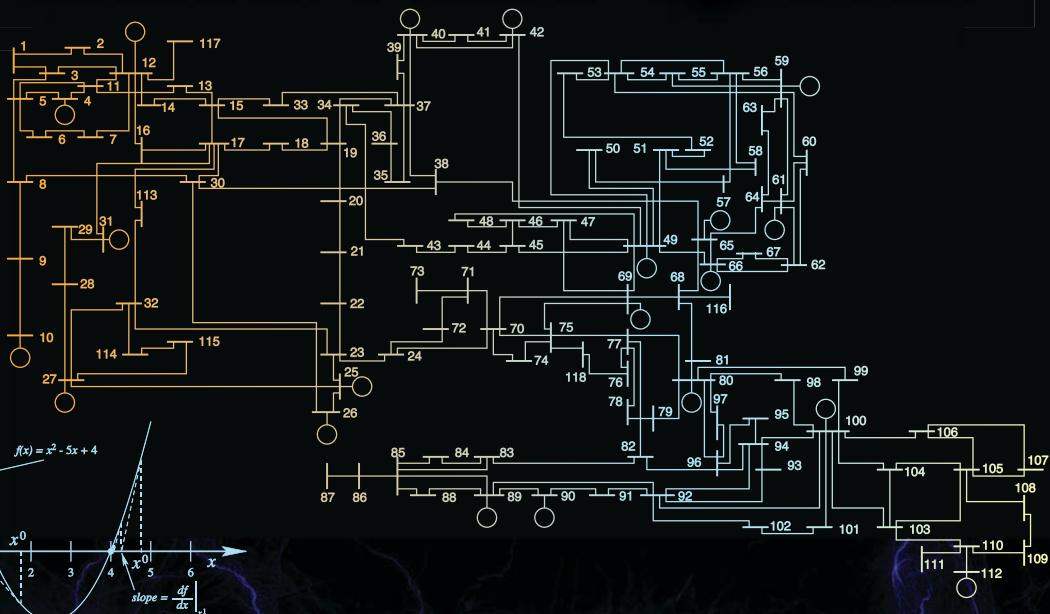


Third Edition

Computational Methods for Electric Power Systems



Mariesa L. Crow



CRC Press
Taylor & Francis Group

Third Edition

Computational Methods for Electric Power Systems

The ELECTRIC POWER ENGINEERING Series

Series Editor Leo L. Grigsby

Published Titles

Computational Methods for Electric Power Systems, Third Edition, Mariesa L. Crow

Electric Power Generation, Transmission, and Distribution, Third Edition, Leonard L. Grigsby

Linear Synchronous Motors: Transportation and Automation Systems, Second Edition,

Jacek Gieras, Zbigniew J. Piech, and Bronislaw Tomczuk

The Induction Machines Design Handbook, Second Edition, Ion Boldea and Syed Nasar

Electric Energy Systems: Analysis and Operation, Antonio Gómez-Expósito,

Antonio J. Conejo, and Claudio Cañizares

Distribution System Modeling and Analysis, Second Edition, William H. Kersting

Electric Machines, Charles A. Gross

Harmonics and Power Systems, Francisco C. De La Rosa

Electric Drives, Second Edition, Ion Boldea and Syed Nasar

Power System Operations and Electricity Markets, Fred I. Denny and David E. Dismukes

Power Quality, C. Sankaran

Electromechanical Systems, Electric Machines, and Applied Mechatronics,

Sergey E. Lyshevski

Electrical Energy Systems, Second Edition, Mohamed E. El-Hawary

Electric Power Substations Engineering, John D. McDonald

Electric Power Transformer Engineering, James H. Harlow

Electric Power Distribution Handbook, Tom Short

Third Edition

Computational Methods for Electric Power Systems

Mariesa L. Crow

Missouri University of Science and Technology, Rolla, USA



CRC Press
Taylor & Francis Group
Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group, an **informa** business

MATLAB® is a trademark of The MathWorks, Inc. and is used with permission. The MathWorks does not warrant the accuracy of the text or exercises in this book. This book's use or discussion of MATLAB® software or related products does not constitute endorsement or sponsorship by The MathWorks of a particular pedagogical approach or particular use of the MATLAB® software.

CRC Press
Taylor & Francis Group
6000 Broken Sound Parkway NW, Suite 300
Boca Raton, FL 33487-2742

© 2016 by Taylor & Francis Group, LLC
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works
Version Date: 20150629

International Standard Book Number-13: 978-1-4987-1160-9 (eBook - PDF)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Visit the Taylor & Francis Web site at
<http://www.taylorandfrancis.com>

and the CRC Press Web site at
<http://www.crcpress.com>

To
Jim, David, and Jacob

Contents

Preface to the third edition	xi
1 Introduction	1
2 The Solution of Linear Systems	3
2.1 Gaussian Elimination	4
2.2 LU Factorization	9
2.2.1 LU Factorization with Partial Pivoting	16
2.2.2 LU Factorization with Complete Pivoting	20
2.3 Condition Numbers and Error Propagation	23
2.4 Stationary Iterative Methods	24
2.5 Conjugate Gradient Methods	30
2.6 Generalized Minimal Residual Algorithm	36
2.7 Preconditioners for Iterative Methods	42
2.7.1 Jacobi	42
2.7.2 Symmetric Successive Overrelaxation	43
2.7.3 Symmetric Gauss–Seidel	43
2.7.4 Incomplete LU Factorization	43
2.7.5 Graph Based	44
2.8 Problems	47
3 Systems of Nonlinear Equations	53
3.1 Fixed-Point Iteration	54
3.2 Newton–Raphson Iteration	61
3.2.1 Convergence Properties	64
3.2.2 The Newton–Raphson for Systems of Nonlinear Equations	65
3.3 Quasi-Newton Methods	68
3.3.1 Secant Method	69
3.3.2 Broyden’s Method	72
3.3.3 Modifications to the Newton–Raphson Method	74
3.3.4 Numerical Differentiation	75
3.3.5 Newton–GMRES	79
3.4 Continuation Methods	83
3.5 Power System Applications	86
3.5.1 Power Flow	87

3.5.2	Regulating Transformers	95
3.5.3	Decoupled Power Flow	99
3.5.4	Fast Decoupled Power Flow	101
3.5.5	PV Curves and Continuation Power Flow	105
3.5.6	Three-Phase Power Flow	112
3.6	Problems	113
4	Sparse Matrix Solution Techniques	117
4.1	Storage Methods	118
4.2	Sparse Matrix Representation	127
4.3	Ordering Schemes	127
4.3.1	Scheme 0	141
4.3.2	Scheme I	142
4.3.3	Scheme II	148
4.3.4	Other Schemes	151
4.4	Power System Applications	152
4.5	Problems	156
5	Numerical Integration	163
5.1	One-Step Methods	164
5.1.1	Taylor Series-Based Methods	164
5.1.2	Forward Euler Method	165
5.1.3	Runge–Kutta Methods	165
5.2	Multistep Methods	166
5.2.1	Adams Methods	172
5.2.2	Gear’s Methods	175
5.3	Accuracy and Error Analysis	176
5.4	Numerical Stability Analysis	180
5.5	Stiff Systems	187
5.6	Step Size Selection	191
5.7	Differential-Algebraic Equations	198
5.8	Power System Applications	200
5.8.1	Transient Stability Analysis	200
5.8.2	Midterm Stability Analysis	208
5.9	Problems	211
6	Optimization	219
6.1	Least Squares State Estimation	220
6.1.1	Weighted Least Squares Estimation	223
6.1.2	Bad Data Detection	226
6.1.3	Nonlinear Least Squares State Estimation	229
6.2	Linear Programming	230
6.2.1	Simplex Method	231
6.2.2	Interior Point Method	235
6.3	Nonlinear Programming	240

6.3.1	Quadratic Programming	241
6.3.2	Steepest Descent Algorithm	243
6.3.3	Sequential Quadratic Programming Algorithm	248
6.4	Power System Applications	251
6.4.1	Optimal Power Flow	251
6.4.2	State Estimation	262
6.5	Problems	266
7	Eigenvalue Problems	273
7.1	The Power Method	274
7.2	The QR Algorithm	276
7.2.1	Deflation	283
7.2.2	Shifted QR	283
7.2.3	Double Shifted QR	284
7.3	Arnoldi Methods	286
7.4	Singular Value Decomposition	293
7.5	Modal Identification	296
7.5.1	Prony Method	298
7.5.2	The Matrix Pencil Method	301
7.5.3	The Levenberg–Marquardt Method	302
7.5.4	Eigensystem Realization Algorithm	305
7.5.5	Examples	306
7.6	Power System Applications	311
7.6.1	Participation Factors	311
7.7	Problems	312
References		315
Index		321

Preface to the third edition

This book is intended for a graduate level course. The material is first presented in a general algorithmic manner followed by power system applications. Users do not necessarily have to have a power systems background to find this book useful, but many of the comprehensive exercises do require a working knowledge of power system problems and notation.

This new edition has been updated to include new material. Specifically, this new edition has added the following material:

- Updated examples on sparse LU factorization
- Preconditioners for linear iterative methods
- Broyden's method
- Jacobian free Newton–Krylov methods
- Double-shift method for computing complex eigenvalues
- Eigensystem Realization Algorithm

and additional problems and examples.

A course structure would typically include the following chapters in sequence: Chapters 1, 2, and 3. Chapter 2 provides a basic background in linear system solution (both direct and iterative) followed by a discussion of nonlinear system solution in Chapter 3. Chapter 2 can be directly followed by Chapter 4, which covers sparse storage and computation and follows directly from LU factorization. Chapters 5, 6, and 7 can be covered in any order after Chapter 3 depending on the interest of the reader.

Many of the methods presented in this book have commercial software packages that will accomplish their solution far more rigorously with many failsafe attributes included (such as accounting for ill conditioning, etc.). It is not my intent to make students experts in each topic, but rather to develop an appreciation for the methods behind the packages. Many commercial packages provide default settings or choices of parameters for the user; through better understanding of the methods driving the solution, informed users can make better choices and have a better understanding of the situations in which the methods may fail. If this book provides any reader with more confidence in using commercial packages, I have succeeded in my intent.

Mariesa L. Crow
Rolla, Missouri

Lecture notes are available from the CRC Web site:
<http://www.crcpress.com/product/isbn/9781498711593>

MATLAB is a registered trademark of The MathWorks, Inc. For product information, please contact:

The MathWorks, Inc.
3 Apple Hill Drive
Natick, MA 01760-2098 USA
Tel: 508 647 7000
Fax: 508 647 7001
E-mail: info@mathworks.com
Web: www.mathworks.com

1

Introduction

In today's deregulated environment, the nation's electric power network is being forced to operate in a manner for which it was not intentionally designed. Therefore, system analysis is very important to predict and continually update the operating status of the network. This includes estimating the current power flows and bus voltages (power flow analysis and state estimation), determining the stability limits of the system (continuation power flow, numerical integration for transient stability, and eigenvalue analysis), and minimizing costs (optimal power flow). This book provides an introductory study of the various computational methods that form the basis of many analytical studies in power systems and other engineering and science fields. This book provides the analytical background of the algorithms used in numerous commercial packages. By understanding the theory behind many of the algorithms, the reader/user can better use the software and make more informed decisions (i.e., choice of integration method and step size in simulation packages).

Due to the sheer size of the power grid, hand-based calculations are nearly impossible, and computers offer the only truly viable means for system analysis. The power industry is one of the largest users of computer technology and one of the first industries to embrace the potential of computer analysis when mainframes first became available. Although the first algorithms for power system analysis were developed in the 1940s, it wasn't until the 1960s that computer usage became widespread within the power industry. Many of the analytical techniques and algorithms used today for the simulation and analysis of large systems were originally developed for power system applications.

As power systems increasingly operate under stressed conditions, computer simulation will play a large role in control and security assessment. Commercial packages routinely fail or give erroneous results when used to simulate stressed systems. Understanding of the underlying numerical algorithms is imperative to correctly interpret the results of commercial packages. For example, will the system really exhibit the simulated behavior or is the simulation simply an artifact of a numerical inaccuracy? The educated user can make better judgments about how to compensate for numerical shortcomings in such packages, either by better choice of simulation parameters or by posing the problem in a more numerically tractable manner. This book will provide the background for a number of widely used numerical algorithms that underlie many commercial packages for power system analysis and design.

This book is intended to be used as a text in conjunction with a semester-long graduate level course in computational algorithms. While the majority of examples in this text are based on power system applications, the theory is presented in a general manner so as to be applicable to a wide range of engineering systems. Although some knowledge of power system engineering may be required to fully appreciate the subtleties of some of the illustrations, such knowledge is not a prerequisite for understanding the algorithms themselves. The text and examples are used to provide an introduction to a wide range of numerical methods without being an exhaustive reference. Many of the algorithms presented in this book have been the subject of numerous modifications and are still the object of on-going research. As this text is intended to provide a foundation, many of these new advances are not explicitly covered, but are rather given as references for the interested reader. The examples in this text are intended to be simple and thorough enough to be reproduced easily. Most “real world” problems are much larger in size and scope, but the methodologies presented in this text should sufficiently prepare the reader to cope with any difficulties he/she may encounter.

Most of the examples in this text were produced using code written in MATLAB[®]. Although this was the platform used by the author, in practice, any computer language may be used for implementation. There is no practical reason for a preference for any particular platform or language.

2

The Solution of Linear Systems

In many branches of engineering and science it is desirable to be able to mathematically determine the state of a system based on a set of physical relationships. These physical relationships may be determined from characteristics such as circuit topology, mass, weight, or force, to name a few. For example, the injected currents, network topology, and branch impedances govern the voltages at each node of a circuit. In many cases, the relationship between the known, or input, quantities and the unknown, or output, states is a linear relationship. Therefore, a linear system may be generically modeled as

$$Ax = b \quad (2.1)$$

where b is the $n \times 1$ vector of known quantities, x is the $n \times 1$ unknown state vector, and A is the $n \times n$ matrix that relates x to b . For the time being, it will be assumed that the matrix A is invertible, or non-singular; thus each vector b will yield a unique corresponding vector x . Thus the matrix A^{-1} exists and

$$x^* = A^{-1}b \quad (2.2)$$

is the unique solution to Equation (2.1).

The natural approach to solving Equation (2.1) is to directly calculate the inverse of A and multiply it by the vector b . One method to calculate A^{-1} is to use *Cramer's rule*:

$$A^{-1}(i,j) = \frac{1}{\det(A)} (A_{ij})^T \quad \text{for } i = 1, \dots, n, j = 1, \dots, n \quad (2.3)$$

where $A^{-1}(i,j)$ is the ij th entry of A^{-1} and A_{ij} is the cofactor of each entry a_{ij} of A . This method requires the calculation of $(n+1)$ determinants, which results in $2(n+1)!$ multiplications to find $A^{-1}!$. For large values of n , the calculation requirement grows too rapidly for computational tractability; thus alternative approaches have been developed.

Basically, there are two approaches to solving Equation (2.1):

- *Direct methods*, or elimination methods, find the exact solution (within the accuracy of the computer) through a finite number of arithmetic operations. The solution x of a direct method would be completely accurate were it not for computer roundoff errors.

- *Iterative methods*, on the other hand, generate a sequence of (hopefully) progressively improving approximations to the solution based on the application of the same computational procedure at each step. The iteration is terminated when an approximate solution is obtained having some prespecified accuracy or when it is determined that the iterates are not improving.

The choice of solution methodology usually relies on the structure of the system under consideration. Certain systems lend themselves more amenable to one type of solution method versus the other. In general, direct methods are best for full matrices, whereas iterative methods are better for matrices that are large and sparse. But, as with most generalizations, there are notable exceptions to this rule of thumb.

2.1 Gaussian Elimination

An alternate method for solving Equation (2.1) is to solve for x without calculating A^{-1} explicitly. This approach is a *direct method* of linear system solution, since x is found directly. One common direct method is the method of *Gaussian elimination*. The basic idea behind Gaussian elimination is to use the first equation to eliminate the first unknown from the remaining equations. This process is repeated sequentially for the second unknown, the third unknown, etc., until the elimination process is completed. The n th unknown is then calculated directly from the input vector b . The unknowns are then recursively substituted back into the equations until all unknowns have been calculated.

Gaussian elimination is the process by which the augmented $n \times (n + 1)$ matrix

$$[A \mid b]$$

is converted to the $n \times (n + 1)$ matrix

$$[I \mid b^*]$$

through a series of elementary row operations, where

$$\begin{aligned} Ax &= b \\ A^{-1}Ax &= A^{-1}b \\ Ix &= A^{-1}b = b^* \\ x^* &= b^* \end{aligned}$$

Thus, if a series of elementary row operations exist that can transform the matrix A into the identity matrix I , then the application of the same set of

elementary row operations will also transform the vector b into the solution vector x^* .

An elementary row operation consists of one of three possible actions that can be applied to a matrix:

- interchange any two rows of the matrix
- multiply any row by a constant
- take a linear combination of rows and add it to another row

The elementary row operations are chosen to transform the matrix A into an upper triangular matrix that has ones on the diagonal and zeros in the subdiagonal positions. This process is known as the *forward elimination* step. Each step in the forward elimination can be obtained by successively premultiplying the matrix A by an elementary matrix ξ , where ξ is the matrix obtained by performing an elementary row operation on the identity matrix.

Example 2.1

Find a sequence of elementary matrices that, when applied to the following matrix, will produce an upper triangular matrix.

$$A = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 2 & 1 & 2 & 3 \\ 4 & 3 & 5 & 8 \\ 9 & 2 & 7 & 4 \end{bmatrix}$$

Solution 2.1 To upper triangularize the matrix, the elementary row operations will need to systematically zero out each column below the diagonal. This can be achieved by replacing each row of the matrix below the diagonal with the difference of the row itself and a constant times the diagonal row, where the constant is chosen to result in a zero sum in the column under the diagonal. Therefore, row 2 of A is replaced by (row 2 minus 2(row 1)) and the elementary matrix is

$$\xi_1 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ -2 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

and

$$\xi_1 A = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 0 & -5 & -6 & -13 \\ 4 & 3 & 5 & 8 \\ 9 & 2 & 7 & 4 \end{bmatrix}$$

Note that all rows except row 2 remain the same and row 2 now has a 0 in the column under the first diagonal. Similarly, the two elementary matrices that complete the elimination of the first column are

$$\xi_2 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ -4 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\xi_3 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ -9 & 0 & 0 & 1 \end{bmatrix}$$

and

$$\xi_3 \xi_2 \xi_1 A = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 0 & -5 & -6 & -13 \\ 0 & -9 & -11 & -24 \\ 0 & -25 & -29 & -68 \end{bmatrix} \quad (2.4)$$

The process is now applied to the second column to zero out everything below the second diagonal and scale the diagonal to one. Therefore,

$$\xi_4 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & -\frac{9}{5} & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\xi_5 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & -\frac{25}{5} & 0 & 1 \end{bmatrix}$$

$$\xi_6 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & -\frac{1}{5} & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Similarly,

$$\xi_6 \xi_5 \xi_4 \xi_3 \xi_2 \xi_1 A = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 0 & 1 & \frac{6}{5} & \frac{13}{5} \\ 0 & 0 & -\frac{5}{5} & -\frac{3}{5} \\ 0 & 0 & 1 & -3 \end{bmatrix} \quad (2.5)$$

Similarly,

$$\xi_7 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 5 & 1 \end{bmatrix}$$

$$\xi_8 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & -5 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

yielding

$$\xi_8 \xi_7 \xi_6 \xi_5 \xi_4 \xi_3 \xi_2 \xi_1 A = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 0 & 1 & \frac{6}{5} & \frac{13}{5} \\ 0 & 0 & 1 & 3 \\ 0 & 0 & 0 & -6 \end{bmatrix} \quad (2.6)$$

Finally

$$\xi_9 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -\frac{1}{6} \end{bmatrix}$$

and

$$\xi_9 \xi_8 \xi_7 \xi_6 \xi_5 \xi_4 \xi_3 \xi_2 \xi_1 A = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 0 & 1 & \frac{6}{5} & \frac{13}{5} \\ 0 & 0 & 1 & 3 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (2.7)$$

which completes the upper triangularization process. ■

Once an upper triangular matrix has been achieved, the solution vector x^* can be found by successive substitution (or *back substitution*) of the states.

