the Simpson's rule instead of the Trapezoidal rule)
 - Using Richardson's extrapolation (i.e., using two less accurate estimates to obtain a more accurate estimate).
- The **degree of precision**, also referred to as the *order of accuracy*, of a quadrature formula is p if and only if the associated truncation error is zero for all polynomials of degree less than or equal to p , and the error is not zero for some polynomial of degree greater than p . Note that the Trapezoidal rule is based on the interpolating polynomial of degree 1 (linear polynomial). Therefore, it produces the exact result while integrating a polynomial of degree 1. Hence it has the degree of precision as 1. The Simpson's 1/3 rule might be expected to have a degree of precision as 2 because it is based on interpolating polynomial of degree 2 (quadratic polynomial). However, it produces the exact result while integrating a polynomial of degree 2, as well as degree 3. Hence, it has the degree of precision as 3. This fact is also evident while deriving the error term for the Simpson's 1/3 rule. This property, together with certain other reasons, makes the Composite Simpson's 1/3 rule often the best choice among the Newton-Cotes integration formulas.
- A concise description of the error terms associated with the Newton-Cotes formulas and relevant degrees of precision can be found under Question 23 (page 276).
- The Gaussian Quadrature is an advanced numerical integration technique in which the quadrature nodes are selected in the interval of integration using the roots of some special polynomial to obtain an optimal approximation of the integral.



Chapter Exercises

Exercise 01: Approximate the integral $\int_a^b f(x)dx$ for the following functions over the interval $[0, 1]$ using the Trapezoidal, Simpson's 1/3 and Simpson's 3/8 rules.

$$(i) \quad f(x) = x^2 + x - 1 \qquad (ii) \quad f(x) = \ln(1 + x) \qquad (iii) \quad f(x) = \frac{1}{1 + x^2}$$

$$(iv) \quad f(x) = \cos\left(\frac{x}{\pi}\right) \qquad (v) \quad f(x) = \frac{1}{\sqrt{x^2 + 4}}$$

Exercise 02: Approximate the integral $\int_a^b f(x)dx$ for the following functions over the interval $[0, 1]$ using the Composite Trapezoidal, Simpson's 1/3, and Simpson's 3/8 rules with $h = 0.1$.

$$(i) \quad f(x) = x^2 + x - 1 \qquad (ii) \quad f(x) = \ln(1 + x) \qquad (iii) \quad f(x) = \frac{1}{1 + x^2}$$

$$(iv) \quad f(x) = \cos\left(\frac{x}{\pi}\right) \qquad (v) \quad f(x) = \frac{1}{\sqrt{x^2 + 4}}$$

Exercise 03: Approximate the integral

$$\int_4^{16} \sin\left(\frac{\pi\sqrt{x}}{4}\right) dx$$

using the Composite Trapezoidal rule with $h = 1$ and five-digit rounding arithmetic.

Exercise 04: Find an approximate value of the integral $\int_0^2 (2 + \sin(2\sqrt{x}))dx$ using the Composite Trapezoidal rule for $n = 10$ and five-digit rounding arithmetic.

Exercise 05: Approximate the arc length of the following functions over the interval $[0, \pi]$

$$(i) \quad f(x) = \sin^2 x \qquad (ii) \quad f(x) = \ln\left(\frac{4 + x}{\pi}\right)$$

using the Composite Simpson's 1/3 rule for $h = \frac{\pi}{6}$ and four-digit rounding arithmetic.

Exercise 06: Find the approximate value of the integral $\int_3^8 (f(x))^2 dx$ using the Composite Simpson's 1/3 rule, given that

x_j	3	4	5	6	7	8	9
$f(x_j)$	0.205	0.240	0.259	0.262	0.250	0.224	0.220

Exercise 07: Approximate the area of a surface of revolution of the following curves:

$$(i) \quad x = 4y, \qquad (ii) \quad x = \tan y$$

about the y – axis from $0 \leq y \leq 1$ using the Composite Simpson's $3/8$ rule for $n = 10$ and four-digit rounding arithmetic.

Exercise 08: Find the approximate value of the integral

$$f(x) = \int_0^3 \frac{x}{x^2 + 3} dx$$

using the Composite Boole's rule with step size $h = 0.25$ and five-digit rounding arithmetic.

Exercise 09: Find the approximate value of the integral

$$f(x) = \int_2^5 \ln(x - 1) dx$$

using the Composite Six-Point rule with step size $h = 0.3$ and five-digit rounding arithmetic.

Exercise 10: Find the approximate value of the integral

$$f(x) = \int_1^4 \sinh(x^2) dx$$

using the Composite Weddle's rule with step size $h = 0.25$ and five-digit rounding arithmetic.

Exercise 11: Suppose that $f(0) = 1$, $f(0.5) = 2.5$, $f(1) = 2$ and $f(0.25) = f(0.75) = \alpha$. Find α if the Composite Trapezoidal rule with $n = 4$ gives the value 1.75 for $\int_0^1 f(x) dx$.

Exercise 12: Suppose that $f(4) = 0.240$, $f(6) = 0.262$, $f(8) = 0.224$, $f(3) = f(5) = f(7) = \alpha$, and $f(9) = 0.220$. Find α if the Composite Simpson's $1/3$ Rule gives the value 1.473 for

$$I = \int_3^9 f(x) dx$$

Exercise 13: Suppose that $f(0.2) = 1.56$, $f(0.4) = 2.00$, $f(0.6) = 3.01$, $f(0.1) = f(0.3) = f(0.5) = \alpha$, and $f(0.7) = 3.32$. Find α if the Composite Simpson's $3/8$ rule gives the value 1.30312 for

$$I = \int_{0.1}^{0.7} f(x) dx$$

Exercise 14: To approximate the integral of $f(x)$ over the interval $[0, 1]$ with an absolute error less than $\frac{1}{2} \times 10^{-4}$, how many subintervals are needed, in case of (a) the Composite Trapezoidal rule, (b) the Composite Simpson's $1/3$ rule, and (c) the Composite Simpson's $3/8$ rule? Given that,

$$(i) \quad f(x) = x^2 + x - 1 \qquad (ii) \quad f(x) = \ln(1 + x) \qquad (iii) \quad f(x) = \frac{1}{1 + x^2}$$

$$(iv) \quad f(x) = \cos\left(\frac{x}{\pi}\right) \qquad (v) \quad f(x) = \frac{1}{\sqrt{x^2 + 4}}$$

Exercise 15: Suppose we wish to evaluate the integral

$$f(x) = \int_0^\pi \sin(\sqrt{x}) dx$$

numerically, with an error of magnitude less than 10^{-5} . How many subintervals are needed if we wish to use the Composite Trapezoidal and Composite Simpson 1/3 rules?

Exercise 16: Find the number of subintervals n or step length h so that the error E_{TC} for the Composite Trapezoidal rule and error E_{SC} for the Composite Simpson's 1/3 rule is less than 5×10^{-4} for numerically integrating the Legendre polynomial,

$$P_4(x) = x^4 - \frac{6}{7}x^2 + \frac{3}{35}$$

over the interval $[-1, 1]$.

Exercise 17: Obtain an upper bound on the absolute error when the *Chebyshev* polynomial of degree four,

$$T_4(x) = 8x^4 - 8x^2 + 1$$

is integrated over the interval $[-1, 1]$ by means of the Composite Simpson's 3/8 rule.

Exercise 18: Obtain an upper bound on the absolute error when the *Laquerre* polynomial of degree four

$$L_4(x) = x^4 - 16x^3 + 72x^2 - 96x + 24$$

is integrated over the interval $[-1, 1]$, by means of the Composite Simpson's 3/8 rule.

Exercise 19: A car travels the loop of the racing track in 65 seconds. The speed of the car in meter/second is recorded after every 5 seconds as shown in the following table:

Time	0	5	10	15	20	25	30	35	40	45	50	55	60	65
Speed	0	40	62	70	72	65	71	79	75	72	68	63	75	82

Estimate the length of the loop of the racing track?

Hint for the Solution:

Clearly, the speed say S is shown to be a function of time, say t , and its values $S(t)$ for different time instants t are given. Obtain the estimate of the integral distance $= \int_0^{65} S(t) dt$ using any appropriate numerical integration rule with the data given in the Table.

Exercise 20: The prime number theorem states that the number of primes in an interval $a \leq x \leq b$ is approximately

$$\int_a^b \frac{1}{\ln x} dx$$

Estimate the number of primes existing in $[50,150]$.

Hint for the Solution: Numerically evaluate the given integral for $a = 50$ and $b = 150$ using different values of $f(x) = \frac{1}{\ln x}$ at equispaced nodes in $[50,150]$, separated by step length $h = 10$ or 20 .

Exercise 21: The depths D (in meters) of a 80 meters wide river at different horizontal distances s from the bank is given in the following table.

s	0	10	20	30	40	50	60	70	80
D	0	3.5	6	12	10	15	9	5	0

Estimate the area of the cross-section of the river.

Hint for the Solution: Clearly, D is shown to be a function of s and its values $D(s)$ for different points s are given. Obtain the estimate of the integral, $area = \int_0^{80} D(s) ds$ using any appropriate numerical integration rule with the data given in the Table.

Exercise 22: A rectangular swimming pool is 35 feet wide and 60 feet long. At different positions P in feet along the length of the pool, the depths D in feet are shown in the following Table. Estimate the volume of the pool using numerical integration.

P	0	6	12	18	24	30	36	42	48	54	60
D	3	3.5	4	4.5	5	5.5	6	6.5	7	7.5	8

Hint for the Solution:

Clearly, D is shown to be a function of P and its values $D(P)$ for different points P are given. Obtain the estimate of the integral $w = \int_0^{60} D(P) dP$ using any appropriate numerical integration rule with the data given in the Table. Note that w is the estimated area of one side-wall of the pool along the length. Multiplying it with the width of 35 feet will give the volume of the pool.

Exercise 23: We know that

$$\int_0^1 \frac{1}{1+x^2} dx = \tan^{-1} x \Big|_0^1 = \tan^{-1} 1 = \frac{\pi}{4}$$

This means that the value of π can be obtained evaluating the above integral and then multiplying the answer by 4. Suppose that we want to approximate π to four decimal places. This means absolute error must

be less than 5.0×10^{-5} . This means that the error in approximating the integral must be less than $\frac{1}{4} \times (5.0 \times 10^{-5}) = 1.25 \times 10^{-5}$. Use the Composite Simpson's 1/3 rule to approximate the value of π . For this, first determine that what should be the minimum number of subintervals that would keep the error less than the tolerance.

Exercise 24: The number of subintervals required to apply the Composite Simpson's 1/3 rule should be

- (A) Multiple of 1 (B) Multiple of 2
(C) Multiple of 3 (D) unconditionally many (E) None of above

Exercise 25: The Simpson's 1/3 rule is based on the integration of interpolating polynomial of degree 2. The Simpson's 1/3 rule can accurately integrate the polynomials of degree

- (A) up to 1 (B) up to 2
(C) up to 3 (D) up to any (E) None of above

Exercise 26: The Gaussian quadrature is different from the Newton's Cotes Integration in regards to

- (A) selection of polynomial degree (B) selection of quadrature nodes
(C) problem dependence (D) None of above



Numerical Differentiation

Corridor I: BASICS

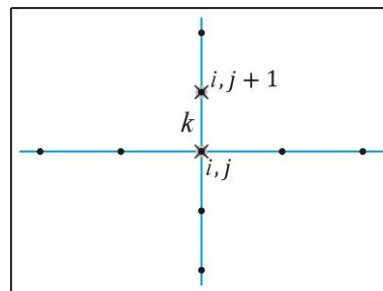
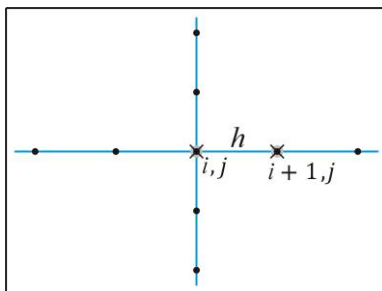
Let's plan it

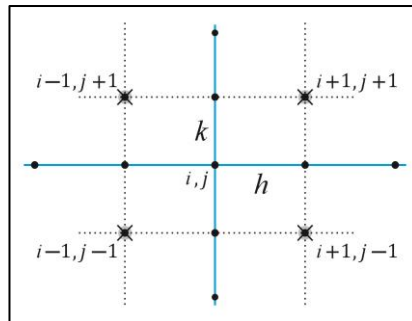
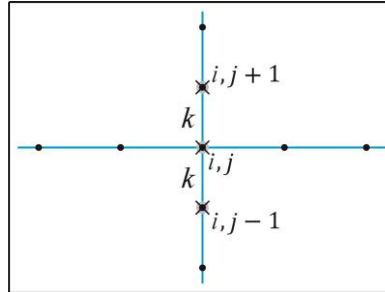
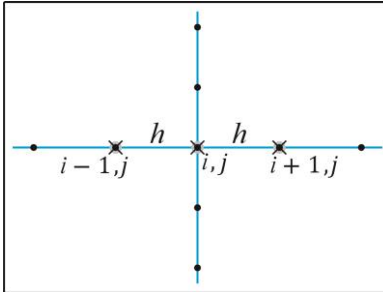
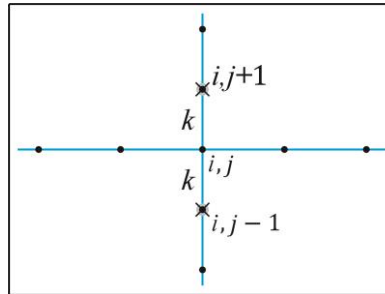
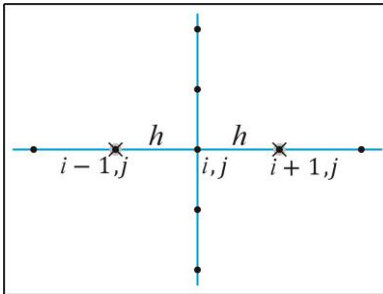
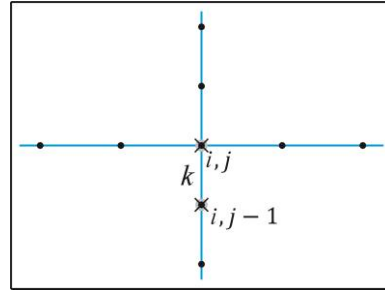
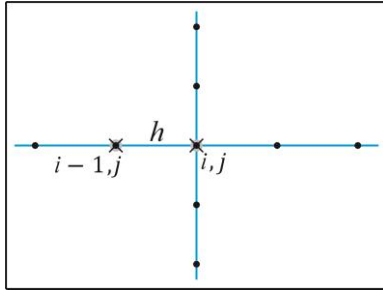
- 5.1 Introduction
- 5.2 Finite Difference Approximations of Derivatives using the Taylor Series
 - 5.2.1 First Order Derivatives
 - 5.2.2 Second Order Derivatives
- 5.3 Listing of the Derivative Formulas

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)

■■■





Direct Linear Solvers

Corridor I: BASICS

Let's plan it

- 6.1 Introduction to Linear Systems
- 6.2 Solving Linear Systems using the Gaussian Elimination Method
- 6.3 Pivoting Strategies
 - Partial Pivoting
 - Scaled Partial Pivoting
 - Complete Pivoting
- 6.4 The Gauss-Jordan Method
- 6.5 Solving Linear Systems using the LU Factorization Method
 - 6.5.1 The Doolittle's Method
 - 6.5.2 The Crout's Method
 - 6.5.3 The Cholesky's Method

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



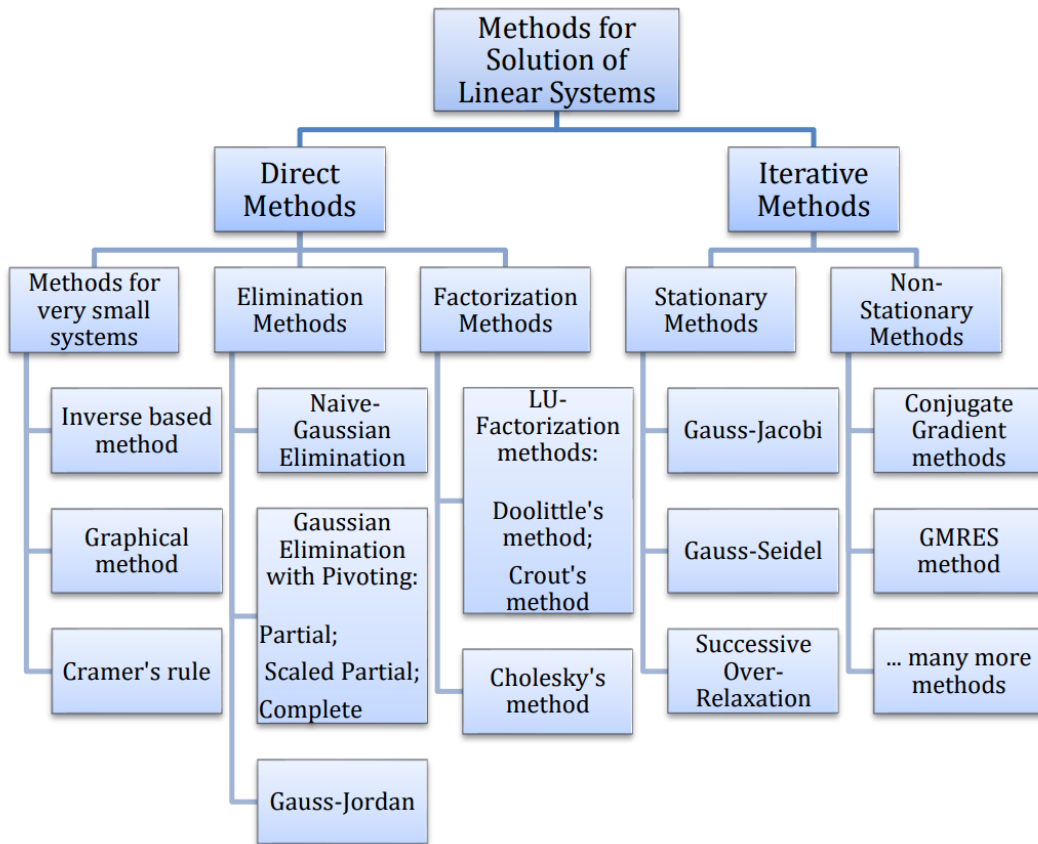


Fig. (6.3): A classification chart of linear solvers.

Corridor II: ANALYSIS

Let's think deep

- 6.6 Operation Count Analysis
- 6.7 Matrix Inversion

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



Corridor III: PROGRAMMING ARCADE

Let's do it

People have been communicating and interconnecting since the beginning, but in this era the communications and interconnections without modern technologies (like phones, networks, internet, radio, and television) stand nowhere in regards to possibility or survival. Likewise, people have been doing mathematics since early ages, but in this modern era the mathematical applicability without making use of the computers stands nowhere. Let's modernize "OUR" culture of doing mathematics so that it can be useful for all the disciplines of science and engineering. It's time to lead the frontiers of the knowledge and its applicability,

Remark: Suggestion: Before this Section, study, Corridor III of Chapter 07 to cope the difficulty level.

6.8 Algorithms and Implementations

The Gaussian Elimination Method with Partial Pivoting

Solving $AX = B$ using the Doolittle's Method

Solving $AX = B$ using the Crout's Method

Solving $AX = B$ using the Cholesky's Method

Performing Operation Count Analysis

To cross-check the results/output of the computer programs you would execute, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



6.8 Algorithms and Implementations

Question 20: Write down an algorithm (pseudo code) to solve a linear system $AX = B$ using the Gaussian Elimination method with partial pivoting.

Algorithm: To solve $AX = B$.

INPUTS: $\left\{ \begin{array}{l} n: \text{an integer as the number of equations and unknowns} \\ A = (a_{ij}), 1 \leq i, j \leq n: \text{a real valued square matrix as the coefficient matrix} \\ B = [b_1, b_2, \dots, b_n]^T: \text{a real valued vector as the vector of right hand side constants} \end{array} \right.$

OUTPUTS: $\left\{ \begin{array}{l} X = [x_1, x_2, \dots, x_n]^T: \text{a real valued vector as the solution vector} \\ \text{or a message that the given system has no unique solution} \end{array} \right.$

Step 1 Receive the inputs as stated above

Step 2 (Forward Elimination Phase)

for $i = 1, 2, \dots, n - 1$

$\left. \begin{array}{l} \text{Set } r = i \\ \text{for } j = i + 1, \dots, n \\ \text{if } (|a_{ri}| < |a_{ji}|) \quad r = j \end{array} \right\} \begin{array}{l} \text{(Searching largest absolute coefficient)} \\ \text{in } i\text{th column for partial pivoting} \end{array}$

if $(a_{ri} = 0)$ OUTPUT ("The given system has no unique solution') and **STOP**

else

if $(r \neq i)$, then interchange the i th row with r th row, and b_i with b_r

for $k = i + 1, i + 2, \dots, n$

$\left. \begin{array}{l} \text{multiplier} = \frac{a_{ki}}{a_{ii}} \\ \text{for } j = i + 1, i + 2, \dots, n \\ \quad a_{kj} = a_{kj} - \text{multiplier} \times a_{ij} \\ b_k = b_k - \text{multiplier} \times b_i \end{array} \right\} \begin{array}{l} \text{(row replacement in the} \\ \text{augmented matrix for} \\ \text{eliminating the coefficients} \\ \text{below the pivot)} \end{array}$

Step 3 if $(a_{nn} = 0)$ OUTPUT ("The given system has no unique solution') and **STOP**

else go to step 4

Step 4 (Back Substitution Phase)

$$x_n = \frac{b_n}{a_{nn}}$$

for $i = n - 1, \dots, 2, 1$

$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } j = i + 1, i + 2, \dots, n \\ \quad \text{sum} = \text{sum} + a_{ij} \times x_j \\ x_i = \frac{[b_i - \text{sum}]}{a_{ii}} \end{array} \right\} \left(x_i = \frac{1}{a_{ii}} \left[b_i - \sum_{j=i+1}^n a_{ij} x_j \right] \right)$

Step 5 Print the output: $X = [x_1, x_2, \dots, x_n]^T$ and **STOP**.

Problem 16: Write a C++ program to solve the following linear system using the Gaussian Elimination method with partial pivoting. For simplification, specify the linear system within the program.

$$1.7x_1 + 2.3x_2 - 1.5x_3 = 2.35$$

$$1.1x_1 + 1.6x_2 - 1.9x_3 = -0.94$$

$$2.7x_1 - 2.2x_2 + 1.5x_3 = 2.70$$

```

1  #include<iostream>
2  #include<cmath>
3  using namespace std;
4  # define n 3                               // number of unknowns
5
6  int main ( )
7  {
8      char quit ;
9      int i , j , k , n , r ;
10     double a[n][n] = {{1.7, 2.3, -1.5} , {1.1, 1.6, -1.9} , {2.7, -2.2, 1.5}} ;
11     double b[n] = {2.35, -0.94, 2.7};      // right-hand side constants
12     double x[n] ;                          // solution vector
13     double sum , temp , temp1 , multiplier ;
14
15     //----- Processing Section -----//
16
17     // Forward Elimination Phase
18     // Searching largest absolute coefficient in the ith column for partial pivoting
19
20     for ( i=0 ; i<n-1 ; i++ )
21     {
22         r = i ;
23         for ( j=i+1 ; j<n ; j++ )
24         {
25             if ( abs( a[r][i] )<abs( a[j][i] ) )
26                 r = j ;
27         }
28         if ( a[r][i] == 0 )
29         {
30             cout<< "The given system has no unique solution" ;
31             break ;
32         }
33         else if ( r != i )
34             for( j=0 ; j<n ; j++ )

```

```

35     {
36         temp = a[i][j] ;
37         a[i][j] = a[r][j] ;
38         a[r][j] = temp ;
39     }
40     temp1 = b[i] ;
41     b[i] = b[r] ;
42     b[r] = temp1 ;
43     for ( k=i+1 ; k<n ; k++ )
44     {
45         multiplier = a[k][i] / a[i][i] ;
46         for ( j=i+1 ; j<n ; j++ )
47             a[k][j] = a[k][j] - multiplier * a[i][j] ;
48
49         b[k] = b[k] - multiplier * b[i] ;
50     }
51 }
52
53 if ( a[n-1][n-1]==0 ) {
54     cout<< "The system has no unique solution. " ;
55     break ; }
56 else
57     x[n-1] = b[n-1] / a[n-1][n-1] ;
58
59 for ( i=n-2 ; i>=0 ; i-- )
60     {
61         sum = 0.0 ;
62         for( j=i+1 ; j<n ; j++ )
63             sum = sum + a[i][j] * x[j] ;
64
65         x[i] = ( b[i] - sum ) / a[i][i] ;
66     }
67 //----- Output Section -----//
68
69 cout<< "The solution of the given system is " <<endl ;
70 for ( i=0; i<n ; i++ )
71     cout<< x[i] << " " ;
72
73 cin>> quit ;
74 return 0 ;
75 }

```

row replacement in the augmented matrix for eliminating the coefficient below the pivot

$$\left(x_i = \frac{1}{a_{ii}} \left[b_i - \sum_{j=i+1}^n a_{ij} x_j \right] \right)$$

Remark: The C++ program in Problem 16 can be modified to receive the linear system at the execution time (instead of fixing in the code). For this, lines 10 and 11 in the solution of Problem 16 should be replaced by the following code segment:

```
//----- Input Section -----//

cout<<"\nEnter the coefficient matrix row-wise: "<<n<<" unknowns<<endl;
for ( i=0 ; i<n ; i++)
    for ( j=0 ; j<n ; j++)
        cin >> a[i][j] ;

cout<<"Enter the elements of constant vector B "<<endl ;
for ( i=0 ; i<n ; i++)
    cin >> b[i] ;
```

Question 21: Write down an algorithm (pseudo code) to solve a linear system using the Doolittle's method.

Algorithm: To solve a linear system $\mathbf{AX} = \mathbf{B}$, for which the factorization $\mathbf{A} = \mathbf{LU}$ is possible.