Example 2.2

Using the upper triangular matrix of Example 2.1, find the solution to

$$\begin{bmatrix} 1 & 3 & 4 & 8 \\ 2 & 1 & 2 & 3 \\ 4 & 3 & 5 & 8 \\ 9 & 2 & 7 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

Solution 2.2 Note that the product of a series of lower triangular matrices is lower triangular; therefore, the product

$$W = \xi_9 \xi_8 \xi_7 \xi_6 \xi_5 \xi_4 \xi_3 \xi_2 \xi_1 \quad (2.8)$$

is lower triangular. Since the application of the elementary matrices to the matrix A results in an upper triangular matrix, then

$$WA = U \quad (2.9)$$

where U is the upper triangular matrix that results from the forward elimination process. Premultiplying Equation (2.1) by W yields

$$WAx = Wb \quad (2.10)$$

$$Ux = Wb \quad (2.11)$$

$$= b' \quad (2.12)$$

where $Wb = b'$.

From Example 2.1:

$$W = \begin{bmatrix} 1 & 0 & 0 & 0 \\ \frac{2}{5} & -\frac{1}{5} & 0 & 0 \\ 2 & 9 & -5 & 0 \\ \frac{1}{6} & \frac{14}{6} & -\frac{5}{6} & -\frac{1}{6} \end{bmatrix}$$

and

$$b' = W \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ \frac{1}{5} \\ \frac{6}{5} \\ \frac{3}{2} \end{bmatrix}$$

Thus

$$\begin{bmatrix} 1 & 3 & 4 & 8 \\ 0 & 1 & \frac{6}{5} & \frac{13}{5} \\ 0 & 0 & 1 & 3 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 \\ \frac{1}{5} \\ \frac{6}{5} \\ \frac{3}{2} \end{bmatrix} \quad (2.13)$$

By inspection, $x_4 = \frac{3}{2}$. The third row yields

$$x_3 = 6 - 3x_4 \quad (2.14)$$

Substituting the value of x_4 into Equation (2.14) yields $x_3 = \frac{3}{2}$. Similarly,

$$x_2 = \frac{1}{5} - \frac{6}{5}x_3 - \frac{13}{5}x_4 \quad (2.15)$$

and substituting x_3 and x_4 into Equation (2.15) yields $x_2 = -\frac{11}{2}$. Solving for x_1 in a similar manner produces

$$x_1 = 1 - 3x_2 - 4x_3 - 8x_4 \quad (2.16)$$

$$= -\frac{1}{2} \quad (2.17)$$

Thus

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} -1 \\ -11 \\ 3 \\ 3 \end{bmatrix} \blacksquare$$

The solution methodology of successively substituting values of x back into the equation as they are found gives rise to the name of *back substitution* for

this step of Gaussian elimination. Therefore, Gaussian elimination consists of two main steps: forward elimination and back substitution. Forward elimination is the process of transforming the matrix A into triangular factors. Back substitution is the process by which the unknown vector x is found from the input vector b and the factors of A . Gaussian elimination also provides the framework under which the LU factorization process is developed.

2.2 LU Factorization

The forward elimination step of Gaussian elimination produces a series of upper and lower triangular matrices that are related to the A matrix as given in Equation (2.9). The matrix W is a lower triangular matrix and U is an upper triangular matrix with ones on the diagonal. Recall that the inverse of a lower triangular matrix is also a lower triangular matrix; therefore, if

$$L \triangleq W^{-1}$$

then

$$A = LU$$

The matrices L and U give rise to the name of the factorization/elimination algorithm known as “LU factorization.” In fact, given any nonsingular matrix A , there exists some permutation matrix P (possibly $P = I$) such that

$$LU = PA \quad (2.18)$$

where U is upper triangular with unit diagonals, L is lower triangular with nonzero diagonals, and P is a matrix of ones and zeros obtained by rearranging the rows and columns of the identity matrix. Once a proper matrix P is chosen, this factorization is unique [8]. Once P, L , and U are determined, then the system

$$Ax = b \quad (2.19)$$

can be solved expeditiously. Premultiplying Equation (2.19) by the matrix P yields

$$PAx = Pb = b' \quad (2.20)$$

$$LUx = b' \quad (2.21)$$

where b' is just a rearrangement of the vector b . Introducing a “dummy” vector y such that

$$Ux = y \quad (2.22)$$

thus

$$Ly = b' \quad (2.23)$$

Consider the structure of Equation (2.23):

$$\begin{bmatrix} l_{11} & 0 & 0 & \cdots & 0 \\ l_{21} & l_{22} & 0 & \cdots & 0 \\ l_{31} & l_{32} & l_{33} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ l_{n1} & l_{n2} & l_{n3} & \cdots & l_{nn} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} b'_1 \\ b'_2 \\ b'_3 \\ \vdots \\ b'_n \end{bmatrix}$$

The elements of the vector y can be found by straightforward substitution:

$$\begin{aligned} y_1 &= \frac{b'_1}{l_{11}} \\ y_2 &= \frac{1}{l_{22}} (b'_2 - l_{21}y_1) \\ y_3 &= \frac{1}{l_{33}} (b'_3 - l_{31}y_1 - l_{32}y_2) \\ &\vdots \\ y_n &= \frac{1}{l_{nn}} \left(b'_n - \sum_{j=1}^{n-1} l_{nj}y_j \right) \end{aligned}$$

After the vector y has been found, then x can be easily found from

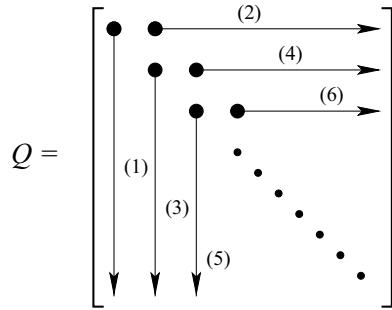
$$\begin{bmatrix} 1 & u_{12} & u_{13} & \cdots & u_{1n} \\ 0 & 1 & u_{23} & \cdots & u_{2n} \\ 0 & 0 & 1 & \cdots & u_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \vdots & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix}$$

Similarly, the solution vector x can be found by backward substitution:

$$\begin{aligned} x_n &= y_n \\ x_{n-1} &= y_{n-1} - u_{n-1,n}x_n \\ x_{n-2} &= y_{n-2} - u_{n-2,n}x_n - u_{n-2,n-1}x_{n-1} \\ &\vdots \\ x_1 &= y_1 - \sum_{j=2}^n u_{1j}x_j \end{aligned}$$

The value of LU factorization is that, once A is factored into the upper and lower triangular matrices, the solution for the solution vector x is straightforward. Note that the inverse to A is never explicitly found.

Several methods for computing the LU factors exist and each method has its advantages and disadvantages. One common factorization approach is

**FIGURE 2.1**Order of calculating columns and rows of Q

known as *Crout's* algorithm for finding the LU factors [8]. Let the matrix Q be defined as

$$Q \stackrel{\triangle}{=} L + U - I = \begin{bmatrix} l_{11} & u_{12} & u_{13} & \cdots & u_{1n} \\ l_{21} & l_{22} & u_{23} & \cdots & u_{2n} \\ l_{31} & l_{32} & l_{33} & \cdots & u_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ l_{n1} & l_{n2} & l_{n3} & \cdots & l_{nn} \end{bmatrix} \quad (2.24)$$

Crout's algorithm computes the elements of Q first by column and then row, as shown in Figure 2.1. Each element q_{ij} of Q depends only on the a_{ij} entry of A and previously computed values of Q .

Crout's Algorithm for Computing LU from A

1. Initialize Q to the zero matrix. Let $j = 1$.
2. Complete the j th column of Q (j th column of L) as

$$q_{kj} = a_{kj} - \sum_{i=1}^{j-1} q_{ki}q_{ij} \quad \text{for } k = j, \dots, n \quad (2.25)$$

3. If $j = n$, then stop.
4. Assuming that $q_{jj} \neq 0$, complete the j th row of Q (j th row of U) as

$$q_{jk} = \frac{1}{q_{jj}} \left(a_{jk} - \sum_{i=1}^{j-1} q_{ji}q_{ik} \right) \quad \text{for } k = j+1, \dots, n \quad (2.26)$$

5. Set $j = j + 1$. Go to step 2.

Once the LU factors are found, then the dummy vector y can be found by forward substitution.

Forward Substitution

$$y_k = \frac{1}{q_{kk}} \left(b_k - \sum_{j=1}^{k-1} q_{kj} y_j \right) \text{ for } k = 1, \dots, n \quad (2.27)$$

Similarly, the solution vector x can be found by backward substitution.

Backward Substitution:

$$x_k = y_k - \sum_{j=k+1}^n q_{kj} x_j \text{ for } k = n, n-1, \dots, 1 \quad (2.28)$$

One measure of the computation involved in the LU factorization process is to count the number of multiplications and divisions required to find the solution, since these are both floating point operations. Computing the j th column of Q (j th column of L) requires

$$\sum_{j=1}^n \sum_{k=j}^n (j-1)$$

multiplications and divisions. Similarly, computing the j th row of Q (j th row of U) requires

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

multiplications and divisions. The forward substitution step requires

$$\sum_{j=1}^n j$$

and the backward substitution step requires

$$\sum_{j=1}^n (n-j)$$

multiplications and divisions. Taken together, the LU factorization procedure requires

$$\frac{1}{3} (n^3 - n)$$

and the substitution steps require n^2 multiplications and divisions. Therefore, the whole process of solving the linear system of Equation (2.1) requires a total of

$$\frac{1}{3} (n^3 - n) + n^2 \quad (2.29)$$

multiplications and divisions. Compare this to the requirements of Cramer's rule, which requires $2(n + 1)!$ multiplications and divisions. Obviously, for a system of any significant size, it is far more computationally efficient to use LU factorization and forward/backward substitution to find the solution x .

Example 2.3

Using LU factorization with forward and backward substitution, find the solution to the system of Example 2.2.

Solution 2.3 The first step is to find the LU factors of the A matrix:

$$A = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 2 & 1 & 2 & 3 \\ 4 & 3 & 5 & 8 \\ 9 & 2 & 7 & 4 \end{bmatrix}$$

Starting with $j = 1$, Equation (2.25) indicates that the elements of the first column of Q are identical to the elements of the first column of A . Similarly, according to Equation (2.26), the first row of Q becomes

$$\begin{aligned} q_{12} &= \frac{a_{12}}{q_{11}} = \frac{3}{1} = 3 \\ q_{13} &= \frac{a_{13}}{q_{11}} = \frac{4}{1} = 4 \\ q_{14} &= \frac{a_{14}}{q_{11}} = \frac{8}{1} = 8 \end{aligned}$$

Thus, for $j = 1$, the Q matrix becomes

$$Q = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 2 & & & \\ 4 & & & \\ 9 & & & \end{bmatrix}$$

For $j = 2$, the second column and row of Q below and to the right of the diagonal, respectively, will be calculated. For the second column of Q :

$$\begin{aligned} q_{22} &= a_{22} - q_{21}q_{12} = 1 - (2)(3) = -5 \\ q_{32} &= a_{32} - q_{31}q_{12} = 3 - (4)(3) = -9 \\ q_{42} &= a_{42} - q_{41}q_{12} = 2 - (9)(3) = -25 \end{aligned}$$

Each element of Q uses the corresponding element of A and elements of Q that have been previously computed. Note also that the inner indices of the products are always the same and the outer indices are the same as the indices of the element being computed. This holds true for both column and row calculations. The second row of Q is computed

$$\begin{aligned} q_{23} &= \frac{1}{q_{22}} (a_{23} - q_{21}q_{13}) = \frac{1}{-5} (2 - (2)(4)) = \frac{6}{5} \\ q_{24} &= \frac{1}{q_{22}} (a_{24} - q_{21}q_{14}) = \frac{1}{-5} (3 - (2)(8)) = \frac{13}{5} \end{aligned}$$

After $j = 2$, the Q matrix becomes

$$Q = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 2 & -5 & \frac{6}{5} & \frac{13}{5} \\ 4 & -9 & & \\ 9 & -25 & & \end{bmatrix}$$

Continuing on for $j = 3$, the third column of Q is calculated

$$\begin{aligned} q_{33} &= a_{33} - (q_{31}q_{13} + q_{32}q_{23}) = 5 - \left((4)(4) + (-9)\frac{6}{5} \right) = -\frac{1}{5} \\ q_{43} &= a_{43} - (q_{41}q_{13} + q_{42}q_{23}) = 7 - \left((9)(4) + (-25)\frac{6}{5} \right) = 1 \end{aligned}$$

and the third row of Q becomes

$$\begin{aligned} q_{34} &= \frac{1}{q_{33}} (a_{34} - (q_{31}q_{14} + q_{32}q_{24})) \\ &= (-5) \left(8 - \left((4)(8) + (-9) \left(\frac{13}{5} \right) \right) \right) = 3 \end{aligned}$$

yielding

$$Q = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 2 & -5 & \frac{6}{5} & \frac{13}{5} \\ 4 & -9 & -\frac{1}{5} & 3 \\ 9 & -25 & 1 & \end{bmatrix}$$

Finally, for $j = 4$, the final diagonal element is found:

$$\begin{aligned} q_{44} &= a_{44} - (q_{41}q_{14} + q_{42}q_{24} + q_{43}q_{34}) \\ &= 4 - \left((9)(8) + (-25) \left(\frac{13}{5} \right) + (3)(1) \right) = -6 \end{aligned}$$

Thus

$$Q = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 2 & -5 & \frac{6}{5} & \frac{13}{5} \\ 4 & -9 & -\frac{1}{5} & 3 \\ 9 & -25 & 1 & -6 \end{bmatrix}$$

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 2 & -5 & 0 & 0 \\ 4 & -9 & -\frac{1}{5} & 0 \\ 9 & -25 & 1 & -6 \end{bmatrix}$$

$$U = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 0 & 1 & \frac{6}{5} & \frac{13}{5} \\ 0 & 0 & 1 & 3 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

One method of checking the correctness of the solution is to check if $LU = A$, which in this case it does.

Once the LU factors have been found, then the next step in the solution process is forward elimination using the L matrix and the b vector to find the dummy vector y . Using forward substitution to solve $Ly = b$ for y :

$$y_1 = \frac{b_1}{L_{11}} = \frac{1}{1} = 1$$

$$y_2 = \frac{(b_2 - L_{21}y_1)}{L_{22}} = \frac{(1 - (2)(1))}{-5} = \frac{1}{5}$$

$$y_3 = \frac{(b_3 - (L_{31}y_1 + L_{32}y_2))}{L_{33}} = (-5) \left(1 - \left((4)(1) + (-9)\frac{1}{5} \right) \right) = 6$$

$$y_4 = \frac{(b_4 - (L_{41}y_1 + L_{42}y_2 + L_{43}y_3))}{L_{44}}$$

$$= \frac{\left(1 - \left((9)(1) + (-25)\left(\frac{1}{5}\right) + (1)(6) \right) \right)}{-6} = \frac{3}{2}$$

Thus

$$y = \begin{bmatrix} 1 \\ \frac{1}{5} \\ 6 \\ \frac{3}{2} \end{bmatrix}$$

Similarly, backward substitution is then applied to $Ux = y$ to find the solution vector x :

$$x_4 = y_4 = \frac{3}{2}$$

$$x_3 = y_3 - U_{34}x_4 = 6 - (3) \left(\frac{3}{2} \right) = \frac{3}{2}$$

$$x_2 = y_2 - (U_{24}x_4 + U_{23}x_3) = \frac{1}{5} - \left(\left(\frac{13}{5} \right) \left(\frac{3}{2} \right) + \left(\frac{6}{5} \right) \left(\frac{3}{2} \right) \right) = -\frac{11}{2}$$

$$x_1 = y_1 - (U_{14}x_4 + U_{13}x_3 + U_{12}x_2)$$

$$= 1 - \left((8) \left(\frac{3}{2} \right) + (4) \left(\frac{3}{2} \right) + (3) \left(-\frac{11}{2} \right) \right) = -\frac{1}{2}$$

yielding the final solution vector

$$x = \frac{1}{2} \begin{bmatrix} -1 \\ -11 \\ 3 \\ 3 \end{bmatrix}$$

which is the same solution found by Gaussian elimination and backward substitution in Example 2.2. A quick check to verify the correctness of the solution is to substitute the solution vector x back into the linear system $Ax = b$. ■

2.2.1 LU Factorization with Partial Pivoting

The LU factorization process presented assumes that the diagonal element is nonzero. Not only must the diagonal element be nonzero, it must be of the same order of magnitude as the other nonzero elements. Consider the solution of the following linear system

$$\begin{bmatrix} 10^{-10} & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 5 \end{bmatrix} \quad (2.30)$$

By inspection, the solution to this linear system is

$$\begin{aligned} x_1 &\approx 2 \\ x_2 &\approx 1 \end{aligned}$$

The LU factors for A are

$$\begin{aligned} L &= \begin{bmatrix} 10^{-10} & 0 \\ 2 & (1 - 2 \times 10^{10}) \end{bmatrix} \\ U &= \begin{bmatrix} 1 & 10^{10} \\ 0 & 1 \end{bmatrix} \end{aligned}$$

Applying forward elimination to solve for the dummy vector y yields

$$\begin{aligned} y_1 &= 10^{10} \\ y_2 &= \frac{(5 - 2 \times 10^{10})}{(1 - 2 \times 10^{10})} \approx 1 \end{aligned}$$

Back substituting y into $Ux = y$ yields

$$\begin{aligned} x_2 &= y_2 \approx 1 \\ x_1 &= 10^{10} - 10^{10}x_2 \approx 0 \end{aligned}$$

The solution for x_2 is correct, but the solution for x_1 is considerably off. Why did this happen? The problem with the equations arranged the way they are

in Equation (2.30) is that 10^{-10} is too near zero for most computers. However, if the equations are rearranged such that

$$\begin{bmatrix} 2 & 1 \\ 10^{-10} & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 5 \\ 1 \end{bmatrix} \quad (2.31)$$

then the LU factors become

$$L = \begin{bmatrix} 2 & 0 \\ 10^{-10} & (1 - \frac{1}{2} \times 10^{-10}) \end{bmatrix}$$

$$U = \begin{bmatrix} 1 & \frac{1}{2} \\ 0 & 1 \end{bmatrix}$$

The dummy vector y becomes

$$y_1 = \frac{5}{2}$$

$$y_2 = \frac{(1 - \frac{5}{2} \times 10^{-10})}{(1 - \frac{1}{2} \times 10^{-10})} \approx 1$$

and by back substitution, x becomes

$$x_2 \approx 1$$

$$x_1 \approx \frac{5}{2} - \frac{1}{2}(1) = 2$$

which is the solution obtained by inspection of the equations. Therefore, even though the diagonal entry may not be exactly zero, it is still good practice to rearrange the equations such that the largest magnitude element lies on the diagonal. This process is known as *pivoting* and gives rise to the permutation matrix P of Equation (2.18).

Since Crout's algorithm computes the Q matrix by column and row with increasing index, only *partial pivoting* can be used, that is, only the rows of Q (and correspondingly A) can be exchanged. The columns must remain static. To choose the best pivot, the column beneath the j th diagonal (at the j th step in the LU factorization) is searched for the element with the largest absolute value. The corresponding row and the j th row are then exchanged. The pivoting strategy may be succinctly expressed as

Partial Pivoting Strategy

- At the j th step of LU factorization, choose the k th row as the exchange row such that

$$|q_{jj}| = \max |q_{kj}| \text{ for } k = j, \dots, n \quad (2.32)$$

- Exchange rows and update A , P , and Q correspondingly.

The permutation matrix P is composed of ones and zeros and is obtained as the product of a series of elementary permutation matrices $P^{j,k}$ which represent the exchange of rows j and k . The elementary permutation matrix $P^{j,k}$, shown in Figure 2.2, is obtained from the identity matrix by interchanging rows j and k . A pivot is achieved by the premultiplication of a properly chosen $P^{j,k}$. Since this is only an interchange of rows, the order of the unknown vector does not change.

FIGURE 2.2
Elementary permutation matrix $P^{j,k}$

Example 2.4

Repeat Example 2.3 using partial pivoting.

Solution 2.4 The A matrix is repeated here for convenience.

$$A = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 2 & 1 & 2 & 3 \\ 4 & 3 & 5 & 8 \\ 9 & 2 & 7 & 4 \end{bmatrix}$$

For $j = 1$, the first column of Q is exactly the first column of A . Applying the pivoting strategy of Equation (2.32), the q_{41} element has the largest magnitude of the first column; therefore, rows four and one are exchanged. The

elementary permutation matrix $P^{1,4}$ is

$$P^{1,4} = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

The corresponding A matrix becomes

$$A = \begin{bmatrix} 9 & 2 & 7 & 4 \\ 2 & 1 & 2 & 3 \\ 4 & 3 & 5 & 8 \\ 1 & 3 & 4 & 8 \end{bmatrix}$$

and Q at the $j = 1$ step

$$Q = \begin{bmatrix} 9 & \frac{2}{9} & \frac{7}{9} & \frac{4}{9} \\ 2 & \\ 4 & \\ 1 & \end{bmatrix}$$

At $j = 2$, the calculation of the second column of Q yields

$$Q = \begin{bmatrix} 9 & \frac{2}{9} & \frac{7}{9} & \frac{4}{9} \\ 2 & \frac{5}{9} & \\ 4 & \frac{19}{9} & \\ 1 & \frac{25}{9} & \end{bmatrix}$$

Searching the elements in the j th column below the diagonal, the fourth row of the j th (i.e., second) column once again yields the largest magnitude. Therefore, rows two and four must be exchanged, yielding the elementary permutation matrix $P^{2,4}$

$$P^{2,4} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

Similarly, the updated A is

$$\begin{bmatrix} 9 & 2 & 7 & 4 \\ 1 & 3 & 4 & 8 \\ 4 & 3 & 5 & 8 \\ 2 & 1 & 2 & 3 \end{bmatrix}$$

which yields the following Q

$$Q = \begin{bmatrix} 9 & \frac{2}{9} & \frac{7}{9} & \frac{4}{9} \\ 1 & \frac{25}{9} & \frac{29}{25} & \frac{68}{25} \\ 4 & \frac{19}{9} & \\ 2 & \frac{5}{9} & \end{bmatrix}$$

For $j = 3$, the calculation of the third column of Q yields

$$Q = \begin{bmatrix} 9 & \frac{2}{9} & \frac{7}{9} & \frac{4}{9} \\ 1 & \frac{25}{9} & \frac{29}{9} & \frac{68}{25} \\ 4 & \frac{19}{9} & -\frac{14}{25} & \frac{25}{25} \\ 2 & \frac{5}{9} & -\frac{1}{5} & \end{bmatrix}$$

In this case, the diagonal element has the largest magnitude, so no pivoting is required. Continuing with the calculation of the 3rd row of Q yields

$$Q = \begin{bmatrix} 9 & \frac{2}{9} & \frac{7}{9} & \frac{4}{9} \\ 1 & \frac{25}{9} & \frac{29}{9} & \frac{68}{25} \\ 4 & \frac{19}{9} & -\frac{14}{25} & -\frac{13}{14} \\ 2 & \frac{5}{9} & -\frac{1}{5} & \end{bmatrix}$$

Finally, calculating q_{44} yields the final Q matrix

$$Q = \begin{bmatrix} 9 & \frac{2}{9} & \frac{7}{9} & \frac{4}{9} \\ 1 & \frac{25}{9} & \frac{29}{9} & \frac{68}{25} \\ 4 & \frac{19}{9} & -\frac{14}{25} & -\frac{13}{14} \\ 2 & \frac{5}{9} & -\frac{1}{5} & \frac{3}{7} \end{bmatrix}$$

The permutation matrix P is found by multiplying together the two elementary permutation matrices:

$$\begin{aligned} P &= P^{2,4}P^{1,4}I \\ &= \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \end{aligned}$$

The results can be checked to verify that $PA = LU$. The forward and backward substitution steps are carried out on the modified vector $b' = Pb$. ■

2.2.2 LU Factorization with Complete Pivoting

An alternate LU factorization that allows complete pivoting is *Gauss'* method. In this approach, two permutation matrices are developed: one for row exchange as in partial pivoting, and a second matrix for column exchange. In this approach, the LU factors are found such that

$$P_1AP_2 = LU \quad (2.33)$$

Therefore, to solve the linear system of equations $Ax = b$ requires that a slightly different approach be used. As with partial pivoting, the permutation matrix P_1 premultiplies the linear system:

$$P_1Ax = P_1b = b' \quad (2.34)$$

Now, define a new vector z such that

$$x = P_2 z \quad (2.35)$$

Then substituting Equation (2.35) into Equation (2.34) yields

$$\begin{aligned} P_1 A P_2 z &= P_1 b = b' \\ L U z &= b' \end{aligned} \quad (2.36)$$

where Equation (2.36) can be solved using forward and backward substitution for z . Once z is obtained, then the solution vector x follows from Equation (2.35).

In complete pivoting, both rows and columns may be interchanged to place the largest element (in magnitude) on the diagonal at each step in the LU factorization process. The pivot element is chosen from the remaining elements below and to the right of the diagonal.

Complete Pivoting Strategy

- At the j th step of LU factorization, choose the pivot element such that

$$|q_{jj}| = \max |q_{kl}| \text{ for } k = j, \dots, n, \text{ and } l = j, \dots, n \quad (2.37)$$

- Exchange rows and update A , P , and Q correspondingly.

Gauss' Algorithm for Computing LU from A

- Initialize Q to the zero matrix. Let $j = 1$.
- Set the j th column of Q (j th column of L) to the j th column of the reduced matrix $A^{(j)}$, where $A^{(1)} = A$, and

$$q_{kj} = a_{kj}^{(j)} \text{ for } k = j, \dots, n \quad (2.38)$$

- If $j = n$, then stop.
- Assuming that $q_{jj} \neq 0$, set the j th row of Q (j th row of U) as

$$q_{jk} = \frac{a_{jk}^{(j)}}{q_{jj}} \text{ for } k = j + 1, \dots, n \quad (2.39)$$

- Update $A^{(j+1)}$ from $A^{(j)}$ as

$$a_{ik}^{(j+1)} = a_{ik}^{(j)} - q_{ij} q_{jk} \text{ for } i = j + 1, \dots, n, \text{ and } k = j + 1, \dots, n \quad (2.40)$$

- Set $j = j + 1$. Go to step 2.

This factorization algorithm gives rise to the same number of multiplications and divisions as Crout's algorithm for LU factorization. Crout's algorithm uses each entry of the A matrix only once, whereas Gauss' algorithm updates the A matrix each time. One advantage of Crout's algorithm over Gauss' algorithm is each element of the A matrix is used only once. Since each q_{jk} is a function of a_{jk} and then a_{jk} is never used again, the element q_{jk} can be written *over* the a_{jk} element. Therefore, rather than having to store two $n \times n$ matrices in memory (A and Q), only one matrix is required.

Crout's and Gauss' algorithms are only two of numerous algorithms for LU factorization. Other methods include Doolittle and bifactorization algorithms [21], [26], [54]. Most of these algorithms require similar numbers of multiplications and divisions and only differ slightly in performance when implemented on traditional serial computers. However, these algorithms differ considerably when factors such as memory access, storage, and parallelization are considered. Consequently, it is wise to choose the factorization algorithm to fit the application and the computer architecture upon which it will be implemented.

Note that, as Q is constructed, there are now entries that may appear in the matrix structure of Q that were not originally present in the structure of A . Entries that change from an initial zero to a nonzero value during the LU factorization are called “fill-ins” or *fills* for short. The number and placement of fills play an important role in multiple applications, including preconditioning and sparse computation.

Example 2.5

For the matrix A given below, find the Q matrix and identify the fills.

$$A = \begin{bmatrix} 10 & 1 & 0 & 3 & 0 & 0 & 0 & 5 & 0 & 0 \\ 2 & 9 & 0 & 0 & 0 & 0 & 5 & 0 & 0 & 2 \\ 0 & 0 & 21 & 5 & 7 & 0 & 0 & 0 & 0 & 4 \\ 4 & 0 & 1 & 18 & 8 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 4 & 7 & 25 & 4 & 1 & 0 & 0 & 2 \\ 0 & 0 & 0 & 0 & 3 & 14 & 9 & 0 & 0 & 0 \\ 0 & 1 & 4 & 0 & 2 & 3 & 12 & 1 & 1 & 0 \\ 1 & 0 & 5 & 0 & 0 & 0 & 5 & 10 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 6 & 0 & 20 & 0 \\ 0 & 2 & 3 & 0 & 4 & 0 & 0 & 0 & 0 & 35 \end{bmatrix}$$

Solution 2.5

The resulting Q is given below. The fills are indicated by the boxed numbers. By comparing Q with A , it can be seen that there are 24 fills, or entries in Q that were previously zero in A .

$$Q = \begin{bmatrix} 10.0000 & 0.1000 & 0 & 0.3000 & 0 & 0 & 0 & 0.5000 & 0 \\ 2.0000 & 8.8000 & 0 & -0.0682 & 0 & 0.5682 & -0.1136 & 0 & 0.2273 \\ 0 & 0 & 21.0000 & 0.2381 & 0.3333 & 0 & 0 & 0 & 0 \\ 4.0000 & -0.4000 & 1.0000 & 16.5346 & 0.4637 & 0 & 0.0137 & -0.1237 & 0 \\ 0 & 0 & 4.0000 & 6.0476 & 20.8625 & 0.1917 & 0.0439 & 0.0359 & 0 \\ 0 & 0 & 0 & 0 & 3.0000 & 13.4248 & 0.6606 & -0.0080 & 0 \\ 0 & 1.0000 & 4.0000 & -0.8842 & 1.0766 & 2.7936 & 9.5513 & 0.1034 & 0.1047 \\ 1.0000 & -0.1000 & 5.0000 & -1.4973 & -0.9724 & 0.1864 & 4.9970 & 8.8229 & -0.0593 \\ 0 & 0 & 0 & 0 & 0 & 0 & 6.0000 & -0.6207 & 19.3350 \\ 0 & 2.0000 & 3.0000 & -0.5779 & 3.2680 & -0.6266 & -0.8581 & 0.1223 & 0.0971 \end{bmatrix}$$

■

2.3 Condition Numbers and Error Propagation

The Gaussian elimination and LU factorization algorithms are considered direct methods because they calculate the solution vector $x^* = A^{-1}b$ in a finite number of steps without an iterative refinement. On a computer with infinite precision, direct methods would yield the exact solution x^* . However, since computers have finite precision, the solution obtained has limited accuracy. The *condition number* of a matrix is a useful measure for determining the level of accuracy of a solution. The condition number of the matrix A is generally defined as

$$\kappa(A) = \sqrt{\frac{\lambda_{\max}}{\lambda_{\min}}} \quad (2.41)$$

where λ_{\max} and λ_{\min} denote the largest and smallest eigenvalues of the matrix $A^T A$. These eigenvalues are real and nonnegative regardless of whether the eigenvalues of A are real or complex.

The condition number of a matrix is a measure of the linear independence of the eigenvectors of the matrix. A singular matrix has at least one zero eigenvalue and contains at least one degenerate row (i.e., the row can be expressed as a linear combination of other rows). The identity matrix, which gives rise to the most linearly independent eigenvectors possible and has every eigenvalue equal to one, has a condition number of 1. If the condition number of a matrix is much greater than one, then the matrix is said to be *ill conditioned*. The larger the condition number, the more sensitive the solution process is to slight perturbations in the elements of A and the more numerical error likely to be contained in the solution.

Because of numerical error introduced into the solution process, the computed solution \tilde{x} of Equation (2.1) will differ from the exact solution x^* by a finite amount Δx . Other errors, such as approximation, measurement, or round-off error, may be introduced into the matrix A and vector b . Gaussian elimination produces a solution that has roughly

$$t \log_{10} \beta - \log_{10} \kappa(A) \quad (2.42)$$

correct decimal places in the solution, where t is the bit length of the mantissa ($t = 24$ for a typical 32-bit binary word), β is the base ($\beta = 2$ for binary operations), and κ is the condition number of the matrix A . One interpretation of Equation (2.42) is that the solution will lose about $\log_{10} \kappa$ digits of accuracy during Gaussian elimination (and consequently LU factorization). Based upon the known accuracy of the matrix entries, the condition number, and the machine precision, the accuracy of the numerical solution \tilde{x} can be predicted [39].

2.4 Stationary Iterative Methods

Stationary iterative methods, also known as relaxation methods, are iterative in nature and produce a sequence of vectors that ideally converge to the solution $x^* = A^{-1}b$. Relaxation methods can be incorporated into the solution of Equation (2.1) in several ways. In all cases, the principal advantage of using an iterative method stems from not requiring a direct solution of a large system of linear equations and from the fact that the relaxation methods permit the simulator to exploit the latent portions of the system (those portions which are relatively unchanging at the present time) effectively. In addition, with the advent of parallel-processing technology, relaxation methods lend themselves more readily to parallel implementation than do direct methods. The two most common stationary methods are the Jacobi and the Gauss–Seidel methods [61].

These relaxation methods may be applied for the solution of the linear system

$$Ax = b \quad (2.43)$$

A general approach to relaxation methods is to define a *splitting matrix* M such that Equation (2.43) can be rewritten in equivalent form as

$$Mx = (M - A)x + b \quad (2.44)$$

This splitting leads to the iterative process

$$Mx^{k+1} = (M - A)x^k + b \quad k = 1, \dots, \infty \quad (2.45)$$

where k is the iteration index. This iteration produces a sequence of vectors x^1, x^2, \dots for a given initial guess x^0 . Various iterative methods can be developed by different choices of the matrix M . The objective of a relaxation method is to choose the splitting matrix M such that the sequence is easily computed and the sequence converges rapidly to a solution.

Let A be split into $L + D + U$, where L is strictly lower triangular, D is a diagonal matrix, and U is strictly upper triangular. Note that these matrices

are different from the L and U obtained from LU factorization. The vector x can then be solved for in an iterative manner using the Jacobi relaxation method,

$$x^{k+1} = -D^{-1} ((L + U) x^k - b) \quad (2.46)$$

or identically in scalar form,

$$x_i^{k+1} = - \sum_{j \neq i}^n \left(\frac{a_{ij}}{a_{ii}} \right) x_j^k + \frac{b_i}{a_{ii}} \quad 1 \leq i \leq n, k \geq 0 \quad (2.47)$$

In the Jacobi relaxation method, all of the updates of the approximation vector x^{k+1} are obtained by using only the components of the previous approximation vector x^k . Therefore, this method is also sometimes called the method of simultaneous displacements.

The Gauss–Seidel relaxation method is similar:

$$x^{k+1} = - (L + D)^{-1} (U x^k - b) \quad (2.48)$$

or in scalar form

$$x_i^{k+1} = - \sum_{j=1}^{i-1} \left(\frac{a_{ij}}{a_{ii}} \right) x_j^{k+1} - \sum_{j=i+1}^n \left(\frac{a_{ij}}{a_{ii}} \right) x_j^k + \frac{b_i}{a_{ii}} \quad 1 \leq i \leq n, k \geq 0 \quad (2.49)$$

The Gauss–Seidel method has the advantage that each new update x_i^{k+1} relies only on previously computed values at that iteration: $x_1^{k+1}, x_2^{k+1}, \dots, x_{i-1}^{k+1}$. Since the states are updated one by one, the new values can be stored in the same locations held by the old values, thus reducing the storage requirements.

Since relaxation methods are iterative, it is essential to determine under what conditions they are guaranteed to converge to the exact solution

$$x^* = A^{-1} b \quad (2.50)$$

It is well known that a necessary and sufficient condition for the Jacobi relaxation method to converge given any initial guess x_0 is that all eigenvalues of

$$M_J \stackrel{\Delta}{=} -D^{-1} (L + U) \quad (2.51)$$

must lie within the unit circle in the complex plane [61]. Similarly, the eigenvalues of

$$M_{GS} \stackrel{\Delta}{=} - (L + D)^{-1} U \quad (2.52)$$

must lie within the unit circle in the complex plane for the Gauss–Seidel relaxation algorithm to converge for any initial guess x_0 . In practice, these conditions are difficult to confirm. There are several more general conditions that are easily confirmed under which convergence is guaranteed. In particular, if A is strictly diagonally dominant, then both the Jacobi and Gauss–Seidel methods are guaranteed to converge to the exact solution.

The initial vector x_0 can be arbitrary; however, if a good guess of the solution is available, it should be used for x_0 to produce more rapid convergence to within some predefined tolerance.

In general, the Gauss–Seidel method converges faster than the Jacobi for most classes of problems. If A is lower-triangular, the Gauss–Seidel method will converge in one iteration to the exact solution, whereas the Jacobi method will take n iterations. The Jacobi method has the advantage, however, that, at each iteration, each x_i^{k+1} is independent of all other x_j^{k+1} for $j \neq i$. Thus the computation of all x_i^{k+1} can proceed in parallel. This method is therefore well suited to parallel processing [40].

Both the Jacobi and Gauss–Seidel methods can be generalized to the block-Jacobi and block-Gauss–Seidel methods where A is split into block matrices $L + D + U$, where D is block diagonal and L and U are lower- and upper-block triangular, respectively. The same necessary and sufficient convergence conditions exist for the block case as for the scalar case, that is, the eigenvalues of M_J and M_{GS} must lie within the unit circle in the complex plane.

Example 2.6

Solve

$$\begin{bmatrix} -10 & 2 & 3 & 6 \\ 0 & -9 & 1 & 4 \\ 2 & 6 & -12 & 2 \\ 3 & 1 & 0 & -8 \end{bmatrix} x = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \quad (2.53)$$

for x using (1) the Gauss–Seidel method, and (2) the Jacobi method.

Solution 2.6 The Gauss–Seidel method given in Equation (2.49) with the initial vector $x = [0 \ 0 \ 0 \ 0]$ leads to the following updates:

k	x_1	x_2	x_3	x_4
1	0.0000	0.0000	0.0000	0.0000
2	-0.1000	-0.2222	-0.3778	-0.5653
3	-0.5969	-0.5154	-0.7014	-0.7883
4	-0.8865	-0.6505	-0.8544	-0.9137
5	-1.0347	-0.7233	-0.9364	-0.9784
6	-1.1126	-0.7611	-0.9791	-1.0124
7	-1.1534	-0.7809	-1.0014	-1.0301
8	-1.1747	-0.7913	-1.0131	-1.0394
9	-1.1859	-0.7968	-1.0193	-1.0443
10	-1.1917	-0.7996	-1.0225	-1.0468
11	-1.1948	-0.8011	-1.0241	-1.0482
12	-1.1964	-0.8019	-1.0250	-1.0489
13	-1.1972	-0.8023	-1.0255	-1.0492
14	-1.1976	-0.8025	-1.0257	-1.0494
15	-1.1979	-0.8026	-1.0259	-1.0495
16	-1.1980	-0.8027	-1.0259	-1.0496

The Gauss–Seidel iterates have converged to the solution

$$x = [-1.1980 \ -0.8027 \ -1.0259 \ -1.0496]^T$$

From Equation (2.47) and using the initial vector $x = [0 \ 0 \ 0 \ 0]$, the following updates are obtained for the Jacobi method:

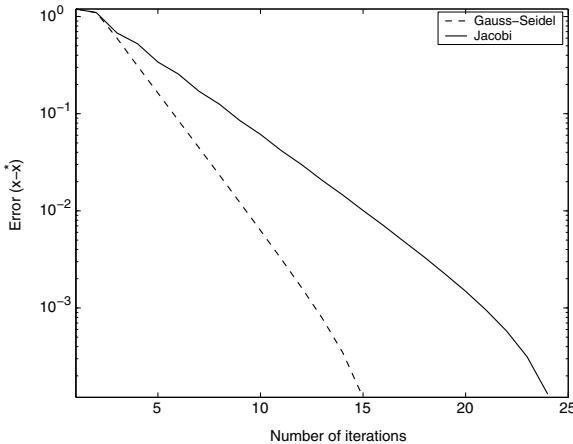
k	x_1	x_2	x_3	x_4
1	0.0000	0.0000	0.0000	0.0000
2	-0.1000	-0.2222	-0.2500	-0.5000
3	-0.5194	-0.4722	-0.4611	-0.5653
4	-0.6719	-0.5247	-0.6669	-0.7538
5	-0.8573	-0.6314	-0.7500	-0.8176
6	-0.9418	-0.6689	-0.8448	-0.9004
7	-1.0275	-0.7163	-0.8915	-0.9368
8	-1.0728	-0.7376	-0.9355	-0.9748
9	-1.1131	-0.7594	-0.9601	-0.9945
10	-1.1366	-0.7709	-0.9810	-1.0123
11	-1.1559	-0.7811	-0.9936	-1.0226
12	-1.1679	-0.7871	-1.0037	-1.0311
13	-1.1772	-0.7920	-1.0100	-1.0363
14	-1.1832	-0.7950	-1.0149	-1.0404
15	-1.1877	-0.7974	-1.0181	-1.0431
16	-1.1908	-0.7989	-1.0205	-1.0451
17	-1.1930	-0.8001	-1.0221	-1.0464
18	-1.1945	-0.8009	-1.0233	-1.0474
19	-1.1956	-0.8014	-1.0241	-1.0480
20	-1.1963	-0.8018	-1.0247	-1.0485
21	-1.1969	-0.8021	-1.0250	-1.0489
22	-1.1972	-0.8023	-1.0253	-1.0491
23	-1.1975	-0.8024	-1.0255	-1.0492
24	-1.1977	-0.8025	-1.0257	-1.0494
25	-1.1978	-0.8026	-1.0258	-1.0494

The Jacobi iterates have converged to the same solution as the Gauss–Seidel method. The error in the iterates is shown in Figure 2.3 on a semilog scale, where the error is defined as the maximum $|x_i^k - x_i^*|$ for all $i = 1, \dots, 4$. Both the Gauss–Seidel and the Jacobi methods exhibit *linear convergence*, but the Gauss–Seidel converges with a steeper slope and will therefore reach the convergence tolerance sooner for the same initial condition. ■

Example 2.7

Repeat Example 2.2 using the Jacobi iterative method.