INPUTS: $\left\{ \begin{array}{l} \mathbf{n}: \text{an integer as the number of equations and unknowns} \\ \mathbf{A} = (\mathbf{a}_{ij}), 1 \leq i, j \leq \mathbf{n}: \text{a real valued square matrix as the coefficient matrix} \\ \mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n]^T: \text{a real valued vector as the vector of right hand side constants} \end{array} \right.$

OUTPUTS: $\left\{ \begin{array}{l} \mathbf{L} = (\mathbf{l}_{ij}), 1 \leq i, j \leq \mathbf{n}: \text{a real valued square matrix as the lower triangular matrix} \\ \mathbf{U} = (\mathbf{u}_{ij}), 1 \leq i, j \leq \mathbf{n}: \text{a real valued square matrix as the upper triangular matrix} \\ \mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T: \text{a real valued vector as the solution vector} \end{array} \right.$

Step 1 (Formation of \mathbf{L} and \mathbf{U} as factors of \mathbf{A} , i.e., $\mathbf{A} = \mathbf{LU}$)

for $i = 1, 2, \dots, \mathbf{n}$

Set $\mathbf{l}_{ii} = 1$

For $j = i, i + 1, \dots, \mathbf{n}$

$sum = 0$

for $s = 1, 2, \dots, i - 1$

$sum = sum + \mathbf{l}_{is} \times \mathbf{u}_{sj}$

$\mathbf{u}_{ij} = \mathbf{a}_{ij} - sum$

$$\left(\mathbf{u}_{ij} = \mathbf{a}_{ij} - \sum_{s=1}^{i-1} \mathbf{l}_{is} \mathbf{u}_{sj} \right)$$

for $j = i + 1, i + 2, \dots, \mathbf{n}$

$sum = 0$

for $s = 1, 2, \dots, i - 1$

$sum = sum + \mathbf{l}_{js} \times \mathbf{u}_{si}$

$\mathbf{l}_{ji} = \frac{[\mathbf{a}_{ji} - sum]}{\mathbf{u}_{ii}}$

$$\left(\mathbf{l}_{ji} = \frac{1}{\mathbf{u}_{ii}} \left[\mathbf{a}_{ji} - \sum_{s=1}^{i-1} \mathbf{l}_{js} \mathbf{u}_{si} \right] \right)$$

Step 2 (Forward substitution phase for solving $LY = B$)

$$y_1 = b_1$$

for $i = 2, 3, \dots, n$

$$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } j = 1, 2, \dots, i - 1 \\ \quad \text{sum} = \text{sum} + l_{ij} \times y_j \\ y_i = b_i - \text{sum} \end{array} \right\}$$

$$\left(y_i = b_i - \sum_{j=1}^{i-1} l_{ij} y_j \right)$$

Step 3 (Back Substitution Phase for solving $UX = Y$)

$$x_n = \frac{y_n}{u_{nn}}$$

for $i = n - 1, \dots, 2, 1$

$$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } j = i + 1, i + 2, \dots, n \\ \quad \text{sum} = \text{sum} + u_{ij} \times x_j \\ x_i = \frac{[y_i - \text{sum}]}{u_{ii}} \end{array} \right\}$$

$$\left(x_i = \frac{1}{u_{ii}} \left[y_i - \sum_{j=i+1}^n u_{ij} x_j \right] \right)$$

STOP.

Problem 18: Write a C++ program to solve the following linear system using the Doolittle's method. For simplification, specify the linear system within the program.

$$1.7x_1 + 2.3x_2 - 1.5x_3 = 2.35$$

$$1.1x_1 + 1.6x_2 - 1.9x_3 = -0.94$$

$$2.7x_1 - 2.2x_2 + 1.5x_3 = 2.70$$

```

1  #include<iostream>
2  using namespace std ;
3  # define n 3
4
5  int main ( )
6  {
7      char quit ;
8      int i , j , s ;
9      double a[n][n] = { {1.7, 2.3, -1.5} , {1.1, 1.6, -1.9} , {2.7, -2.2, 1.5} } ;
10     double b[n] = {2.35, -0.94, 2.7} ;           // right-hand side constant vector
11     double x[n] ;
12     double u[n][n] , l[n][n] , y[n] , sum ;
13
14     //----- Processing Section -----//
15     // Formation of L and U as factors of A, i.e., A=LU
16

```



```

17   for ( i=0 ; i<n ; i++ )
18       for ( j=0 ; j<n ; j++ )
19           l[i][j] = u[i][j] = 0.0 ;
20
21   for ( j=0 ; j<n ; j++ )
22       {
23           u[0][j] = a[0][j] ;
24           l[j][0] = a[j][0] / u[0][0] ;
25       }
26   for ( i=1 ; i<n ; i++ )
27       {
28           l[i][i] = 1 ;
29           for ( j=i ; j<n ; j++ )
30               {
31                   sum = 0.0 ;
32                   for ( s=0 ; s<=i-1 ; s++ )
33                       sum = sum + l[i][s] * u[s][j] ;
34                   u[i][j] = a[i][j] - sum ;
35               }
36
37           for ( j=i+1 ; j<n ; j++ )
38               {
39                   sum = 0.0 ;
40                   for ( s=0 ; s<=i-1 ; s++ )
41                       sum = sum + l[j][s] * u[s][i] ;
42                   l[j][i] = ( a[j][i] - sum ) / u[i][i] ;
43               }
44       }
45
46   // Forward substitution phase for solving LY=B
47
48   y[0] = b[0] ;
49   for ( i=1 ; i<n ; i++ )
50       {
51           sum = 0.0 ;
52           for ( j=0 ; j<=i-1 ; j++ )
53               sum = sum + l[i][j] * y[j] ;
54           y[i] = b[i] - sum ;
55       }
56
57   // Back Substitution Phase for solving UX=Y

```

$$\left(u_{ij} = a_{ij} - \sum_{s=1}^{i-1} l_{is} u_{sj} \right)$$

$$\left(l_{ji} = \frac{1}{u_{ii}} \left[a_{ji} - \sum_{s=1}^{i-1} l_{js} u_{si} \right] \right)$$

$$\left(y_i = b_i - \sum_{j=1}^{i-1} l_{ij} y_j \right)$$

```

58
59 x[n-1] = y[n-1] / u[n-1][n-1] ;
60 for ( i=n-2 ; i>=0 ; i-- )
61 {
62     sum = 0.0 ;
63     for ( j=i+1 ; j<n ; j++ )
64         sum = sum + u[i][j] * x[j] ;
65     x[i] = ( y[i] - sum ) / u[i][i] ;
66 }
67
68 //----- Output Section -----//
69
70 cout<< "The L matrix is" << endl ;
71 for ( i=0 ; i<n ; i++ )
72 {
73     for ( j = 0 ; j<n ; j++ )
74     {
75         cout<< l[i][j]<< " " ;
76     }
77     cout<< endl ;
78 }
79 cout<< endl ;
80 cout<< "The U matrix is" << endl ;
81 for ( i=0 ; i<n ; i++ )
82 {
83     for ( j = 0 ; j<n ; j++ )
84     {
85         cout<< u[i][j] << " " ;
86     }
87     cout<< endl ;
88 }
89 cout<< endl ;
90 cout<< "The required solution is" << endl ;
91 for ( i=0 ; i<n ; i++ )
92     cout<< x[i] << " " ;
93
94 cin>> quit ;
95 return 0 ;
96 }

```

$$x_i = \frac{1}{u_{ii}} \left[y_i - \sum_{j=i+1}^n u_{ij} x_j \right]$$



Remark: Replace the lines 9 and 10 in the solution of Problem 18 with the following code segment to receive the linear system at the execution time (instead of fixing in the code):

```
//----- Input Section -----//

cout<<"\nEnter the coefficient matrix row-wise: "<<n<<" unknowns<<endl;
for ( i=0 ; i<n ; i++)
    for ( j=0 ; j<n ; j++)
        cin >> a[i][j] ;

cout<<"Enter the elements of constant vector B "<<endl ;
for ( i=0 ; i<n ; i++)
    cin >> b[i] ;
```

Question 22: Write down an algorithm (pseudo code) to solve a linear system using the Crout's method.

Algorithm: To solve a linear system $AX = B$, for which the factorization $A = LU$ is possible.

INPUTS: $\left\{ \begin{array}{l} n: \text{an integer as the number of equations and unknowns} \\ A = (a_{ij}), 1 \leq i, j \leq n: \text{a real valued square matrix as the coefficient matrix} \\ B = [b_1, b_2, \dots, b_n]^T: \text{a real valued vector as the vector of right hand side constants} \end{array} \right.$

OUTPUTS: $\left\{ \begin{array}{l} L = (l_{ij}), 1 \leq i, j \leq n: \text{a real valued square matrix as the lower triangular matrix} \\ U = (u_{ij}), 1 \leq i, j \leq n: \text{a real valued square matrix as the upper triangular matrix} \\ X = [x_1, x_2, \dots, x_n]^T: \text{a real valued vector as the solution vector} \end{array} \right.$

Step 1 (Formation of L and U as factors of A , i.e., $A = LU$)

for $i = 1, 2, \dots, n$

Set $u_{ii} = 1$

for $j = i, i + 1, \dots, n$

$sum = 0$

for $s = 1, 2, \dots, i - 1$

$sum = sum + l_{js} \times u_{si}$

$l_{ji} = a_{ji} - sum$

$$\left(l_{ji} = a_{ji} - \sum_{s=1}^{i-1} l_{js} u_{si} \right)$$

for $j = i + 1, i + 2, \dots, n$

$sum = 0$

for $s = 1, 2, \dots, i - 1$

$sum = sum + l_{is} \times u_{sj}$

$u_{ij} = \frac{[a_{ij} - sum]}{l_{ii}}$

$$\left(u_{ij} = \frac{1}{l_{ii}} \left[a_{ij} - \sum_{s=1}^{i-1} l_{is} u_{sj} \right] \right)$$

Step 2 (Forward Substitution Phase for solving $LY = B$)

$$y_1 = \frac{b_1}{l_{11}}$$

for $i = 2, 3, \dots, n$

$$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } j = 1, 2, \dots, i - 1 \\ \quad \text{sum} = \text{sum} + l_{ij} \times y_j \\ y_i = \frac{[b_i - \text{sum}]}{l_{ii}} \end{array} \right\}$$

$$\left(y_i = \frac{1}{l_{ii}} \left[b_i - \sum_{j=1}^{i-1} l_{ij} y_j \right] \right)$$

Step 3 (Back Substitution Phase for solving $UX = Y$)

$$x_n = y_n$$

for $i = n - 1, \dots, 2, 1$

$$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } j = i + 1, i + 2, \dots, n \\ \quad \text{sum} = \text{sum} + u_{ij} \times x_j \\ x_i = y_i - \text{sum} \end{array} \right\}$$

$$\left(x_i = y_i - \sum_{j=i+1}^n u_{ij} x_j \right)$$

STOP.

Problem 20: Write a C++ program to solve the following linear system using the Crout's method. For simplification, specify the linear system within the program.

$$1.7x_1 + 2.3x_2 - 1.5x_3 = 2.35$$

$$1.1x_1 + 1.6x_2 - 1.9x_3 = -0.94$$

$$2.7x_1 - 2.2x_2 + 1.5x_3 = 2.70$$

```

1 #include<iostream>
2 using namespace std ;
3 # define n 3 // number of unknowns
4
5 int main ( )
6 {
7     char quit ;
8     int i , j , s ;
9     double a[n][n] = {{1.7, 2.3, -1.5} , {1.1, 1.6, -1.9} , {2.7, -2.2, 1.5}} ;
10    double b[n] = {2.35, -0.94, 2.7} ; // right hand side constant vector
11    double x[n] ;
12    double u[n][n] , l[n][n] , y[n] , sum ;
13
14    //----- Processing Section -----//
15    // Formation of L and U as factors of A, i.e., A=LU
16
```

```

17   for ( i=0 ; i<n ; i++ )
18       for ( j=0 ; j<n ; j++ )
19           l[i][j] = u[i][j] = 0.0 ;
20
21   for ( j=0 ; j<n ; j++ )
22   {
23       l[j][0] = a[j][0] ;
24       u[0][j] = a[0][j] / l[0][0] ;
25   }
26
27   for ( i=1 ; i<n ; i++ )
28   {
29       u[i][i] = 1 ;
30       for ( j=i ; j<n ; j++ )
31       {
32           sum = 0.0 ;
33           for ( s=0 ; s<=i-1 ; s++ )
34               sum = sum + l[j][s] * u[s][i];
35           l[j][i] = a[j][i] - sum ;
36       }
37
38       for ( j=i+1 ; j<n ; j++ )
39       {
40           sum = 0.0 ;
41           for ( s=0 ; s<=i-1 ; s++ )
42               sum = sum + l[i][s] * u[s][j] ;
43           u[i][j] = ( a[i][j] - sum ) / l[i][i] ;
44       }
45   }
46
47   // Forward substitution phase for solving LY=B
48
49   y[0] = b[0] / l[0][0] ;
50   for ( i=1 ; i<n ; i++ )
51   {
52       sum = 0.0 ;
53       for ( j=0 ; j<=i-1 ; j++ )
54           sum = sum + l[i][j] * y[j] ;
55       y[i] = ( b[i] - sum ) / l[i][i] ;
56   }
57

```

$$\left(l_{ji} = a_{ji} - \sum_{s=1}^{i-1} l_{js} u_{si} \right)$$

$$\left(u_{ij} = \frac{1}{l_{ii}} \left[a_{ij} - \sum_{s=1}^{i-1} l_{is} u_{sj} \right] \right)$$

$$\left(y_i = \frac{1}{l_{ii}} \left[b_i - \sum_{j=1}^{i-1} l_{ij} y_j \right] \right)$$

```

58 // Back Substitution Phase for solving UX=Y
59
60 x[n-1] = y[n-1] ;
61 for ( i=n-2 ; i>=0 ; i-- )
62 {
63     sum = 0.0 ;
64     for ( j=i+1 ; j<n ; j++ )
65         sum = sum + u[i][j] * x[j] ;
66     x[i] = y[i] - sum ;
67 }
68
69 //----- Output Section -----//
70
71 cout<< "The L matrix is" << endl ;
72 for ( i=0 ; i<n ; i++ )
73 {
74     for ( j = 0 ; j<n ; j++ )
75     {
76         cout<< l[i][j]<< " " ;
77     }
78     cout<< endl ;
79 }
80 cout<< endl ;
81 cout<< "The U matrix is" << endl ;
82 for ( i=0 ; i<n ; i++ )
83 {
84     for ( j = 0 ; j<n ; j++ )
85     {
86         cout<< u[i][j] << " " ;
87     }
88     cout<< endl ;
89 }
90 cout<< endl ;
91 cout<< "The required solution is" << endl ;
92 for ( i=0 ; i<n ; i++ )
93     cout<< x[i] << " " ;
94
95 cin>> quit ;
96 return 0 ;
97 }

```

$$\left(x_i = y_i - \sum_{j=i+1}^n u_{ij}x_j \right)$$

Remark: Replace the lines 9 and 10 in the solution of Problem 20 with the following code segment to receive the linear system at the execution time (instead of fixing in the code):

```
//----- Input Section -----//  
  
cout<<"\nEnter the coefficient matrix row-wise: "<<n<<" unknowns<<endl;  
for ( i=0 ; i<n ; i++)  
    for ( j=0 ; j<n ; j++)  
        cin >> a[i][j] ;  
  
cout<<"Enter the elements of constant vector B "<<endl ;  
for ( i=0 ; i<n ; i++)  
    cin >> b[i] ;
```

■

Question 23: Write down an algorithm (pseudo code) to solve a linear system using the Cholesky's method.

Algorithm: To solve a linear system $AX = B$, for which the factorization $A = LL^T$ is possible.

INPUTS: $\left\{ \begin{array}{l} n: \text{an integer as the number of equations and unknowns} \\ A = (a_{ij}), 1 \leq i, j \leq n: \text{a real valued square matrix as the coefficient matrix} \\ B = [b_1, b_2, \dots, b_n]^T: \text{a real valued vector as the vector of right hand side constants} \end{array} \right.$

OUTPUTS: $\left\{ \begin{array}{l} L = (l_{ij}), 1 \leq i, j \leq n: \text{a real valued square matrix as the lower triangular matrix} \\ X = [x_1, x_2, \dots, x_n]^T: \text{a real valued vector as the solution vector} \end{array} \right.$

Step 1 (Formation of L as factors of A , i.e., $A = LL^T$)

for $i = 1, 2, \dots, n$

$$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } k = 1, 2, \dots, i-1 \\ \quad \text{sum} = \text{sum} + l_{ik} \times l_{ik} \\ l_{ii} = \text{sqr}t(a_{ii} - \text{sum}) \end{array} \right\}$$

$$\left(l_{ii} = \left[a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2 \right]^{\frac{1}{2}} \right)$$

for $j = i + 1, i + 2, \dots, n$

$$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } k = 1, 2, \dots, i-1 \\ \quad \text{sum} = \text{sum} + l_{ik} \times l_{jk} \\ l_{ji} = \frac{[a_{ji} - \text{sum}]}{l_{ii}} \end{array} \right\}$$

$$\left(l_{ji} = \frac{1}{l_{ii}} \left[a_{ji} - \sum_{k=1}^{i-1} l_{ik} l_{jk} \right] \right)$$

Step 2 (Forward Substitution Phase for solving $LY = B$)

$$y_1 = \frac{b_1}{l_{11}}$$

for $i = 2, 3, \dots, n$

$$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } j = 1, 2, \dots, i-1 \\ \quad \text{sum} = \text{sum} + l_{ij} \times y_j \\ y_i = \frac{[b_i - \text{sum}]}{l_{ii}} \end{array} \right\}$$

$$\left(y_i = \frac{1}{l_{ii}} \left[b_i - \sum_{j=1}^{i-1} l_{ij} y_j \right] \right)$$

Step 3 (Back Substitution Phase for solving $L^T X = Y$)

$$x_n = \frac{y_n}{l_{nn}}$$

for $i = n - 1, \dots, 2, 1$

$$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } j = i + 1, i + 2, \dots, n \\ \quad \text{sum} = \text{sum} + l_{ji} \times x_j \\ x_i = \frac{[y_i - \text{sum}]}{l_{ii}} \end{array} \right\}$$

$$\left(x_i = \frac{1}{l_{ii}} \left[y_i - \sum_{j=i+1}^n l_{ji} x_j \right] \right)$$

STOP.

Problem 22: Write a C++ program to solve the following positive definite linear system using the Cholesky's method. For simplification, specify the linear system within the program.

$$\begin{aligned} 0.4x_1 & & + 0.12x_3 & = 1.4 \\ & 0.64x_2 & + 0.32x_3 & = 1.6 \\ -0.12x_1 & + 0.32x_2 & + 0.56x_3 & = 5.4 \end{aligned}$$

```

1  #include <iostream>
2  #include <cmath>
3  using namespace std ;
4  # define n 3                               // number of unknowns
5
6  int main ( )
7  {
8      char quit ;
9      int i , j , k ;
10     double a[n][n] = {{0.4, 0, 0.12} , {0, 0.64, 0.32} , {-0.12, 0.32, 0.56}} ;
11     double b[n] = {1.4, 1.6, 5.4} ;          // right hand side constant vector
12     double x[n] ;
13     double l[n][n] , y[n] , sum ;
14
15     //----- Processing Section -----//
16     // Formation of L and U as factors of A, i.e., A=LU
17
18     for ( i=0 ; i<n ; i++ )
19         for ( j=0 ; j<n ; j++ )
20             l[i][j] = 0.0 ;
21     l[0][0] = sqrt( a[0][0] ) ;
22
23     for ( j=1 ; j<n ; j++ )
24         l[j][0] = a[j][0] / l[0][0] ;
25
26     for ( i=1 ; i<n ; i++ )
27     {
28         sum = 0.0 ;
29         for ( k=0 ; k<=i-1 ; k++ )
30             sum = sum + l[i][k] * l[i][k] ;
31         l[i][i] = sqrt( a[i][i] - sum ) ;
32

```

$$l_{ii} = \left[a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2 \right]^{\frac{1}{2}}$$

```

33     for ( j=i+1 ; j<n ; j++ )
34     {
35         sum = 0.0 ;
36         for ( k=0 ; k<=i-1 ; k++ )
37             sum = sum + l[i][k] * l[j][k] ;
38         l[j][i] = (a[j][i] - sum) / l[i][i] ;
39     }
40 }
41
42 // Forward substitution phase for solving LY=B
43
44 y[0] = b[0] / l[0][0] ;
45 for ( i=1 ; i<n ; i++ )
46 {
47     sum = 0.0 ;
48     for ( j=0 ; j<=i-1 ; j++ )
49         sum = sum + l[i][j] * y[j] ;
50     y[i] = ( b[i] - sum ) / l[i][i] ;
51 }
52
53 // Back Substitution Phase for solving L^T X=Y
54
55 x[n-1] = y[n-1] / l[n-1][n-1] ;
56 for ( i=n-2 ; i>=0 ; i-- )
57 {
58     sum = 0.0 ;
59     for ( j=i+1 ; j<n ; j++ )
60         sum = sum + l[j][i] * x[j] ;
61     x[i] = (y[i] - sum) / l[i][i] ;
62 }
63
64 //----- Output Section -----//
65
66 cout<< "The L matrix is" << endl ;
67 for ( i=0 ; i<n ; i++ )
68 {
69     for ( j = 0 ; j<n ; j++ )
70     {
71         cout<< l[i][j]<< " " ;
72     }
73     cout<< endl ;

```

$$l_{ji} = \frac{1}{l_{ii}} \left[a_{ji} - \sum_{k=1}^{i-1} l_{ik} l_{jk} \right]$$

$$y_i = \frac{1}{l_{ii}} \left[b_i - \sum_{j=1}^{i-1} l_{ij} y_j \right]$$

$$x_i = \frac{1}{l_{ii}} \left[y_i - \sum_{j=i+1}^n l_{ji} x_j \right]$$

```

74     }
75     cout<< endl ;
76     cout<< "The required solution is" << endl ;
77     for ( i=0 ; i<n ; i++ )
78         cout<< x[i] << " " ;
79
80     cin>> quit ;
81     return 0 ;
82 }

```

■

Remark: Following are some notations and formulas that might be useful in carrying out operation count analysis of the algorithms.

$$\sum_{p=1}^n cf(p) = c \sum_{p=1}^n f(p)$$

$$\sum_{p=1}^n [f(p) + g(p)] = \sum_{p=1}^n f(p) + \sum_{p=1}^n g(p)$$

$$\sum_{p=1}^n 1 = 1 + 1 + \dots + 1 = n$$

$$\sum_{p=k}^n 1 = n - k + 1$$

$$\sum_{p=1}^n p = 1 + 2 + 3 + \dots + n = \frac{n(n+1)}{2} = \frac{n^2}{2} + \mathcal{O}(n)$$

$$\sum_{p=1}^n p^2 = 1^2 + 2^2 + \dots + n^2 = \frac{n(n+1)(2n+1)}{6} = \frac{n^3}{3} + \mathcal{O}(n^2)$$

Question 24: Perform the operation count analysis of the algorithm that involves the following phases to solve an $n \times n$ linear system:

- (1) Forward elimination to obtain the upper triangular form using the Gauss Elimination method.
- (2) Back substitution to solve the upper triangular system.

(1) The forward elimination phase occurs just after setting the inputs in the algorithm. This phase contains three nested loops. The first loop, say i -loop (which ranges from $i = 1$ to $n - 1$), corresponds to the $n - 1$ elimination stages of the method. For each row i , the i th element is considered a pivot element. The second loop, say k -loop (which ranges from $k = i + 1$ to n), corresponds to the elements below the pivot element to make them zero. The third loop, say j -loop (which ranges from $j = i + 1$ to n) corresponds to the columns after the pivot element.

Note that, for any loop with index ranging from $i + 1$ to n , the number of passes/iterations will be $n - (i + 1) + 1$ (or simply $(n - i)$ passes). Therefore, each of the k -loop and j -loop has $(n - i)$ passes.

Each pass of k -loop will perform one division to obtain the multiplier, and one multiplication and subtraction to update the right-hand side constant, b_k . Moreover, in each pass of k -loop, $(n - i)$ multiplications and $(n - i)$ subtractions will be performed in j -loop to update the relevant entries of the coefficient matrix, a_{kj} . Thus, in each pass of k -loop, the total number of multiplications/divisions will be $(1 + 1 + (n - i))$ or $(n - i + 2)$ and the total number of additions/subtractions will be $(1 + n - i)$.

As there are $(n - i)$ passes of k -loop in each pass of i -loop, therefore there will be $(n - i) \times (n - i + 2)$ multiplications/divisions and $(n - i) \times (n - i + 1)$ additions/subtractions in each pass of i -loop.

Hence, the total number of multiplications/divisions in $n - 1$ passes of i -loop of the forward elimination phase will be

$$\begin{aligned}
 \sum_{i=1}^{n-1} (n-i)(n-i+2) &= \sum_{i=1}^{n-1} (n-i)((n+2)-i) \\
 &= \sum_{i=1}^{n-1} [n(n+2) - ni - i(n+2) + i^2] = \sum_{i=1}^{n-1} [n(n+2) - 2i(n+1) + i^2] \\
 &= n(n+2) \sum_{i=1}^{n-1} 1 - 2(n+1) \sum_{i=1}^{n-1} i + \sum_{i=1}^{n-1} i^2 \\
 &= n(n+2)(n-1) - 2(n+1) \left[\frac{(n-1)n}{2} \right] + \left[\frac{(n-1)n(2n-1)}{6} \right] \\
 &= n(n-1) \left[n+2 - n-1 + \frac{n}{3} - \frac{1}{6} \right] \\
 &= (n^2 - n) \left[\frac{n}{3} + \frac{5}{6} \right] = \frac{n^3}{3} + \frac{n^2}{2} - \frac{5n}{6} = \frac{n^3}{3} + \mathcal{O}(n^2)
 \end{aligned}$$

Similarly, the total number of additions/subtractions in $n - 1$ passes of i -loop of the forward elimination phase will be

$$\begin{aligned}
\sum_{i=1}^{n-1} (n-i)(n-i+1) &= \sum_{i=1}^{n-1} (n-i)((n+1)-i) \\
&= \sum_{i=1}^{n-1} [n(n+1) - ni - i(n+1) + i^2] \\
&= \sum_{i=1}^{n-1} [n(n+1) - i(2n+1) + i^2] \\
&= n(n+1) \sum_{i=1}^{n-1} 1 - (2n+1) \sum_{i=1}^{n-1} i + \sum_{i=1}^{n-1} i^2 \\
&= n(n+1)(n-1) - (2n+1) \left[\frac{(n-1)n}{2} \right] + \left[\frac{(n-1)n(2n-1)}{6} \right] \\
&= n(n-1) \left[n+1 - n - \frac{1}{2} + \frac{n}{3} - \frac{1}{6} \right] \\
&= (n^2 - n) \left[\frac{n}{3} + \frac{1}{3} \right] = \frac{n^3}{3} - \frac{n}{3} = \frac{n^3}{3} + \mathcal{O}(n)
\end{aligned}$$

The summary of the operation count of the Gaussian Elimination phase is given as:

	Operations		Total flops
	Multiplications/divisions	Additions/subtractions	
Forward Elimination	$\frac{n^3}{3} + \frac{n^2}{2} - \frac{5n}{6}$	$\frac{n^3}{3} - \frac{n}{3}$	$\frac{2n^3}{3} + \mathcal{O}(n^2)$

(2) The back substitution phase occurs after the forward elimination phase. This phase contains two nested loops. The first loop, say i -loop (which ranges from $i = n - 1$ to 1), corresponds to $n - 1$ of the components of the solution vector. The second loop, say j -loop (which ranges from $j = i + 1$ to n), corresponds to the columns after the diagonal elements.

Each pass of i -loop will perform one subtraction and one division to obtain the value of x_i . Moreover, in each pass of i -loop, the number of both of the multiplications and additions will be $n - (i + 1) + 1$ (or simply $(n - i)$) in j -loop. Thus, in each pass of i -loop, the total number of both of the multiplications/divisions and additions/subtractions will be $(n - i + 1)$.

Hence, the total number of the multiplications/divisions in the back substitution phase will be

$$\begin{aligned} 1 + \sum_{i=1}^{n-1} (n+1-i) &= 1 + (n+1) \sum_{i=1}^{n-1} 1 - \sum_{i=1}^{n-1} i = 1 + (n+1)(n-1) - \left[\frac{(n-1)n}{2} \right] \\ &= 1 + n^2 - 1 - \frac{n^2}{2} + \frac{n}{2} = \frac{n^2}{2} + \frac{n}{2} = \frac{n^2}{2} + \mathcal{O}(n) \end{aligned}$$

Similarly, the total number of the additions/subtractions in the back substitution phase will be

$$\sum_{i=1}^{n-1} (n+1-i) = \frac{n^2}{2} + \frac{n}{2} - 1 = \frac{n^2}{2} + \mathcal{O}(n)$$

Finally, the summary of the operation count of the complete algorithm (including the two phases) is given as:

	Operations		Total flops
	Multiplications/divisions	Additions/subtractions	
Forward elimination	$\frac{n^3}{3} + \frac{n^2}{2} - \frac{5n}{6}$	$\frac{n^3}{3} - \frac{n}{3}$	$\frac{2n^3}{3} + \frac{n^2}{2} - \frac{7n}{6}$
Back Substitution	$\frac{n^2}{2} + \frac{n}{2}$	$\frac{n^2}{2} - \frac{n}{2}$	n^2
Totals	$\frac{n^3}{3} + n^2 - \frac{n}{3}$	$\frac{n^3}{3} + \frac{n^2}{2} - \frac{5n}{6}$	$\frac{2n^3}{3} + \frac{3n^2}{2} - \frac{7n}{6}$

Question 25: Perform the operation count analysis of the algorithm that involves the following phases to solve an $n \times n$ linear system:

- (1) Factorization of the coefficient matrix using the Doolittle's method.
- (2) Forward substitution to solve the lower triangular system.
- (3) Back substitution to solve the upper triangular system.

(1) The factorization of the coefficient matrix A into the product of the unit lower triangular L and the upper triangular U matrices occurs just after setting the inputs in the algorithm. The formulation of L and U as the factors of A contains three nested loops. The first loop, say i -loop (which ranges from $i = 1$ to n), corresponds to the i th row and column of U and L respectively. The second loop, say j -loop (ranges from $j = i$ to n), corresponds to the column j of U and (ranges

from $j = i + 1$ to n), corresponds to the row j of L . The third loop, say s -loop (which ranges from $s = 1$ to $i - 1$), corresponds to the multiplication of the rows of L and columns of U .

Note that, the j -loop corresponding to column j of U ranging from i to n , the number of passes/iterations will be $n - i + 1$. Similarly, the number of passes/iterations in j -loop, corresponds to row j of L ranging from $i + 1$ to n , will be $n - (i + 1) + 1$ (or simply $(n - i)$).

Each pass of j -loop will perform one subtraction to obtain the entry u_{ij} of U . Moreover, in each pass of j -loop, the number of both of the multiplications and additions will be $(i - 1) - 1 + 1$ (or simply $(i - 1)$) in s -loop. Thus, in each pass of j -loop, the total number of multiplications/divisions will be $(i - 1)$ and the total number of additions/subtractions will be $(1 + i - 1)$ or (i) . Similarly, in each pass of j -loop, to obtain the entry l_{ji} of L , the total number of both of the multiplications and additions will be $(i - 1) + 1$ (or simply (i)).

As there are $(n - i + 1)$ passes of j -loop in each pass of i -loop, therefore there will be $(n - i + 1) \times (i - 1)$ multiplications/divisions and $(n - i + 1) \times (i)$ additions/subtractions in each pass of i -loop for the formulation of row i of U . Similarly, in each pass of i -loop, there will be $(n - i) \times (i)$ multiplications/divisions and $(n - i) \times (i)$ additions/subtractions in each pass of i -loop for the formulation of column i of L .

Hence, the total number of multiplications/divisions in n passes of i -loop for the formulation of upper triangular matrix U will be

$$\begin{aligned}
\sum_{i=1}^n (n - i + 1)(i - 1) &= \sum_{i=1}^n (n + 1 - i)(i - 1) \\
&= \sum_{i=1}^n [(n + 1)i - (n + 1) - i^2 + i] = \sum_{i=1}^n [(n + 2)i - i^2 - (n + 1)] \\
&= (n + 2) \sum_{i=1}^n i - \sum_{i=1}^n i^2 - (n + 1) \sum_{i=1}^n 1 \\
&= (n + 2) \left[\frac{n(n + 1)}{2} \right] - \left[\frac{n(n + 1)(2n + 1)}{6} \right] - (n + 1)n \\
&= n(n + 1) \left[\frac{n}{2} + 1 - \frac{n}{3} - \frac{1}{6} - 1 \right] = (n^2 + n) \left[\frac{n}{6} - \frac{1}{6} \right] \\
&= \frac{n^3}{6} - \frac{n}{6} = \frac{n^3}{6} + \mathcal{O}(n)
\end{aligned}$$

Similarly, the total number of additions/subtractions in n passes of i -loop for the formulation of upper triangular matrix U will be

$$\sum_{i=1}^n (n-i+1)(i) = \frac{n^3}{6} + \frac{n}{6} = \frac{n^3}{6} + \mathcal{O}(n)$$

Moreover, the total number of multiplications/divisions and additions/subtractions in n passes of i -loop for the formulation of unit lower triangular matrix L will be

$$\begin{aligned} \sum_{i=1}^n (n-i)(i) &= \sum_{i=1}^n (n-i)(i) = \sum_{i=1}^n (ni - i^2) \\ &= n \sum_{i=1}^n i - \sum_{i=1}^n i^2 = n \left[\frac{n(n+1)}{2} \right] - \left[\frac{n(n+1)(2n+1)}{6} \right] \\ &= n(n+1) \left[\frac{n}{2} - \frac{n}{3} - \frac{1}{6} \right] = (n^2 + n) \left[\frac{n}{6} - \frac{1}{6} \right] \\ &= \frac{n^3}{6} - \frac{n}{6} = \frac{n^3}{6} + \mathcal{O}(n) \end{aligned}$$

The summary of the operation count of the LU -factorization is given as:

	Operations		Total flops
	Multiplication/Division	Addition/Subtraction	
Upper Triangular Matrix U	$\frac{n^3}{6} - \frac{n}{6}$	$\frac{n^3}{6} + \frac{n}{6}$	$\frac{n^3}{3}$
Lower Triangular Matrix L	$\frac{n^3}{6} - \frac{n}{6}$	$\frac{n^3}{6} - \frac{n}{6}$	$\frac{n^3}{3} - \frac{n}{3}$
LU -factorization	$\frac{n^3}{3} - \frac{n}{3}$	$\frac{n^3}{3}$	$\frac{2n^3}{3} - \frac{n}{3}$

(2) The forward substitution phase occurs after the formulation of L and U as factors of the coefficient matrix for solving the lower triangular system. This phase contains two nested loops. The first loop, say i -loop (which ranges from $i = 2$ to n), corresponds to $n - 1$ of the components

of the intermediate vector Y . The second loop, say j -loop (which ranges from $j = 1$ to $i - 1$), corresponds to the columns before the diagonal elements.

Each pass of i -loop will perform one subtraction to obtain the value of y_i . Moreover, in each pass of i -loop, the number of both of the multiplications and additions will be $(i - 1) - 1 + 1$ (or simply $(i - 1)$) in j -loop. Thus, in each pass of i -loop, the total number of multiplications/divisions will be $(i - 1)$ and the total number of additions/subtractions will be $(1 + i - 1)$ or (i) .

Hence, the total number of the multiplications/divisions in the forward substitution phase will be

$$\begin{aligned} \sum_{i=2}^n (i - 1) &= \sum_{i=2}^n i - \sum_{i=2}^n 1 = \left[\frac{n(n+1)}{2} - 1 \right] - (n - 2 + 1) \\ &= \frac{n^2}{2} + \frac{n}{2} - 1 - n + 1 = \frac{n^2}{2} - \frac{n}{2} = \frac{n^2}{2} + \mathcal{O}(n) \end{aligned}$$

Similarly, the total number of the additions/subtractions in the forward substitution phase will be

$$\sum_{i=2}^n (i) = \frac{n^2}{2} + \frac{n}{2} - 1 = \frac{n^2}{2} + \mathcal{O}(n)$$

The summary of the operation count of the Unit Lower triangular system $LY = B$ is given as:

	Operations		Total flops
	Multiplication/Division	Addition/Subtraction	
Unit Lower triangular system $LY = B$	$\frac{n^2}{2} - \frac{n}{2}$	$\frac{n^2}{2} + \frac{n}{2} - 1$	$n^2 - 1$

(3) The back substitution phase occurs after the solution of the lower triangular system. This phase contains two nested loops. The first loop, say i -loop (which ranges from $i = n - 1$ to 1), corresponds to $n - 1$ of the components of solution vector X . The second loop, say j -loop (which ranges from $j = i + 1$ to n), corresponds to the columns after the diagonal elements.

Each pass of i -loop will perform one subtraction and one division to obtain the value of x_i . Moreover, in each pass of i -loop, the number of both of the multiplications and additions will be $n - (i + 1) + 1$ (or simply $(n - i)$) in j -loop. Thus, in each pass of i -loop, the total number of both of the multiplications/divisions and additions/subtractions will be $(n - i + 1)$.

Hence, the total number of the multiplications/divisions in the back substitution phase will be

$$\begin{aligned} 1 + \sum_{i=1}^{n-1} (n+1-i) &= 1 + (n+1) \sum_{i=1}^{n-1} 1 - \sum_{i=1}^{n-1} i \\ &= 1 + (n+1)(n-1) - \left[\frac{(n-1)n}{2} \right] \\ &= 1 + n^2 - 1 - \frac{n^2}{2} + \frac{n}{2} = \frac{n^2}{2} + \frac{n}{2} = \frac{n^2}{2} + \mathcal{O}(n) \end{aligned}$$

Similarly, the total number of the additions/subtractions in the back substitution phase will be

$$\sum_{i=1}^{n-1} (n+1-i) = \frac{n^2}{2} + \frac{n}{2} - 1 = \frac{n^2}{2} + \mathcal{O}(n)$$

The summary of the operation count of the Upper triangular system $UX = Y$ is given as:

	Operations		Total flops
	Multiplications/Divisions	Additions/Subtractions	
Upper triangular system $UX = Y$	$\frac{n^2}{2} + \frac{n}{2}$	$\frac{n^2}{2} + \frac{n}{2} - 1$	$n^2 + n - 1$

Question 26: Perform the operation count analysis of the algorithm that involves the following phases to solve an $n \times n$ linear system:

- (1) Factorization of the coefficient matrix using the Doolittle's method
- (2) Forward substitution to solve the lower triangular system.
- (3) Back substitution to solve the upper triangular system.

(1) The factorization of the coefficient matrix A into the product of the lower triangular L and the unit upper triangular U matrices occur just after setting the inputs in the algorithm. The formulation of L and U as the factors of A contains three nested loops. The first loop, say i -loop (which ranges from $i = 1$ to n), corresponds to the i th column of L and i th row of U respectively. The second loop, say j -loop (ranges from $j = i$ to n), corresponds to the j th row of L and (ranges from $j = i + 1$ to n), corresponds to the j th column of U . The third loop, say s -loop (which ranges from $s = 1$ to $i - 1$), corresponds to the multiplication of the rows of L and columns of U .

Note that, the j -loop corresponding to row j of L ranging from i to n , the number of passes/iterations will be $n - i + 1$. Similarly, the number of passes/iterations in j -loop, corresponds to column j of U ranging from $i + 1$ to n , will be $n - (i + 1) + 1$ (or simply $(n - i)$).

Each pass of j -loop will perform one subtraction to obtain the entry l_{ji} of L . Moreover, in each pass of j -loop, the number of both of the multiplications and additions will be $(i - 1) - 1 + 1$ (or simply $(i - 1)$) in s -loop. Thus, in each pass of j -loop, the total number of multiplications/divisions will be $(i - 1)$ and the total number of additions/subtractions will be $(1 + i - 1)$ or (i) . Similarly, in each pass of j -loop, to obtain the entry u_{ij} of U , the total number of both of the multiplications and additions will be $(i - 1) + 1$ (or simply (i)).

As there are $(n - i + 1)$ passes of j -loop in each pass of i -loop, therefore there will be $(n - i + 1) \times (i - 1)$ multiplications/divisions and $(n - i + 1) \times (i)$ additions/subtractions in each pass of i -loop for the formulation of column i of L . Similarly, in each pass of i -loop, there will be $(n - i) \times (i)$ multiplications/divisions and $(n - i) \times (i)$ additions/subtractions in each pass of i -loop for the formulation of row i of U .

Hence, the total number of multiplications/divisions in n passes of i -loop for the formulation of the lower triangular matrix L will be

$$\begin{aligned}
 \sum_{i=1}^n (n - i + 1)(i - 1) &= \sum_{i=1}^n (n + 1 - i)(i - 1) \\
 &= \sum_{i=1}^n [(n + 1)i - (n + 1) - i^2 + i] = \sum_{i=1}^n [(n + 2)i - i^2 - (n + 1)] \\
 &= (n + 2) \sum_{i=1}^n i - \sum_{i=1}^n i^2 - (n + 1) \sum_{i=1}^n 1 \\
 &= (n + 2) \left[\frac{n(n + 1)}{2} \right] - \left[\frac{n(n + 1)(2n + 1)}{6} \right] - (n + 1)n \\
 &= n(n + 1) \left[\frac{n}{2} + 1 - \frac{n}{3} - \frac{1}{6} - 1 \right] = (n^2 + n) \left[\frac{n}{6} - \frac{1}{6} \right] \\
 &= \frac{n^3}{6} - \frac{n}{6} = \frac{n^3}{6} + \mathcal{O}(n)
 \end{aligned}$$

Similarly, the total number of additions/subtractions in n passes of i -loop for the formulation of the lower triangular matrix L will be

$$\sum_{i=1}^n (n-i+1)(i) = \frac{n^3}{6} + \frac{n}{6} = \frac{n^3}{6} + \mathcal{O}(n)$$

Moreover, the total number of multiplications/divisions and additions/subtractions in n passes of i -loop for the formulation of the unit upper triangular matrix U will be

$$\begin{aligned} \sum_{i=1}^n (n-i)(i) &= \sum_{i=1}^n (n-i)(i) = \sum_{i=1}^n (ni - i^2) \\ &= n \sum_{i=1}^n i - \sum_{i=1}^n i^2 = n \left[\frac{n(n+1)}{2} \right] - \left[\frac{n(n+1)(2n+1)}{6} \right] \\ &= n(n+1) \left[\frac{n}{2} - \frac{n}{3} - \frac{1}{6} \right] = (n^2 + n) \left[\frac{n}{6} - \frac{1}{6} \right] \\ &= \frac{n^3}{6} - \frac{n}{6} = \frac{n^3}{6} + \mathcal{O}(n) \end{aligned}$$

The summary of the operation count of the LU -factorization is given as:

	Operations		Total flops
	Multiplication/Division	Addition/Subtraction	
Lower Triangular Matrix L	$\frac{n^3}{6} - \frac{n}{6}$	$\frac{n^3}{6} + \frac{n}{6}$	$\frac{n^3}{3}$
Upper Triangular Matrix U	$\frac{n^3}{6} - \frac{n}{6}$	$\frac{n^3}{6} - \frac{n}{6}$	$\frac{n^3}{3} - \frac{n}{3}$
LU -factorization	$\frac{n^3}{3} - \frac{n}{3}$	$\frac{n^3}{3}$	$\frac{2n^3}{3} - \frac{n}{3}$

(2) The forward substitution phase occurs after the formulation of L and U as factors of the coefficient matrix for solving the lower triangular system. This phase contains two nested loops. The first loop, say i -loop (which ranges from $i = 2$ to n), corresponds to $n - 1$ of the components of the intermediate vector Y . The second loop, say j -loop (which ranges from $j = 1$ to $i - 1$), corresponds to the columns before the diagonal elements.

Each pass of i -loop will perform one subtraction and one division to obtain the value of y_i . Moreover, in each pass of i -loop, the number of both of the multiplications and additions will be $(i - 1) - 1 + 1$ (or simply $(i - 1)$) in j -loop. Thus, in each pass of i -loop, the total number of both of the multiplications/divisions and additions/subtractions will be $(1 + i - 1)$ or simply (i) .

Hence, the total number of the multiplications/divisions in the forward substitution phase will be

$$\begin{aligned} 1 + \sum_{i=2}^n (i) &= 1 + \sum_{i=2}^n i = 1 + \left[\frac{n(n+1)}{2} - 1 \right] \\ &= \frac{n^2}{2} + \frac{n}{2} = \frac{n^2}{2} + \frac{n}{2} = \frac{n^2}{2} + \mathcal{O}(n) \end{aligned}$$

Similarly, the total number of the additions/subtractions in the forward substitution phase will be

$$\sum_{i=2}^n (i) = \frac{n^2}{2} + \frac{n}{2} - 1 = \frac{n^2}{2} + \mathcal{O}(n)$$

The summary of the operation count of the Unit Lower triangular system $LY = B$ is given as:

	Operations		Total flops
	Multiplication/Division	Addition/Subtraction	
Lower triangular system $LY = B$	$\frac{n^2}{2} + \frac{n}{2}$	$\frac{n^2}{2} + \frac{n}{2} - 1$	$n^2 + n - 1$

(3) The back substitution phase occurs after the solution of the lower triangular system. This phase contains two nested loops. The first loop, say i -loop (which ranges from $i = n - 1$ to 1), corresponds to $n - 1$ of the components of solution vector X . The second loop, say j -loop (which ranges from $j = i + 1$ to n), corresponds to the columns after the diagonal elements.

Each pass of i -loop will perform one subtraction to obtain the value of x_i . Moreover, in each pass of i -loop, the number of both of the multiplications and additions will be $n - (i + 1) + 1$ (or simply $(n - i)$) in j -loop. Thus, in each pass of i -loop, the total number of multiplications/divisions will be $(n - i)$ and the total number of additions/subtractions will be $(n - i + 1)$.

Hence, the total number of the multiplications/divisions in the back substitution phase will be

$$\begin{aligned}\sum_{i=1}^{n-1} (n-i) &= n \sum_{i=1}^{n-1} 1 - \sum_{i=1}^{n-1} i = n(n-1) - \left[\frac{(n-1)n}{2} \right] \\ &= n^2 - n - \frac{n^2}{2} + \frac{n}{2} = \frac{n^2}{2} - \frac{n}{2} = \frac{n^2}{2} + \mathcal{O}(n)\end{aligned}$$

Similarly, the total number of the additions/subtractions in the back substitution phase will be

$$\sum_{i=1}^{n-1} (n+1-i) = \frac{n^2}{2} + \frac{n}{2} - 1 = \frac{n^2}{2} + \mathcal{O}(n)$$

The summary of the operation count of the Upper triangular system $UX = Y$ is given as:

	Operations		Total flops
	Multiplications/Divisions	Additions/Subtractions	
Upper triangular system $UX = Y$	$\frac{n^2}{2} - \frac{n}{2}$	$\frac{n^2}{2} + \frac{n}{2} - 1$	$n^2 - 1$



Chapter Summary

- A **system of linear equations** (simply called as **linear system**) is a set or collection of two or more linear equations with the same set of variables whose simultaneous solution satisfies all the equations. Precisely, a linear system can be referred to as a set of **simultaneous linear algebraic equations**.
- If $m > n$, where m is the number of equations and n is the number of unknowns, then the linear system is called **over-determined**. If $m < n$, then the linear system is called **under-determined**.
- A linear system $AX = B$ is called **homogenous** if B is a zero vector (i.e., $B = \bar{0}$), and **non-homogeneous** or **inhomogeneous** if $B \neq \bar{0}$.
- A non-homogeneous linear system $AX = B$ is called **consistent** if it has a unique solution or infinitely many solutions, and it is called **inconsistent** if it has no solution.
- If A^{-1} does not exist, then matrix A is called **singular** or **non-invertible**. If A^{-1} exists then A is called **non-singular** matrix and is **invertible**.
- If $\det(A) = 0$, then A^{-1} does not exist and the system $AX = B$ does not have a unique solution; the system either has no solution or infinitely many solutions.
- Although the steps of the algorithms for the solution of a linear system are elementary in nature, there might be certain pitfalls. This raises the need of skillful selection and use of an appropriate algorithm for obtaining the solution.
- In general, methods for the solution of linear systems (also called **linear solvers**) are evaluated based on their *accuracy*, *speed of convergence*, and computer *resource requirements* (CPU-requirements, memory requirements).
- A linear equation in two variables, say x and y , represents a **line** in xy -plane. If there exists a unique solution of the system then it is the point where the two lines intersect.
- A linear equation in three variables, say x , y , and z , represents a **plane** in xyz -space. If there exists a unique solution of such a system then it is the point where the three planes intersect.
- There are two broad categories of methods to solve linear systems: the **direct** (also called **exact**) methods and **iterative** methods. The prominent features of these two categories can be found in Question 5 (Section 6.1).
- An $n \times n$ square matrix $U = (u_{ij})$ is called the **upper triangular matrix** if $u_{ij} = 0$ whenever $i > j$. A linear system $UX = Y$ is said to be **upper triangular system** if its coefficient matrix is an upper triangular one. It has a unique solution if no diagonal element is zero (i.e., $|u_{ii}| \neq 0$, for $i = 1, 2, \dots, n$), otherwise it has either no solution or infinitely many solutions. If there is a unique solution of an upper

triangular system then the solution can easily be obtained by a so-called **back substitution** process. In analogy, the said propositions also hold for a lower triangular matrix $L = (l_{ij})$ for which $l_{ij} = 0$ whenever $i < j$. The solution of a lower-triangular system can be obtained by a similar so-called **forward substitution** process.

- To solve a linear system $AX = B$, the **Gaussian Elimination** method aims at obtaining an upper triangular system $UX = Y$, equivalent to $AX = B$. This process may be termed as **forward elimination**. The upper triangular system can then be solved by back substitution.
- To guard against the pitfalls of the **Gaussian Elimination** method, the process of **pivoting** is performed while using the method. The pivoting could be any of **partial, scaled** or **complete**.
- **Pivoting** refers to the interchanging of two rows of the augmented matrix so that the diagonal coefficient (to be used as the pivot element) is of greatest magnitude among the possible ones for the row under consideration.
- Pivoting **must be** performed if the main diagonal coefficient is zero (to make the triangular system non-singular). Pivoting **should be** performed if the magnitude of the main diagonal element is a smaller one (to prevent the propagation of the round-off error).
- The **Gauss-Jordan method** is a variant of the Gaussian Elimination method. It is based on the same elementary row operations; however, it eliminates all the elements below as well as above the pivot element (in the same column). Thus it does not produce an upper-triangular system for back-substitution; rather it obtains a diagonal matrix in which the solution vector is almost readily available.
- The **LU Factorization** or **LU Decomposition** method is another direct solver. A concise description of this method (and its variants) can be found in Question 12 (Section 6.5).
- The operation count analysis of an algorithm usually refers to the counting of the arithmetic operations involved. This is useful in determining the execution time required by the algorithm. For numerical computations, the operation count analysis is mostly considered as the counting of the **floating-point operations** (simply called as **flops**) involved in the algorithm.
- The additions/subtractions are considered to be requiring less CPU-time (being lighter operations) as compared to the multiplications/divisions. Therefore, it might be appropriate to count the two types of operations separately for the operation count analysis.



Chapter Exercises

Exercise 01: Solve the following system using the Gaussian Elimination method with back substitution.

$$\begin{aligned} 2x_1 - 3x_2 + x_3 &= -1 \\ 4x_1 + 4x_2 - 3x_3 &= 3 \\ -2x_1 + 3x_2 + x_3 &= 7 \end{aligned}$$

Exercise 02: Solve the following system using the Gaussian Elimination method with partial pivoting.

$$\begin{aligned} x_1 + x_2 + x_3 &= 6 \\ 3x_1 + 3x_2 + x_3 &= 12 \\ 2x_1 + x_2 + 5x_3 &= 20 \end{aligned}$$

Exercise 03: Solve the following system using the Gaussian Elimination method with partial pivoting and three-digit rounding arithmetic.

$$\begin{aligned} 2.5x_1 - 3x_2 + 4.6x_3 &= -1.05 \\ -3.5x_1 + 2.6x_2 + 1.5x_3 &= -14.46 \\ -6.5x_1 + -3.5x_2 + 7.3x_3 &= -17.735 \end{aligned}$$

Exercise 04: Solve the following system using the Gaussian Elimination method with scaled partial pivoting.

$$\begin{aligned} x_1 + x_2 - 2x_3 &= 3 \\ 4x_1 - 2x_2 + x_3 &= 5 \\ 3x_1 - x_2 + 3x_3 &= 8 \end{aligned}$$

Exercise 05: Solve the following system using the Gaussian Elimination method with scaled partial pivoting and four-digit rounding arithmetic.

$$\begin{aligned} 3.03x_1 - 12.1x_2 + 14x_3 &= -119 \\ -3.03x_1 + 12.1x_2 - 7x_3 &= 120 \\ 6.11x_1 - 14.2x_2 + 21x_3 &= -139 \end{aligned}$$

Exercise 06: Solve the following system using the Gaussian Elimination method with complete pivoting.

$$\begin{aligned} x_1 + 2x_2 + 2x_3 &= 1 \\ 2x_1 + 6x_2 + 10x_3 &= -2 \\ 3x_1 + 14x_2 + 28x_3 &= -11 \end{aligned}$$

Exercise 07: Solve the following system using the Gaussian Elimination method with complete pivoting and three-digit rounding arithmetic.

$$\begin{aligned} 1.012x_1 - 2.132x_2 + 3.104x_3 &= 1.984 \\ -2.132x_1 + 4.096x_2 - 7.013x_3 &= -5.049 \\ 3.104x_1 - 7.013x_2 + 0.014x_3 &= -3.895 \end{aligned}$$

Exercise 08: Solve the following system using the Gauss-Jordan method

$$\begin{aligned}x_1 + 2x_2 + x_3 &= 6 \\2x_1 + 3x_2 + 4x_3 &= 12 \\4x_1 + 3x_2 + 2x_3 &= 12\end{aligned}$$

Exercise 09: Solve the following system using the Gauss-Jordan method and three-digit rounding arithmetic.

$$\begin{aligned}0.125x_1 + 0.201x_2 + 0.401x_3 &= 2.306 \\0.375x_1 + 0.501x_2 + 0.601x_3 &= 4.806 \\0.501x_1 + 0.301x_2 + 0.001x_3 &= 2.91\end{aligned}$$

Exercise 10: Solve the following linear system $AX = B$ using the Doolittle's method.

$$\begin{aligned}x_1 + x_2 + x_3 &= 3 \\2x_1 - x_2 + 2x_3 &= 16 \\3x_1 + x_2 + x_3 &= -3\end{aligned}$$

Exercise 11: Solve the following linear system $AX = B$ using the Doolittle's method.

$$\begin{aligned}x_1 + x_2 + 2x_3 + 2x_4 &= 9 \\2x_1 + 4x_2 + 7x_3 + 3x_4 &= 25 \\-x_1 - 5x_2 - 6x_3 + 2x_4 &= -17 \\x_1 - x_2 + 3x_3 + 8x_4 &= 15\end{aligned}$$

Exercise 12: Solve the following linear system $AX = B$ using the Crout's method.

$$\begin{aligned}8x_1 + x_2 - x_3 &= 8 \\2x_1 + x_2 + 9x_3 &= 12 \\x_1 - 7x_2 + 2x_3 &= -4\end{aligned}$$

Exercise 13: Solve the following linear system $AX = B$ using the Crout's method.

$$\begin{aligned}x_1 + x_2 + 0x_3 + 3x_4 &= 9 \\2x_1 + x_2 - x_3 + x_4 &= 5 \\3x_1 - x_2 + x_3 + 2x_4 &= 6 \\-x_1 + 2x_2 + 3x_3 - x_4 &= 4\end{aligned}$$

Exercise 14: Solve the following linear system $AX = B$ using the Cholesky's method.

$$\begin{aligned}2x_1 + 3x_2 + 4x_3 &= 1 \\3x_1 + 8x_2 + 5x_3 &= 6 \\4x_1 + 5x_2 + 10x_3 &= -1\end{aligned}$$

Exercise 15: Solve the given linear system $AX = B$ using the Cholesky's method

$$\begin{aligned}4x_1 + x_2 + x_3 + x_4 &= 9 \\x_1 + 3x_2 - x_3 + x_4 &= 4 \\x_1 - x_2 + 2x_3 + 0x_4 &= 4 \\x_1 + x_2 + 0x_3 + 2x_4 &= 6\end{aligned}$$

Exercise 16: The upward velocity of a rocket at three different times after its launching are given as follows:

Time, t in (s)	Velocity, v in (m/s)
6	115.7
9	182.5
12	295.6

The velocity data is approximated by a polynomial as

$$v(t) = a_1 t^2 + a_2 t + a_3, \quad 5 \leq t \leq 12$$

Thus, the coefficients a_1 , a_2 and a_3 for the above expression are given by

$$\begin{bmatrix} 36 & 6 & 1 \\ 81 & 9 & 1 \\ 144 & 12 & 1 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \begin{bmatrix} 115.7 \\ 182.5 \\ 295.6 \end{bmatrix}$$

Find the values of a_1 , a_2 and a_3 using a linear solver. Then, calculate the velocity at $t = 7, 8, 10$, and 11 seconds.