Solution 2.7 Repeating the solution procedure of Example 2.6 yields the

**FIGURE 2.3**

Convergence rates of the Gauss–Seidel and Jacobi methods

following iterations for the Jacobi method:

k	x_1	x_2	x_3	x_4
1	0	0	0	0
2	1.0000	1.0000	0.2000	0.2500
3	-4.8000	-2.1500	-1.6000	-2.8500
4	36.6500	22.3500	9.8900	14.9250
5	-225.0100	-136.8550	-66.4100	-110.6950

Obviously these iterates are not converging. To understand why they are diverging, consider the iterative matrix for the Jacobi matrix:

$$\begin{aligned} M_J &= -D^{-1}(L + U) \\ &= \begin{bmatrix} 0.00 & -3.00 & -4.00 & -8.00 \\ -2.00 & 0.00 & -2.00 & -3.00 \\ -0.80 & -0.60 & 0.00 & -1.60 \\ -2.25 & -0.50 & -1.75 & 0.00 \end{bmatrix} \end{aligned}$$

The eigenvalues of M_J are

$$\begin{bmatrix} -6.6212 \\ 4.3574 \\ 1.2072 \\ 1.0566 \end{bmatrix}$$

which are all greater than one and lie outside the unit circle. Therefore, the Jacobi method will not converge to the solution regardless of choice of initial condition and cannot be used to solve the system of Example 2.2. ■

If the largest eigenvalue of the iterative matrix M_J or M_{GS} is less than,

but almost, unity, then the convergence may proceed very slowly. In this case it is desirable to introduce a weighting factor ω that will improve the rate of convergence. From

$$x^{k+1} = -(L + D)^{-1} (Ux^k - b) \quad (2.54)$$

it follows that

$$x^{k+1} = x^k - D^{-1} (Lx^{k+1} + (D + U)x^k - b) \quad (2.55)$$

A new iterative method can be defined with the weighting factor ω such that

$$x^{k+1} = x^k - \omega D^{-1} (Lx^{k+1} + (D + U)x^k - b) \quad (2.56)$$

This method is known as the *successive overrelaxation (SOR)* method with relaxation coefficient $\omega > 0$. This method takes the form of a weighted average between the previous iterate and the computed Gauss–Seidel iterate successively for each component. Note that, if the relaxation iterates converge, they converge to the solution $x^* = A^{-1}b$. One necessary condition for the SOR method to be convergent is that $0 < \omega < 2$ [29]. The idea is to choose ω such that the convergence of the iterates is accelerated. The calculation of the optimal value for ω is difficult, except in a few simple cases. The optimal value is usually determined through trial and error, but analysis shows that, for systems larger than $n = 30$, the optimal SOR can be more than forty times faster than the Jacobi method [29]. The improvement in the speed of convergence often increases as n increases.

If $\omega = 1$, the SOR method simplifies to the Gauss–Seidel method. In principle, given the spectral radius ρ of the Jacobi iteration matrix, ω can be determined

$$\omega_{opt} = \frac{2}{1 + \sqrt{1 - \rho^2}} \quad (2.57)$$

but this calculation is typically not computationally efficient.

One further modification can be added to the SOR method to achieve a *symmetric successive overrelaxation (SSOR)* method. This is essentially the combination of a forward SOR step combined with a backward SOR step:

$$x^{k+1} = B_1 B_2 x^k + \omega(2 - \omega)(D - \omega U)^{-1} D (D - \omega L)^{-1} b \quad (2.58)$$

where

$$B_1 = (D - \omega U)^{-1} (\omega L + (1 - \omega)D) \quad (2.59)$$

$$B_2 = (D - \omega L)^{-1} (\omega U + (1 - \omega)D) \quad (2.60)$$

This method is analogous to a two-step SOR applied to a symmetric matrix. In the first step, the unknowns are updated in forward order, and in the second step, or backward step, the unknowns are updated in reverse order.

2.5 Conjugate Gradient Methods

In this section, the Krylov subspace method is described for solving linear systems. The conjugate gradient method is a nonstationary iterative method. Nonstationary methods differ from stationary methods in that the computations do not use an iteration matrix and typically involve information that changes at each iteration. Krylov subspace methods define a space such that the k th Krylov subspace \mathcal{K}_k is

$$\mathcal{K}_k = \text{span} (r_0, Ar_0, \dots, A^{k-1}r_0) \text{ for } k \geq 1 \quad (2.61)$$

where the residual is defined as

$$r_k = b - Ax_k \quad (2.62)$$

Krylov subspace methods generate vector sequences of iterates (i.e., successive approximations to the solution), residuals corresponding to the iterates, and search directions used in updating the iterates and residuals.

A common Krylov iterative method for solving $Ax = b$ is the *conjugate gradient* method. The conjugate gradient (CG) method was originally intended as a direct method, but has been adopted as an iterative method and has generally superseded the Jacobi–Gauss–Seidel–SOR family of methods [27].

The conjugate gradient method is intended to solve symmetric positive definite systems of equations. Recall that the matrix A is symmetric positive definite if $A = A^T$ and

$$x^T Ax > 0 \text{ for all } x \neq 0$$

The CG method can be considered a minimization method for the function

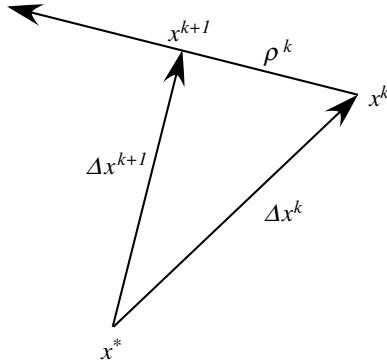
$$E(x) = \|Ax - b\|^2 \quad (2.63)$$

along a succession of rays. Therefore, the k th iterate x_k of the CG iteration minimizes

$$\phi(x) = \frac{1}{2}x^T Ax - x^T b \quad (2.64)$$

over $x_0 + \mathcal{K}_k$.

One attractive feature of this method is that it is guaranteed to converge in at most n steps (neglecting round-off error) if the A matrix is positive definite. The conjugate gradient method is most frequently used instead of Gaussian elimination if the A matrix is very large and sparse, in which case the solution may be obtained in less than n steps. This is especially true if the A matrix is well conditioned. If the matrix is ill conditioned, then round-off errors may prevent the algorithm from obtaining a sufficiently accurate solution after n steps.

**FIGURE 2.4**

The conjugate gradient method

In the conjugate gradient method, a succession of search directions ρ_k is employed and a parameter α_k is computed such that $f(x^k - \alpha_k \rho_k)$ is minimized along the ρ_k direction. Upon setting x^{k+1} equal to $x^k - \alpha_k \rho_k$, the new search direction is found. As the conjugate gradient method progresses, each error function is associated with a specific ray, or orthogonal expansion. Therefore, the conjugate gradient method is reduced to the process of generating the orthogonal vectors and finding the proper coefficients to represent the desired solution. The conjugate gradient method is illustrated in Figure 2.4. Let x^* denote the exact (but unknown) solution, x^k an approximate solution, and $\Delta x^k = x^k - x^*$. Given any search direction ρ^k , the minimal distance from the line to x^* is found by constructing Δx^{k+1} perpendicular to x^k . Since the exact solution is unknown, the residual is made to be perpendicular to ρ^k . Regardless of how the new search direction is chosen, the norm of the residual will not increase.

All Krylov iterative methods for solving $Ax = b$ define an iterative process such that

$$x^{k+1} = x^k + \alpha_{k+1} \rho_{k+1} \quad (2.65)$$

where x^{k+1} is the updated value, α_k is the step length, and ρ_k defines the direction $\in R^n$ in which the algorithm moves to update the estimate.

Let the residual, or mismatch, vector at step k be given by

$$r_k = Ax^k - b \quad (2.66)$$

and the error function given by

$$E_k(x^k) = \|Ax^k - b\|^2 \quad (2.67)$$

Once the search direction ρ_{k+1} is determined, α_{k+1} can be computed from

the minimization property of the iteration, in which

$$\frac{d\phi(x_k + \alpha\rho_{k+1})}{d\alpha} = 0 \quad (2.68)$$

for $\alpha = \alpha_{k+1}$. Equation (2.68) can be written as

$$\rho_{k+1}^T A x_k + \alpha \rho_{k+1}^T A \rho_{k+1} - \rho_{k+1}^T b = 0 \quad (2.69)$$

Then the coefficient that minimizes the error function at step $k+1$ is

$$\begin{aligned} \alpha_{k+1} &= \frac{\rho_{k+1}^T (b - Ax_k)}{\rho_{k+1}^T A \rho_{k+1}} \\ &= \frac{\rho_{k+1}^T r_k}{\rho_{k+1}^T A \rho_{k+1}} \\ &= \frac{\|A^T r_k\|^2}{\|A \rho_{k+1}\|^2} \end{aligned}$$

This has the geometric interpretation of minimizing E_{k+1} along the ray defined by ρ_{k+1} . Further, an improved algorithm is one that seeks the minimum of E_{k+1} in a plane spanned by two direction vectors, such that

$$x^{k+1} = x^k + \alpha_{k+1} (\rho_{k+1} + \beta_{k+1} \sigma_{k+1}) \quad (2.70)$$

where the rays ρ_{k+1} and σ_{k+1} span a plane in R^n . The process of selecting direction vectors and coefficients to minimize the error function E_{k+1} is optimized when the chosen vectors are orthogonal, such that

$$\langle A \rho_{k+1}, A \sigma_{k+1} \rangle = 0 \quad (2.71)$$

where $\langle \cdot \rangle$ denotes inner product. Vectors that satisfy the orthogonality condition of Equation (2.71) are said to be mutually conjugate with respect to the operator $A^T A$, where A^T is the conjugate transpose of A . One method of choosing appropriate vectors is to choose σ_{k+1} as a vector orthogonal to ρ_k , thus eliminating the need to specify two orthogonal vectors at each step. While this simplifies the procedure, there is now an implicit recursive dependence for generating the ρ vectors.

Conjugate Gradient Algorithm for Solving $Ax = b$

Initialization: Let $k = 1$, and

$$r_0 = Ax^0 - b \quad (2.72)$$

$$\rho_0 = \|r_0\|^2 \quad (2.73)$$

While $\|r_k\| \geq \varepsilon$,

$$\sigma = r_{k-1} \text{ if } k = 0, \text{ else } \beta = \frac{\rho_{k-1}}{\rho_{k-2}} \text{ and } \sigma = r_{k-1} + \beta \sigma \quad (2.74)$$

$$w = A\sigma \quad (2.75)$$

$$\alpha = \frac{\rho_{k-1}}{\sigma^T w} \quad (2.76)$$

$$x = x + \alpha\sigma \quad (2.77)$$

$$r_k = r_{k-1} - \alpha w \quad (2.78)$$

$$\rho_k = \|r_k\|^2 \quad (2.79)$$

$$k = k + 1 \quad (2.80)$$

For an arbitrary symmetric positive definite matrix A , the conjugate gradient method will produce a solution in at most n steps (neglecting round-off error). This is a direct consequence of the fact that the n direction vectors ρ_0, ρ_1, \dots span the solution space. Finite step termination is a significant advantage of the conjugate gradient method over other iterative methods such as relaxation methods.

Note that the matrix A itself need not be formed or stored. The only requirement is a routine for producing matrix-vector products. Krylov space methods are often called *matrix-free* methods for this reason. Each iteration requires only a single matrix-vector product (to compute $w = A\sigma$) and two scalar products ($\sigma^T w$ and $\|r_k\|^2$).

Example 2.8

Use the conjugate gradient method to solve the following system of equations:

$$\begin{bmatrix} 22 & 2 & 8 & 5 & 3 \\ 2 & 19 & 6 & 6 & 2 \\ 8 & 6 & 27 & 7 & 4 \\ 5 & 6 & 7 & 24 & 0 \\ 3 & 2 & 4 & 0 & 9 \end{bmatrix} x = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}$$

with

$$x^0 = [0 \ 0 \ 0 \ 0 \ 0]$$

Solution 2.8 It can be easily verified that the matrix A is symmetric and positive definite.

Following the order of the conjugate gradient algorithm yields

Initialization

$$r_0 = [1 \ 2 \ 3 \ 4 \ 5]^T$$

$$\rho_0 = 55$$

k=1

$$\sigma = [1 \ 2 \ 3 \ 4 \ 5]^T$$

$$w = \begin{bmatrix} 85 \\ 92 \\ 149 \\ 134 \\ 64 \end{bmatrix}$$

$$\alpha = 0.0350$$

$$x^1 = \begin{bmatrix} 0.0350 \\ 0.0700 \\ 0.1050 \\ 0.1399 \\ 0.1749 \end{bmatrix}$$

$$r_1 = \begin{bmatrix} -1.9739 \\ -1.2188 \\ -2.2131 \\ -0.6883 \\ 2.7608 \end{bmatrix}$$

$$\rho_1 = 18.3756$$

$$\|r_1\| = 4.2867$$

k=2

$$\beta = 0.3341$$

$$\sigma = \begin{bmatrix} -1.6398 \\ -0.5506 \\ -1.2108 \\ 0.6481 \\ 4.4313 \end{bmatrix}$$

$$w = \begin{bmatrix} -30.3291 \\ -8.2550 \\ -26.8518 \\ -4.4238 \\ 29.0180 \end{bmatrix}$$

$$\alpha = 0.0865$$

$$x^2 = \begin{bmatrix} -0.1068 \\ 0.0224 \\ 0.0003 \\ 0.1960 \\ 0.5581 \end{bmatrix}$$

$$r_2 = \begin{bmatrix} 0.6486 \\ -0.5050 \\ 0.1087 \\ -0.3058 \\ 0.2517 \end{bmatrix}$$

$$\rho_2 = 0.8444$$

$$\|r_2\| = 0.9189$$

Similarly

k	3	4	5
β_k	0.0460	0.0041	0.0001
σ_k	$\begin{bmatrix} 0.5732 \\ -0.5303 \\ 0.0531 \\ -0.2760 \\ 0.4553 \end{bmatrix}$	$\begin{bmatrix} 0.0212 \\ -0.0146 \\ -0.0346 \\ 0.0389 \\ -0.0081 \end{bmatrix}$	$\begin{bmatrix} 0.0002 \\ 0.0003 \\ -0.0002 \\ -0.0001 \\ 0.0001 \end{bmatrix}$
w	$\begin{bmatrix} 11.9610 \\ -9.3567 \\ 2.7264 \\ -6.5682 \\ 4.9692 \end{bmatrix}$	$\begin{bmatrix} 0.3306 \\ -0.2250 \\ -0.6121 \\ 0.7098 \\ -0.1767 \end{bmatrix}$	$\begin{bmatrix} 0.0024 \\ 0.0042 \\ -0.0026 \\ -0.0019 \\ 0.0010 \end{bmatrix}$
α	0.0526	0.0566	0.0713
x^k	$\begin{bmatrix} -0.0766 \\ -0.0056 \\ 0.0031 \\ 0.1815 \\ 0.5821 \end{bmatrix}$	$\begin{bmatrix} -0.0754 \\ -0.0064 \\ 0.0011 \\ 0.1837 \\ 0.5816 \end{bmatrix}$	$\begin{bmatrix} -0.0754 \\ -0.0064 \\ 0.0011 \\ 0.1837 \\ 0.5816 \end{bmatrix}$
ρ_k	$\begin{bmatrix} 0.0034 \\ 0.0189 \\ -0.0124 \\ -0.0348 \\ 0.0400 \\ -0.0099 \end{bmatrix}$	$\begin{bmatrix} 0.0000 \\ 0.0002 \\ 0.0003 \\ -0.0002 \\ -0.0001 \\ 0.0001 \end{bmatrix}$	$\begin{bmatrix} 0.0000 \\ -0.0000 \\ -0.0000 \\ -0.0000 \\ -0.0000 \\ -0.0000 \end{bmatrix}$
$\ r_k\ $	0.0585	0.0004	0.0000

The iterations converged in five iterations, as the algorithm guarantees. Note that, depending on the size of the convergence criterion ϵ , it is possible to have convergence after four iterations, as the norm of $r_4 = 0.0004$ is relatively small at the conclusion of the fourth iteration. ■

The conjugate gradient method is more numerically competitive for matrices that are very large and sparse or that have a special structure that cannot be easily handled by LU factorization. In some cases, the speed of convergence of the conjugate gradient method can be improved by *preconditioning*. As seen with the Gauss–Seidel and Jacobi iteration, the convergence rate of iterative algorithms is closely related to the eigenvalue spectrum of the iterative matrix. Consequently, a scaling or matrix transformation that converts

the original system of equations into one with a better eigenvalue spectrum may significantly improve the rate of convergence. This procedure is known as preconditioning and is discussed in Section 2.7.

2.6 Generalized Minimal Residual Algorithm

If the matrix A is neither symmetric nor positive definite, then the term

$$\langle A\rho_{k+1}, A\sigma_{k+1} \rangle$$

is not guaranteed to be zero and the search vectors are not mutually orthogonal. Mutual orthogonality is required to generate a basis of the solution space. Hence this basis must be explicitly constructed. The extension of the conjugate gradient method, called the generalized minimal residual algorithm (GMRES), minimizes the norm of the residual in a subspace spanned by the set of vectors

$$r^0, Ar^0, A^2r^0, \dots, A^{k-1}r^0$$

where vector r^0 is the initial residual $r^0 = \|b - Ax^0\|$, and the k th approximation to the solution is chosen from this space. This subspace, a *Krylov subspace*, is made orthogonal by the well-known Gram–Schmidt procedure, known as the Arnoldi process when applied to a Krylov subspace [41]. At each step k , the GMRES algorithm applies the Arnoldi process to a set of k orthonormal basis vectors for the k th Krylov subspace to generate the next basis vector. Arnoldi methods are described in greater detail in Section 7.3.

The k th iteration of the GMRES method is the solution to the least squares problem

$$\text{minimize}_{x \in x_0 + \mathcal{K}_k} \|b - Ax\| \quad (2.81)$$

At each step, the algorithm multiplies the previous Arnoldi vector v_j by A and then orthonormalizes the resulting vector w_j against all previous v_i 's. The columns $V = [v_1, v_2, \dots, v_k]$ form an orthonormal basis for the Krylov subspace and H is the orthogonal projection of A onto this space.

An orthogonal matrix triangularization such as the Arnoldi method consists in determining an $n \times n$ orthogonal matrix Q such that

$$Q^T = \begin{bmatrix} R \\ 0 \end{bmatrix} \quad (2.82)$$

where R is an $m \times m$ upper triangular matrix R . Then the solution process reduces to solving the triangular system $Rx = Py$, where P consists of the first m rows of Q .

To clear one element at a time to upper triangularize a matrix, Given's rotation can be applied. Given's rotation is a transformation based on the matrix

$$G_{jk} = \begin{bmatrix} 1 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & & & & \vdots \\ 0 & \dots & \text{cs} & \dots & \text{sn} & \dots & 0 \\ \vdots & & \vdots & \ddots & & & \vdots \\ 0 & \dots & -\text{sn} & \dots & \text{cs} & \dots & 0 \\ \vdots & & \vdots & & \ddots & & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 1 \end{bmatrix} \quad (2.83)$$

where properly chosen $\text{cs} = \cos(\phi)$ and $\text{sn} = \sin(\phi)$ (for some rotation angle ϕ) can be used to zero out the element A_{kj} . The orthogonal matrix G_{jk} rotates the vector $(c, -s)$, which makes an angle of $-\phi$ with the x -axis, such that it overlaps the x -axis. Given's rotations are used to remove single nonzero elements of matrices in reduction to triangular form.