Exercise 17: A factory produces three products, say Prod1, Prod2, and Prod3, by using three kinds of raw materials, say Raw1, Raw2, and Raw3. The units of each of the raw materials needed to produce one unit of each of the products are shown the table below.

Sectors	Raw1	Raw2	Raw3
Prod1	5	3	1
Prod2	4	4	3
Prod3	2	1	3

If 335 units of Raw1, 532 units of Raw2, and 440 units of Raw3 are available, then how much each of the three products can be produced.

Hint for the Solution:

Assume that x_1 , x_2 and x_3 represent the quantities of the products: Prod1, Prod2, and Prod3, respectively. The problem can be represented by a linear system whose solution would provide the required values.

$$5x_1 + 4x_2 + 2x_3 = 335$$

$$3x_1 + 4x_2 + x_3 = 532$$

$$x_1 + 3x_2 + 3x_3 = 440$$

Exercise 18: Assume that the economy of a country depends on the three sectors: Food, Cloth, and House. The production of one unit of each of these needs certain units of each of these sectors, as shown in the following table:

Sectors	Food Units Needed	Cloth Units Needed	House Units Needed
Food	0.45	0.18	0.15
Cloth	0.25	0.27	0.07
House	0.30	0.40	0.45

The consumer demand is as in the table below:

Sector	worth in billion rupees
Food	220
Cloth	185
House	550

For satisfying the above demands, what total output is required from each of the sectors.

Hint for the Solution in MATLAB: Assume that x_1 , x_2 and x_3 represent the total outputs in units from Food, Cloth and House sectors, respectively. The problem can be represented by a linear system whose solution would provide the required values.

Exercise 19: A bakery produces three products: Cake, Pastry, and Muffin. It uses three kinds of materials: Flour, Milk, and Sugar. The units of each of the raw materials needed to produce one unit of each of the bakery products are shown the table below.

Product ->	Cake	Pastry	Muffin
Flour	6	5	3
Milk	4	5	2
Sugar	2	3	3

If 347 units of Flour, 604 units of Milk, and 502 units of Sugar are available, then how much each of the three products can be produced.

Exercise 20: Pivoting is necessary with the Gaussian elimination if

- (A) the coefficient matrix is singular (B) the linear system is homogenous
 (C) the linear system is ill conditioned (D) None of above

Exercise 21: Cholesky decomposition for a linear system is not possible, if

- (A) the linear system is ill conditioned (B) the linear system is homogenous
 (C) the coefficient matrix is asymmetric (D) None of above

Iterative Linear Solvers

Corridor I: BASICS

Let's plan it

- 7.1 Vector Norms and Distances
- 7.2 Convergence Criteria for Linear Solvers
- 7.3 Basic Methods
 - 7.3.1 The Jacobi Method
 - 7.3.2 The Gauss-Seidel Method
 - 7.3.3 The SOR Method

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)

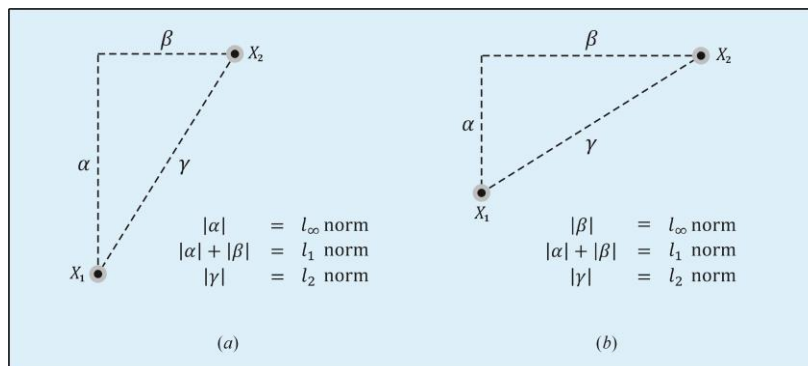


Fig. (7.4): Explanation of the different types of distances between the two vectors.

Corridor II: ANALYSIS

Let's think deep

- 7.4 Matrix Norms and Conditioning
- 7.5 Iteration Matrix and Matrix Form of a Solver

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



Corridor III: PROGRAMMING ARCADE

Let's do it

- 7.6 Algorithms and Implementations
 - The Jacobi Method
 - Modification in the Jacobi Method's algorithm for the Gauss-Seidel Method
 - Modification in the Jacobi Method's algorithm for the SOR Method

To cross-check the results/output of the computer programs you would execute, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



7.6 Algorithms and Implementations

Question 22: Write down an algorithm (pseudo code) to solve a linear system using the Jacobi method.

Algorithm: To solve $AX = B$, given an initial approximation $X^{(0)}$.

INPUTS: $\left\{ \begin{array}{l} n: \text{an integer as the number of equations and unknowns} \\ A = (a_{ij}), 1 \leq i, j \leq n: \text{a real valued square matrix as the coefficient matrix} \\ B = [b_1, b_2, \dots, b_n]^T: \text{a real valued vector as the vector of right hand side constants} \\ X = [x_1, x_2, \dots, x_n]^T: \text{a real valued vector (having initial approximation, } X^{(0)}) \\ TOL: \text{a real value as the error tolerance} \\ N: \text{an integer as the maximum number of iterations} \end{array} \right.$

OUTPUT: $\left\{ \begin{array}{l} X = [x_1, x_2, \dots, x_n]^T: \text{a real valued vector as the approximate solution} \\ \text{(either on convergence, or on completing } N \text{ iterations – which ever happens first)} \end{array} \right.$

Step 1 Receive the inputs as stated above

Step 2 for $k = 1, 2, 3, \dots, N$ perform steps 3-6

Step 3 for $i = 1, 2, \dots, n$ Set $x_{p_i} = x_i$ $\left\{ \begin{array}{l} XP = [xp_1, xp_2, \dots, xp_n]^T \text{ is to keep a copy of present} \\ \text{approximation } X, \text{ because } X \text{ is going to be updated} \end{array} \right.$

Step 4 for $i = 1, 2, \dots, n$ (compute the components of solution vector X)

$$\left. \begin{array}{l} sum = 0 \\ \text{for } j = 1, 2, \dots, n \\ \quad \text{if } (j \neq i) \quad sum = sum + a_{ij} \times XP_j \\ x_i = \frac{[b_i - sum]}{a_{ii}} \end{array} \right\} \left(x_i^{(k)} = \frac{1}{a_{ii}} \left[b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(k-1)} \right] \right)$$

Step 5 Compute $err = \|X - XP\|$ (or $err = \|X - XP\|/\|X\|$) Here $\|\cdot\|$ is any suitable norm.

Step 6

$\left. \begin{array}{l} \text{if } (err < TOL) \text{ then} \\ \quad \text{Exit/Break the loop} \end{array} \right\}$ This means that the consecutive approximations are nearly the same. Therefore, stop iterations.

end for loop of Step 2 (Go to Step 3)

Step 7 Print the output: $X = [x_1, x_2, \dots, x_n]^T$

if $(err < TOL)$ OUTPUT ('The desired accuracy achieved; Solution converged.')
 else OUTPUT ('The number of iterations exceeded the maximum limit.')

STOP.

Question 23: What modification a programmer needs to make in the algorithm (pseudo code) of the Jacobi method (as given in the answer of Question 22) to convert it into the Gauss-Seidel method for solving a linear system.

The algorithm (pseudo code) of the Jacobi method (as given in the answer of Question 22) can be converted into the algorithm of the Gauss-Seidel method simply by replacing its Step 4 with the following:

Step 4 for $i = 1, 2, \dots, n$ (compute the components of solution vector \mathbf{X})

$$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } j = 1, 2, \dots, n \\ \quad \text{if } (j \neq i) \quad \text{sum} = \text{sum} + a_{ij} \times x_j \\ x_i = \frac{[b_i - \text{sum}]}{a_{ii}} \end{array} \right\} \left(\begin{array}{l} \text{sum} = \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} + \sum_{j=i+1}^n a_{ij}x_j^{(k-1)} \\ x_i^{(k)} = \frac{1}{a_{ii}} [b_i - \text{sum}] \end{array} \right)$$

■

Question 24: What modification a programmer needs to make in the algorithm (pseudo code) of the Jacobi method (as given in the answer of Question 22) to convert it into the Gauss-Seidel method with over-relaxation (i.e., the SOR method) for solving a linear system.

The algorithm (pseudo code) of the Jacobi method (as given in the answer of Question 22) can be converted into the algorithm of the SOR method simply by taking one more input:

$WF = 1.3$: a real value as the over – relaxation / weighting factor

And then replacing Step 4 with the following:

Step 4

for $i = 1, 2, \dots, n$ (compute the components of solution vector \mathbf{X})

$$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } j = 1, 2, \dots, n \\ \quad \text{if } (j \neq i) \quad \text{sum} = \text{sum} + a_{ij} \times x_j \\ x_i = \frac{[b_i - \text{sum}]}{a_{ii}} \end{array} \right\} \left(\begin{array}{l} \text{sum} = \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} + \sum_{j=i+1}^n a_{ij}x_j^{(k-1)} \\ x_i^{(k)} = \frac{1}{a_{ii}} [b_i - \text{sum}] \end{array} \right)$$

$$x_i = WF \times x_i + (1 - WF)XP_i \quad (\text{apply over – relaxation})$$

■

Problem 17: Write a C++ program to solve the following linear system using the Jacobi method. Take initial approximate solution as: $X^{(0)} = [0, 0, 0]^T$. The iterations of the method should stop when either the approximation is accurate within 10^{-6} , or the number of iterations exceeds 200, whichever happens first.

$$5x_1 + 3x_2 + 2x_3 = 17, \quad 3x_1 + 4x_2 - x_3 = 8, \quad -x_1 + x_2 - 3x_3 = -8$$

```

1  # include <iostream>
2  # include <cmath>
3  using namespace std ;
4
5  # define n 3                                // number of unknowns
6  # define TOL 0.000001                       // error tolerance
7  # define N 200                              // maximum number of iterations
8
9  int main ( )
10 {
11
12     int i , j , k ;
13     double a[n][n] = {{5, 3, 2}, {3, 4, -1}, {-1, 1,-3}} ; // coefficient matrix
14     double b[n] = {17, 8, -8} ;                // right-hand side constants
15     double x[n] ;                             // the solution vector
16     double xp[n] , sum , err ;
17
18     cout<<"The Gauss-Jacobi Method for solving a system of "<<n<<" unknowns." ;
19
20     for ( i=0 ; i<n ; i++)
21         x[i] = 0.0 ;                          // setting initial approximation as zero vector
22
23     //----- Processing Section -----//
24
25     for ( k=1 ; k<=N ; k++ )
26     {
27         for ( i=0 ; i<n ; i++)
28             xp[i] = x[i] ;
29
30         for ( i=0 ; i<n ; i++)
31         {
32             sum = 0 ;
33             for ( j=0 ; j<n ; j++)
34                 if ( j != i ) sum = sum + a[i][j] * xp[j] ;
35
36             x[i] = ( b[i] - sum ) / a[i][i] ;
37         }
38
39         sum = 0.0 ;
40         for ( i=0 ; i<n ; i++)
41             sum = sum + ( ( x[i] - xp[i] ) * ( x[i] - xp[i] ) ) ;
42         err = sqrt(sum) ;
43
44         if ( err < TOL )    break ;
45     }

```

$$\sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(k-1)}$$

$$x_i^{(k)} = \frac{1}{a_{ii}} [b_i - sum]$$

Computing
 $l_2 - norm$

```

46
47 //----- Output Section -----//
48
49 cout << " The latest approximate solution vector is given by : " << endl;
50 for ( i=0 ; i<n ; i++)
51     cout << x[i] << "\t" ;
52
53 if ( err < TOL)
54     cout<< "\nThe desired accuracy achieved; Solution converged." ;
55 else
56     cout<< "\nThe number of iterations exceeded the maximum limit." ;
57
58 }

```

Remark: Replacing $xp[j]$ by $x[j]$ in line 34 of the C++ program in Problem 17 would convert the program for the Gauss-Seidel method, because it would then correspond to computing

$$sum = \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} + \sum_{j=i+1}^n a_{ij}x_j^{(k-1)}$$

Remark: In the program of Problem 17, the code segment of lines 45-47 can be placed just before line 45 to print the latest result on completion of each of the iterations.

Problem 19: Write a C++ program to solve the following linear system using the Gauss-Seidel method with over-relaxation (the SOR method). Take initial approximate solution as: $X^{(0)} = [0, 0, 0]^T$ and over-relaxation factor as 1.2. The iterations of the method should stop when either the approximation is accurate within 10^{-6} , or the number of iterations exceeds 200, whichever happens first.

$$5x_1 + 3x_2 + 2x_3 = 17, \quad 3x_1 + 4x_2 - x_3 = 8, \quad -x_1 + x_2 - 3x_3 = -8$$

```

1 # include <iostream>
2 # include <cmath>
3 using namespace std ;
4
5 # define n 3 // number of unknowns
6 # define TOL 0.000001 // error tolerance
7 # define N 200 // maximum number of iterations
8 # define WF 1.2 // over-relaxation factor
9
10 int main ( )
11 {
12

```

```

13  int i , j , k ;
14  double a[n][n] = {{5, 3, 2}, {3, 4, -1}, {-1, 1,-3}} ; // coefficient matrix
15  double b[n] = {17, 8, -8} ; // right-hand side constants
16  double x[n] ; // the solution vector
17  double xp[n] , sum , err ;
18
19  cout<<"The Gauss-Seidel method with over-relaxation for solving a system.\n" ;
20
21  for ( i=0 ; i<n ; i++)
22      x[i] = 0.0 ; // setting initial approximation as zero vector
23
24  //----- Processing Section -----//
25
26  for ( k=1 ; k<=N ; k++ )
27  {
28      for ( i=0 ; i<n ; i++)
29          xp[i] = x[i] ;
30
31      for ( i=0 ; i<n ; i++)
32      {
33          sum = 0 ;
34          for ( j=0 ; j<n ; j++)
35              if ( j != i ) sum = sum + a[i][j] * x[j] ;
36
37
38          x[i] = ( b[i] - sum ) / a[i][i] ;
39
40
41          x[i] = WF * x[i] + ( 1 - WF ) * xp[i] ;
42
43      }
44
45      sum = 0.0 ;
46      for ( i=0 ; i<n ; i++)
47          sum = sum + ( ( x[i] - xp[i] ) * ( x[i] - xp[i] ) ) ;
48      err = sqrt(sum) ;
49
50      if ( err < TOL ) break ;
51  }
52
53  //----- Output Section -----//
54
55  cout << " The latest approximate solution vector is given by : " << endl;
56  for ( i=0 ; i<n ; i++)
57      cout << x[i] << "\t" ;
58
59  if ( err < TOL )
60      cout<< "\nThe desired accuracy achieved; Solution converged." ;
61  else
62      cout<< "\nThe number of iterations exceeded the maximum limit." ;

```

$$\sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} \times (\text{latest value of } x_j)$$

$$x_i^{(k)} = \frac{1}{a_{ii}} [b_i - \text{sum}]$$

$$x_i = WF \times x_i + (1 - WF)XP_i$$

Computing
 l_2 - norm

```
63 }
64 }
```

Remark: Note that setting the weighting factor of over-relaxation (WF) as 1.0 in the solutions of Problem 19 would make the programs for the Gauss-Seidel method. ■

Remark: In the program of Problem 19, the code segment of lines 55-57 can be placed just before line 51 to print the latest result on completion of each of the iterations. ■

Remark: The C++ programs in Problem 17 and 19 can be modified to receive the linear system at the execution time (instead of fixing in the code). For this, the lines 13-14 in the program of Problem 17 and the lines 14-15 in the program of Problem 19 should be replaced by the following code segment:

```
cout<<"\nEnter the coefficient matrix row-wise: "<<n<<" unknowns<<endl;
for ( i=0 ; i<n ; i++)
    for ( j=0 ; j<n ; j++)
        cin >> a[i][j] ;

cout<<"Enter the elements of constant vector B "<<endl ;
for ( i=0 ; i<n ; i++)
    cin >> b[i] ;
for i = 1:n
    b(i) = input('Enter the element of constant vector: ');
end
```

■ ■ ■

Chapter Summary

- The **norm of a vector** is a real-valued function that provides a measure of “size”, “length”, or “magnitude” of the vector. Let \mathbb{R} denotes the set of real numbers, and \mathbb{R}^n denotes the space of n -dimensional real-valued column vectors. A norm of a vector on \mathbb{R}^n is a function, $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$, with the following properties,
 1. $\|X\| \geq 0$, for all $X \in \mathbb{R}^n$
 2. $\|X\| = 0$, if and only if $X = \mathbf{0}$ in \mathbb{R}^n
 3. $\|\alpha X\| = |\alpha| \|X\|$, for all $\alpha \in \mathbb{R}$ and $X \in \mathbb{R}^n$
 4. $\|X + Y\| \leq \|X\| + \|Y\|$, for all $X, Y \in \mathbb{R}^n$

- The vector norm definitions, as well as the concerning illustrations, can be found in Question 01 (Section 7.1).
- The norm of a vector gives a measure for the distance between an arbitrary vector and the zero vector, just as the absolute value of a real number is its distance from 0.
- The **distance between two vectors** is defined as the norm of the “difference vector” of the two vectors, just as the distance between two real numbers is the absolute value of their difference. The definitions of different vector distances, as well as the concerning illustrations, can be found in Question 02 (Section 7.1).
- To determine the convergence of an iterative solution, the norm of the difference vector of every two consecutive approximations is ensured to be smaller than a pre-specified error tolerance τ , i.e.,

$$\|X^{(k)} - X^{(k-1)}\| < \tau$$

- A square matrix, say $A = (a_{ij})_{n \times n}$, is said to be **diagonally dominant** if, for $i = 1, 2, \dots, n$

$$|a_{ii}| \geq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|,$$

- A linear system is said to be diagonally dominant if its coefficient matrix is **diagonally dominant** (i.e., the magnitude of the diagonal entry in a row is greater than or equal to the sum of the magnitudes of all other entries in that row).
- If “ \geq ” is replaced by “ $>$ ” in the above equation, then A is said to be **strictly diagonally dominant**. A strictly diagonally dominant matrix is always non-singular.
- If a linear system is not diagonally dominant, then a rearrangement of its rows might make it diagonally dominant.
- The Gauss-Jacobi, Gauss-Seidel, and SOR methods must converge if the linear system to be solved is diagonally dominant.
- Suppose that $AX = B$ is a $n \times n$ linear system to be solved such that $A = (a_{ij})_{n \times n}$ is the coefficient matrix, $B = (b_i)_{n \times 1}$ is the vector of right-hand side constants, and $X = (x_i)_{n \times 1}$ is the vector of unknowns.

➤ The **Jacobi method** can be written in a compact form as

$$x_i^{(k)} = \frac{1}{a_{ii}} \left[b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(k-1)} \right], \quad \text{for } i = 1, 2, \dots, n$$

- The **Gauss-Seidel method** can be written in a compact form as

$$x_i^{(k)} = \frac{1}{a_{ii}} \left[b_i - \left(\sum_{j=1}^{i-1} a_{ij} x_j^{(k)} + \sum_{j=i+1}^n a_{ij} x_j^{(k-1)} \right) \right], \quad \text{for } i = 1, 2, \dots, n$$

- The **successive over-relaxation (SOR) method** can be written in a compact form as

$$\bar{x}_i^{(k)} = \frac{1}{a_{ii}} \left[b_i - \left(\sum_{j=1}^{i-1} a_{ij} x_j^{(k)} + \sum_{j=i+1}^n a_{ij} x_j^{(k-1)} \right) \right], \quad \text{for } i = 1, 2, \dots, n$$

$$x_i^{(k)} = \omega \bar{x}_i^{(k)} + (1 - \omega) x_i^{(k-1)} \quad (\text{for } 1 \leq \omega \leq 2, \text{ usually the best is around } 1.2)$$

Here $k = 1, 2, 3, \dots$, represents the iterations and $x_i^{(k)}$ represents the k th approximation of the i th unknown. The iterative procedure is started with an initial approximation vector $X^{(0)} = [x_1^{(0)}, x_2^{(0)}, x_3^{(0)}, \dots, x_n^{(0)}]^T$ and produces a sequence of successive approximations $\{X^{(k)}\}_{k=1}^{\infty}$, such that $X^{(k)} = [x_1^{(k)}, x_2^{(k)}, x_3^{(k)}, \dots, x_n^{(k)}]^T$. The sequence is anticipated to refine/improve the approximate solution gradually and ultimately converge to the exact solution vector (theoretically). In practice, the iterations of the method are stopped when a sufficient level of accuracy is achieved.

- The **norm of a matrix** is a real-valued function that provides a measure of “size”, “length”, or “magnitude” of the matrix. Let \mathbb{R} denotes the set of real numbers, and \mathbb{M}^n denotes the set of $n \times n$ real-valued matrices. The norm of a matrix on \mathbb{M}^n is a function, $\|\cdot\| : \mathbb{M}^n \rightarrow \mathbb{R}$, with the following properties,
 1. $\|A\| \geq 0$, for all $A \in \mathbb{M}^n$
 2. $\|A\| = 0$, if and only if $A = \mathbf{0}$ in \mathbb{M}^n
 3. $\|\alpha A\| = |\alpha| \|A\|$, for all $\alpha \in \mathbb{R}$ and $A \in \mathbb{M}^n$
 4. $\|A + B\| \leq \|A\| + \|B\|$, for all $A, B \in \mathbb{M}^n$
 5. $\|AB\| \leq \|A\| \|B\|$, for all $A, B \in \mathbb{M}^n$
- The matrix norm definitions can be found in Question 11 (Section 7.4).
- The **distance between two matrices** A and B with respect to a certain norm $\|\cdot\|$ is defined as the norm of the “difference matrix” of the two matrices, i.e., $\|A - B\|$.
- The **condition number** of a non-singular matrix A with respect to a matrix norm $\|\cdot\|$ is defined as

$$\mathcal{K}(A) = \|A\| \|A^{-1}\|, \quad (\text{and } \mathcal{K}(A) \geq 1)$$

- The **condition number of a linear system** is the condition number of its coefficient matrix.

- A computational problem is called **ill-conditioned** (or ill-posed) if small changes in the data (the input) cause large changes in the solution (the output). On the other hand, a problem is called **well-conditioned** (or well-posed) if small changes in the data cause only small changes in the solution.
- The main issue while solving an ill-conditioned problem is that the round-off errors can cause production of wide range worthless solutions (which appear to be original ones because they approximately satisfy the given problem). Therefore, minimizing the round-off errors becomes more relevant for the ill-conditioned problems.
- If $AX = B$ is an **ill-conditioned linear system** then the solution of its perturbed system (the one which is obtained by making small changes in the original system, either through small changes in A , or in B) is much different from that of the original linear system. In that case, the matrix A is said to be an ill-conditioned matrix. The determinant of an ill-conditioned matrix A is usually close to zero (NOT the zero). Remind that if the determinant is exactly zero then a relevant linear system $AX = B$ has either no solution, or an infinite number of solutions.
- There is no strict line between the well-conditioning and ill-conditioning of a system, as these concepts are qualitative. A linear system whose condition number (i.e., the condition number of its coefficient matrix) is close to 1 is well-conditioned, whereas a condition number significantly larger than 1 indicates that the linear system is ill-conditioned. If the condition number is below 100, it is usually not a reason for concern. However, a condition number of more than 100 calls for caution. It may be noted that a coefficient matrix, having magnitudes of diagonal elements larger than that of other elements in each of the rows, indicates well-conditioning of the linear system.
- In general, an iterative linear solver involves a process that converts an $n \times n$ system $AX = B$ into an equivalent system of the form $X = TX + C$ for some fixed matrix T and vector C . After the initial vector $X^{(0)}$ is selected, the sequence of approximate solution vectors, $X^{(1)}$, $X^{(2)}$, $X^{(3)}$, \dots , is generated by computing

$$X^{(k)} = TX^{(k-1)} + C, \quad \text{for } k = 1, 2, 3, \dots$$

The matrix T is called the **iteration matrix** of the iterative method, and the relation is called the **matrix form** of the iterative method.

- The iterative linear solvers for which the iteration matrix remains unchanged (or fixed) during the iterative process are said to be **stationary solvers**, whereas the iterative linear solvers for which the iteration matrix changes from iteration to iteration are referred to as **non-stationary solvers**.
- Examples of stationary solvers include simple methods like the Jacobi, Gauss-Seidel, and SOR methods. Examples of the non-stationary solvers include more sophisticated methods like the

Krylov subspace methods: especially, Conjugate Gradient (CG) methods, Minimal Residual methods (especially GMRES), and many more.