One of the difficulties with the GMRES methods is that, as k increases, the number of vectors requiring storage increases as k and the number of multiplications as $\frac{1}{2}k^2n$ (for an $n \times n$ matrix). To remedy this difficulty, the algorithm can be applied iteratively, i.e., it can be restarted every m steps, where m is some fixed integer parameter.

GMRES Algorithm for Solving $Ax = b$

Initialization:

$$\begin{aligned} r_0 &= b - Ax^0 \\ e_1 &= [1 \ 0 \ 0 \ \dots \ 0]^T \\ v_1 &= r_0 / \|r_0\| \\ s &= \|r_0\|e_1 \\ k &= 1 \\ \text{cs} &= [0 \ 0 \ 0 \ \dots \ 0]^T \\ \text{sn} &= [0 \ 0 \ 0 \ \dots \ 0]^T \end{aligned}$$

While $\|r_k\| \geq \varepsilon$ and $k \leq k_{\max}$, set

1. $H(j, k) = (Av_k)^T v_j, \ j = 1, \dots, k$
2. $v_{k+1} = Av_k - \sum_{j=1}^k H(j, k)v_j$
3. $H(k + 1, k) = \|v_{k+1}\|$
4. $v_{k+1} = v_{k+1} / \|v_{k+1}\|$
5. Givens rotation:

(a)

$$\begin{bmatrix} H(j, k) \\ H(j+1, k) \end{bmatrix} = \begin{bmatrix} \text{cs}(j) & \text{sn}(j) \\ -\text{sn}(j) & \text{cs}(j) \end{bmatrix} \begin{bmatrix} H(j, k) \\ H(j+1, k) \end{bmatrix}, \quad j = 1, \dots, k-1$$

(b)

$$\begin{aligned} \text{cs}(k) &= \frac{H(k, k)}{\sqrt{H(k+1, k)^2 + H(k, k)^2}} \\ \text{sn}(k) &= \frac{H(k+1, k)}{\sqrt{H(k+1, k)^2 + H(k, k)^2}} \end{aligned}$$

(c) Approximate residual norm

$$\begin{aligned} \alpha &= \text{cs}(k)s(k) \\ s(k+1) &= -\text{sn}(k)s(k) \\ s(k) &= \alpha \\ \text{error} &= |s(k+1)| \end{aligned}$$

(d) Set

$$\begin{aligned} H(k, k) &= \text{cs}(k)H(k, k) + \text{sn}(k)H(k+1, k) \\ H(k+1, k) &= 0 \end{aligned}$$

6. If error $\leq \varepsilon$

- (a) Solve $Hy = s$ for y
- (b) Calculate $x = x - Vy$
- (c) Method has converged. Return.

Otherwise $k = k + 1$ **Example 2.9**

Repeat Example 2.6 using the GMRES method.

Solution 2.9 The problem of Example 2.6 is repeated here for convenience.
Solve

$$\begin{bmatrix} -10 & 2 & 3 & 6 \\ 0 & -9 & 1 & 4 \\ 2 & 6 & -12 & 2 \\ 3 & 1 & 0 & -8 \end{bmatrix} x = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \quad (2.84)$$

with $x^0 = [0 \ 0 \ 0 \ 0]^T$. Let $\varepsilon = 10^{-3}$.

k=1 Solving the Arnoldi process yields

$$H = \begin{bmatrix} -4.0333 & 0 & 0 & 0 \\ 6.2369 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$V = \begin{bmatrix} 0.1826 & 0.9084 & 0 & 0 \\ 0.3651 & 0.2654 & 0 & 0 \\ 0.5477 & -0.0556 & 0 & 0 \\ 0.7303 & -0.3181 & 0 & 0 \end{bmatrix}$$

Applying Given's rotation:

At $k = 1$, the H matrix is not updated; therefore,

$$\text{cs}(k = 1) = -0.5430$$

$$\text{sn}(k = 1) = 0.8397$$

The new H matrix becomes

$$H = \begin{bmatrix} 7.4274 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$s = [-2.9743 \quad -4.5993 \quad 0 \quad 0]^T$$

Since error ($= |s(2)| = 4.5993$) is greater than ε , $k = k + 1$ and repeat.

k=2 Solving the Arnoldi process yields

$$H = \begin{bmatrix} 7.4274 & 2.6293 & 0 & 0 \\ 0 & -12.5947 & 0 & 0 \\ 0 & 1.9321 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$V = \begin{bmatrix} 0.1826 & 0.9084 & 0.1721 & 0 \\ 0.3651 & 0.2654 & -0.6905 & 0 \\ 0.5477 & -0.0556 & 0.6728 & 0 \\ 0.7303 & -0.3181 & -0.2024 & 0 \end{bmatrix}$$

Applying Given's rotation yields

$$H = \begin{bmatrix} 7.4274 & -12.0037 & 0 & 0 \\ 0 & 4.6314 & 0 & 0 \\ 0 & 1.9321 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{aligned}\text{cs}(k=2) &= 0.9229 \\ \text{sn}(k=2) &= 0.3850\end{aligned}$$

Updating H

$$H = \begin{bmatrix} 7.4274 & -12.0037 & 0 & 0 \\ 0 & 5.0183 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$s = [-2.9743 \ -4.2447 \ 1.7708 \ 0]^T$$

Since error ($= |s(3)| = 1.7708$) is greater than ε , $k = k + 1$ and repeat.

k=3 Solving the Arnoldi process yields

$$H = \begin{bmatrix} 7.4274 & -12.0037 & -3.8697 & 0 \\ 0 & 5.0183 & -0.2507 & 0 \\ 0 & 0 & -13.1444 & 0 \\ 0 & 0 & 2.6872 & 0 \end{bmatrix}$$

$$V = \begin{bmatrix} 0.1826 & 0.9084 & 0.1721 & 0.3343 \\ 0.3651 & 0.2654 & -0.6905 & -0.5652 \\ 0.5477 & -0.0556 & 0.6728 & -0.4942 \\ 0.7303 & -0.3181 & -0.2024 & 0.5697 \end{bmatrix}$$

Applying Given's rotation yields

$$H = \begin{bmatrix} 7.4274 & -12.0037 & 1.8908 & 0 \\ 0 & 5.0183 & -1.9362 & 0 \\ 0 & 0 & -13.4346 & 0 \\ 0 & 0 & 2.6872 & 0 \end{bmatrix}$$

$$\begin{aligned}\text{cs}(k=3) &= -0.9806 \\ \text{sn}(k=3) &= 0.1961\end{aligned}$$

$$H = \begin{bmatrix} 7.4274 & -12.0037 & 1.8908 & 0 \\ 0 & 5.0183 & -1.9362 & 0 \\ 0 & 0 & 13.7007 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$s = [-2.9743 \ -4.2447 \ -1.7364 \ -0.3473]^T$$

Since error ($= |s(4)| = 0.3473$) is greater than ε , $k = k + 1$ and repeat.

k=4 Solving the Arnoldi process yields

$$H = \begin{bmatrix} 7.4274 & -12.0037 & 1.8908 & 1.4182 \\ 0 & 5.0183 & -1.9362 & 0.5863 \\ 0 & 0 & 13.7007 & -1.4228 \\ 0 & 0 & 0 & -9.2276 \end{bmatrix}$$

$$V = \begin{bmatrix} -0.1826 & -0.9084 & -0.1721 & -0.3343 & 0.7404 \\ -0.3651 & -0.2654 & 0.6905 & 0.5652 & 0.2468 \\ -0.5477 & 0.0556 & -0.6728 & 0.4942 & 0.6032 \\ -0.7303 & 0.3181 & 0.2024 & -0.5697 & 0.1645 \end{bmatrix}$$

Applying Given's rotation yields

$$\begin{aligned} \text{cs}(k=4) &= 1.0000 \\ \text{sn}(k=4) &= 0.0000 \end{aligned}$$

$$H = \begin{bmatrix} 7.4274 & -12.0037 & 1.8908 & -0.2778 \\ 0 & 5.0183 & -1.9362 & -1.9407 \\ 0 & 0 & 13.7007 & -1.0920 \\ 0 & 0 & 0 & 9.1919 \end{bmatrix}$$

$$s = [-2.9743 \ -4.2447 \ -1.7364 \ -0.3473 \ 0.0000]^T$$

Since error ($= |s(5)| = 0$), the iteration has converged.

Solving for y from $Hy = s$ yields

$$y = \begin{bmatrix} -1.8404 \\ -0.9105 \\ -0.1297 \\ -0.0378 \end{bmatrix} \quad (2.85)$$

Note that, since H is upper triangular, y can be found quickly using forward elimination.

Solving for x from

$$x = x - Vy \quad (2.86)$$

yields

$$x = \begin{bmatrix} -1.1981 \\ -0.8027 \\ -1.0260 \\ -1.0496 \end{bmatrix} \quad (2.87)$$

which is the same as the previous example. ■

2.7 Preconditioners for Iterative Methods

As noted previously, the rate of convergence of iterative methods depends on the condition number of the iteration matrix. Therefore, it may be advantageous to transform the linear system of equations into an equivalent system (i.e., one that has the same solution) that has more favorable properties. The matrix, or matrices, that perform this transformation are called *preconditioners*.

For example, for any nonsingular preconditioner matrix M , the transformed system

$$M^{-1}Ax = M^{-1}b \quad (2.88)$$

has the same solution as the original system $Ax = b$. Furthermore, if M is chosen such that it approximates A , then the resulting condition number of $M^{-1}A$ is much improved. There is, of course, a trade-off between the reduction in computation achieved by the improved convergence properties and the amount of computation required to find the matrix M^{-1} .

In addition, note that, even if A is symmetric, $M^{-1}A$ will most likely be nonsymmetric, and methods that require symmetry (such as the conjugate gradient method) will fail to converge. For this reason, the matrix M is often split such that

$$M = M_1 M_2 \quad (2.89)$$

The preconditioned system becomes

$$M_1^{-1} A M_2^{-1} (M_2 x) = M_1^{-1} b \quad (2.90)$$

The matrix M_1 is called the *left preconditioner* and M_2 is the *right preconditioner*. If $M_1^T = M_2$, then $M_1^{-1} A M_2^{-1}$ is symmetric. This leads directly to the ability to transform any iterative method into an equivalent iterative method by replacing the initial residual r_0 by $M_1^{-1} r_0$ and replacing the final solution x_k by $M_2^{-1} x_k$.

2.7.1 Jacobi

Just as stationary methods are based on a splitting of the A matrix such that $A = L + U + D$, preconditioners for stationary methods can be based on the splitting matrices. The simplest preconditioner consists of just the diagonal D of the A matrix:

$$m_{i,j} \begin{cases} a_{i,i} & \text{if } i = j \\ 0 & \text{otherwise} \end{cases} \quad (2.91)$$

This is known as the (point) Jacobi preconditioner. It can be split such that $m_1(i, i) = m_2(i, i) = \sqrt{a_{i,i}}$. The inverses of the matrices are straightforward to calculate. In the interest of computational efficiency, the inverse of these matrices is usually stored so that repeated divisions are not required.

The Jacobi preconditioner can be easily generalized to the *block-Jacobi* preconditioner in which small blocks along the diagonal are kept rather than individual points. Block preconditioning is useful in cases where the A matrix has blocks of coupled variables with sparse interconnections between blocks.

2.7.2 Symmetric Successive Overrelaxation

A preconditioner matrix based on the SSOR method can be defined

$$M = \frac{1}{\omega(2-\omega)}(D + \omega L)D^{-1}(D + \omega U) \quad (2.92)$$

The constant coefficient $\frac{1}{\omega(2-\omega)}$ only has the effect of scaling the equations of the preconditioned system. Furthermore, if this M matrix is used as a preconditioner, it is not necessary to choose ω as carefully as for the underlying fixed-point iteration.

2.7.3 Symmetric Gauss-Seidel

Taking $\omega = 1$ in the SSOR preconditioner leads to the symmetric Gauss-Seidel iteration

$$M = (D + L)D^{-1}(D + U) \quad (2.93)$$

Note that the preconditioned systems (2.88) or (2.90) may be a full system even if the original A matrix was sparse due to M^{-1} . Even if M is sparse, M^{-1} may be full. Trying to find M^{-1} directly may adversely impact the computational efficiency of the solution method. This must be taken into account when selecting a particular solution technique. Quite often, recasting the problem as a series of matrix-vector products can replace finding the inverse directly.

2.7.4 Incomplete LU Factorization

Since it is desired to find a preconditioner matrix M^{-1} that approximates A^{-1} , it may be reasonable to combine direct methods and iterative methods to derive the preconditioner matrix. LU factorization produces matrices L and U from which

$$U^{-1}L^{-1} = A^{-1} \quad (2.94)$$

Since L and U are lower and upper triangular, respectively, they may be computationally more efficient. An LU factorization is called *incomplete* if, during the factorization process, one or more fill elements are ignored. This has the impact of maintaining the underlying sparsity of the original matrix.

If all of the fills that result from the factorization process are neglected, this is referred to as the ILU(0) case, corresponding to incomplete LU factorization of level 0.

Example 2.10

Repeat Example 2.5 using the ILU(0) method. Compare the resulting Q matrix for both cases.

Solution 2.10 The incomplete LU factorization method results in the factors shown below. Note their similarity to the LU factors of Example 2.5.

$$Q = \begin{bmatrix} 10.000 & 0.1000 & 0 & 0.3000 & 0 & 0 & 0 & 0.5000 & 0 & 0 \\ 2.000 & 8.8000 & 0 & 0 & 0 & 0 & 0.5682 & 0 & 0 & 0.2273 \\ 0 & 0 & 21.0000 & 0.2381 & 0.3333 & 0 & 0 & 0 & 0 & 0.1905 \\ 4.0000 & 0 & 1.0000 & 16.5619 & 0.4629 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 4.0000 & 6.0476 & 20.8672 & 0.1917 & 0.0479 & 0 & 0 & 0.0593 \\ 0 & 0 & 0 & 0 & 3.0000 & 13.4249 & 0.6597 & 0 & 0 & 0 \\ 0 & 1.0000 & 4.0000 & 0 & 0.6667 & 2.8722 & 9.5051 & 0.1052 & 0.1052 & 0 \\ 1.0000 & 0 & 5.0000 & 0 & 0 & 0 & 5.0000 & 8.9740 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 6.0000 & 0 & 19.3688 & 0 \\ 0 & 2.0000 & 3.0000 & 0 & 3.0000 & 0 & 0 & 0 & 0 & 33.7960 \end{bmatrix}$$

■

Modifications to the ILU method focus on different methods of identifying which fills to include in the LU factors. Several methods are presented in [42].

2.7.5 Graph Based

A recent development in preconditioners is to use the underlying graph structure of the matrix to drive the structure of the preconditioner M . A graph-based preconditioning method is designed to be applied to symmetric systems that can be represented by graphs [30]. In this approach, matrix A is represented by a graph $G = (V, E, \omega)$ with each nonzero entry of the matrix representing a connection in the power system. To increase the speed of computation, it is beneficial to reduce the density of A . Graph theory has a natural approach to accomplishing this by removing trivial edges in graphs through graph sparsification [30] [34].

The preconditioner is derived from the low-stretch spanning tree of the matrix. To find the low-stretch spanning tree, the system matrix A is represented by a graph with m representing the number of edges in the graph or nonzero entries in the matrix. The concept of stretch is critical for constructing a good spanning tree preconditioner [6]. In order to have a spanning tree that serves as an effective preconditioner, the off-tree edges must have an average stretch δ over a spanning tree in order for the spanning tree to be an $O(\delta m)$ -approximation of the graph [30]. There exists a unique “detour” path in the tree between vertices u and v , for every edge $e(u, v)$. Stretch is defined as the distortion caused by the detour required by taking the tree path. The stretch of the edges in the tree is given by

$$\text{stretch}(e) = \frac{\sum_{i=1}^k w'(e_i)}{w'(e)} \quad (2.95)$$

The denominator of the stretch expression contains the distance between the vertices (u, v) in the graph, and the numerator sums the distances of the edges

along the unique path in the tree connecting the vertices (u, v) , with distance equaling the inverse of the edge weight $w' = 1/e$.

The spanning tree generates a preconditioner with a provably small condition number between the original graph G and the approximate graph \tilde{G} , preserving the spectral properties of the system. The approximate graph \tilde{G} , contains a nearly linear number of edges that approximates the given graph derived from A . Most of the current progress in this area is focused on various methods to construct the spanning trees [51] [28].

Example 2.11

Solve the system of Example 2.5 using the following preconditioners with the GMRES method.

1. Jacobi
 2. SSOR with $\omega = 1.0$
 3. ILU(0)

Let $b = [1 \ 1 \ \dots \ 1]^T$ and x^0 be the zero vector. Compare the number of iterations required in each method. Use a convergence tolerance of 10^{-5} .

Solution 2.5 The matrix A is repeated here for convenience.

$$A = \begin{bmatrix} 10 & 1 & 0 & 3 & 0 & 0 & 0 & 5 & 0 & 0 \\ 2 & 9 & 0 & 0 & 0 & 0 & 5 & 0 & 0 & 2 \\ 0 & 0 & 21 & 5 & 7 & 0 & 0 & 0 & 0 & 4 \\ 4 & 0 & 1 & 18 & 8 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 4 & 7 & 25 & 4 & 1 & 0 & 0 & 2 \\ 0 & 0 & 0 & 0 & 3 & 14 & 9 & 0 & 0 & 0 \\ 0 & 1 & 4 & 0 & 2 & 3 & 12 & 1 & 1 & 0 \\ 1 & 0 & 5 & 0 & 0 & 0 & 5 & 10 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 6 & 0 & 20 & 0 \\ 0 & 2 & 3 & 0 & 4 & 0 & 0 & 0 & 0 & 35 \end{bmatrix}$$

- For the Jacobi preconditioner, M will be the diagonal of the A matrix.

2. For the SSOR preconditioner with $\omega = 1$,

$$M = (D + L) D^{-1} (D + U)$$

where

with

$$M = \begin{bmatrix} 10.0000 & 1.0000 & 0 & 3.0000 & 0 & 0 & 0 & 5.0000 & 0 & 0 \\ 2.0000 & 9.2000 & 0 & 0.6000 & 0 & 0 & 5.0000 & 1.0000 & 0 & 2.0000 \\ 0 & 0 & 21.0000 & 5.0000 & 7.0000 & 0 & 0 & 0 & 0 & 4.0000 \\ 4.0000 & 0.4000 & 1.0000 & 19.4381 & 8.3333 & 0 & 0 & 2.0000 & 0 & 0.1905 \\ 0 & 0 & 4.0000 & 7.9524 & 29.4444 & 4.0000 & 1.0000 & 0 & 0 & 2.7619 \\ 0 & 0 & 0 & 0 & 3.0000 & 14.4800 & 9.1200 & 0 & 0 & 0.2400 \\ 0 & 1.0000 & 4.0000 & 0.9524 & 3.3333 & 3.3200 & 14.5641 & 1.0000 & 1.0000 & 1.1441 \\ 1.0000 & 0.1000 & 5.0000 & 1.4905 & 1.6667 & 0 & 5.0000 & 10.9167 & 0.4167 & 0.9524 \\ 0 & 0 & 0 & 0 & 0 & 0 & 6.0000 & 0.5000 & 20.5000 & 0 \\ 0 & 2.0000 & 3.0000 & 0.7143 & 5.0000 & 0.6400 & 1.2711 & 0 & 0 & 36.3359 \end{bmatrix}$$

3. The incomplete LU factorization method results in the pseudo LU factors:

and $M = LU$:

$$M = \begin{bmatrix} 10.0000 & 1.0000 & 0 & 3.0000 & 0 & 0 & 0 & 5.0000 & 0 & 0 \\ 2.0000 & 9.0000 & 0 & 0.6000 & 0 & 0 & 5.0000 & 1.0000 & 0 & 2.0000 \\ 0 & 0 & 21.0000 & 5.0000 & 7.0000 & 0 & 0 & 0 & 0 & 4.0000 \\ 4.0000 & 0.4000 & 1.0000 & 18.0000 & 8.0000 & 0 & 0 & 2.0000 & 0 & 0.1905 \\ 0 & 0 & 4.0000 & 7.0000 & 25.0000 & 4.0000 & 1.0000 & 0 & 0 & 2.0000 \\ 0 & 0 & 0 & 0 & 3.0000 & 14.0000 & 9.0000 & 0 & 0 & 0.1780 \\ 0 & 1.0000 & 4.0000 & 0.9524 & 2.0000 & 3.0000 & 12.0000 & 1.0000 & 1.0000 & 1.0287 \\ 1.0000 & 0.1000 & 5.0000 & 1.4905 & 1.6667 & 0 & 5.0000 & 10.0000 & 0.5260 & 0.9524 \\ 0 & 0 & 0 & 0 & 0 & 0 & 6.0000 & 0.6312 & 20.0000 & 0 \\ 0 & 2.0000 & 3.0000 & 0.7143 & 4.0000 & 0.5751 & 1.2801 & 0 & 0 & 35.0000 \end{bmatrix}$$

The results of each method are summarized.