■■■

Chapter Exercises

Exercise 01: Workout first three iterations of (i) the Jacobi method, (ii) the Gauss-Seidel method, and (iii) the Gauss-Seidel method with successive over-relaxation factor $\omega = 1.2$ and $\omega = 1.5$ for solving the following systems for any initial approximation. Perform computations with a precision of 4 decimal digits, at least. Assume the error tolerance as 0.0001.

(a)

$$\begin{aligned} x_1 - 0.25x_2 - 0.25x_3 &= 9 \\ -0.25x_1 + x_2 - 0.25x_3 &= 4 \\ -0.25x_1 - 0.25x_2 + x_3 &= -1 \end{aligned}$$

(b)

$$\begin{aligned} 4x_1 + x_2 - x_3 + x_4 &= 2.5 \\ x_1 + 4x_2 - x_3 - x_4 &= 0.5 \\ -x_1 - x_2 + 5x_3 + x_4 &= 5 \\ x_1 - x_2 + x_3 + 3x_4 &= 4 \end{aligned}$$

(c)

$$\begin{aligned} 2x_1 - x_2 + x_3 &= -3 \\ 2x_1 + 4x_2 + 2x_3 &= 8 \\ -x_1 - x_2 + 2x_3 &= 1 \end{aligned}$$

(d)

$$\begin{aligned} x_1 - 0.25x_2 - 0.25x_3 + 0x_4 &= 11 \\ -0.25x_1 + x_2 + 0x_3 - 0.25x_4 &= 7 \\ -0.25x_1 + 0x_2 + x_3 - 0.25x_4 &= 3 \\ 0x_1 - 0.25x_2 - 0.25x_3 + x_4 &= -1 \end{aligned}$$

(e)

$$\begin{aligned} 0.2x_1 + 0.3x_2 + 0x_3 &= 0.1 \\ 0.3x_1 + 0x_2 + 0.2x_3 &= 0.1 \\ 0x_1 + 0.2x_2 + 0.3x_3 &= 0.8 \end{aligned}$$

(f)

$$\begin{aligned} 8x_1 + 4x_2 + 0x_3 + 0x_4 &= 10 \\ 4x_1 + 12x_2 + 2x_3 + 0x_4 &= 12 \\ 0x_1 + 2x_2 + 7x_3 + 2.5x_4 &= 9.25 \\ 0x_1 + 0x_2 + 2.5x_3 + 4.5x_4 &= 4.75 \end{aligned}$$

■■■

Eigenvalues and Eigenvectors

Corridor I: BASICS

Let's plan it

- 8.1 Basic Definitions and Concepts
- 8.2 General Approach of Finding Eigenvalues and Eigenvectors
- 8.3 Some Numerical Methods for Eigenvalues
 - The Power Method
 - The Householder Method
 - The QR Factorization Method
 - The Sturm Method

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)

Corridor II: ANALYSIS

Let's think deep

- 8.4 Further Discussions
 - The Power Theorem
 - The Gerschgorin Circle Theorems
 - The Singular Value Decomposition (SVD)

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)

Corridor III: PROGRAMMING ARCADE

Let's do it

8.5 Algorithms and Implementations The Power Method

To cross-check the results/output of the computer programs you would execute, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



8.5 Algorithms and Implementations

Question 12: Write down an algorithm (pseudo code) to find dominant eigenvalue and a corresponding eigenvector of a matrix using the Power method.

Algorithm: To approximate the dominant eigenvalue and associated eigenvector of an $n \times n$ matrix A , given a nonzero normalized vector X (i.e., having 1 as the largest component) as the initial approximation.

INPUTS: $\left\{ \begin{array}{l} n: \text{an integer as the length of the vector } X \\ X = [x_1, x_2, \dots, x_n]^T: \text{a real valued vector (as a normalised initial approximation)} \\ A = (a_{ij}), 1 \leq i, j \leq n: \text{a real valued square matrix whose eigenvalue is to be obtained} \\ \text{TOL: a real value as the tolerance} \\ N: \text{an integer as the maximum number of iterations} \end{array} \right.$

OUTPUT: $\left\{ \begin{array}{l} B: \text{a real value as the approximate eigenvalue} \\ X = [x_1, x_2, \dots, x_n]^T: \text{a normalized vector as the eigenvector corresponding to } B \end{array} \right.$

Step 1 Receive the inputs as stated above

Step 2 for $k = 1, 2, 3, \dots, N$ perform steps 3-6

Step 3 for $i = 1, 2, \dots, n$ Set $xp_i = x_i$ $\left\{ \begin{array}{l} XP = [xp_1, xp_2, \dots, xp_n]^T \text{ is to keep a copy of present} \\ \text{(approximation } X, \text{ because } X \text{ is going to be updated)} \end{array} \right.$

Step 4 (Compute the vector such that $X^{(k)} = AX^{(k-1)}$)

for $i = 1, 2, \dots, n$

$\left. \begin{array}{l} \text{sum} = 0 \\ \text{for } j = 1, 2, \dots, n \\ \quad \text{sum} = \text{sum} + a_{ij} \times xp_j \\ x_i = \text{sum} \end{array} \right\} \left(x_i^{(k)} = \sum_{j=1}^n a_{ij} x_j^{(k-1)} \right)$

Step 5 (Approximate the eigenvalue B and normalize the vector X)

$\left. \begin{array}{l} \text{set } r = 1 \\ \text{for } i = 1, 2, \dots, n \\ \quad \text{if } (|x_i| > |x_r|) \text{ } r = i \\ \text{set } B = x_r \end{array} \right\} \left(\begin{array}{l} \text{Finding the element of } X \text{ with} \\ \text{the largest absolute value} \\ \text{and then setting it as } B \end{array} \right)$

for $i = 1, 2, \dots, n$

$x_i = x_i / B$ (Normalizing the vector X)

Step 6

$\left. \begin{array}{l} \text{if } (err < TOL) \text{ then} \\ \quad \text{Exit/Break the loop} \end{array} \right\} \left. \begin{array}{l} \text{This means that the consecutive} \\ \text{approximations are nearly the same.} \\ \text{Therefore, stop iterations.} \end{array} \right.$

end for loop of Step 2 (Go to Step 3)

Step 9 Print the output: eigenvalue B , and eigenvector $X = [x_1, x_2, \dots, x_n]^T$

if $(err < TOL)$ OUTPUT ('The desired accuracy achieved; Solution converged.')

else OUTPUT ('The number of iterations exceeded the maximum limit.')

STOP.

Problem 10: Write a C++ program to find the dominant eigenvalue of the following matrix using the Power method. For simplification, specify the matrix within the program. Take $X^{(0)} = [1, 1, 1]^T$ as the initial approximation. The iterations of the method should stop when either the approximation is accurate within 10^{-5} , or the number of iterations exceeds 100, whichever happens first.

$$A = \begin{bmatrix} 4 & 1 & 0 \\ 2 & 5 & 0 \\ 7 & 2 & 1 \end{bmatrix}$$

```

1  # include <iostream>
2  # include <cmath>
3  using namespace std ;
4
5  # define n 3                               // number of components
6  # define TOL 0.00001                       // error tolerance
7  # define N 100                             // maximum number of iterations
8
9  int main ( )
10 {
11     char quit ;
12     int i , j , r , k ;
13     double a[n][n] = { {4,1,0} , {2,5,0} , {7,2,1} } ;           // the problem matrix
14     double x[n] = { 1, 1, 1 } ;           // initial approx. to the dominant eigenvector
15     double xp[n] , sum , sum1 , B , err ;
16
17 //----- Processing Section -----//
18
19     for ( k=1 ; k<=N ; k++ )
20     {
21
22         for ( i=0 ; i<n ; i++ )
23             xp[i] = x[i] ;
24
25         // Computing the vector X^(k) = A * X^(k-1)
26         for ( i=0 ; i<n ; i++ )
27         {
28             sum = 0.0 ;
29             for ( j=0 ; j<n ; j++ )
30                 sum = sum + a[i][j] * xp[j] ;
31             x[i] = sum ;
32         }
33
34         // Approximating the eigenvalue B and normalizing the vector X
35         r = 0 ;
36         for ( i= 1 ; i<n ; i++ )
37             if ( abs(x[i]) > abs(x[r]) )    r = i ;
38
39         B = x[r] ;
40
41         for ( i=0 ; i<n ; i++ )
42             x[i] = x[i] / B ;

```

$$\left(x_i^{(k)} = \sum_{j=1}^n a_{ij} x_j^{(k-1)} \right)$$

```

43
44 // Computing the error as L2-norm
45 sum1 = 0.0 ;
46 for (i=0 ; i<n ; i++ )
47     sum1 = sum1 + ( x[i] -xp[i] ) * ( x[i] - xp[i] ) ;           Computing
48     err = sqrt( sum1 ) ;                                       l2 - norm
49
50     if ( err < TOL )    break ;
51
52 }
53
54 //----- Output Section -----//
55
56 cout<< "\nThe approximate dominant eigenvalue is " << B ;
57
58 cout<< "\nThe approximate corresponding eigenvector is " << endl ;
59 for ( i=0 ; i<n ; i++ )    cout<< x[i] << " " ;
60
61 if( err<TOL )
62     cout<< "\nThe desired accuracy achieved; Solution converged" ;
63 else
64     cout<< "\nThe number of iterations exceeded the maximum limit." ;
65
66 return 0 ;
67
68 }

```

Remark: In the program of Problem 10, the code segment of lines 56-59 can be placed just before line 52 to print the latest results on completion of each of the iterations.

Remark: The C++ program in Problem 10 can be modified to receive the square matrix and the initial approximation of the Eigenvector at the execution time (instead of fixing in the code). For this, the code segment at lines 13 and 14 in the program of Problem 10 should be replaced by the following code segment:

```

cout<<"\nEnter the matrix row-wise: "<<endl;
for ( i=0 ; i<n ; i++)
    for ( j=0 ; j<n ; j++)
        cin >> a[i][j] ;

cout<<"Enter the initial approximation:"<<endl ;
for ( i=0 ; i<n ; i++)
    cin >> b[i] ;

```

Chapter Summary

- An **eigenvalue** of a square matrix $A = (a_{ij})_{n \times n}$ is a number λ such that the vector equation

$$AX = \lambda X$$

- has a non-zero solution vector X . The solution vector X is then called an **eigenvector** of the matrix A corresponding to the eigenvalue λ . The set of all eigenvalues of a matrix is called the **spectrum** of the matrix. An eigenvalue is also called a *characteristic value* or *latent root*. Likewise, an eigenvector is also called a *characteristic vector* or *latent vector*.

- A concise account of the results and techniques relevant to the eigenvalues and eigenvectors is given in Section 8.1.

- The theorem of the Power method: Suppose that an $n \times n$ matrix A has n eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$ and associated n linearly independent eigenvectors, V_1, V_2, \dots, V_n . Further, suppose that $X^{(0)}$ is a normalized vector (i.e., a vector having maximum absolute value as 1) in the space of the said eigenvectors. The sequence of normalized vectors $\{X^{(k)}\}_{k=1}^{\infty}$ and the sequence of scalars $\{\beta_k\}_{k=1}^{\infty}$ generated recursively by

$$X^{(k)} = \frac{1}{\beta_k} Y^{(k)},$$

where $Y^{(k)} = AX^{(k-1)}$, and $\beta_k = y_r^{(k)}$ such that $|y_r^{(k)}| = \|Y^{(k)}\|_{\infty}$,

will converge to the dominant eigenvector and eigenvalue, respectively.

- In the Power method, both the sequences of the scalars $\{\beta_k\}_{k=1}^{\infty}$ and the normalized vectors $\{X^{(k)}\}_{k=1}^{\infty}$ converge linearly to the dominant eigenvalue λ_1 and a corresponding eigenvector V_1 , respectively. Thus, the order of convergence of the Power method is linear.
- Aitken's Δ^2 method offers a technique for accelerating the convergence of any sequence that is linearly convergent. Using a given sequence, say $\{\beta_k\}_{k=1}^{\infty}$, which converges linearly to λ_1 , another sequence $\{\hat{\beta}_k\}_{k=1}^{\infty}$ (that also converges to λ_1 with possibly improved convergence rate) is constructed by using the Aitken's Δ^2 process as:

$$\hat{\beta}_k = \beta_k - \frac{(\beta_{k+1} - \beta_k)^2}{\beta_{k+2} - 2\beta_{k+1} + \beta_k} = \beta_k - \frac{(\Delta\beta_k)^2}{\Delta^2\beta_k}, \quad \text{for } k = 0, 1, 2, \dots$$

- Suppose that λ is a non-zero eigenvalue of a square matrix A and X is an eigenvector corresponding to λ . Then, $1/\lambda$ is an eigenvalue of A^{-1} and the same X is an eigenvector corresponding to $1/\lambda$. Thus, the reciprocal of all the non-zero eigenvalues of a square matrix A are the eigenvalues of A^{-1} (having the same set of eigenvectors). Hence, the largest of the absolute eigenvalues of A is the smallest of the

eigenvalues of A^{-1} (and vice-versa). Thus, the Power method can be used to obtain the largest eigenvalue of A^{-1} and then taking its reciprocal gives the smallest eigenvalue of A .

■ ■ ■

Chapter Exercises

Exercise 01: Find all the eigenvalues and eigenvectors of the following matrices using the characteristic equations. Also find the spectrum, spectral radius, trace, and determinant of the given matrix.

$$(i) \begin{bmatrix} 3 & 2 & -1 \\ 2 & 6 & 4 \\ -1 & 4 & 5 \end{bmatrix}$$

$$(ii) \begin{bmatrix} 3 & -2 & 0.5 \\ -1 & -2 & 1.5 \\ -4 & 0 & 4 \end{bmatrix}$$

$$(iii) \begin{bmatrix} 2 & 0 & 0 \\ -6 & 8 & -14 \\ 0 & 0 & -6 \end{bmatrix}$$

$$(iv) \begin{bmatrix} -15.5 & -10 & 10 \\ 3 & 4.5 & -3 \\ -17 & -10 & 11.5 \end{bmatrix}$$

$$(v) \begin{bmatrix} 4.5 & 0 & 1.5 \\ -6 & 9 & 6 \\ 1.5 & 0 & 4.5 \end{bmatrix}$$

Exercise 02: Apply the Power method to find the dominant eigenvalue and corresponding eigenvector of the given matrices.

$$(i) \begin{bmatrix} 3 & 2 & -1 \\ 2 & 6 & 4 \\ -1 & 4 & 5 \end{bmatrix}$$

$$(ii) \begin{bmatrix} 3 & -2 & 0.5 \\ -1 & -2 & 1.5 \\ -4 & 0 & 4 \end{bmatrix}$$

$$(iii) \begin{bmatrix} 2 & 0 & 0 \\ -6 & 8 & -14 \\ 0 & 0 & -6 \end{bmatrix}$$

$$(iv) \begin{bmatrix} -15.5 & -10 & 10 \\ 3 & 4.5 & -3 \\ -17 & -10 & 11.5 \end{bmatrix}$$

$$(v) \begin{bmatrix} 4.5 & 0 & 1.5 \\ -6 & 9 & 6 \\ 1.5 & 0 & 4.5 \end{bmatrix}$$

Exercise 03: Apply the Power method to find the dominant eigenvalue and corresponding eigenvector of the given matrices.

$$(i) \begin{bmatrix} 8 & 1 & 0 & 0 \\ 0 & 7 & 0 & 0 \\ -2 & 1 & 10 & 0 \\ -4 & -1 & 4 & 6 \end{bmatrix}$$

$$(ii) \begin{bmatrix} 1 & 10 & 6 & -6 \\ 0 & -9 & 0 & 0 \\ -0.5 & 16.5 & 7.5 & 0.5 \\ -6.5 & 10.5 & 6.5 & 1.5 \end{bmatrix}$$

Exercise 04: Use Householder's method to place the following matrices in tridiagonal form.

$$(i) \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$

$$(ii) \begin{bmatrix} 2 & -1 & -1 \\ -1 & 2 & -1 \\ -1 & -1 & 2 \end{bmatrix}$$

$$(iii) \begin{bmatrix} 5 & -2 & -0.5 & 1.5 \\ -2 & 5 & 1.5 & -0.5 \\ -0.5 & 1.5 & 5 & -2 \\ 1.5 & -0.5 & -2 & 5 \end{bmatrix}$$

$$(iv) \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & 0 & 0 \\ -1 & 0 & 4 & 1 \\ 0 & -2 & 2 & 3 \end{bmatrix}$$

Exercise 05: Apply two iterations of the QR Factorization method without shifting the following matrices.

$$(i) \begin{bmatrix} 4 & -1 & 0 \\ -1 & 3 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$

$$(ii) \begin{bmatrix} 3 & 1 & 0 \\ 1 & 4 & 2 \\ 0 & 2 & 1 \end{bmatrix}$$

$$(iii) \begin{bmatrix} 4 & 2 & 0 & 0 \\ 2 & 4 & 2 & 0 \\ 0 & 2 & 4 & 2 \\ 0 & 0 & 2 & 4 \end{bmatrix}$$

$$(iv) \begin{bmatrix} 0.5 & 0.25 & 0 & 0 \\ 0.25 & 0.8 & 0.4 & 0 \\ 0 & 0.4 & 0.6 & 0.1 \\ 0 & 0 & 0.1 & 1 \end{bmatrix}$$

Exercise 06: Determine a singular value decomposition for the following matrices.

$$(i) \begin{bmatrix} 1 & 1 & 0 \\ -1 & 0 & 1 \\ 0 & 1 & -1 \end{bmatrix}$$

$$(ii) \begin{bmatrix} 2 & 1 \\ -1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$(iii) \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}$$

$$(iv) \begin{bmatrix} 2 & 1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$

■■■

Numerical Solution of Ordinary Differential Equations (ODEs)

Corridor I: BASICS

Let's plan it

- 9.1 Introduction
- 9.2 Solving IVPs using Single Step Methods and Multistep Methods
 - The Euler Method
 - The Mid-point Method (an RK2 method of Order 2)
 - The Modified/Improved Euler Method (an RK2 method of Order 2)
 - The RK Method of order 4 (RK4)
- 9.3 Solving IVPs using Predictor-Corrector Methods
 - The Adams-Bashforth-Moulton Method of Order 4
- 9.4 Solving Systems of ODEs and Higher Order ODEs
 - Using the Classical RK4 Method
- 9.5 Solving Linear BVPs using the Finite Difference Method

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



Corridor II: ANALYSIS

let's think deep

9.6 Some Theoretical Concepts and Error Analysis

To unleash the topics of this Corridor, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)

■■■

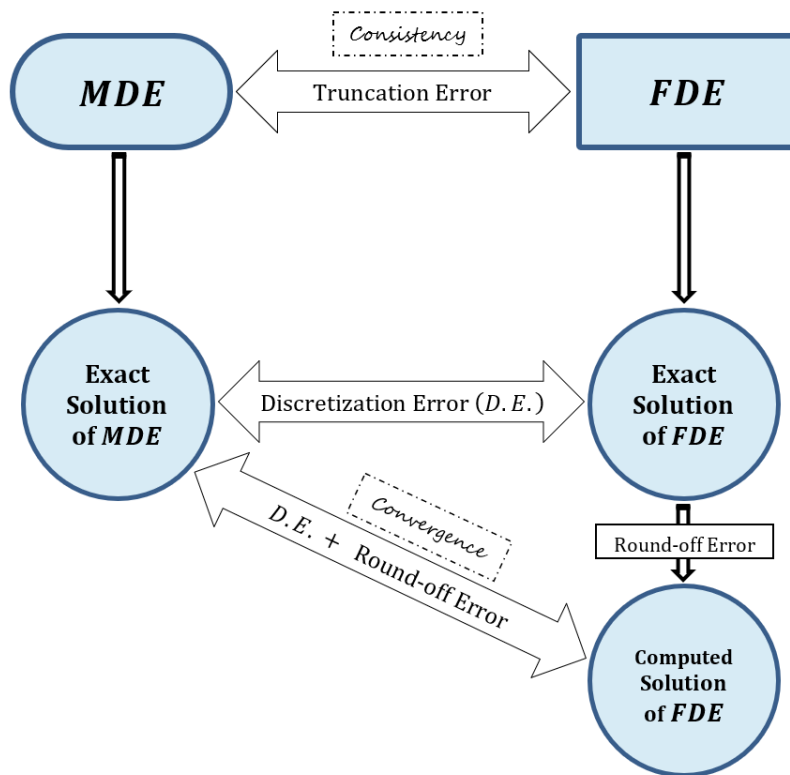


Figure: The connection between various terms related to MDE (Model Differential Equation/s - ODE/PDE) and the related FDE (Finite Difference Equation/s).

Corridor III: PROGRAMMING ARCADE

Let's think deep

9.7 Algorithms and Implementations

- Euler method
- Mid-point method
- Modified/Improved Euler method
- RK method of order 4 (RK4)
- Adams-Bashforth method of order 4
- Adams-Bashforth-Moulton method of order 4
- RK4 method for a system of two ODEs
- RK4 method for a system of three ODEs
- RK4 method for Second Order ODE
- RK4 method for Third Order ODE
- Linear FDM for BVP

To cross-check the results/output of the computer programs you would execute, please delve into the principal book:

Simplified Numerical Analysis (Fourth Edition; by Dr. Amjad Ali; www.TimeRenders.com.pk)



9.7 Algorithms and Implementations

Question 16: Write down an algorithm (pseudo code) to solve a first-order ODE using the Explicit Euler's method (the Taylor method of order 1).

Algorithm: To solve $y' = f(x, y)$, for $a \leq x \leq b$ and $y(a) = \alpha$ by approximating $y = y(x)$ at $(m + 1)$ equispaced nodes $x_0, x_1, x_2, \dots, x_m$, such that $a = x_0 < x_1 < x_2 < \dots < x_m = b$, $h = (b - a)/m$ and $y(x_i) = y_i$ using the Explicit Euler's method (the Taylor method of order 1): For $i = 1, 2, 3, \dots, m$

$$w_i = w_{i-1} + h \times f(x_{i-1}, y_{i-1})$$

INPUTS: $\begin{cases} a, b: \text{real values as the endpoints of the interval: } x \in [a, b] \\ m: \text{an integer as the number of nodes (other than } a) \text{ in the } x \text{ direction} \\ \text{alpha: a real value as the initial condition } y(a) \\ \text{A definition of the function } f(x, y) \text{ in an appropriate way} \end{cases}$

OUTPUT: $\begin{cases} W = [w_0, w_1, \dots, w_m]^T: \text{a real valued vector as the approximate solution } w(x_i) \\ \text{at the nodes } x_0, x_1, x_2, \dots, x_m \end{cases}$

Auxiliary Variables: $\begin{cases} h: \text{a real value as the step length in } x \text{ direction such that } h = (b - a)/m \\ x = (x_i) \text{ for } i = 0, 1, \dots, m: \text{a real valued vector to represent } x_i \text{s} \end{cases}$

Step 1 Receive the inputs as stated above

Step 2 Set $h = (b - a)/m$
Set $x(0) = a$
Set $x(m) = b$

Step 3 for $i = 1, 2, \dots, m - 1$
Set $x(i) = x(0) + i \times h$ (Constructing interior mesh points, x_i)
end for

Step 4 Set $w(0) = \text{alpha}$ (Setting the initial condition)

Step 5 for $i = 1, 2, \dots, m$
 $fval = f(x(i - 1), w(i - 1))$ (Computing the value $f(x_{i-1}, y_{i-1})$)
 $w(i) = w(i - 1) + h \times fval$
end for

Step 6 Print the output: $W = [w_0, w_1, \dots, w_m]^T$;

STOP.

Problem 09: Write down a C++ program to solve the IVP, $y' = 4y + 4x^2 + 3x$, for $0 \leq x \leq 1$, with initial condition $y(0) = \alpha = 1/2$, using the Explicit Euler's method (the Taylor method of order 1). Compute the solution for 10 steps. At each step, compare the approximate solution with the exact solution, to be obtained by $y(x) = -x^2 - \frac{5}{4}x - \frac{5}{16} + \frac{13}{16}e^{4x}$, by finding the relative error between the two solutions.

```

1  #include <iostream>
2  #include <cmath>
3  #include <iomanip>
4  using namespace std ;
5
6  #define a      0.0           // starting point of domain
7  #define b      1.0           // ending point of domain
8  #define alpha  0.5           // initial condition
9  #define m      10           // number of steps
10
11 #define fval(x,y) 4*y+4*x*x+3*x
12 #define fexact(x) -x*x-1.2*x-(5.0/16.0)+(13.0/16.0)*exp(4*x)
13
14 int main()
15 {
16     double h ;               // step size
17     double fv ;              // RHS function value
18     int i ;
19     double x[m+1] ;          // vector of nodes
20     double w[m+1] ;          // vector of solution values
21     double sol , err ;
22
23     h = (b-a)/m ;
24
25     x[0] = a ;
26
27     for (i=1 ; i<=m ; i++)
28     {
29         x[i] = x[i-1] + h ;    // x[i] = x[0] + i*h ;
30     }
31
32     w[0]= alpha ;            // setting initial condition
33

```

```

34 // ----- Computing solutions with the Euler method ----- //
35     for (i=1 ; i<=m ; i++)
36     {
37         fv = fval(x[i-1], w[i-1]) ;
38         w[i] = w[i-1] + h*fv ;
39     }
40
41 // ----- Printing Solutions ----- //
42     for (i=0 ; i<=m ; i++)
43     {
44         cout << "Node= " << setw(2) << i ;
45         cout << "\tx= " << setw(8) << x[i] ;
46         cout << "\tw= " << setw(8) << w[i] ;
47         sol = fexact(x[i]) ;
48         cout << "\tExact sol= " << setw(8) << sol ;
49         err = fabs(sol-w[i])/fabs(sol) ;
50         cout << "\tRelative Error= " << err << endl ;
51     }
52
53     return 0;
54 }

```

The above program can be written in a better way that a C++ function for the Euler method is formed to compute the solution. This makes the program better manageable and modular. The new program is given as follows.

```

1  #include <iostream>
2  #include <cmath>
3  #include <iomanip>
4  using namespace std ;
5
6  #define a      0.0           // starting point of domain
7  #define b      1.0           // ending point of domain
8  #define alpha  0.5           // initial condition
9  #define m      10           // number of steps
10
11 #define fval(x,y) 4*y+4*x*x+3*x
12 #define fexact(x) -x*x-1.2*x-(5.0/16.0)+(13.0/16.0)*exp(4*x)
13

```

```
14 void euler(double [], double [], double) ;
15
16 int main()
17 {
18     double h ;                // step size
19     double sol , err ;
20
21     int i ;
22     double x[m+1] ;          // vector of nodes
23     double w[m+1] ;          // vector of solution values
24
25     h = (b-a)/m ;
26
27     // forming vector of x-values or nodes
28     x[0] = a ;
29     x[m] = b ;
30
31     for (i=1 ; i<=m ; i++)
32     {
33         x[i] = x[i-1] + h ;    // x[i] = x[0] + i*h ;
34     }
35
36     w[0]= alpha ;            // setting initial condition
37
38     euler(x, w, h) ;         // Call-by-reference to the function
39
40 // ----- Printing Solutions ----- //
41     for (i=0 ; i<=m ; i++)
42     {
43         cout << "Node= " << setw(2) << i ;
44         cout << "\tx= " << setw(8) << x[i] ;
45         cout << "\tw= " << setw(8) << w[i] ;
46         sol = fexact(x[i]) ;
47         cout << "\tExact sol= " << setw(8) << sol ;
48         err = fabs(sol-w[i])/fabs(sol) ;
49         cout << "\tRelative Error= " << err << endl ;
50     }
51
52     return 0;
53 }
54
```

```

55 // ----- User-defined function for the Euler method ----- //
56
57 void euler(double x[m+1], double w[m+1], double h)
58 {
59     double fv ;           // RHS function value
60     int i ;
61
62     for (i=1 ; i<=m ; i++)
63     {
64         fv = fval(x[i-1], w[i-1]) ;
65         w[i] = w[i-1] + h*fv ;    // computing next solution value
66     }
67 }

```

The results are shown in the following table.