Method	No pre-conditioning	Jacobi	SSOR	ILU(0)
iterations	10	7	4	4
error	3.1623	0.2338	0.1319	0.1300
	0.9439	0.0124	0.0190	0.0165
	0.2788	0.0032	0.0022	0.0012
	0.0948	0.0011	0.0002	0.0001
	0.0332	0.0003		
	0.0056	0.0001		
	0.0018	0.0000		
	0.0005			
	0.0003			
	0.0000			

Note that, without preconditioning, the GMRES method requires the full $n = 10$ iterations to converge, but with preconditioning the number of iterations is reduced. The ILU(0) is the best method, as indicated by the smallest convergence error, but only slightly over the SSOR method. In this particular example, the choice of $\omega = 1$ led to the fewest number of SSOR iterations, but this is not always the case. ■

2.8 Problems

1. Show that the number of multiplications and divisions required in the LU factorization of an $n \times n$ square matrix is $n(n^2 - 1)/3$.
2. Consider the system $Ax = b$, where

$$a_{ij} = \frac{1}{i + j - 1} \quad i, j = 1, \dots, 4$$

and

$$b_i = \frac{1}{3} \sum_{j=1}^4 a_{ij}$$

Using only four decimal places of accuracy, solve this system using LU factorization with

- (a) no pivoting
- (b) partial pivoting

Comment on the differences in solutions (if any).

3. Prove that the matrix

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

does not have an LU factorization.

4. Assuming that an LU factorization of A is available, write an algorithm to solve the equation $x^T A = b^T$.

5. For the following matrix, find $A = LU$ (no pivoting) and $PA = LU$ (with partial pivoting)

(a)

$$A = \begin{bmatrix} 6 & -2 & 2 & 4 \\ 12 & -8 & 4 & 10 \\ 3 & -13 & 3 & 3 \\ -6 & 4 & 2 & -18 \end{bmatrix}$$

(b)

$$A = \begin{bmatrix} -2 & 1 & 2 & 5 \\ 2 & -1 & 4 & 1 \\ 1 & 4 & -3 & 2 \\ 8 & 2 & 3 & -6 \end{bmatrix}$$

6. Write an LU factorization-based algorithm to find the inverse of any nonsingular matrix A .

7. Solve the system of Problem 5(a) with the vector

$$b = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

- (a) Using LU factorization and forward/backward substitution
- (b) Using a Jacobi iteration. How many iterations are required?
- (c) Using a Gauss–Seidel iteration. How many iterations are required?
- (d) Using the conjugate gradient method. How many iterations are required?
- (e) Using the GMRES method. How many iterations are required?

Use a starting vector of

$$x = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

and a convergence error of 10^{-5} for the iterative methods.

8. Apply the Gauss–Seidel iteration to the system

$$A = \begin{bmatrix} 0.96326 & 0.81321 \\ 0.81321 & 0.68654 \end{bmatrix}$$

$$b = \begin{bmatrix} 0.88824 \\ 0.74988 \end{bmatrix}$$

Use $x^0 = [0.33116 \ 0.70000]^T$ and explain what happens.

9. Solve the system of equations in Problem 2 using the conjugate gradient method.

10. Solve the system of equations in Problem 2 using the GMRES method.

11. Consider an $n \times n$ tridiagonal matrix of the form

$$T_a = \begin{bmatrix} a & -1 & & & \\ -1 & a & -1 & & \\ & -1 & a & -1 & \\ & & -1 & a & -1 \\ & & & -1 & a & -1 \\ & & & & -1 & a \end{bmatrix}$$

where a is a real number.

- (a) Verify that the eigenvalues of T_a are given by

$$\lambda_j = a - 2 \cos(j\theta) \quad j = 1, \dots, n$$

where

$$\theta = \frac{\pi}{n+1}$$

- (b) Let $a = 2$.
- Will the Jacobi iteration converge for this matrix?
 - Will the Gauss–Seidel iteration converge for this matrix?
12. An alternative conjugate gradient algorithm for solving $Ax = b$ may be based on the error functional $E_k(x^k) = \langle x^k - x, x^k - x \rangle$ where $\langle \cdot \rangle$ denotes inner product. The solution is given as

$$x^{k+1} = x^k + \alpha_k \sigma_k$$

Using $\sigma_1 = -A^T r_0$ and $\sigma_{k+1} = -A^T r_k + \beta_k \sigma_k$, derive this conjugate gradient algorithm. The coefficients α_k and β_k can be expressed as

$$\alpha_{k+1} = \frac{\|r_k\|^2}{\|\sigma_{k+1}\|^2}$$

$$\beta_{k+1} = \frac{\|r_{k+1}\|^2}{\|r_k\|^2}$$

Repeat Example 2.8 using this conjugate gradient algorithm.

13. Write a subroutine with two inputs (A , flag) that will generate, for any nonsingular matrix A , the outputs (Q, P) such that if
- flag=0, $A = LU, P = I$
 - flag=1, $PA = LU$

where

$$L = \begin{bmatrix} l_{11} & 0 & 0 & \dots & 0 \\ l_{21} & l_{22} & 0 & \dots & 0 \\ l_{31} & l_{32} & l_{33} & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ l_{n1} & l_{n2} & l_{n3} & \dots & l_{nn} \end{bmatrix} \quad \text{and} \quad U = \begin{bmatrix} 1 & u_{12} & u_{13} & \dots & u_{1n} \\ 0 & 1 & u_{23} & \dots & u_{2n} \\ 0 & 0 & 1 & \dots & u_{3n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}$$

and

$$Q = L + U - I$$

14. For the following nonsingular matrices, use the subroutine of Problem 13 and obtain matrices P and Q in each of the following cases:
- (a)

$$\begin{bmatrix} 0 & 0 & 1 \\ 3 & 1 & 4 \\ 2 & 1 & 0 \end{bmatrix}$$

(b)

$$\begin{bmatrix} 10^{-10} & 0 & 0 & 1 \\ 0 & 0 & 1 & 4 \\ 0 & 2 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

15. Write a subroutine with two inputs (A, b) that will generate, for any nonsingular matrix A , the output (x) such that

$$Ax = b$$

using forward and backward substitution. This subroutine should incorporate the subroutine developed in Problem 13.

16. Using the subroutines of Problems 13 and 15, solve the following system of equations:

$$\begin{bmatrix} 2 & 5 & 6 & 11 \\ 4 & 6 & 8 & 2 \\ 4 & 3 & 7 & 0 \\ 1 & 26 & 3 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

17. For the matrix given below, using the following preconditioners with the GMRES method.

- (a) Jacobi
- (b) SSOR with $\omega = 1.0$
- (c) ILU(0)

Let $b = [1 \ 1 \ \dots \ 1]^T$ and x^0 be the zero vector. Compare the number of iterations required in each method. Use a convergence tolerance of 10^{-5} .

$$A = \begin{bmatrix} 10 & 0 & 43 & 0 & 23 & 0 & 0 & 0 & 0 & 0 \\ 60 & 41 & 10 & 0 & 0 & 47 & 0 & 0 & 64 & 0 \\ 48 & 0 & 27 & 96 & 85 & 0 & 0 & 29 & 95 & 36 \\ 0 & 0 & 0 & 25 & 0 & 0 & 0 & 0 & 0 & 67 \\ 0 & 97 & 29 & 0 & 13 & 78 & 61 & 0 & 71 & 0 \\ 0 & 54 & 0 & 29 & 0 & 45 & 0 & 69 & 0 & 85 \\ 0 & 0 & 0 & 0 & 0 & 0 & 3 & 55 & 12 & 84 \\ 0 & 0 & 0 & 0 & 0 & 48 & 25 & 43 & 61 & 26 \\ 32 & 0 & 0 & 7 & 0 & 0 & 92 & 0 & 46 & 62 \\ 0 & 78 & 0 & 0 & 0 & 0 & 0 & 65 & 0 & 59 \end{bmatrix}$$

3

Systems of Nonlinear Equations

Many systems can be modeled generically as

$$F(x) = 0 \quad (3.1)$$

where x is an n -vector and F represents a nonlinear mapping with both its domain and range in the n -dimensional real linear space R^n . The mapping F can also be interpreted as being an n -vector of functions

$$F(x) = \begin{bmatrix} f_1(x_1, x_2, \dots, x_n) \\ f_2(x_1, x_2, \dots, x_n) \\ \vdots \\ f_n(x_1, x_2, \dots, x_n) \end{bmatrix} = 0 \quad (3.2)$$

where at least one of the functions is nonlinear. Each function may or may not involve all n states x_i , but it is assumed that every state appears at least once in the set of functions. The solution x^* of the nonlinear system cannot, in general, be expressed in closed form. Thus nonlinear systems are usually solved numerically. In many cases, it is possible to find an approximate solution \hat{x} arbitrarily close to the actual solution x^* , by replacing each approximation with successively better (more accurate) approximations until

$$F(\hat{x}) \approx 0$$

Such methods are usually *iterative*. An iterative solution is one in which an initial guess (x^0) to the solution is used to create a sequence x^0, x^1, x^2, \dots that (hopefully) converges arbitrarily close to the desired solution x^* .

Three principal issues arise with the use of iterative methods, namely,

1. Is the iterative process well defined? That is, can it be successively applied without numerical difficulties?
2. Do the iterates (i.e., the sequence of updates) converge to a solution of Equation (3.1)? Is the solution the desired solution?
3. How economical is the entire solution process?

The complete (or partial) answers to these issues are enough to fill several volumes, and as such cannot be discussed in complete detail in this chapter.

These issues, however, are central to the solution of nonlinear systems and cannot be fully ignored. Therefore, this chapter will endeavor to provide sufficient detail for the reader to be aware of the advantages (and disadvantages) of different types of iterative methods without providing exhaustive coverage.

3.1 Fixed-Point Iteration

Solving a system of nonlinear equations is a complex problem. To better understand the mechanisms involved in a large-scale system, it is instructive to first consider the one-dimensional, or scalar, nonlinear system

$$f(x) = 0 \quad (3.3)$$

One approach to solving any nonlinear equation is the tried-and-true “trial and error” method that most engineering and science students have used at one time or another in their careers.

Example 3.1

Find the solution to

$$f(x) = x^2 - 5x + 4 = 0 \quad (3.4)$$

Solution 3.1 This is a quadratic equation that has a closed form solution. The two solutions are

$$x_1^*, x_2^* = \frac{5 \pm \sqrt{(-5)^2 - (4)(4)}}{2} = 1, 4$$

If a closed form solution did not exist, however, one approach would be to use a trial and error approach. Since the solution occurs when $f(x) = 0$, the value of $f(x)$ can be monitored and used to refine the estimates to x^* .

k	x	$f(x)$
0	0	$0 - 0 + 4 = 4 > 0$
1	2	$4 - 10 + 4 = -2 < 0$
2	0.5	$0.25 - 2.5 + 4 = 1.75 > 0$
3	1.5	$2.25 - 7.5 + 4 = -1.25 < 0$

By noting the sign of the function and whether or not it changes sign, the interval in which the solution lies can be successively narrowed. If a function $f(x)$ is continuous and $f(a) \cdot f(b) < 0$, then the equation $f(x) = 0$ has at least one solution in the interval (a, b) . Since $f(0.5) > 0$ and $f(1.5) < 0$, it can be concluded that one of the solutions lies in the interval $(0.5, 1.5)$. ■

This process, however, tends to be tedious, and there is no guidance to determine what the next guess should be other than the bounds established by the change in sign of $f(x)$. A better method would be to write the sequence of updates in terms of the previous guesses. Thus an iterative function can be defined as

$$I : \quad x^{k+1} = g(x^k), \quad k = 1, \dots, \infty \quad (3.5)$$

This is known as a fixed-point iteration because at the solution

$$x^* = g(x^*) \quad (3.6)$$

Example 3.2

Find the solution to Equation (3.4) using a fixed-point iteration.

Solution 3.2 Equation (3.4) can be rewritten as

$$x = \frac{x^2 + 4}{5} \quad (3.7)$$

Adopting the notation of Equation (3.5), the iterative function becomes

$$x^{k+1} = g(x^k) = \frac{(x^k)^2 + 4}{5} \quad (3.8)$$

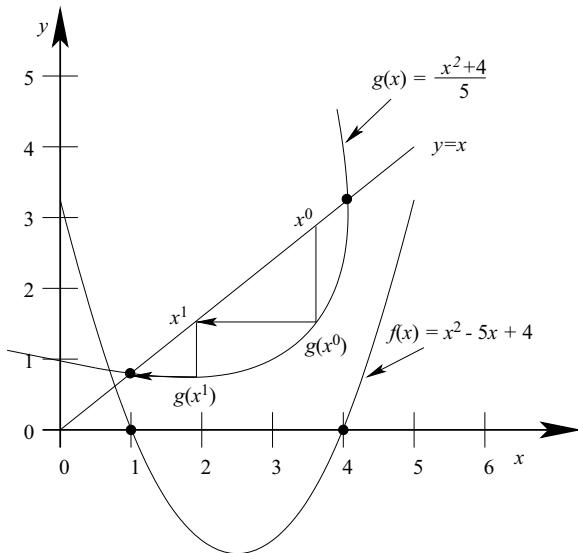
Using this iterative function, the estimates to x^* are

k	x^k	x^{k+1}
0	0	$\frac{0+4}{5} = 0.8$
1	0.8	$\frac{0.64+4}{5} = 0.928$
2	0.928	$\frac{0.856+4}{5} = 0.972$
3	0.971	$\frac{0.943+4}{5} = 0.989$

It is obvious that this sequence is converging to the solution $x^* = 1$.

Now consider the same example, except with a different initial guess:

k	x^k	x^{k+1}
0	5	$\frac{25+4}{5} = 5.8$
1	5.8	$\frac{33.64+4}{5} = 7.528$
2	7.528	$\frac{56.67+4}{5} = 12.134$

**FIGURE 3.1**

Graphical interpretation of the fixed-point iteration

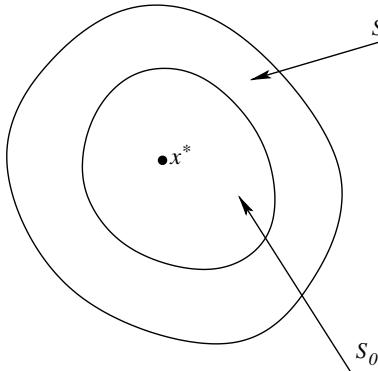
In this case, the iterates are increasing rapidly and after a few more iterations would approach infinity. In this case, it is said that the iteration is *diverging*. ■

This example brings up two very important points: will a sequence of iterates converge, and, if so, to what solution will they converge? In order to address these questions, consider first a graphical interpretation of Example 3.2. Plotting both sides of the function in Equation (3.7) yields the two lines shown in Figure 3.1.

These two lines intersect at the same two points in which the original function $f(x) = 0$. The fixed-point iteration works by finding this intersection. Consider the initial guess x^0 shown in Figure 3.1. The function $g(x)$ evaluated at x^0 gives the updated iterate x^1 . Thus a vertical line projected from x^0 points to $g(x^0)$ and a horizontal line projected from $g(x^0)$ gives x^1 .

The projection of the function $g(x^1)$ yields x^2 . Similar vertical and horizontal projections will eventually lead directly to the point at which the two lines intersect. In this way, the solution to the original function $f(x)$ can be obtained.

In this example, the solution $x^* = 1$ is the *point of attraction* of the fixed-point iteration. A point x^* is said to be a point of attraction of an iterative function I if there exists an open neighborhood S_0 of x^* such that, for all initial guesses x^0 in the subset S_0 of S , the iterates will remain in S and

**FIGURE 3.2**Domain of attraction of x^*

$$\lim_{k \rightarrow \infty} x^k = x^* \quad (3.9)$$

The neighborhood S_0 is called the *domain of attraction* of x^* [38]. This concept is illustrated in Figure 3.2 and implies that the iterates of I will converge to x^* whenever x^0 is sufficiently close to x^* . In Example 3.2, the fixed point $x^* = 1$ is a point of attraction of

$$I : x^{k+1} = \frac{(x^k)^2 + 4}{5}$$

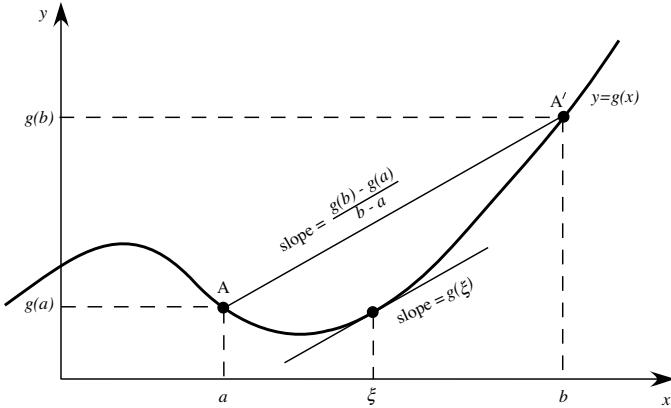
whereas $x^* = 4$ is not. The domain of attraction of $x^* = 1$ is all x in the domain $-\infty < x < 4$.

It is often difficult to determine a priori whether or not an iteration will converge. In some cases, a series of iterates will appear to be converging, but will not approach x^* even as $k \rightarrow \infty$. However, there are a number of theorems that provide insight as to whether an iteration of $x = g(x)$ will converge.

Mean Value Theorem: [52] Suppose a function $g(x)$ and its derivative $g'(x)$ are both continuous in the interval $a \leq x \leq b$. Then there exists at least one ξ , $a < \xi < b$ such that

$$g'(\xi) = \frac{g(b) - g(a)}{b - a} \quad (3.10)$$

The meaning of this theorem is shown in Figure 3.3. If a function $g(x)$ is defined in the region between $x = a$ and $x = b$ and is both differentiable (smooth) and continuous, then a secant line can be drawn between points A

**FIGURE 3.3**

Meaning of the mean value theorem

and A' . The slope of this secant line is

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

The mean value theorem states that there is at least one point on the curve, at $x = \xi$, where the tangent to the curve has the same slope as the line AA' .

Rewriting the equation of the mean value theorem as

$$g(b) - g(a) = g'(\xi)(b - a)$$

then for successive iterates in which $x^{k+1} = b$ and $x^k = a$, then

$$g(x^{k+1}) - g(x^k) = g'(\xi^k)(x^{k+1} - x^k) \quad (3.11)$$

or taking the absolute values

$$|g(x^{k+1}) - g(x^k)| = |g'(\xi^k)| |(x^{k+1} - x^k)| \quad (3.12)$$

As long as the appropriate ξ^k is used, the mean value theorem can be successively applied to each iteration of the sequence. If the entire region includes x^* as well as all of the x^k , then the derivative $g'(x)$ is bounded. Therefore

$$|g'(\xi^k)| \leq M \quad (3.13)$$

for any k where M is the positive upper bound. Then, starting from the initial guess x^0 ,

$$|x^2 - x^1| \leq M |x^1 - x^0| \quad (3.14)$$

$$|x^3 - x^2| \leq M |x^2 - x^1| \quad (3.15)$$

$$\vdots \quad (3.16)$$

$$|x^{k+1} - x^k| \leq M |x^k - x^{k-1}| \quad (3.17)$$

and by combining yields

$$|x^{k+1} - x^k| \leq M^k |x^1 - x^0| \quad (3.18)$$

Thus, for any initial guess x^0 ,

$$|g'(x)| \leq M < 1 \quad (3.19)$$

the iterates will converge.

A similar, but slightly different method of determining whether or not an iterative process will converge is given by the **Ostrowski** theorem [38]. This theorem states that, if the iterative process

$$I : \quad x^{k+1} = g(x^k), \quad k = 1, \dots, \infty$$

has a fixed-point x^* and is continuous and differentiable at x^* , and if $\left| \frac{\partial g(x^*)}{\partial x} \right| < 1$, then x^* is a point of attraction of I .