Steps i	Node $x(i)$	Numerical solution $w_i =$ $w(x_i)$	Exact solution $y_i = y(x_i)$	Error of computer program solution
0	0	0.5	0.5	0
1	0.1	0.7	0.769608	0.0904455
2	0.2	1.014	1.21575	0.165948
3	0.3	1.4956	1.93509	0.227118
4	0.4	2.21984	3.07184	0.277358
5	0.5	3.29178	4.84111	0.320037
6	0.6	4.85849	7.56383	0.357669
7	0.7	7.12588	11.7188	0.391926
8	0.8	10.3844	18.0202	0.423855
9	0.9	15.0311	27.5336	0.45408
10	1.0	21.6376	41.8485	0.482954

Question 17: Write down an algorithm (pseudo code) to solve a first-order ODE using the Midpoint method (which is an RK method of order 2).

Algorithm: To solve $y' = f(x, y)$, for $a \leq x \leq b$ and $y(a) = \alpha$ by approximating $y = y(x)$ at $(m + 1)$ equispaced nodes $x_0, x_1, x_2, \dots, x_m$, such that $a = x_0 < x_1 < x_2 < \dots < x_m = b$, $h = (b - a)/m$ and $y(x_i) = y_i$ using the Midpoint method: For $i = 1, 2, 3, \dots, m$

$$\bar{y}_i = y_{i-1} + \frac{h}{2} \times f(x_{i-1}, y_{i-1})$$

$$y_i = y_{i-1} + h \times f\left(x_{i-1} + \frac{h}{2}, \bar{y}_i\right)$$

INPUTS: $\left\{ \begin{array}{l} a, b: \text{real values as the endpoints of the interval: } x \in [a, b] \\ m: \text{an integer as the number of nodes (other than } a) \text{ in the } x \text{ direction} \\ \text{alpha: a real value as the initial condition } y(a) \\ \text{A definition of the function } f(x, y) \text{ in an appropriate way} \end{array} \right.$

OUTPUT: $\left\{ \begin{array}{l} W = [w_0, w_1, \dots, w_m]^T: \text{a real valued vector as the approximate solution } w(x_i) \\ \text{at the nodes } x_0, x_1, x_2, \dots, x_m \end{array} \right.$

Auxiliary Variables: $\left\{ \begin{array}{l} h: \text{a real value as the step length in } x \text{ direction such that } h = (b - a)/m \\ X = (x_i) \text{ for } i = 0, 1, \dots, m: \text{a real valued vector to represent } x_i \text{'s} \end{array} \right.$

Step 1 Receive the inputs as stated above

Step 2 Set $h = (b - a)/m$
Set $x(0) = a$
Set $x(m) = b$

Step 3 for $i = 1, 2, \dots, m - 1$
 Set $x(i) = x(0) + i \times h$ (Constructing interior mesh points, x_i)
end for

Step 4 Set $w(0) = \text{alpha}$ (Setting the initial condition)

Step 5 for $i = 1, 2, \dots, m$
 $fval1 = f(x(i - 1), w(i - 1))$ (Computing the value $f(x_{i-1}, y_{i-1})$)
 $aux = w(i - 1) + (h/2) \times fval1$
 $fval2 = f(x(i - 1) + (h/2), aux)$ (Computing $f(x_{i-1} + h/2, aux)$)
 $w(i) = w(i - 1) + h \times fval2$
end for

Step 6 Print the output: $W = [w_0, w_1, \dots, w_m]^T$

STOP.

Problem 10: Write down a C++ program to solve the IVP, $y' = 4y + 4x^2 + 3x$, for $0 \leq x \leq 1$, with initial condition $y(0) = \alpha = 1/2$, using the Midpoint method (which is an RK method of order 2). Compute the solution for 10 steps. At each step, compare the approximate solution with the exact solution, to be obtained by $y(x) = -x^2 - \frac{5}{4}x - \frac{5}{16} + \frac{13}{16}e^{4x}$, by finding the relative error between the two solutions.

```

1  #include <iostream>
2  #include <cmath>
3  #include <iomanip>
4  using namespace std ;
5
6  #define a      0.0           // starting point of domain
7  #define b      1.0           // ending point of domain
8  #define alpha  0.5           // initial condition
9  #define m      10           // number of steps
10
11 #define fval(x,y) 4*y+4*x*x+3*x
12 #define fexact(x) -x*x-1.2*x-(5.0/16.0)+(13.0/16.0)*exp(4*x)
13
14 void eulermid(double [], double [], double) ;
15
16 int main()
17 {
18     double h ;                // step size
19     double sol , err ;
20
21     int i ;
22     double x[m+1] ;           // vector of nodes
23     double w[m+1] ;           // vector of solution values
24
25     h = (b-a)/m ;
26
27     // forming vector of x-values or nodes
28     x[0] = a ;
29     x[m] = b ;
30
31     for (i=1 ; i<=m ; i++)
32     {
33         x[i] = x[i-1] + h ;    // x[i] = x[0] + i*h ;
34     }
35

```

```
36     w[0]= alpha ;           // setting initial condition
37
38     eulermid(x, w, h) ;     // Call-by-reference to the function
39
40 // ----- Printing Solutions ----- //
41     for (i=0 ; i<=m ; i++)
42     {
43         cout << "Node= " << setw(2) << i ;
44         cout << "\tx= " << setw(8) << x[i] ;
45         cout << "\tw= " << setw(8) << w[i] ;
46         sol = fexact(x[i]) ;
47         cout << "\tExact sol= " << setw(8) << sol ;
48         err = fabs(sol-w[i])/fabs(sol) ;
49         cout << "\tRelative Error= " << err << endl ;
50     }
51
52     return 0;
53 }
54
55 // ----- User-defined function for the Euler Mid-point ----- //
56
57 void eulermid(double xa[m+1], double wa[m+1], double h)
58 {
59     double fv ;           // RHS function value
60     int i ;
61     double whalf , xhalf ;
62
63     for (i=1 ; i<=m ; i++)
64     {
65         fv = fval(xa[i-1], wa[i-1]) ;
66         whalf = wa[i-1] + (h/2.0)*fv ;
67
68         xhalf = xa[i-1] + (h/2.0) ;
69         fv = fval(xhalf, whalf) ;
70
71         wa[i] = wa[i-1] + h*fv ;
72     }
73 }
```

The results are shown in the following table.

Steps i	Node $x(i)$	Numerical solution $w_i =$ $w(x_i)$	Exact solution $y_i = y(x_i)$	Error of computer program solution
0	0	0.5	0.5	0
1	0.1	0.756	0.769608	0.0176812
2	0.2	1.17968	1.21575	0.0296705
3	0.3	1.86113	1.93509	0.0382248
4	0.4	2.93367	3.07184	0.0449802
5	0.5	4.59463	4.84111	0.0509141
6	0.6	7.13605	7.56383	0.0565563
7	0.7	10.9902	11.7188	0.0621758
8	0.8	16.7966	18.0202	0.0678992
9	0.9	25.5022	27.5336	0.0737776
10	1.0	38.5081	41.8485	0.0798221

Question 18: Write down an algorithm (pseudo code) to solve a first-order ODE using the RK method of order 2 (also known as the Modified or Improved Euler's method).

Algorithm: To solve $y' = f(x, y)$, for $a \leq x \leq b$ and $y(a) = \alpha$ by approximating $y = y(x)$ at $(m + 1)$ equispaced nodes $x_0, x_1, x_2, \dots, x_m$, such that $a = x_0 < x_1 < x_2 < \dots < x_m = b$, $h = (b - a)/m$ and $y(x_i) = y_i$ using the Modified Euler's method of order 2: For $i = 1, 2, 3, \dots, m$

$$\begin{aligned} K_1 &= h \times f(x_{i-1}, y_{i-1}) \\ K_2 &= h \times f(x_i, y_{i-1} + K_1) \\ y_i &= y_{i-1} + \frac{1}{2} \times [K_1 + K_2] \end{aligned}$$

INPUTS: $\begin{cases} a, b: \text{real values as the endpoints of the interval: } x \in [a, b] \\ m: \text{an integer as the number of nodes (other than } a) \text{ in the } x \text{ direction} \\ \text{alpha: a real value as the initial condition } y(a) \\ \text{A definition of the function } f(x, y) \text{ in an appropriate way} \end{cases}$

OUTPUT: $\begin{cases} W = [w_0, w_1, \dots, w_m]^T: \text{a real valued vector as the approximate solution } w(x_i) \\ \text{at the nodes } x_0, x_1, x_2, \dots, x_m \end{cases}$

Auxiliary Variables: $\begin{cases} h: \text{a real value as the step length in } x \text{ direction such that } h = (b - a)/m \\ X = (x_i) \text{ for } i = 0, 1, \dots, m: \text{a real valued vector to represent } x_i \text{'s} \end{cases}$

- Step 1** Receive the inputs as stated above
- Step 2** Set $h = (b - a)/m$
Set $x(0) = a$
Set $x(m) = b$
- Step 3** for $i = 1, 2, \dots, m - 1$
Set $x(i) = x(0) + i \times h$ (Constructing interior mesh points, x_i)
end for
- Step 4** Set $w(0) = \alpha$ (Setting the initial condition)
- Step 5** for $i = 1, 2, \dots, m$
 $k1 = h \times f(x(i - 1), w(i - 1))$
 $k2 = h \times f(x(i), w(i - 1) + k1)$
 $w(i) = w(i - 1) + 0.5 \times (k1 + k2)$
 end for
- Step 6** Print the output: $W = [w_0, w_1, \dots, w_m]^T$
- STOP.**

Problem 11: Write down a C++ program to solve the IVP, $y' = 4y + 4x^2 + 3x$, for $0 \leq x \leq 1$, with initial condition $y(0) = \alpha = 1/2$, using the RK method of order 2 (also known as the Modified or Improved Euler's method). Compute the solution for 10 steps. At each step, compare the approximate solution with the exact solution, to be obtained by $y(x) = -x^2 - \frac{5}{4}x - \frac{5}{16} + \frac{13}{16}e^{4x}$, by finding the relative error between the two solutions.

```

1  #include <iostream>
2  #include <cmath>
3  #include <iomanip>
4  using namespace std ;
5
6  #define a      0.0          // starting point of domain
7  #define b      1.0          // ending point of domain
8  #define alpha  0.5          // initial condition
9  #define m      10          // number of steps
10
11 #define fval(x,y) 4*y+4*x*x+3*x
12 #define fexact(x) -x*x-1.2*x-(5.0/16.0)+(13.0/16.0)*exp(4*x)
13

```

```

14 void eulerimp(double [], double [], double) ;
15
16 int main()
17 {
18     double h ;                // step size
19     double sol , err ;
20
21     int i ;
22     double x[m+1] ;          // vector of nodes
23     double w[m+1] ;          // vector of solution values
24
25     h = (b-a)/m ;
26
27     // forming vector of x-values or nodes
28     x[0] = a ;
29     x[m] = b ;
30
31     for (i=1 ; i<=m ; i++)
32     {
33         x[i] = x[i-1] + h ;    // x[i] = x[0] + i*h ;
34     }
35
36     w[0]= alpha ;            // setting initial condition
37
38     eulerimp(x, w, h) ;      // Call-by-reference to the function
39
40 // ----- Printing Solutions ----- //
41     for (i=0 ; i<=m ; i++)
42     {
43         cout << "Node= " << setw(2) << i ;
44         cout << "\tx= " << setw(8) << x[i] ;
45         cout << "\tw= " << setw(8) << w[i] ;
46         sol = fexact(x[i]) ;
47         cout << "\tExact sol= " << setw(8) << sol ;
48         err = fabs(sol-w[i])/fabs(sol) ;
49         cout << "\tRelative Error= " << err << endl ;
50     }
51
52     return 0;
53 }
54

```

```

55 // ----- User-defined function for the RK 2----- //
56
57 void eulerimp(double xa[m+1], double wa[m+1], double h)
58 {
59     double fv, fvnext ;           // RHS function value
60     int i ;
61     double wnext ;
62
63     for (i=1 ; i<=m ; i++)
64     {
65         fv = fval(xa[i-1], wa[i-1]) ;
66         wnext = wa[i-1] + h*fv ;
67
68         fvnext = fval(xa[i], wnext) ;
69
70         wa[i] = wa[i-1] + h*(fv+fvnext)/2.0 ;
71     }
72 }

```

The results are shown in the following table.

Steps i	Node $x(i)$	Numerical solution $w_i =$ $w(x_i)$	Exact solution $y_i = y(x_i)$	Error of computer program solution
0	0	0.5	0.5	0
1	0.1	0.757	0.769608	0.0163818
2	0.2	1.18216	1.21575	0.0276306
3	0.3	1.8658	1.93509	0.0358113
4	0.4	2.94158	3.07184	0.0424044
5	0.5	4.60734	4.84111	0.0482887
6	0.6	7.15586	7.56383	0.0539372
7	0.7	11.0205	11.7188	0.0595885
8	0.8	16.8425	18.0202	0.0653535
9	0.9	25.5711	27.5336	0.0712755
10	1.0	38.611	41.8485	0.0773618



Question 19: Write down an algorithm (pseudo code) to solve a first-order ODE using the RK method of order 4.

Algorithm: To solve $y' = f(x, y)$, for $a \leq x \leq b$ and $y(a) = \alpha$ by approximating $y = y(x)$ at $(m + 1)$ equispaced nodes $x_0, x_1, x_2, \dots, x_m$, such that $a = x_0 < x_1 < x_2 < \dots < x_m = b$, $h = (b - a)/m$ and $y(x_i) = y_i$ using the RK method of order 4: For $i = 1, 2, 3, \dots, m$

$$\begin{aligned} K_1 &= h \times f(x_{i-1}, y_{i-1}) \\ K_2 &= h \times f(x_{i-1} + 0.5h, y_{i-1} + 0.5K_1) \\ K_3 &= h \times f(x_{i-1} + 0.5h, y_{i-1} + 0.5K_2) \\ K_4 &= h \times f(x_i, y_{i-1} + K_3) \\ y_i &= y_{i-1} + \frac{1}{6} \times [K_1 + 2K_2 + 2K_3 + K_4] \end{aligned}$$

INPUTS: $\begin{cases} a, b: \text{real values as the endpoints of the interval: } x \in [a, b] \\ m: \text{an integer as the number of nodes (other than } a) \text{ in the } x \text{ direction} \\ \text{alpha: a real value as the initial condition } y(a) \\ \text{A definition of the function } f(x, y) \text{ in an appropriate way} \end{cases}$

OUTPUT: $\begin{cases} W = [w_0, w_1, \dots, w_m]^T: \text{a real valued vector as the approximate solution } y(x_i) \\ \text{at the nodes } x_0, x_1, x_2, \dots, x_m \end{cases}$

Auxiliary Variables: $\begin{cases} h: \text{a real value as the step length in } x \text{ direction such that } h = (b - a)/m \\ X = (x_i) \text{ for } i = 0, 1, \dots, m: \text{a real valued vector to represent } x_i \end{cases}$

Step 1 Receive the inputs as stated above

Step 2 Set $h = (b - a)/m$

Set $x(0) = a$

Set $x(m) = b$

Step 3 for $i = 1, 2, \dots, m - 1$

Set $x(i) = x(0) + i \times h$

(Constructing interior mesh points, x_i)

end for

Step 4 Set $w(0) = \text{alpha}$

(Setting the initial condition)

Step 5 for $i = 1, 2, \dots, m$

$k1 = h \times f(x(i - 1), w(i - 1))$

$k2 = h \times f(x(i - 1) + 0.5 \times h, w(i - 1) + 0.5 \times k1)$

$k3 = h \times f(x(i - 1) + 0.5 \times h, w(i - 1) + 0.5 \times k2)$

$k4 = h \times f(x(i), w(i - 1) + k3)$

$w(i) = w(i - 1) + (k1 + 2 \times k2 + 2 \times k3 + k4)/6$

end for

Step 6 Print the output: $W = [w_0, w_1, \dots, w_m]^T$

STOP.

Problem 12: Write down a C++ program to solve the IVP, $y' = 4y + 4x^2 + 3x$, for $0 \leq x \leq 1$, with initial condition $y(0) = \alpha = 1/2$, using the RK method of order 4. Computer the solution for 10 steps. At each step, compare the approximate solution with the exact solution, to be obtained by $y(x) = -x^2 - \frac{5}{4}x - \frac{5}{16} + \frac{13}{16}e^{4x}$, by finding the relative error between the two solutions.

```
1  #include <iostream>
2  #include <cmath>
3  #include <iomanip>
4  using namespace std ;
5
6  #define a      0.0           // starting point of domain
7  #define b      1.0           // ending point of domain
8  #define alpha  0.5           // initial condition
9  #define m      10           // number of steps
10
11 #define fval(x,y) 4*y+4*x*x+3*x
12 #define fexact(x) -x*x-1.2*x-(5.0/16.0)+(13.0/16.0)*exp(4*x)
13
14 void rk4(double [], double [], double) ;
15
16 int main()
17 {
18     double h ;                // step size
19     double sol , err ;
20
21     int i ;
22     double x[m+1] ;           // vector of nodes
23     double w[m+1] ;           // vector of solution values
24
25     h = (b-a)/m ;
26
27     // forming vector of x-values or nodes
28     x[0] = a ;
29     x[m] = b ;
30
31     for (i=1 ; i<=m ; i++)
32     {
33         x[i] = x[i-1] + h ;    // x[i] = x[0] + i*h ;
34     }
35
```

```

36     w[0]= alpha ;           // setting initial condition
37
38     rk4(x, w, h) ;         // Call-by-reference to the function
39
40 // ----- Printing Solutions ----- //
41     for (i=0 ; i<=m ; i++)
42     {
43         cout << "Node= " << setw(2) << i ;
44         cout << "\tx= " << setw(8) << x[i] ;
45         cout << "\tw= " << setw(8) << w[i] ;
46         sol = fexact(x[i]) ;
47         cout << "\tExact sol= " << setw(8) << sol ;
48         err = fabs(sol-w[i])/fabs(sol) ;
49         cout << "\tRelative Error= " << err << endl ;
50     }
51
52     return 0;
53 }
54
55 // ----- User-defined function for the RK 4----- //
56
57 void rk4(double x[m+1], double w[m+1], double h)
58 {
59     double k1 , k2 , k3 , k4 ;
60     int i ;
61     double p , q ;
62
63     for (i=1 ; i<=m ; i++)
64     {
65         k1 = h*(fval( x[i-1] , w[i-1] ));
66         k2 = h*(fval( (x[i-1]+(h/2.0)) , (w[i-1]+(k1/2.0)) ));
67         k3 = h*(fval( (x[i-1]+(h/2.0)) , (w[i-1]+(k2/2.0)) ));
68         k4 = h*(fval( x[i] , (w[i-1]+k3) ));
69
70         w[i] = w[i-1] + (k1+2*k2+2*k3+k4)/6.0 ;
71     }
72 }

```

The results are shown in the following table.

Steps i	Node $x(i)$	Numerical solution $w_i =$ $w(x_i)$	Exact solution $y_i = y(x_i)$	Error of computer program solution
0	0	0.5	0.5	0
1	0.1	0.764547	0.769608	0.00657595
2	0.2	1.20556	1.21575	0.00838021
3	0.3	1.91966	1.93509	0.00797517
4	0.4	3.05096	3.07184	0.00679678
5	0.5	4.81444	4.84111	0.00550804
6	0.6	7.53081	7.56383	0.00436537
7	0.7	11.6785	11.7188	0.00343964
8	0.8	17.9711	18.0202	0.00272617
9	0.9	27.4731	27.5336	0.00219447
10	1.0	41.7728	41.8485	0.0018091

Question 20: Write down an algorithm (pseudo code) to solve a first-order ODE using the Adams-Bashforth method of order 4.

Algorithm: To solve $y' = f(x, y)$, for $a \leq x \leq b$ and $y(a) = \alpha$ by approximating $y = y(x)$ at $(m + 1)$ equispaced nodes $x_0, x_1, x_2, \dots, x_m$, such that $a = x_0 < x_1 < x_2 < \dots < x_m = b$, $h = (b - a)/m$ and $y(x_i) = y_i$. Having $y(x_0) = \alpha_0$, $y(x_1) = \alpha_1$, $y(x_2) = \alpha_2$, and $y(x_3) = \alpha_3$, compute y_i using the 4-step explicit Adams-Bashforth method of order 4: For $i = 4, 5, 6, \dots, m$

$$y_i = y_{i-1} + \frac{h}{24} \times [55f(x_{i-1}, y_{i-1}) - 59f(x_{i-2}, y_{i-2}) + 37f(x_{i-3}, y_{i-3}) - 9f(x_{i-4}, y_{i-4})]$$

INPUTS: $\begin{cases} a, b: \text{real values as the endpoints of the interval: } x \in [a, b] \\ m: \text{an integer as the number of nodes (other than } a) \text{ in the } x \text{ direction} \\ \text{alpha: a real value as the initial condition } y(a) \\ \text{A definition of the function } f(x, y) \text{ in an appropriate way} \end{cases}$

OUTPUT: $\begin{cases} W = [w_0, w_1, \dots, w_m]^T: \text{a real valued vector as the approximate solution } w(x_i) \\ \text{at the nodes } x_0, x_1, x_2, \dots, x_m \end{cases}$

Auxiliary Variables: $\begin{cases} h: \text{a real value as the step length in } x \text{ direction such that } h = (b - a)/m \\ X = (x_i) \text{ for } i = 0, 1, \dots, m: \text{a real valued vector to represent } x_i \text{s} \end{cases}$

Step 1 Receive the inputs as stated above


```
8 #define alpha 0.5 // initial condition
9 #define m 10 // number of steps
10
11 #define fval(x,y) 4*y+4*x*x+3*x
12 #define fexact(x) -x*x-1.2*x-(5.0/16.0)+(13.0/16.0)*exp(4*x)
13
14 void adamsb4(double [], double [], double) ;
15 void rk4(double [], double [], double) ;
16
17 int main()
18 {
19     double h ; // step size
20     double sol , err ;
21
22     int i ;
23     double x[m+1] ; // vector of nodes
24     double w[m+1] ; // vector of solution values
25
26     h = (b-a)/m ;
27
28     // forming vector of x-values or nodes
29     x[0] = a ;
30     x[m] = b ;
31
32     for (i=1 ; i<=m ; i++)
33     {
34         x[i] = x[i-1] + h ; // x[i] = x[0] + i*h ;
35     }
36
37     w[0]= alpha ; // setting initial condition
38
39     rk4(x, w, h) ; //Call-by-reference to function RK4
40     adamsb4(x, w, h) ; // Call-by-reference to function adamsb4
41
42 // ----- Printing Solutions ----- //
43     for (i=0 ; i<=m ; i++)
44     {
45         cout << "Node= " << setw(2) << i ;
46         cout << "\tx= " << setw(8) << x[i] ;
47         cout << "\tw= " << setw(8) << w[i] ;
48         sol = fexact(x[i]) ;
```

```

49         cout << "\tExact sol= " << setw(8) << sol ;
50         err = fabs(sol-w[i])/fabs(sol) ;
51         cout << "\tRelative Error= " << err << endl ;
52     }
53
54     return 0;
55 }
56
57 // ---- User-defined function for Adams-Bashforth method ---- //
58
59 void adamsb4(double x[m+1], double w[m+1], double h)
60 {
61     double k1 , k2 , k3 , k4 , fv ;
62     double p , q ;
63     int i ;
64
65     for (i=4 ; i<=m ; i++)
66     {
67         k1 = fval( x[i-1] , w[i-1] );
68         k2 = fval( x[i-2] , w[i-2] );
69         k3 = fval( x[i-3] , w[i-3] );
70         k4 = fval( x[i-4] , w[i-4] );
71
72         w[i] = w[i-1] + (h/24.0)*(55*k1 - 59*k2 + 37*k3 - 9*k4) ;
73     }
74 }
75
76 // ----- User-defined function for the RK 4----- //
77
78 void rk4(double x[m+1], double w[m+1], double h)
79 {
80     double k1 , k2 , k3 , k4 , fv ;
81     int i ;
82     double p , q ;
83
84     for (i=1 ; i<=3 ; i++)
85     {
86         k1 = h*(fval( x[i-1] , w[i-1] ));
87         k2 = h*(fval( (x[i-1]+(h/2.0)) , (w[i-1]+(k1/2.0)) ));
88         k3 = h*(fval( (x[i-1]+(h/2.0)) , (w[i-1]+(k2/2.0)) ));
89         k4 = h*(fval( x[i] , (w[i-1]+k3) ));

```

```

90
91     w[i] = w[i-1] + (k1+2*k2+2*k3+k4)/6.0 ;
92 }
93 }

```

The results are shown in the following table.

Steps i	Node $x(i)$	Numerical solution $w_i =$ $w(x_i)$	Exact solution $y_i = y(x_i)$	Error of computer program solution
0	0	0.5	0.5	0
1	0.1	0.764547	0.769608	0.00657595
2	0.2	1.20556	1.21575	0.00838021
3	0.3	1.91966	1.93509	0.00797517
4	0.4	3.04469	3.07184	0.00883945
5	0.5	4.79306	4.84111	0.00992469
6	0.6	7.48205	7.56383	0.0108119
7	0.7	11.5814	11.7188	0.011726
8	0.8	17.7896	18.0202	0.0127952
9	0.9	27.1477	27.5336	0.0140137
10	1.0	41.2059	41.8485	0.0153558

Question 21: Write down an algorithm (pseudo code) to solve a first-order ODE using the Adams-Bashforth-Moulton method of order 4.

Algorithm: To solve $y' = f(x, y)$, for $a \leq x \leq b$ and $y(a) = \alpha$ by approximating $y = y(x)$ at $(m + 1)$ equispaced nodes $x_0, x_1, x_2, \dots, x_m$, such that $a = x_0 < x_1 < x_2 < \dots < x_m = b$, $h = (b - a)/m$ and $y(x_i) = y_i$. Having $y(x_0) = \alpha_0$, $y(x_1) = \alpha_1$, $y(x_2) = \alpha_2$, and $y(x_3) = \alpha_3$, compute y_i using

(1) the 4-step explicit Adams-Bashforth method of order 4 as the predictor:

$$y_i = y_{i-1} + \frac{h}{24} \times [55f(x_{i-1}, y_{i-1}) - 59f(x_{i-2}, y_{i-2}) + 37f(x_{i-3}, y_{i-3}) - 9f(x_{i-4}, y_{i-4})]$$

(2) the 3-step implicit Adams-Moulton method of order 4 as the corrector:

$$y_i = y_{i-1} + \frac{h}{24} \times [9f(x_i, y_i) + 19f(x_{i-1}, y_{i-1}) - 5f(x_{i-2}, y_{i-2}) + f(x_{i-3}, y_{i-3})]$$

for $i = 4, 5, 6, \dots, m$.

INPUTS: $\begin{cases} a, b: \text{real values as the endpoints of the interval: } x \in [a, b] \\ m: \text{an integer as the number of nodes (other than } a) \text{ in the } x \text{ direction} \\ \text{alpha: a real value as the initial condition } y(a) \\ \text{A definition of the function } f(x, y) \text{ in an appropriate way} \end{cases}$

OUTPUT: $\begin{cases} W = [w_0, w_1, \dots, w_m]^T: \text{a real valued vector as the approximate solution } y(x_i) \\ \text{at the nodes } x_0, x_1, x_2, \dots, x_m \end{cases}$

Auxiliary Variables: $\begin{cases} h: \text{a real value as the step length in } x \text{ direction such that } h = (b - a)/m \\ X = (x_i) \text{ for } i = 0, 1, \dots, m: \text{a real valued vector to represent } x_i \end{cases}$

Step 1 Receive the inputs as stated above

Step 2 Set $h = (b - a)/m$
Set $x(0) = a$
Set $x(m) = b$

Step 3 for $i = 1, 2, \dots, m - 1$
Set $x(i) = x(0) + i \times h$ (Constructing interior mesh points, x_i)
end for

Step 4 Set $w(0) = \text{alpha}$ (Setting the initial condition)

Step 5 Obtain or compute (using some other basic method for ODEs) the following:

$w(1) = \text{alpha1}$
 $w(2) = \text{alpha2}$
 $w(3) = \text{alpha3}$

Step 5 for $i = 4, 5, 6, \dots, m$
 $fv1 = f(x(i - 1), w(i - 1))$ (Computing the value $f(x_{i-1}, y_{i-1})$)
 $fv2 = f(x(i - 2), w(i - 2))$ (Computing the value $f(x_{i-2}, y_{i-2})$)
 $fv3 = f(x(i - 3), w(i - 3))$ (Computing the value $f(x_{i-3}, y_{i-3})$)
 $fv4 = f(x(i - 4), w(i - 4))$ (Computing the value $f(x_{i-4}, y_{i-4})$)

$$w(i) = w(i - 1) + \left(\frac{h}{24}\right) \times (55fv1 - 59fv2 + 37fv3 - 9fv4)$$

$fv = f(x(i), w(i))$ (Computing the value $f(x_{i-4}, y_{i-4})$)

$$w(i) = w(i - 1) + \left(\frac{h}{24}\right) \times (9fv + 19fv1 - 5fv2 + fv3)$$

end for

Step 6 Print the output: $W = [w_0, w_1, \dots, w_m]^T$
STOP.

Problem 14: Write down a C++ program to solve the IVP, $y' = 4y + 4x^2 + 3x$, for $0 \leq x \leq 1$, with initial condition $y(0) = \alpha = 1/2$, using the Adams-Bashforth-Moulton method of order 4. Compute the solution for 10 steps. For computing the approximate solution at the first three steps, use the RK4 method. At each step, compare the approximate solution with the exact solution, to be obtained by $y(x) = -x^2 - \frac{5}{4}x - \frac{5}{16} + \frac{13}{16}e^{4x}$, by finding the relative error between the two solutions.

```

1  #include <iostream>
2  #include <cmath>
3  #include <iomanip>
4  using namespace std ;
5
6  #define a      0.0          // starting point of domain
7  #define b      1.0          // ending point of domain
8  #define alpha  0.5          // initial condition
9  #define m      10          // number of steps
10
11 #define fval(x,y) 4*y+4*x*x+3*x
12 #define fexact(x) -x*x-1.2*x-(5.0/16.0)+(13.0/16.0)*exp(4*x)
13
14 void adamsb4m3(double [], double [], double) ;
15 void rk4(double [], double [], double) ;
16
17 int main()
18 {
19     double h ;                // step size
20     double sol , err ;
21
22     int i ;
23     double x[m+1] ;          // vector of nodes
24     double w[m+1] ;          // vector of solution values
25
26     h = (b-a)/m ;
27
28     // forming vector of x-values or nodes
29     x[0] = a ;
30     x[m] = b ;
31
32     for (i=1 ; i<=m ; i++)
33     {
34         x[i] = x[i-1] + h ;    // x[i] = x[0] + i*h ;
35     }

```

```

36
37     w[0]= alpha ;                // setting initial condition
38
39     rk4(x, w, h) ;                //Call-by-reference to function RK4
40     adamsb4m3(x, w, h) ;         // Call-by-reference to function adamsb4
41
42 // ----- Printing Solutions ----- //
43     for (i=0 ; i<=m ; i++)
44     {
45         cout << "Node= " << setw(2) << i ;
46         cout << "\tx= " << setw(8) << x[i] ;
47         cout << "\tw= " << setw(8) << w[i] ;
48         sol = fexact(x[i]) ;
49         cout << "\tExact sol= " << setw(8) << sol ;
50         err = fabs(sol-w[i])/fabs(sol) ;
51         cout << "\tRelative Error= " << err << endl ;
52     }
53
54     return 0;
55 }
56
57 // ---- User-defined function for Adams-Bashforth-Moulten method ---- //
58
59 void adamsb4m3(double x[m+1], double w[m+1], double h)
60 {
61     double fv1 , fv2 , fv3 , fv4 , fv ;
62     double p , q ;
63     int i ;
64
65     for (i=4 ; i<=m ; i++)
66     {
67         fv1 = fval( x[i-1] , w[i-1] );
68         fv2 = fval( x[i-2] , w[i-2] );
69         fv3 = fval( x[i-3] , w[i-3] );
70         fv4 = fval( x[i-4] , w[i-4] );
71
72         w[i] = w[i-1] + (h/24.0)*(55*fv1 - 59*fv2 + 37*fv3 - 9*fv4) ;
73
74         fv = fval( x[i] , w[i]);
75
76         w[i] = w[i-1] + (h/24.0)*(9*fv + 19*fv1 - 5*fv2 + fv3);

```

```

77     }
78 }
79
80 // ----- User-defined function for the RK 4----- //
81
82 void rk4(double x[m+1], double w[m+1], double h)
83 {
84     double k1 , k2 , k3 , k4, fv ;
85     int i ;
86     double p , q ;
87
88     for (i=1 ; i<=3 ; i++)
89     {
90         k1 = h*(fval( x[i-1] , w[i-1] ));
91         k2 = h*(fval( (x[i-1]+(h/2.0)) , (w[i-1]+(k1/2.0)) ));
92         k3 = h*(fval( (x[i-1]+(h/2.0)) , (w[i-1]+(k2/2.0)) ));
93         k4 = h*(fval( x[i] , (w[i-1]+k3) ));
94
95         w[i] = w[i-1] + (k1+2*k2+2*k3+k4)/6.0 ;
96     }
97 }

```

The results are shown in the following table.

Steps i	Node $x(i)$	Numerical solution $w_i =$ $w(x_i)$	Exact solution $y_i = y(x_i)$	Error of computer program solution
0	0	0.5	0.5	0
1	0.1	0.764547	0.769608	0.00657595
2	0.2	1.20556	1.21575	0.00838021
3	0.3	1.91966	1.93509	0.00797517
4	0.4	3.05087	3.07184	0.00682605
5	0.5	4.81417	4.84111	0.00556428
6	0.6	7.53021	7.56383	0.0044448
7	0.7	11.6773	11.7188	0.00354036
8	0.8	17.9689	18.0202	0.00284706
9	0.9	27.4693	27.5336	0.00233489
10	1.0	41.7661	41.8485	0.00196874

■

Problem 15: Write a C++ program to solve the following system of two ODEs for the functions $y_1 = y_1(x)$ and $y_2 = y_2(x)$, where $x \in [0,1]$:

$$\begin{aligned}y_1' &= y_1 y_2 - 2 \\ y_2' &= 2y_1 - y_2^3\end{aligned}$$

With initial conditions:

$$\begin{aligned}y_1(0) &= 2.0 \\ y_2(0) &= 0.3\end{aligned}$$

Use the RK4 method of order 4 for 5 steps.

For 5 steps the domain is discretized as

$$h = \frac{b - a}{m} = \frac{1.0 - 0.0}{5} = 0.2$$

$$x_0 = 0, \quad x_1 = 0.2, \quad x_2 = 0.4, \quad x_3 = 0.6, \quad x_4 = 0.8, \quad x_5 = 1.0.$$

According to the initial conditions:

$$\begin{aligned}w_{10} = y_{10} &= y_1(x_0) = y_1(0) = 2.0 \\ w_{20} = y_{20} &= y_2(x_0) = y_2(0) = 0.3\end{aligned}$$

The problem is to find approximations w_{1i} to $y_{1i} = y_1(x_i)$ and w_{2i} to $y_{2i} = y_2(x_i)$, for $i = 1, 2, 3, 4, 5$.

The C++ program for the solution is as follows.

```

1  #include <iostream>
2  #include <cmath>
3  #include <iomanip>
4  using namespace std ;
5
6  #define a      0.0           // starting point of domain
7  #define b      1.0           // ending point of domain
8  #define alpha1 2.0           //initial condition for first variable y1=y1(x)
9  #define alpha2 0.3          // initial condition for second variable y2=y2(x)
10 #define m      5            // number of steps
11
12 #define f1(x,y1,y2) y1*y2 - 2
13 #define f2(x,y1,y2) 2*y1 - y2*y2*y2
14
16 void rk4system2(double [], double [], double[], double) ;
17
18 int main()
19 {
20     double h ;                // step size
21     double sol , err ;

```

```
22
23     int i ;
24     double x[m+1] ;           // vector of nodes
25     double w1[m+1] ;         // vector of solution values for y1=y1(x)
26     double w2[m+1] ;         // vector of solution values for y2=y2(x)
27
28     h = (b-a)/m ;
29
30     // forming vector of x-values or nodes
31     x[0] = a ;
32     x[m] = b ;
33
34     for (i=1 ; i<=m ; i++)
35     {
36         x[i] = x[i-1] + h ;    // x[i] = x[0] + i*h ;
37     }
38
39     w1[0] = alpha1 ;          // setting initial condition for y1
40     w2[0] = alpha2 ;          // setting initial condition for y2
41
42     // Call-by-reference to the solver of ODEs system of 2 equations
43     rk4system2(x , w1 , w2 , h) ;
44
45     // ----- Printing Solutions ----- //
46     for (i=0 ; i<=m ; i++)
47     {
48         cout<< "Node=" << setw(3) << i ;
49         cout << "\tx= " << setw(8) << x[i] ;
50         cout<< "\tw1=" << setw(8) << w1[i] ;
51         cout<< "\tw2=" << setw(8) << w2[i] << endl ;
52     }
53
54     return 0;
55 }
56
57 // --- User-defined function for ODE system of 2 equations using RK4 --- //
58
59 void rk4system2(double x[m+1], double w1[m+1], double w2[m+1], double h)
60 {
61     double k11, k12, k13, k14 ;
62     double k21, k22, k23, k24 ;
63     int i ;
64
65     for (i=1 ; i<=m ; i++)
66     {
```

```

67     k11 = h * (f1(x[i-1], w1[i-1], w2[i-1])) ;
68     k21 = h * (f2(x[i-1], w1[i-1], w2[i-1])) ;
69
70     k12 = h*(f1((x[i-1]+0.5*h), (w1[i-1]+0.5*k11), (w2[i-1]+0.5*k21))) ;
71     k22 = h * (f2((x[i-1]+0.5*h), (w1[i-1]+0.5*k11), (w2[i-1]+0.5*k21))) ;
72
73     k13 = h * (f1((x[i-1]+0.5*h), (w1[i-1]+0.5*k12), (w2[i-1]+0.5*k22))) ;
74     k23 = h * (f2((x[i-1]+0.5*h), (w1[i-1]+0.5*k12), (w2[i-1]+0.5*k22))) ;
75
76     k14 = h * (f1((x[i-1]+h), (w1[i-1]+k13), (w2[i-1]+k23))) ;
77     k24 = h * (f2((x[i-1]+h), (w1[i-1]+k13), (w2[i-1]+k23))) ;
78
79     w1[i] = w1[i-1] + (k11 + 2*k12 + 2*k13 + k14)/6.0 ;
80     w2[i] = w2[i-1] + (k21 + 2*k22 + 2*k23 + k24)/6.0 ;
81 }
82 }

```

The results are shown in the following table.

Steps i	Node $x(i)$	Numerical solution $w_{1i} = w_1(x_i)$	Numerical solution $w_{2i} = w_2(x_i)$
0	0.0	2	0.3
1	0.2	1.815132	0.985522
2	0.4	1.90079	1.36485
3	0.6	2.08065	1.52571
4	0.8	2.38251	1.62306
5	1.0	2.85383	1.72393

Problem 16: Write a C++ program to solve the following system of two ODEs for the functions $y_1 = y_1(x)$, $y_2 = y_2(x)$, and $y_3 = y_3(x)$, where $x \in [0,1]$:

$$y_1' = y_1 + 3y_2 - 3y_3 + e^{-x}$$

$$y_2' = 2y_2 + y_3 - 3e^{-x}$$

$$y_3' = y_1 + 2y_2 + e^{-x}$$

With initial conditions:

$$y_1(0) = 2.5$$

$$y_2(0) = -1.5$$

$$y_3(0) = -1.0$$

Use the RK4 method of order 4 for 10 steps.

For 10 steps, the domain is discretized as

$$h = \frac{b - a}{m} = \frac{1.0 - 0.0}{10} = 0.1$$

$x_0 = 0, x_1 = 0.1, x_2 = 0.2, x_3 = 0.3, x_4 = 0.4, x_5 = 0.5, x_6 = 0.6, x_7 = 0.7, x_8 = 0.8, x_9 = 0.9, x_{10} = 1.0$.

According to the initial conditions:

$$w_{10} = y_{10} = y_1(x_0) = y_1(0) = 2.5$$

$$w_{20} = y_{20} = y_2(x_0) = y_2(0) = -1.5$$

$$w_{30} = y_{30} = y_3(x_0) = y_3(0) = -1.0$$

The problem is to find approximations w_{1i} to $y_{1i} = y_1(x_i)$, w_{2i} to $y_{2i} = y_2(x_i)$, and w_{3i} to $y_{3i} = y_3(x_i)$, for $i = 1, 2, \dots, 10$.

The C++ program for the solution is as follows.

```

1  #include <iostream>
2  #include <cmath>
3  #include <iomanip>
4  using namespace std ;
5
6  #define a      0.0           // starting point of domain
7  #define b      1.0           // ending point of domain
8  #define alpha1  2.5         //initial condition for first variable y1=y1(x)
9  #define alpha2 -1.5         // initial condition for second variable y2=y2(x)
10 #define alpha3 -1.0        // initial condition for third variable y3=y3(x)
11 #define m      5           // number of steps
12
13 #define f1(x,y1,y2,y3)  y1 + 3*y2 - 3*y3 + exp(-x)
13 #define f2(x,y1,y2,y3)  2*y2 + y3 - 3*exp(-x)
14 #define f3(x,y1,y2,y3)  y1 + 2*y2 + exp(-x)
15
16
17 void rk4system3(double [], double [], double[], double [], double) ;
18
19 int main()
20 {
21     double h ;                // step size
22     double sol , err ;
23
24     int i ;
25     double x[m+1] ;          // vector of nodes

```

```

26     double w1[m+1] ;           // vector of solution values for y1=y1(x)
27     double w2[m+1] ;           // vector of solution values for y2=y2(x)
28     double w3[m+1] ;           // vector of solution values for y3=y3(x)
29
30     h = (b-a)/m ;
31
32     // forming vector of x-values or nodes
33     x[0] = a ;
34     x[m] = b ;
35
36     for (i=1 ; i<=m ; i++)
37     {
38         x[i] = x[i-1] + h ;     // x[i] = x[0] + i*h ;
39     }
40
41     w1[0] = alpha1 ;           // setting initial condition for y1
42     w2[0] = alpha2 ;           // setting initial condition for y2
43     w3[0] = alpha3 ;           // setting initial condition for y3
44
45     // Call-by-reference to the solver of ODEs system of 3 equations
46     rk4system3(x , w1 , w2 , w3 , h) ;
47
48     // ----- Printing Solutions ----- //
49     for (i=0 ; i<=m ; i++)
50     {
51         cout<< "Node=" << setw(3) << i ;
52         cout << "\tx= " << setw(8) << x[i] ;
53         cout<< "\tw1=" << setw(8) << w1[i] ;
54         cout<< "\tw2=" << setw(8) << w2[i] ;
55         cout<< "\tw3=" << setw(8) << w3[i] << endl ;
56     }
57
58     return 0;
59 }
60
61 // --- User-defined function for ODE system of 3 equations using RK4 --- //
62 void rk4system3(double x[m+1], double w1[m+1], double w2[m+1],
63 double w3[m+1], double h)
64 {
65     double k11, k12, k13, k14 ;
66     double k21, k22, k23, k24 ;
67     double k31, k32, k33, k34 ;
68     int i ;

```



```

67   for (i=1 ; i<=m ; i++)
68   {
69       k11 = h * (f1(x[i-1], w1[i-1], w2[i-1], w3[i-1])) ;
70       k21 = h * (f2(x[i-1], w1[i-1], w2[i-1], w3[i-1])) ;
71       k31 = h * (f3(x[i-1], w1[i-1], w2[i-1], w3[i-1])) ;
72
73       k12 = h * (f1((x[i-1]+0.5*h), (w1[i-1]+0.5*k11),
74                   (w2[i-1]+0.5*k21), (w3[i-1]+0.5*k31))) ;
75       k22 = h * (f2((x[i-1]+0.5*h), (w1[i-1]+0.5*k11),
76                   (w2[i-1]+0.5*k21), (w3[i-1]+0.5*k31))) ;
77       k32 = h * (f3((x[i-1]+0.5*h), (w1[i-1]+0.5*k11),
78                   (w2[i-1]+0.5*k21), (w3[i-1]+0.5*k31))) ;
79
80       k13 = h * (f1((x[i-1]+0.5*h), (w1[i-1]+0.5*k12),
81                   (w2[i-1]+0.5*k22), (w3[i-1]+0.5*k32))) ;
82       k23 = h * (f2((x[i-1]+0.5*h), (w1[i-1]+0.5*k12),
83                   (w2[i-1]+0.5*k22), (w3[i-1]+0.5*k32))) ;
84       k33 = h * (f3((x[i-1]+0.5*h), (w1[i-1]+0.5*k12),
85                   (w2[i-1]+0.5*k22), (w3[i-1]+0.5*k32))) ;
86
87       k14 = h * (f1((x[i-1]+h), (w1[i-1]+k13),
88                   (w2[i-1]+k23), (w3[i-1]+k33))) ;
89       k24 = h * (f2((x[i-1]+h), (w1[i-1]+k13),
90                   (w2[i-1]+k23), (w3[i-1]+k33))) ;
91       k34 = h * (f3((x[i-1]+h), (w1[i-1]+k13),
92                   (w2[i-1]+k23), (w3[i-1]+k33))) ;
93
94       w1[i] = w1[i-1] + (k11 + 2*k12 + 2*k13 + k14)/6.0 ;
95       w2[i] = w2[i-1] + (k21 + 2*k22 + 2*k23 + k24)/6.0 ;
96       w3[i] = w3[i-1] + (k31 + 2*k32 + 2*k33 + k34)/6.0 ;
97   }
98 }

```

The results are shown in the following table.

Steps i	Node $x(i)$	Numerical solution $w_{1i} = w_1(x_i)$	Numerical solution $w_{2i} = w_2(x_i)$	Numerical solution $w_{3i} = w_3(x_i)$
0	0.0	2.5	-1.5	-1.0
1	0.2	2.45262	-3.16688	-1.22125
2	0.4	1.43325	-5.68432	-2.39815
3	0.6	-0.713292	-9.79838	-5.20827
4	0.8	-4.26124	-16.8302	-10.7742
5	1.0	-9.74139	-29.1044	-21.007

Problem 17: Write a C++ program to find the numerical solution of the ODE, $xy'' - y' + 8x^3y^3 = 0$ with initial condition $y(1) = 0.5$ and $y'(1) = -0.5$ for $y(1.1)$. Consider the step size of $h = 0.1$, thus only step is required. i.e., $m = 1$. Use the exact solution, $y = 1/(1 + x^2)$, to find the error in the numerical solution.

For the solution, consider

$$y' = z$$

Then, the given ODE becomes

$$z' = \frac{(z - 8x^3y^3)}{x}$$

Thus, the second-order IVP is essentially converted to the problem of a first-order system of ODEs of comprising the two equations subject to the initial conditions:

$$w_{10} = y_0 = y(x_0) = y(1) = 0.5$$

$$w_{20} = z_0 = z(x_0) = z(1) = -0.5$$

The problem is to find approximations w_{11} to $y_1 = y(x_1)$ and w_{21} to $z_1 = z(x_1)$.

The C++ program for the solution is as follows.

```

1  #include <iostream>
2  #include <cmath>
3  #include <iomanip>
4  using namespace std ;
5
6  #define a      1.0           // starting point of domain
7  #define b      1.1           // ending point of domain
8  #define alpha1 0.5           //initial condition for first variable y1=y1(x)
9  #define alpha2 -0.5         // initial condition for second variable y2=y2(x)
10 #define m      5            // number of steps
11
12 #define f1(x,y,z) z
13 #define f2(x,y,z) (z-8*x*x*x*y*y*y)/x
14
15 void rk4system2(double [], double [], double[], double) ;
16
17
18 int main()
19 {
20     double h ;                // step size
21     double sol , err ;

```

```
22
23     int i ;
24     double x[m+1] ;           // vector of nodes
25     double w1[m+1] ;         // vector of solution values for y1=y1(x)
26     double w2[m+1] ;         // vector of solution values for z=y'
27
28     h = (b-a)/m ;
29
30     // forming vector of x-values or nodes
31     x[0] = a ;
32     x[m] = b ;
33
34     for (i=1 ; i<=m ; i++)
35     {
36         x[i] = x[i-1] + h ;    // x[i] = x[0] + i*h ;
37     }
38
39     w1[0] = alpha1 ;          // setting initial condition for y
40     w2[0] = alpha2 ;          // setting initial condition for z
41
42     // Call-by-reference to the solver of ODEs system of 2 equations
43     rk4system2(x , w1 , w2 , h) ;
44
45     // ----- Printing Solutions ----- //
46     for (i=0 ; i<=m ; i++)
47     {
48         cout<< "Node=" << setw(3) << i ;
49         cout << "\tx= " << setw(8) << x[i] ;
50         cout<< "\tw1=" << setprecision(10) << setw(13) << w1[i] ;
51         cout<< "\tw2=" << setprecision(10) << setw(13) << w2[i] <<
52         endl ;
53     }
54     return 0;
55 }
56
57 // --- User-defined function for ODE system of 2 equations using RK4 --- //
58
59 void rk4system2(double x[m+1], double w1[m+1], double w2[m+1],
60 double h)
61 {
```

```

61 double k11, k12, k13, k14 ;
62 double k21, k22, k23, k24 ;
63 int i ;
64
65 for (i=1 ; i<=m ; i++)
66 {
67     k11 = h * (f1(x[i-1], w1[i-1], w2[i-1])) ;
68     k21 = h * (f2(x[i-1], w1[i-1], w2[i-1])) ;
69
70     k12 = h*(f1((x[i-1]+0.5*h), (w1[i-1]+0.5*k11),
71                                     (w2[i-1]+0.5*k21))) ;
72     k22 = h * (f2((x[i-1]+0.5*h), (w1[i-1]+0.5*k11),
73                                     (w2[i-1]+0.5*k21))) ;
74     k13 = h * (f1((x[i-1]+0.5*h), (w1[i-1]+0.5*k12),
75                                     (w2[i-1]+0.5*k22))) ;
76
77     k14 = h * (f1((x[i-1]+h), (w1[i-1]+k13), (w2[i-1]+k23))) ;
78     k24 = h * (f2((x[i-1]+h), (w1[i-1]+k13), (w2[i-1]+k23))) ;
79
80     w1[i] = w1[i-1] + (k11 + 2*k12 + 2*k13 + k14)/6.0 ;
81     w2[i] = w2[i-1] + (k21 + 2*k22 + 2*k23 + k24)/6.0 ;
82 }
83 }

```

The results are shown in the following table.

Steps i	Node $x(i)$	Numerical solution $w_{1i} = w_1(x_i)$	Numerical solution $w_{2i} = w_2(x_i)$
0	1	0.5	-0.5
1	1.02	0.4897000998	-0.5299800347
2	1.04	0.4788015942	-0.5598407786
3	1.06	0.4673080456	-0.5894651108
4	1.08	0.4552253218	-0.618739865
5	1.1	0.4425615	-0.6475575603

■

Problem 18: Solve the ODE $y''' = -y'' + 3y' + 3y$ for $y = y(x)$ in $x \in [0,2]$ with the initial conditions: $y(0) = 2.0$, $y'(0) = -1.0$, and $y''(0) = 8.0$. Solve it for 10 steps.

Given the equation,

$$y''' = -y'' + 3y' + 3y \quad \text{---(1)}$$

For $y = y(x)$ in $x \in [0,2]$ with the initial conditions:

$$y(0) = 2.0$$

$$y'(0) = -1.0$$

$$y''(0) = 8.0$$

consider

$$y' = z_1 \quad \text{---(2)}$$

$$y'' = z_1' = z_2 \quad \text{---(3)}$$

Then, the given third-order Eq. (1) becomes

$$z_2' = -z_2 + 3z_1 + 3y \quad \text{---(4)}$$

Thus, the third-order IVP is essentially converted to the problem of a first-order system of ODEs of comprising the three equations (2) - (4) subject to the initial conditions:

$$y(0) = 2.0$$

$$z_1(0) = -1.0$$

$$z_2(0) = 8.0$$

For 10 steps, the domain is discretized as

$$h = \frac{b-a}{m} = \frac{2.0-0.0}{10} = 0.2$$

$x_0 = 0, x_1 = 0.2, x_2 = 0.4, x_3 = 0.6, x_4 = 0.8, x_5 = 1.0, x_6 = 1.2, x_7 = 1.4, x_8 = 1.6, x_9 = 1.8, x_{10} = 2.0$.

According to the initial conditions:

$$w_{10} = y_0 = y(x_0) = y(0) = 2.0$$

$$w_{20} = z_{10} = z_1(x_0) = z_1(0) = -1.0$$

$$w_{30} = z_{20} = z_2(x_0) = z_2(0) = 8.0$$

The problem is to find approximations w_{1i} to $y_i = y(x_i)$, w_{2i} to $z_{1i} = z_1(x_i)$, and w_{3i} to $z_{2i} = z_2(x_i)$, for $i = 1, 2, \dots, 10$.