Example 3.3

Determine whether $x^* = 1$ and $x^* = 4$ are points of attraction of the iterative function of Equation (3.8).

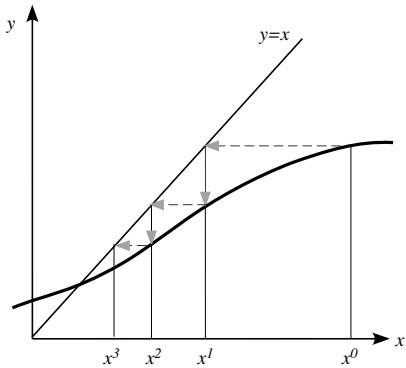
Solution 3.3 The derivative of the iterative process I in Equation (3.8) is

$$\left| \frac{\partial g(x)}{\partial x} \right| = \left| \frac{2}{5}x \right|$$

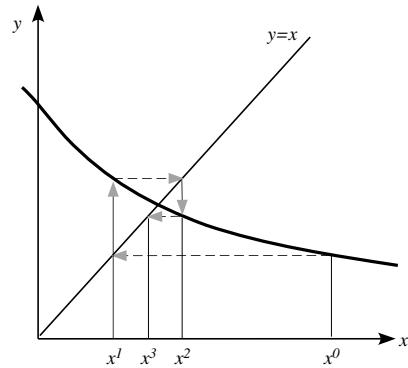
Thus, for $x^* = 1$, $\left| \frac{2}{5}x^* \right| = \frac{2}{5} < 1$ and $x^* = 1$ is a point of attraction of I . For $x^* = 4$, $\left| \frac{2}{5}x^* \right| = \frac{2}{5}(4) = \frac{8}{5} > 1$; thus $x^* = 4$ is not a point of attraction of I . ■

There are four possible convergence types for fixed-point iterations. These are shown graphically in Figure 3.4. Figure 3.4(a) shows what happens if $g'(x)$ is between 0 and 1. Even if the initial guess x_0 is far from x^* , the successive values of x^k approach the solution from one side – this is defined as *monotonic convergence*. Figure 3.4(b) shows the situation when $g'(x)$ is between -1 and 0 . Even if the initial guess x_0 is far from x^* , the successive values of x^k approach the solution from first one side and then the other oscillating around the root. This convergence is *oscillatory convergence*.

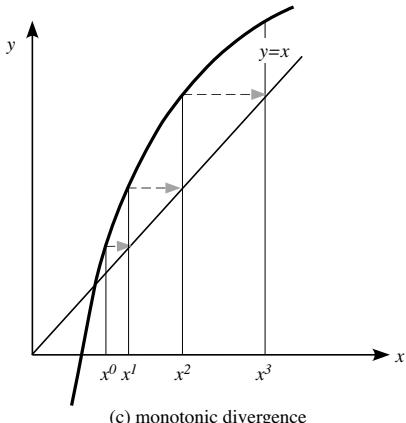
Figure 3.4(c) shows the case when $g'(x)$ is greater than 1, leading to *monotonic divergence*. Figure 3.4(d) illustrates the case when $g'(x) < -1$ and $|g'(x)| > 1$. This is *oscillatory divergence*.



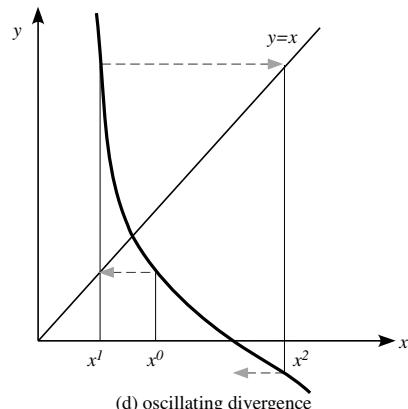
(a) monotonic convergence



(b) oscillating convergence



(c) monotonic divergence



(d) oscillating divergence

FIGURE 3.4

Four possible convergence types in the iteration $x = g(x)$

3.2 Newton–Raphson Iteration

Several iterative methods offer more robust convergence behavior than the simple fixed-point iteration described in the previous section. One of the most widely used iterative methods is the *Newton–Raphson* iterative method. This method can also be described by the iterative process

$$I : \quad x^{k+1} = g(x^k), \quad k = 1, \dots, \infty$$

but frequently offers better convergence properties than the fixed-point iteration.

Consider again the scalar nonlinear function

$$f(x^*) = 0 \quad (3.20)$$

Expanding this function in a Taylor series expansion about the point x^k yields

$$f(x^*) = f(x^k) + \frac{\partial f}{\partial x} \Big|_{x^k} (x^* - x^k) + \frac{1}{2!} \frac{\partial^2 f}{\partial x^2} \Big|_{x^k} (x^* - x^k)^2 + \dots = 0 \quad (3.21)$$

If it is assumed that the iterates will converge to x^* as $k \rightarrow \infty$, then the updated guess x^{k+1} can be substituted for x^* , yielding

$$f(x^{k+1}) = f(x^k) + \frac{\partial f}{\partial x} \Big|_{x^k} (x^{k+1} - x^k) + \frac{1}{2!} \frac{\partial^2 f}{\partial x^2} \Big|_{x^k} (x^{k+1} - x^k)^2 + \dots = 0 \quad (3.22)$$

If the initial guess is “sufficiently close” to x^* and within the domain of attraction of x^* , then the higher-order terms of the expansion can be neglected, yielding

$$f(x^{k+1}) = f(x^k) + \frac{\partial f}{\partial x} \Big|_{x^k} (x^{k+1} - x^k) \approx 0 \quad (3.23)$$

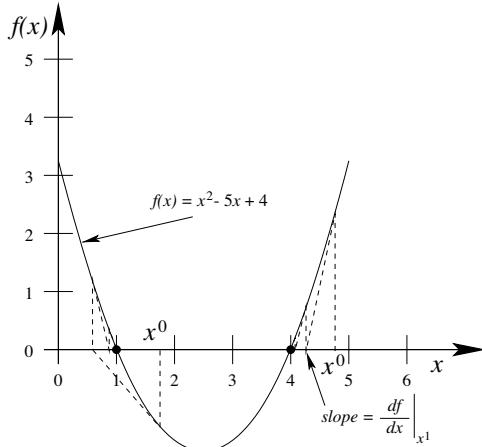
Solving directly for x^{k+1} as a function of x^k yields the following iterative function:

$$I : \quad x^{k+1} = x^k - \left[\frac{\partial f}{\partial x} \Big|_{x^k} \right]^{-1} f(x^k) \quad (3.24)$$

which is the well-known Newton–Raphson iterative method.

The Newton–Raphson method also lends itself to a graphical interpretation. Consider the same function as in Example 3.2 plotted in Figure 3.5. In this method, the slope of the function evaluated at the current iteration is used to produce the next guess. For any guess x^k , there corresponds a point on the function $f(x^k)$ with slope

$$\left. \frac{\partial f}{\partial x} \right|_{x=x^k}$$

**FIGURE 3.5**

Graphical interpretation of the Newton–Raphson method

Therefore, the next guess x^{k+1} is simply the intersection of the slope and the x -axis. This process is repeated until the guesses are sufficiently close to the solution x^* . An iteration is said to have converged at x^k if

$$|f(x^k)| < \varepsilon$$

where ε is some predetermined tolerance.

Example 3.4

Repeat Example 3.2 using a Newton–Raphson iteration.

Solution 3.4 Using the Newton–Raphson method of Equation (3.24), the iterative function is given by

$$I : x^{k+1} = x^k - \frac{(x^k)^2 - 5x^k + 4}{2x^k - 5} \quad (3.25)$$

Using this iterative function, the estimates to x^* from an initial guess of $x^0 = 3$ are

k	x^k	x^{k+1}
0	3	$3 - \frac{9-15+4}{6-5} = 5$
1	5	$5 - \frac{25-25+4}{10-5} = 4.2$
2	4.2	$4.2 - \frac{17.64-21+4}{8.4-5} = 4.012$

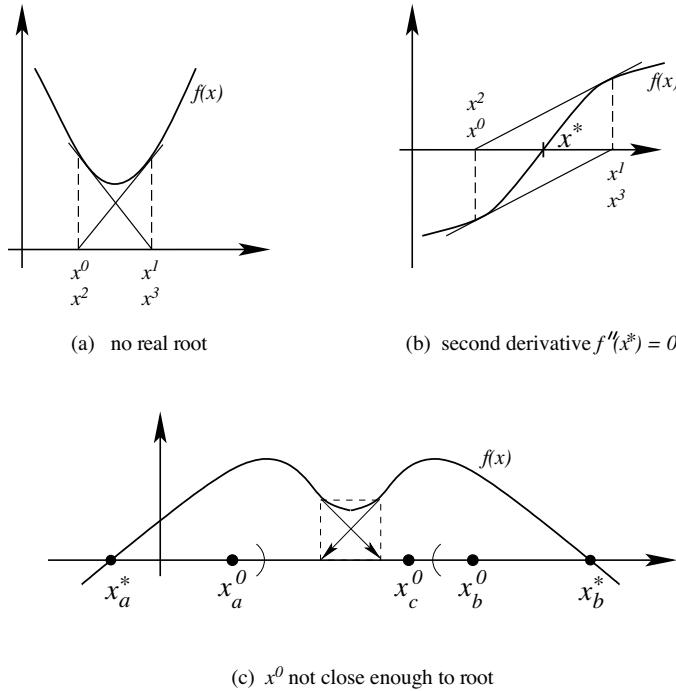


FIGURE 3.6

Newton–Raphson regions of convergence

Similarly, the estimates to x^* from an initial guess of $x^0 = 2$ are

k	x^k	x^{k+1}
0	2	$2 - \frac{4-10+4}{4-5} = 0$
1	0	$0 - \frac{0-0+4}{0-5} = 0.8$
2	0.8	$0.8 - \frac{0.64-4+4}{1.6-5} = 0.988$

In this case, both solutions are points of attraction of the Newton–Raphson iteration. ■

In some cases, however, the Newton–Raphson method will also fail to converge. Consider the functions shown in Figure 3.6. In Figure 3.6(a), the function has no real root. In Figure 3.6(b), the function is symmetric around x^* and the second derivative is zero. In Figure 3.6(c), an initial guess of x_a^0 will converge to the solution x_a^* . An initial guess of x_b^0 will converge to the

solution x_b^* . However, an initial guess of x_c^0 will cause the iterates to get locked in and oscillate in the region denoted by the dashed box without ever converging to a solution. This figure supports the assertion that, if the initial guess is too far away from the actual solution, the iterates may not converge. Or conversely, the initial guess must be sufficiently close to the actual solution for the Newton–Raphson iteration to converge. This supports the initial assumption used to derive the Newton–Raphson algorithm in that, *if the iterates were sufficiently close to the actual solution*, the higher-order terms of the Taylor series expansion could be neglected. If the iterates are not sufficiently close to the actual solution, these higher-order terms are significant and the assumption upon which the Newton–Raphson algorithm is based is not valid.

3.2.1 Convergence Properties

Note that the rate of convergence to the solution in Example 3.4 is much faster than in Example 3.2. This is because the Newton–Raphson method exhibits *quadratic convergence*, whereas the fixed-point iteration exhibits only *linear convergence*. Linear convergence implies that, once the iterates x^k are sufficiently close to the actual solution x^* , then the error

$$\varepsilon^k = |x^k - x^*| \quad (3.26)$$

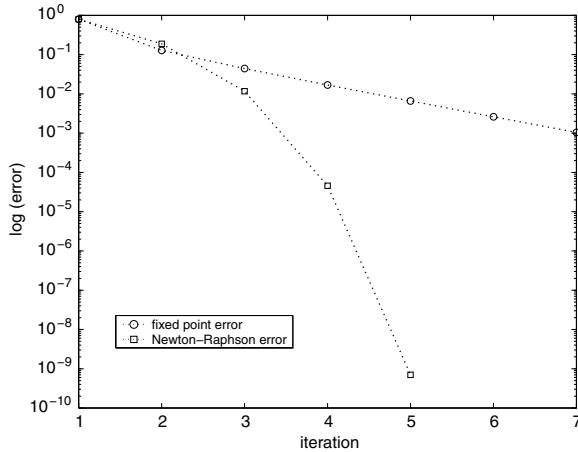
will approach zero in a linear fashion. The convergence of Examples 3.2 and 3.4 is shown in Figure 3.7. Plotted on a log-scale plot, the error for the fixed-point iteration is clearly linear, whereas the Newton–Raphson error exhibits quadratic convergence until it becomes too small to plot. Numerous methods have been proposed to predict the rate of convergence of iterative methods. Let the error of an iterative function be defined as in Equation (3.26). If there exists a number p and a constant $C \neq 0$ such that

$$\lim_{k \rightarrow \infty} \frac{|\varepsilon^{k+1}|}{|\varepsilon^k|^p} = C \quad (3.27)$$

then p is called the *order of convergence* of the iterative sequence and C is the *asymptotic error constant*. If $p = 1$, the convergence is said to be *linear*. If $p = 2$, the convergence is *quadratic*, and if $p = 3$, the order of convergence is *cubic*. The Newton–Raphson method satisfies Equation (3.27) with $p = 2$ if

$$C = \frac{1}{2} \frac{\left| \frac{d^2 f(x^*)}{dx^2} \right|}{\left| \frac{df(x^*)}{dx} \right|}$$

where $C \neq 0$ only if $\frac{d^2 f(x^*)}{dx^2} \neq 0$. Thus, for most functions, the Newton–Raphson method exhibits quadratic convergence.

**FIGURE 3.7**

Nonconverging iteration (fixed-point vs. Newton–Raphson)

3.2.2 The Newton–Raphson for Systems of Nonlinear Equations

In science and engineering, many applications give rise to *systems* of equations such as those in Equation (3.2). With a few modifications, the Newton–Raphson method developed in the previous section can be extended to systems of nonlinear equations. Systems of equations can similarly be represented by Taylor series expansions. By making the assumption once again that the initial guess is sufficiently close to the exact solution, the multidimensional higher-order terms can be neglected, yielding the Newton–Raphson method for n -dimensional systems:

$$x^{k+1} = x^k - [J(x^k)]^{-1} F(x^k) \quad (3.28)$$

where

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$

$$F(x^k) = \begin{bmatrix} f_1(x^k) \\ f_2(x^k) \\ f_3(x^k) \\ \vdots \\ f_n(x^k) \end{bmatrix}$$

and the Jacobian matrix $[J(x^k)]$ is given by

$$[J(x^k)] = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \frac{\partial f_3}{\partial x_1} & \frac{\partial f_3}{\partial x_2} & \frac{\partial f_3}{\partial x_3} & \cdots & \frac{\partial f_3}{\partial x_n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \frac{\partial f_n}{\partial x_3} & \cdots & \frac{\partial f_n}{\partial x_n} \end{bmatrix}$$

Typically, the inverse of the Jacobian $[J(x^k)]$ is not found directly, but rather through LU factorization by posing the Newton–Raphson method as

$$[J(x^k)](x^{k+1} - x^k) = -F(x^k) \quad (3.29)$$

which is now in the form $Ax = b$ where the Jacobian is the matrix A , the function $-F(x^k)$ is the vector b , and the unknown x is the difference vector $(x^{k+1} - x^k)$. Convergence is typically evaluated by considering the norm of the function

$$\|F(x^k)\| < \varepsilon \quad (3.30)$$

Note that the Jacobian is a function of x^k and is therefore updated every iteration along with $F(x^k)$.

Example 3.5

Find the solution to

$$0 = x_1^2 + x_2^2 - 5x_1 + 1 = f_1(x_1, x_2) \quad (3.31)$$

$$0 = x_1^2 - x_2^2 - 3x_2 - 3 = f_2(x_1, x_2) \quad (3.32)$$

with an initial guess of

$$x^{(0)} = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

Solution 3.5 The Jacobian of this system of equations is

$$J(x_1, x_2) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1 - 5 & 2x_2 \\ 2x_1 & -2x_2 - 3 \end{bmatrix}$$

Iteration 1

The Jacobian and the functions f_1 and f_2 are evaluated at the initial condition

$$\begin{bmatrix} 1 & 6 \\ 6 & -9 \end{bmatrix} \begin{bmatrix} x_1^{(1)} - 3 \\ x_2^{(1)} - 3 \end{bmatrix} = \begin{bmatrix} -4 \\ 12 \end{bmatrix} \quad (3.33)$$

Solving this linear system yields

$$\begin{bmatrix} x_1^{(1)} - 3 \\ x_2^{(1)} - 3 \end{bmatrix} = \begin{bmatrix} 0.8 \\ -0.8 \end{bmatrix} \quad (3.34)$$

Thus

$$x_1^{(1)} = 0.8 + x_1^{(0)} = 0.8 + 3 = 3.8 \quad (3.35)$$

$$x_2^{(1)} = -0.8 + x_2^{(0)} = -0.8 + 3 = 2.2 \quad (3.36)$$

The error at iteration 1 is

$$\left\| \begin{bmatrix} f_1(x_1^{(0)}, x_2^{(0)}) \\ f_2(x_1^{(0)}, x_2^{(0)}) \end{bmatrix} \right\|_{\infty} = 12$$

Iteration 2

The Jacobian and the functions f_1 and f_2 are evaluated at $x^{(1)}$

$$\begin{bmatrix} 2.6 & 4.4 \\ 7.6 & -7.4 \end{bmatrix} \begin{bmatrix} x_1^{(2)} - 3.8 \\ x_2^{(2)} - 2.2 \end{bmatrix} = \begin{bmatrix} -1.28 \\ 0.00 \end{bmatrix} \quad (3.37)$$

Solving this linear system yields

$$\begin{bmatrix} x_1^{(2)} - 3.8 \\ x_2^{(2)} - 2.2 \end{bmatrix} = \begin{bmatrix} -0.1798 \\ -0.1847 \end{bmatrix} \quad (3.38)$$

Thus

$$x_1^{(2)} = -0.1798 + x_1^{(1)} = -0.1798 + 3.8 = 3.6202 \quad (3.39)$$

$$x_2^{(2)} = -0.1847 + x_2^{(1)} = -0.1847 + 2.2 = 2.0153 \quad (3.40)$$

The error at iteration 2 is

$$\left\| \begin{bmatrix} f_1(x_1^{(1)}, x_2^{(1)}) \\ f_2(x_1^{(1)}, x_2^{(1)}) \end{bmatrix} \right\|_{\infty} = 1.28$$

Iteration 3

The Jacobian and the functions f_1 and f_2 are evaluated at $x^{(2)}$

$$\begin{bmatrix} 2.2404 & 4.0307 \\ 7.2404 & -7.0307 \end{bmatrix} \begin{bmatrix} x_1^{(3)} - 3.6202 \\ x_2^{(3)} - 2.0153 \end{bmatrix} = \begin{bmatrix} -0.0664 \\ 0.0018 \end{bmatrix} \quad (3.41)$$

Solving this linear system yields

$$\begin{bmatrix} x_1^{(3)} - 3.6202 \\ x_2^{(3)} - 2.0153 \end{bmatrix} = \begin{bmatrix} -0.0102 \\ -0.0108 \end{bmatrix} \quad (3.42)$$

Thus

$$x_1^{(3)} = -0.0102 + x_1^{(2)} = -0.0102 + 3.6202 = 3.6100 \quad (3.43)$$

$$x_2^{(3)} = -0.0108 + x_2^{(2)} = -0.0108 + 2.0153 = 2.0045 \quad (3.44)$$

The error at iteration 3 is

$$\left\| \begin{bmatrix} f_1(x_1^{(2)}, x_2^{(2)}) \\ f_2(x_1^{(2)}, x_2^{(2)}) \end{bmatrix} \right\|_{\infty} = 0.0664$$

At iteration 4, the functions f_1 and f_2 are evaluated at $x^{(3)}$ and yield the following:

$$\begin{bmatrix} f_1(x_1^{(3)}, x_2^{(3)}) \\ f_2(x_1^{(3)}, x_2^{(3)}) \end{bmatrix} = \begin{bmatrix} -0.221 \times 10^{-3} \\ 0.012 \times 10^{-3} \end{bmatrix}$$

Since the norm of this matrix is very small, it can be concluded that the iterates have converged and

$$\begin{bmatrix} x_1^{(3)} \\ x_2^{(3)} \end{bmatrix} = \begin{bmatrix} 3.6100 \\ 2.0045 \end{bmatrix}$$

are within an order of error of 10^{-3} of the actual solution. ■

In the solution to Example 3.5, the error at each iteration is

iteration	error
0	12.0000
1	1.2800
2	0.0664
3	0.0002

Note that, once the solution is sufficiently close to the actual solution, the error at each iteration decreases rapidly. If the iterations were carried far enough, the error at each iteration would become roughly the square of the previous iteration error. This convergence behavior is indicative of the quadratic convergence of the Newton–Raphson method.