The C++ program for the solution is as follows.

```

1  #include <iostream>
2  #include <cmath>
3  #include <iomanip>
4  using namespace std ;
5
6  #define a      0.0          // starting point of domain
7  #define b      2.0          // ending point of domain
8  #define alpha1 2.0          // initial condition: function value
9  #define alpha2 -1.0         // initial condition: first derivative value
10 #define alpha3 8.0          // initial condition: second derivative value
11 #define m      5            // number of steps
12
13 #define f1(x,y,z1,z2)  z1
13 #define f2(x,y,z1,z2)  z2
14 #define f3(x,y,z1,z2)  -z2 + 3*z1 + 3*y
15
16
17 void rk4system3(double [], double [], double[], double [], double) ;
18
19 int main()
20 {
21     double h ;                // step size
22     double sol , err ;
23
24     int i ;
25     double x[m+1] ;           // vector of nodes
26     double w1[m+1] ;          // vector of solution values for y=y(x)
27     double w2[m+1] ;          // vector of solution values for z1 = y'
28     double w3[m+1] ;          // vector of solution values for z2 = y''
29
30     h = (b-a)/m ;
31
32     // forming vector of x-values or nodes
33     x[0] = a ;
34     x[m] = b ;
35
36     for (i=1 ; i<=m ; i++)
37     {
38         x[i] = x[i-1] + h ;    // x[i] = x[0] + i*h ;
39     }
40
41     w1[0] = alpha1 ;           // setting initial condition for y
42     w2[0] = alpha2 ;           // setting initial condition for z1
43     w3[0] = alpha3 ;           // setting initial condition for z2
44
45     // Call-by-reference to the solver of ODEs system of 3 equations

```

```

46     rk4system3(x , w1 , w2 , w3 , h) ;
47
48 // ----- Printing Solutions ----- //
49     for (i=0 ; i<=m ; i++)
50     {
48         cout<< "Node=" << setw(3) << i ;
49         cout << "\tx= " << setw(8) << x[i] ;
50         cout<< "\tw1=" << setprecision(10) << setw(13) << w1[i] ;
51         cout<< "\tw2=" << setprecision(10) << setw(13) << w2[i] ;
52         cout<< "\tw3=" << setprecision(10) << setw(13) << w3[i] <<
           endl ;
53     }
54
55     return 0;
56 }
57
58 // --- User-defined function for ODE system of 3 equations using RK4 --- //
59
60 void rk4system3(double x[m+1], double w1[m+1],
61                double w2[m+1], double w3[m+1], double h)
62 {
63     double k11, k12, k13, k14 ;
64     double k21, k22, k23, k24 ;
65     double k31, k32, k33, k34 ;
66     int i ;
67
68     for (i=1 ; i<=m ; i++)
69     {
70         k11 = h * (f1(x[i-1], w1[i-1], w2[i-1], w3[i-1])) ;
71         k21 = h * (f2(x[i-1], w1[i-1], w2[i-1], w3[i-1])) ;
72         k31 = h * (f3(x[i-1], w1[i-1], w2[i-1], w3[i-1])) ;
73
74         k12 = h * (f1((x[i-1]+0.5*h), (w1[i-1]+0.5*k11),
75                       (w2[i-1]+0.5*k21), (w3[i-1]+0.5*k31))) ;
76
77         k22 = h * (f2((x[i-1]+0.5*h), (w1[i-1]+0.5*k11),
78                       (w2[i-1]+0.5*k21), (w3[i-1]+0.5*k31))) ;
79         k32 = h * (f3((x[i-1]+0.5*h), (w1[i-1]+0.5*k11),
80                       (w2[i-1]+0.5*k21), (w3[i-1]+0.5*k31))) ;
81
82         k13 = h * (f1((x[i-1]+0.5*h), (w1[i-1]+0.5*k12),
83                       (w2[i-1]+0.5*k22), (w3[i-1]+0.5*k32))) ;
84         k23 = h * (f2((x[i-1]+0.5*h), (w1[i-1]+0.5*k12),
85                       (w2[i-1]+0.5*k22), (w3[i-1]+0.5*k32))) ;
86         k33 = h * (f3((x[i-1]+0.5*h), (w1[i-1]+0.5*k12),

```

```

86             (w2[i-1]+0.5*k22), (w3[i-1]+0.5*k32))) ;
87
88     k14 = h * (f1((x[i-1]+h), (w1[i-1]+k13),
89                (w2[i-1]+k23), (w3[i-1]+k33))) ;
90     k24 = h * (f2((x[i-1]+h), (w1[i-1]+k13),
91                (w2[i-1]+k23), (w3[i-1]+k33))) ;
92     k34 = h * (f3((x[i-1]+h), (w1[i-1]+k13),
93                (w2[i-1]+k23), (w3[i-1]+k33))) ;
94
95     w1[i] = w1[i-1] + (k11 + 2*k12 + 2*k13 + k14)/6.0 ;
96     w2[i] = w2[i-1] + (k21 + 2*k22 + 2*k23 + k24)/6.0 ;
97     w3[i] = w3[i-1] + (k31 + 2*k32 + 2*k33 + k34)/6.0 ;
98 }
99 }
```

The results are shown in the following table.

Steps i	Node $x(i)$	Numerical solution $w_{1i} = w_1(x_i)$	Numerical solution $w_{2i} = w_2(x_i)$	Numerical solution $w_{3i} = w_3(x_i)$
0	0.0	2	-1	8
1	0.4	2.2144	2.0592	7.984
2	0.8	3.75537152	5.92358144	12.16498688
3	1.2	7.317856436	12.5913347	22.55617331
4	1.6	14.62815413	25.4339096	44.2884481
5	2.0	29.29652369	50.88508966	88.1604031

Step 4 (Applying the boundary conditions)

Set $w(0) = \text{alpha}$

Set $w(m + 1) = \text{beeta}$

Step 5 (Setting the initial conditions on interior nodes)

for $i = 1, 2, \dots, m$

Set $w(i) = 0$ (Constructing interior mesh points, x_i)

end for

Step 6

for $i = 1, 2, \dots, m$

Set $B(i) = -h \times h \times r(x(i));$ end for

for $i = 1, 2, \dots, m$

Set $D(i) = 2 + h \times h \times q(x(i));$ end for

for $i = 1, 2, \dots, m$

Set $U(i) = -1 + h \times 0.5 \times p(x(i));$ end for

for $i = 1, 2, \dots, m$

Set $L(i) = -1 - h \times 0.5 \times p(x(i));$ end for

Step 7 for $k = 1, 2, 3, \dots, N$ perform steps 8-11

Step 8

for $i = 1, 2, \dots, m$ Set $ZP(i) = W$ (keeping a copy of Z in ZP for taking the norm)

Step 9

for $i = 1, 2, \dots, m$ (compute the components of solution vector Z)

$$w(i) = \frac{B(i) - L(i) \times w(i - 1) - U(i) \times w(i + 1)}{D(i)}$$

end for

Step 10 Compute $err = \|W - ZP\|$

(or $err = \|X - XP\|/\|X\|$) Here $\|\cdot\|$ is any suitable norm.

Step 11

if ($err < TOL$) then } This means that the consecutive
Exit/Break the loop } approximations are nearly the same.
Therefore, stop iterations.

end for loop of Step 7 (Go to Step 8)

Step 12 Print the output: $W = [w_0, w_1, \dots, w_m]^T$; **STOP.**

Problem 19: Write a C++ program that uses a second-order accurate Finite Difference method to solve the following boundary value problem:

$$y'' = y' + 2y + \cos(x), \quad \text{for } y = y(x), \quad \text{where } 0 \leq x \leq \frac{\pi}{2}$$

subject to the following Dirichlet boundary conditions: $y(0) = -0.3$ and $y\left(\frac{\pi}{2}\right) = -0.1$.

For domain discretization, take step sizes as $h = \Delta x = \frac{\pi}{8}$

To form the computational domain, the physical domain $\left[0, \frac{\pi}{2}\right]$ is discretized by considering that it consists of a number of equispaced discrete points or nodes, x_i , for $i = 0, 1, 2, \dots, m + 1$. For the given problem,

$$\text{Number of interior nodes} = m = 3$$

$$p(x) = 1$$

$$q(x) = 2$$

$$r(x) = \cos(x)$$

The target is to obtain the approximations w_i to the function values $y_i = y(x_i)$ at the interior nodes x_i , for $i = 1, 2, 3$. The values of the solution function are known at x_0 and x_4 due to Dirichlet boundary conditions:

$$w_0 = y(x_0) = -0.3$$

$$w_4 = y(x_4) = -0.1$$

A C++ program that uses the Gauss-Seidel approach for the stated solution is as follows.

```

1  #include<iostream>
2  #include<cmath>
3  #include<iomanip>
4  using namespace std;
5
6  #define a  0.0
7  #define b  M_PI_2
8  #define alpha -0.3
9  #define beta  -0.1
10 #define m  3                //number of interior nodes
11 #define N  200              //maximum number of iterations
12 #define TOL  0.0000001
13
14 #define p(x) 1.0
15 #define q(x) 2.0
16 #define r(x) cos(x)

```

```
17
18 void egs(double [], double [], double [], double [], double [], double h) ;
19
20 int main()
21 {
22     double h , err , sum ;
23     double x[m+2] , z[m+2] , zp[m+2] ;
24     double B[m+1] , D[m+1] , U[m+1], L[m+1] ;
25     int i , k ;
26
27     h = (b-a)/(m+1) ;
28     x[0] = a ;
29     x[m+1] = b ;
30
31     for(i=1 ; i<=m ; i++)
32     {
33         x[i] = x[i-1]+h ;
34         cout<< "\tnodes" << x[i] << endl ;
35     }
36
37     z[0] = alpha ;
38     z[m+1] = beta ;
39
40     for (i=1 ; i<=m ;i++)
41         B[i] = (-h)*(h)*(r(x[i]));
42
43     for (i=1 ; i<=m ; i++)
44         D[i] = 2+((h)*(h)*(q(x[i])));
45
46     for (i=1 ; i<=m ; i++)
47         U[i] = -1.0 + ((h*0.5)*(p(x[i])));
48
49     for (i=1 ; i<=m ; i++)
50         L[i] = -1.0-((h*0.5)*(p(x[i])))) ;
51
52     for(i=1 ; i<=m ; i++)
53         z[i] = 0.0 ;
54
55     k = 0 ;
56     cout<< setw(4) << k << ": " ;
57     cout<< "z= " ;
```

```
48     cout<< setprecision(2) << setw(3) << z[0] << " ";
49     for (i=1 ; i<=m ; i++)
50     {
51         cout<< setprecision(8) << setw(9) << z[i] << " ";
52     }
53
54     cout<< setprecision(2) << setw(3) << z[m+1] << " ";
55     cout << endl ;
56
57     egs(z, B, D, U, L, h) ;
58
59     return 0 ;
60 }
61
62 // --- User-defined function for an efficient Gauss-Seidel method --- //
63
64 void egs(double z[m+2], double B[m+1] , double D[m+1] ,
65         double U[m+1], double L[m+1], double h)
66 {
67     int i, k ;
68     double sum , err ;
69     double zp[m+2] ;
70
71
72     for (k=1 ; k<=N ; k++)          // Iterations loop
73     {
74         // Making a copy of the solution vector before updating it
75         for (i=1 ; i<=m ; i++)
76             zp[i] = z[i] ;
77
78         // Updating the solution vector
79         for (i=1 ; i<=m ; i++)
80         {
81             z[i] = (B[i]-(L[i]*z[i-1])-(U[i]*z[i+1]))/D[i] ;
82         }
83
84         // Printing the latest solution vector
85         cout<< setw(4) << k << ": " ;
86         cout<< "z= " ;
87         cout<< setprecision(2) << setw(3) << z[0] << " " ;
```

```

88     for (i=1 ; i<=m ; i++)
89     {
90         cout<< setprecision(8) << setw(9) << z[i] << " " ;
91     }
92     cout<< setprecision(2) << setw(3) << z[m+1] << " " ;
93     cout << endl ;
94
95     // Finding the error as the L2-norm
96     sum = 0.0 ;
97     for (i=1 ; i<=m ; i++)
98         sum = sum + (z[i]-zp[i])*(z[i]-zp[i])/(z[i]*z[i]) ;
99     err = sqrt(sum) ;
100
101     // Testing the convergence
102     if (err<TOL) break ;
103 }
104 }
```

■

Chapter Summary

- The numerical solution of an ODE is not a definition of $y = y(x)$. The numerical solution of the ODE is a set of numbers w_i that are approximations to the function values $y(x_i)$ at some pre-specified discrete values $x_i \in [a, b]$. That is, $w_i \cong y_i = y(x_i)$.
- To solve an initial-value problem consisting of a single first-order ODE in $y = y(x)$ for $a \leq x \leq b$ and an initial-value $y(a) = \alpha$, first the domain $[a, b]$ is discretized by selecting $(m + 1)$ equispaced nodes $x_0, x_1, x_2, \dots, x_m$ in $[a, b]$ such that $a = x_0 < x_1 < x_2 < \dots < x_m = b$, and $h = (b - a)/m$. Then, approximations w_i to the values $y_i = y(x_i)$ for $i = 1, 2, \dots, m$ are obtained with $w_0 = y(a)$. For simplicity, $y(x_i)$ is denoted by y_i .
- There is a wide variety of methods for finding numerical solutions of the ODEs involved in initial value problems (IVPs) and boundary value problems (BVPs).
- Methods for IVPs include single step methods and multi-step methods, each category having explicit and implicit methods. A hybrid method, i.e., predictor-corrector method, involves a combination of explicit and implicit formulas.
- Methods for BVPs are so versatile and involve much richer mathematical constructs.

- The accuracy of the approximate solution can be improved either by using a larger number of steps (a smaller step size), or by using a better numerical method.
- The prime characteristics (or considerations) associated with a finite difference scheme to determine its quality include
 - Stability
 - Local Truncation Error
 - Consistency (Compatibility)
 - Discretization Error
 - Convergence

■ ■ ■

Chapter Exercises

Exercise 01: Find the numerical solution of the ODE, $y' = 3 - 3y - e^{-6x}$, for $0 \leq x \leq 2$, with initial condition $y(0) = 1.0$. Consider the step size of 0.5. Use the exact solution, $y(x) = \frac{1}{3}(e^{-6x} - e^{-3x} + 3)$, to find the error in the numerical solution.

Exercise 02: Find the numerical solution of the ODE, $y' = 1 + (x - y)^2$, for $2 \leq x \leq 3$, with initial condition $y(2) = 1.0$. Consider the step size of 0.5. Use the exact solution, $y(x) = x + 1/(1 - x)$, to find the error in the numerical solution.

Exercise 03: Find the numerical solution of the ODE, $y' = 2 + (x - y)^2$, for $2 \leq x \leq 3$, with initial condition $y(2) = 1.5$. Consider the step size of 0.5. Use the exact solution, $y(x) = x - \tan(-x + 2.463)$, to find the

Exercise 04: Find the numerical solution of the ODE, $y' = (1 + x)/(1 + y)$, for $0 \leq x \leq 1$, with initial condition $y(0) = 2.0$. Consider the step size of 0.5. Use the exact solution, $y(x) = \sqrt{x^2 + 2x + 9} - 1$, to find the error in the numerical solution.

Exercise 05: For the functions $y_1 = y_1(x)$ and $y_2 = y_2(x)$, where $x \in [0,1]$, solve the following system of two ODEs:

$$y_1' = y_1 y_2 - 2$$

$$y_2' = 2y_1 - y_2^3$$

With initial conditions:

$$y_1(0) = 2.0$$

$$y_2(0) = 0.3$$

Use the RK4 method of order 4 for 5 steps.

HINT: For 5 steps the domain is discretized as

$$h = \frac{b-a}{m} = \frac{1.0-0.0}{5} = 0.2$$

$$x_0 = 0, x_1 = 0.2, x_2 = 0.4, x_3 = 0.6, x_4 = 0.8, x_5 = 1.0.$$

According to the initial conditions:

$$w_{10} = y_{10} = y_1(x_0) = y_1(0) = 2.0$$

$$w_{20} = y_{20} = y_2(x_0) = y_2(0) = 0.3$$

The problem is to find approximations w_{1i} to $y_{1i} = y_1(x_i)$ and w_{2i} to $y_{2i} = y_2(x_i)$, for $i = 1, 2, 3, 4, 5$.

Exercise 06: For the functions $y_1 = y_1(x)$, $y_2 = y_2(x)$, and $y_3 = y_3(x)$, where $x \in [0, 1]$, solve the following system of three ODEs:

$$y_1' = y_1 + 3y_2 - 3y_3 + e^{-x}$$

$$y_2' = 2y_2 + y_3 - 3e^{-x}$$

$$y_3' = y_1 + 2y_2 + e^{-x}$$

with initial conditions:

$$y_1(0) = 2.5$$

$$y_2(0) = -1.5$$

$$y_3(0) = -1.0$$

Use the RK4 method of order 4 for 10 steps.

HINT: For 10 steps, the domain is discretized as

$$h = \frac{b-a}{m} = \frac{1.0-0.0}{10} = 0.1$$

$$x_0 = 0, x_1 = 0.1, x_2 = 0.2, x_3 = 0.3, x_4 = 0.4, x_5 = 0.5, x_6 = 0.6, x_7 = 0.7, x_8 = 0.8, x_9 = 0.9, x_{10} = 1.0.$$

According to the initial conditions:

$$w_{10} = y_{10} = y_1(x_0) = y_1(0) = 2.5$$

$$w_{20} = y_{20} = y_2(x_0) = y_2(0) = -1.5$$

$$w_{30} = y_{30} = y_3(x_0) = y_3(0) = -1.0$$

The problem is to find approximations w_{1i} to $y_{1i} = y_1(x_i)$, w_{2i} to $y_{2i} = y_2(x_i)$, and w_{3i} to $y_{3i} = y_3(x_i)$, for $i = 1, 2, \dots, 10$.

Exercise 07: Find the numerical solution of the IVP, $y'' - 8y' + 7y = 16e^{-x}$ for $0 \leq x \leq 1$, with initial condition $y(0) = 4.0$ and $y'(0) = 4.0$. Also find $y(1.1)$. Consider the step size of 0.1. Use the exact solution, $y = (1/3)(e^{7x} + 8e^x + 3e^{-x})$, to find the error in the numerical solution.

HINT: Given the equation,

$$y'' - 8y' + 7y = 16e^{-x} \quad \text{--- (1)}$$

For the solution, consider

$$y' = z \quad \text{--- (2)}$$

Then, the given second-order Eq. (1) becomes

$$z' = 4z - 3y + 7e^{-x} \quad \text{--- (3)}$$

Thus, the second-order IVP is essentially converted to the problem of a first-order system of ODEs of comprising the two equations (2) and (3) subject to the initial conditions:

$$y(0) = 3.0$$

$$z(0) = 3.0$$

For 10 steps, the domain is discretized as

$$h = \frac{b-a}{m} = \frac{1.0-0.0}{10} = 0.1$$

$x_0 = 0, x_1 = 0.1, x_2 = 0.2, x_3 = 0.3, x_4 = 0.4, x_5 = 0.5, x_6 = 0.6, x_7 = 0.7, x_8 = 0.8, x_9 = 0.9, x_{10} = 1.0$.

According to the initial conditions:

$$w_{10} = y_0 = y(x_0) = y(0) = 3.0$$

$$w_{20} = z_0 = z(x_0) = z(0) = 3.0$$

The problem is to find approximations w_{1i} to $y_i = y_i(x_i)$ and w_{2i} to $z_i = z(x_i)$, for $i = 1, 2, \dots, 10$.

Exercise 08: Solve the ODE $y'' = 4y' - 3y + 7e^{-x}$ for $y = y(x)$ in $x \in [0,1]$ with the initial conditions: $y(0) = 3.0$ and $y'(0) = 3.0$. Solve it for 10 steps.

Exercise 09: Find the numerical solution of the BVP, $y'' - 9y' + y = x$ for $0 \leq x \leq 1$, with initial condition $y(0) = 0.0$ and $y'(1) = 6.0$. Consider the step size of 0.1.

Exercise 10: Find the numerical solution of the ODE, $x^2y'' + 3xy' + 3y = 0$, with initial condition $y(1) = 1$ and $y'(1) = -5$ for $y(1.1)$. The exact solution is, $y = \frac{1}{x} \left(\cos(\sqrt{2} \ln x) + \left(\frac{1}{x^2} - 5 \right) \sin(\sqrt{2} \ln x) \right)$.

Exercise 11: Find the numerical solution of the ODE, $y'' - 6y' + 9y = x^2e^{3x}$, with initial condition $y(0) = 2$ and $y'(0) = 6$ for $y(1.1)$. The exact solution is, $y = 2e^{3x} + \frac{1}{12}x^4e^{3x}$.

Exercise 12: Solve the ODE $y''' = -y'' + 3y' + 3y$ for $y = y(x)$ in $x \in [0,2]$ with the initial conditions: $y(0) = 2.0, y'(0) = -1.0$, and $y''(0) = 8.0$. Solve it for 10 steps.

HINT: Given the equation,

$$y''' = -y'' + 3y' + 3y \quad \text{--- (1)}$$

For $y = y(x)$ in $x \in [0,2]$ with the initial conditions:

$$y(0) = 2.0$$

$$y'(0) = -1.0$$

$$y''(0) = 8.0$$

Solve it for 10 steps.

For the solution, consider

$$y' = z_1 \quad \text{---(2)}$$

$$y'' = z'_1 = z_2 \quad \text{---(3)}$$

Then, the given third-order Eq. (1) becomes

$$z'_2 = -z_2 + 3z_1 + 3y \quad \text{---(3)}$$

Thus, the third-order IVP is essentially converted to the problem of a first-order system of ODEs of comprising the three equations (2) - (4) subject to the initial conditions:

$$y(0) = 2.0$$

$$z_1(0) = -1.0$$

$$z_2(0) = 8.0$$

For 10 steps, the domain is discretized as

$$h = \frac{b-a}{m} = \frac{2.0-0.0}{10} = 0.2$$

$$x_0 = 0, x_1 = 0.2, x_2 = 0.4, x_3 = 0.6, x_4 = 0.8, x_5 = 1.0, x_6 = 1.2, x_7 = 1.4, x_8 = 1.6, x_9 = 1.8, x_{10} = 2.0.$$

According to the initial conditions:

$$w_{10} = y_0 = y(x_0) = y(0) = 2.0$$

$$w_{20} = z_{10} = z_1(x_0) = z_1(0) = -1.0$$

$$w_{30} = z_{20} = z_2(x_0) = z_2(0) = 8.0$$

The problem is to find approximations w_{1i} to $y_i = y(x_i)$, w_{2i} to $z_{1i} = z_1(x_i)$, and w_{3i} to $z_{2i} = z_2(x_i)$, for $i = 1, 2, \dots, 10$.

Exercise 13: Using a second-order accurate Finite Difference method, solve the following BVP:

$$y'' = 9y' - y + x, \quad \text{for } y = y(x), \quad \text{where } 0 \leq x \leq 1$$

subject to the following Dirichlet boundary conditions: $y(0) = 0$ and $y(1) = 6$.

For domain discretization, take step sizes as $h = \Delta x = 0.25$.

Exercise 14: Using a second-order accurate Finite Difference method, solve the following BVP:

$$y'' = -5y' - 8y + x^2, \quad \text{for } y = y(x), \quad \text{where } 1 \leq x \leq 2$$

subject to the following Dirichlet boundary conditions: $y(1) = 0$ and $y(2) = 24$. For domain discretization, take step sizes as $h = \Delta x = 0.25$.

Bibliography

1. Richard L. Burden & J. Douglas Faires, (2011), Numerical Analysis, (9th Edition), USA, Brooks/Cole Pub. Co.
2. Steven C. Chapra & Raymond P. Canale, (2006), Numerical Methods for Engineers, (5th Edition), NY, USA, McGraw-Hill Co.
3. David R. Kincaid & E. Ward Cheney, (2002), Numerical Analysis: Mathematics of Scientific Computing, (3rd Edition), USA, Brooks/Cole Pub. Co.
4. E. Ward Cheney & David R. Kincaid, (2013), Numerical Mathematics and Computing, (7th Edition), New-Delhi India, Cengage Learning India Pvt. Ltd.
5. Brian Bradie, (2005), A Friendly Introduction to Numerical Analysis, Pearson.
6. John H. Mathews & Kurtis D. Fink, (2015), Numerical Methods using MATLAB, (4th Edition), India, Pearson India Education Services Pvt. Ltd.
7. M. K. Jain, S. R. K. Iyengar & R. K. Jain, (2012), Numerical Methods for Scientific and Engineering Computation, (6th Edition), New-Delhi India, New Age International Pvt. Ltd.
8. George R. Lindfield & John E. T. Penny, (2013), Numerical Methods using MATLAB, (3rd Edition), USA, Academic Press, An imprint of Elsevier.
9. Amos Gillat, (2011), MATLAB: An Introduction with Applications, (4th Edition), USA, John Wiley & Sons, Inc.
10. Laurene V. Fausett, (2009), Applied Numerical Analysis using MATLAB, (2nd Edition), India, PEARSON Education Inc.
11. Babu ram, (2010), Numerical Methods, India, PEARSON Education Inc.
12. Francis Schied, (1990), 2000 Solved Problems in Numerical Analysis, (International Edition), NY, USA, McGraw-Hill Co.
13. P. Siva Ramakrishna Das & C. Vijayakumari, (2004), Numerical Analysis, (1st Edition), India, Dorling Kindersley Pvt. Ltd.
14. Saeed Akhtar Bhatti & Naveed Akhtar Bhatti, (2008), A First Course in Numerical Analysis with C++, (5th Edition), Lahore, Pakistan, A-ONE Publishers.

15. Mohammad Iqbal, (1990), An Introduction to Numerical Analysis, Urdu Bazar Lahore, Pakistan, Ilmi Kitab Khana.
16. Fiaz Ahmad & Muhammad Afzal Rana, (1995), Elements of Numerical Analysis, Islamabad, Pakistan, National Book Foundation
17. Amjad Pervez, (1996), An Introduction to Numerical Analysis, Urdu Bazar Lahore, Pakistan, A.H. Publishers.
18. Germund Dahlquist & Ake Björck, (2003), Numerical Methods, New Jersey, USA, Prentice-Hall Inc.
19. Erwin Kreyszig, (2011), Advanced Engineering Mathematics, (10th Edition), USA, John Wiley & Sons, Inc.
20. S. S. Sastry, (2019), Introductory Methods of Numerical Analysis, (Fifth Edition), PHI Learning Private Limited.
21. Curtis F. Gerald & Patrick O. Wheatley, (2003), Applied Numerical Analysis, (7th Edition), India, PEARSON Education Inc.
22. K. Sankara Rao, (2009), Numerical Methods for Scientists and Engineers, (Third Edition), PHI Learning Private Limited.
23. Lal Din Baig, (2014), Numerical Analysis, Ilmi Kitab Khana, Lahore.