3.3 Quasi-Newton Methods

Although the full Newton–Raphson method exhibits quadratic convergence and a minimum number of iterations, each iteration may require significant computation. For example, the computation of a full Jacobian matrix requires

n^2 calculations, and each iteration requires on the order of n^3 operations for the LU factorization if the Jacobian is a full matrix. Therefore, most modifications to the Newton–Raphson method propose to reduce either the calculation or the LU factorization of the Jacobian matrix.

Consider once again the iterative statement of the Newton–Raphson method:

$$I : \quad x^{k+1} = x^k - [J(x^k)]^{-1} f(x^k)$$

This iterative statement can be written in a more general form as

$$I : \quad x^{k+1} = x^k - [M(x^k)]^{-1} f(x^k) \quad (3.45)$$

where M is an $n \times n$ matrix that may or may not be a function of x^k . Note that, even if $M \neq J$, this iteration will still converge to a correct solution for x if the function $f(x)$ is driven to zero. So one approach to simplifying the Newton–Raphson method is to find a suitable substitute matrix M that is easier to compute than the system Jacobian. Methods that use a substitute matrix M in place of J are known as *quasi-Newton* methods.

3.3.1 Secant Method

The Newton–Raphson method is based on using the line tangent to the function $y = f(x)$. By using the slope of the tangent line, the update to the iteration can be calculated. The difficulty with this method is the computational complexity required to compute the function derivative, or $f'(x)$. An alternate approach to calculating the slope of the tangent is to take two points close to the desired root and interpolate between them to estimate the slope, as shown in Figure 3.8. This produces the linear function

$$q(x) = a_0 + a_1 x \quad (3.46)$$

where $q(x^0) = f(x^0)$ and $q(x^1) = f(x^1)$. This line is the *secant* line and is given by

$$q(x) = \frac{(x^1 - x) f(x^0) + (x - x^0) f(x^1)}{x^1 - x^0} \quad (3.47)$$

Setting $x_2 = x$ and solving yields

$$x^2 = x^1 - f(x^1) \left[\frac{f(x^1) - f(x^0)}{x^1 - x^0} \right]^{-1} \quad (3.48)$$

The process can now be repeated by using x^2 and x^1 to produce another secant line. By repeatedly updating the secant line, the generalized formula becomes

$$x^{k+1} = x^k - f(x^k) \left[\frac{f(x^k) - f(x^{k-1})}{x^k - x^{k-1}} \right]^{-1} \quad (3.49)$$

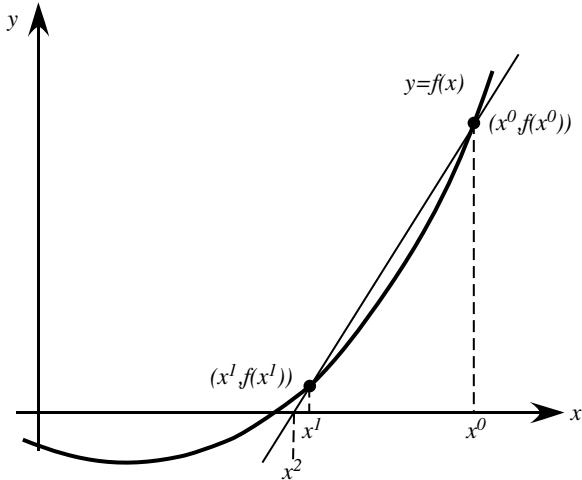
**FIGURE 3.8**

Illustration of secant method

Note that the secant method can be considered an approximation of the Newton–Raphson method

$$x^{k+1} = x^k - \frac{f(x^k)}{f'(x^k)} \quad (3.50)$$

by using the approximation

$$f'(x^k) = \frac{f(x^k) - f(x^{k-1})}{x^k - x^{k-1}} \quad (3.51)$$

The secant method is often faster than the Newton–Raphson method even though it requires a greater number of iterations to converge to the same level of accuracy. This is because the Newton–Raphson method requires two function evaluations ($f(x^k)$ and $f'(x^k)$), whereas the secant method only requires one function evaluation ($f(x^k)$) since $f(x^{k-1})$ can be saved from the previous iteration.

The secant method exhibits *super linear* convergence; its convergence is faster than linear convergence, but not as fast as quadratic convergence (of the Newton–Raphson method). Let the error at iteration k be given by

$$e^k = x^k - x^* \quad (3.52)$$

where x^* is the exact solution. Using the Taylor series expansion

$$f(x^k) = f(x^* + (x^k - x^*)) \quad (3.53)$$

$$= f(x^* + e^k) \quad (3.54)$$

$$= f(x^*) + f'(x^*)e^k + \frac{1}{2}f''(x^*)(e^k)^2 + \dots \quad (3.55)$$

Similarly,

$$f(x^{k-1}) = f(x^*) + f'(x^*)e^{k-1} + \frac{1}{2}f''(x^*)(e^{k-1})^2 + \dots \quad (3.56)$$

Furthermore,

$$x^k - x^{k-1} = (x^k - x^*) - (x^{k-1} - x^*) = e^k - e^{k-1}$$

Subtracting x^* from both sides of Equation (3.49) and recalling that $f(x^*) = 0$ yields

$$e^{k+1} = e^k - \frac{f'(x^*)e^k + \frac{1}{2}f''(x^*)(e^k)^2}{f'(x^*) + \frac{1}{2}f''(x^*)(e^k + e^{k-1})} \quad (3.57)$$

or

$$e^{k+1} = \frac{1}{2} \frac{f''(x^*)}{f'(x^*)} e^k e^{k-1} + O(e^3) \quad (3.58)$$

Let

$$e^k = C_k (e^k)^r$$

where r is the convergence order. If $r > 1$, then the convergence rate is super linear. If the remainder term in Equation (3.58) is negligible, then Equation (3.58) can be rewritten

$$\frac{e^{k+1}}{e^k e^{k-1}} = \frac{1}{2} \frac{f''(x^*)}{f'(x^*)} \quad (3.59)$$

and in the limit

$$\lim_{k \rightarrow \infty} \frac{e^{k+1}}{e^k e^{k-1}} = C' \quad (3.60)$$

For large k ,

$$e^k = C (e^{k-1})^r$$

and

$$e^{k+1} = C (e^k)^r = C \left(C (e^{k-1})^r \right)^r = C^{r+1} (e^{k-1})^{r^2}$$

Substituting this into Equation (3.60) yields

$$\lim_{k \rightarrow \infty} C^r (e^{k+1})^{r^2-r-1} = C' \quad (3.61)$$

Since $\lim_{k \rightarrow \infty} e^k = 0$, this relationship can only be satisfied if $r^2 - r - 1 = 0$, which has the solution

$$r = \frac{1 + \sqrt{5}}{2} > 1 \quad (3.62)$$

and hence super linear convergence.

Example 3.6

Use the secant method to find a solution of

$$0 = e^{x^2 - 2} - 3 \ln(x)$$

starting with $x^0 = 1.5$ and $x^1 = 1.4$.

Solution 3.6 Using Equation (3.49), the following set of values is obtained.

k	x^{k+1}	x^k	x^{k-1}	$f(x^k)$	$f(x^{k+1})$
1	1.4418	1.4000	1.5000	-0.0486	0.0676
2	1.4617	1.4418	1.4000	-0.0157	-0.0486
3	1.4552	1.4617	1.4418	0.0076	-0.0157
4	1.4557	1.4552	1.4617	-0.0006	0.0076
5	1.4557	1.4557	1.4552	-0.0000	-0.0006

■

3.3.2 Broyden's Method

The generalization of the secant method to a system of equations is Broyden's method. Broyden's method is an example of using a secant update. Returning again to the equation

$$x^{k+1} = x^k - [M_k(x^k)]^{-1} f(x^k) \quad (3.63)$$

where the iteration matrix M_k is updated with each subsequent x^{k+1} . The updated matrix is found

$$M_{k+1} = M_k + \frac{(y - M_k s) s^T}{s^T s} \quad (3.64)$$

where

$$y = f(x^{k+1}) - f(x^k) \quad (3.65)$$

and

$$s = x^{k+1} - x^k \quad (3.66)$$

Substituting y from Equation (3.65) into Equation (3.64) yields

$$M_{k+1} = M_k + \frac{f(x^{k+1}) s^T}{s^T s} \quad (3.67)$$

Example 3.7

Repeat Example 3.5 using Broyden's method.

Solution 3.7 As before, the initial condition is

$$x^{(0)} = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

This method also requires an initial guess for the iteration matrix M_0 . One possible initial iteration matrix is the Jacobian matrix evaluated at x^0 .

Initialization ($k = 0$)

$$\begin{aligned} f^0 &= \begin{bmatrix} -4 \\ 12 \end{bmatrix} \\ M_0 &= \begin{bmatrix} 1 & 6 \\ 6 & -9 \end{bmatrix} \text{ (both the same as in Example 3.5)} \end{aligned}$$

Iteration 1: ($k = 1$)

$$\begin{aligned} x^1 &= x^0 - M_0^{-1}f^0 = \begin{bmatrix} 3.8000 \\ 2.2000 \end{bmatrix} \text{ (same as in Example 3.5)} \\ f^1 &= \begin{bmatrix} 1.2800 \\ 0.0000 \end{bmatrix} \\ M_1 &= M_0 + \frac{f^1 (x^1 - x^0)^T}{(x^1 - x^0)^T (x^1 - x^0)} \\ &= \begin{bmatrix} 1.8000 & 5.2000 \\ 6.0000 & -9.0000 \end{bmatrix} \\ \text{error} &= \|f^1\|_\infty = 1.2800 \end{aligned}$$

Iteration 2: ($k = 2$)

$$\begin{aligned} x^2 &= x^1 - M_1^{-1}f^1 = \begin{bmatrix} 3.5570 \\ 2.0380 \end{bmatrix} \text{ (same as in Example 3.5)} \\ f^2 &= \begin{bmatrix} 0.0205 \\ -0.6153 \end{bmatrix} \\ M_2 &= M_1 + \frac{f^2 (x^2 - x^1)^T}{(x^2 - x^1)^T (x^2 - x^1)} \\ &= \begin{bmatrix} 1.7416 & 5.1611 \\ 7.7527 & -7.8315 \end{bmatrix} \\ \text{error} &= \|f^2\|_\infty = 0.6153 \end{aligned}$$

This process continues until the error is driven sufficiently small. In the

solution to Example 3.7, the error at each iteration is

iteration	error
0	12.0000
1	1.2800
2	0.6153
3	0.0506
4	0.0081
5	0.0001
6	0.0000

Note that it takes more iterations for this method to converge than the full Newton–Raphson. However, in most cases, the calculation of the iteration matrix M_k may require far less computation than the full Jacobian. The effort required is a trade-off between the computational effort required in calculating the Jacobian at every iteration and the number of iterations required for convergence. ■

3.3.3 Modifications to the Newton–Raphson Method

Another common simplification is to substitute each of the partial derivative entries $\frac{\partial f_i}{\partial x_j}$ by a difference approximation. For example, a simple approximation might be

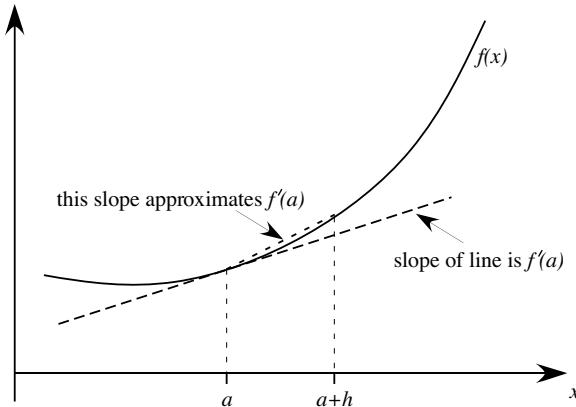
$$\frac{\partial f_i}{\partial x_j} \approx \frac{1}{h_{ij}} [f_i(x + h_{ij}e^j) - f_i(x)] \quad (3.68)$$

where e^j is the j th unit vector:

$$e^j = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

where the 1 occurs in the j th row of the unit vector and all other entries are zero. The scalar h_{ij} can be chosen in numerous ways, but one common choice is to let $h_{ij}^k = x_j^k - x_j^{k-1}$. This choice for h_{ij} leads to a rate of convergence of 1.62, which lies between quadratic and linear convergence rates.

Another common modification to the Newton–Raphson method is to set M equal to the Jacobian matrix at occasional intervals. For example, the matrix M can be reevaluated whenever the convergence slows down, or at more regular intervals, such as every other or every third iteration. This modification is known as the *dishonest Newton* method. An extreme extension

**FIGURE 3.9**

Graphical interpretation of the difference approximation of the slope of $f(a)$

of this method is to set M equal to the initial Jacobian matrix and then to hold it constant throughout the remainder of the iteration. This is commonly called the *very dishonest Newton* method. In addition to the reduction in computation associated with the calculation of the matrix, this method also has the advantage that the M matrix need only be factored into the LU matrices once since it is a constant. This can save considerable computation time in the LU factorization process. Similarly, the matrices of the dishonest Newton method need only be factored when the M matrix is reevaluated.

3.3.4 Numerical Differentiation

Using the Newton–Raphson method or any of its modifications requires the calculation of numerous partial derivatives. In many cases, the analytic computation of the partial derivative may be extremely complex or may be computationally expensive to compute. In these cases, it is desirable to compute the derivative numerically directly from the function $f(x)$ without explicit knowledge of $\frac{\partial f}{\partial x}$.

Consider the scalar function $f(x)$. The derivative of the function f' at the point $x = a$ is equivalent to the slope of the function $f(a)$. A reasonable approximation to the slope of a curve $f(a)$ is to use a nearby point $a + h$ to compute a *difference approximation*, as shown in Figure 3.9.

This has a mathematical basis that can be derived by application of the Taylor series expansion to $f(a + h)$:

$$f(a + h) = f(a) + h \frac{\partial f}{\partial x}(a) + \frac{h^2}{2!} \frac{\partial^2 f}{\partial x^2}(a) + \frac{h^3}{3!} \frac{\partial^3 f}{\partial x^3}(a) + \dots \quad (3.69)$$

By rearranging

$$\frac{f(a+h) - f(a)}{h} = \frac{\partial f}{\partial x}(a) + \frac{h}{2!} \frac{\partial^2 f}{\partial x^2}(a) + \dots \quad (3.70)$$

By neglecting the higher-order terms

$$\frac{\partial f}{\partial x}(a) \approx \frac{f(a+h) - f(a)}{h} \quad (3.71)$$

This approximation becomes increasingly more accurate as h becomes smaller (and is exact in the limit as $h \rightarrow 0$). This approximation is the one-sided difference approximation known as the *forward difference* approximation to the derivative of f . A similar approach can be taken in which the series is expanded about $a-h$ and

$$\frac{\partial f}{\partial x}(a) \approx \frac{f(a) - f(a-h)}{h} \quad (3.72)$$

which is the approximation known as the *backward difference* approximation. Consider now the combination of the two approaches:

$$\frac{\partial f}{\partial x}(a) \approx \frac{f(a+h) - f(a-h)}{2h} \quad (3.73)$$

This combination is often referred to as the *center difference* approximation, and is shown in Figure 3.10. The forward and backward difference approximations both have error on the order of $O(h)$, whereas the center approximation has an error on the order of $O(h^2)$ and will in general have better accuracy than either the forward or backward difference approximations.

Example 3.8

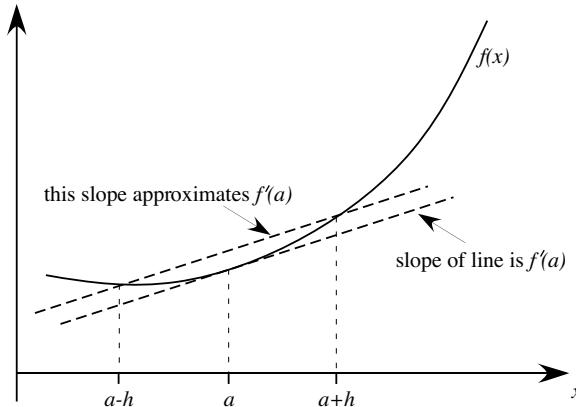
Consider the polynomial

$$f(x) = x^3 + x^2 - \frac{5}{4}x - \frac{3}{4}$$

Approximate the derivative of this polynomial in the range $[-2, 1.5]$ with $h = 0.2$ using the forward, backward, and center difference approximations.

Solution 3.8 The exact derivative of this function is given by

$$f'(x) = 3x^2 + 2x - \frac{5}{4}$$

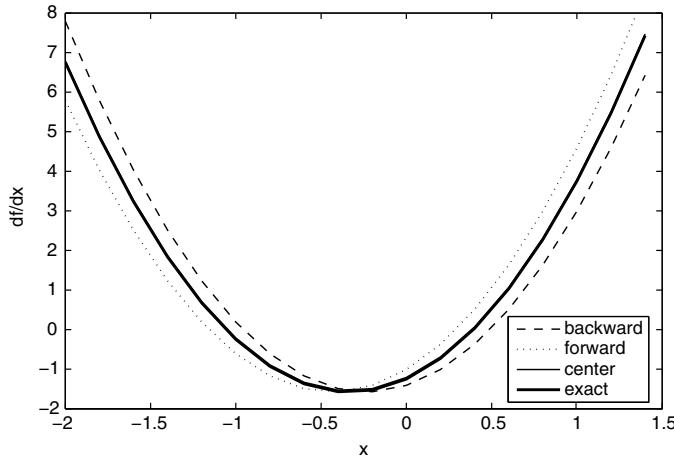
**FIGURE 3.10**

Graphical interpretation of the center difference approximation of the slope of $f(a)$

x	$f(x-h)$	$f(x)$	$f(x+h)$	$f'(x)$ backward	$f'(x)$ forward	$f'(x)$ center	$f'(x)$ exact
-2.0	-3.808	-2.250	-1.092	7.79	5.79	6.79	6.75
-1.8	-2.250	-1.092	-0.286	5.79	4.03	4.91	4.87
-1.6	-1.092	-0.286	0.216	4.03	2.51	3.27	3.23
-1.4	-0.286	0.216	0.462	2.51	1.23	1.87	1.83
-1.2	0.216	0.462	0.500	1.23	0.19	0.71	0.67
-1.0	0.462	0.500	0.378	0.19	-0.61	-0.21	-0.25
-0.8	0.500	0.378	0.144	-0.61	-1.17	-0.89	-0.93
-0.6	0.378	0.144	-0.154	-1.17	-1.49	-1.33	-1.37
-0.4	0.144	-0.154	-0.468	-1.49	-1.57	-1.53	-1.57
-0.2	-0.154	-0.468	-0.750	-1.57	-1.41	-1.49	-1.53
-0.0	-0.468	-0.750	-0.952	-1.41	-1.01	-1.21	-1.25
0.2	-0.750	-0.952	-1.026	-1.01	-0.37	-0.69	-0.73
0.4	-0.952	-1.026	-0.924	-0.37	0.51	0.07	0.03
0.6	-1.026	-0.924	-0.598	0.51	1.63	1.07	1.03
0.8	-0.924	-0.598	-0.000	1.63	2.99	2.31	2.27
1.0	-0.598	-0.000	0.918	2.99	4.59	3.79	3.75
1.2	-0.000	0.918	2.204	4.59	6.43	5.51	5.47
1.4	0.918	2.204	3.906	6.43	8.51	7.47	7.43

Figure 3.11 clearly shows the accuracy levels of the different derivative approximations. ■

By continuing in the same approach of using Taylor series expansions and including additional information, increasingly more accurate approximations can be achieved. One such approximation that is widely used is the *Richardson*

**FIGURE 3.11**

Exact versus approximate derivatives

approximation:

$$f'(x) \approx \frac{f(x-2h) - 8f(x-h) + 8f(x+h) - f(x+2h)}{12h} \quad (3.74)$$

This approximation has an error on the order $O(h^4)$.

Consider once again the Newton–Raphson method, which requires the calculation of the Jacobian. The approximations can be utilized to calculate the derivatives in the Jacobian rather than a direct analytic calculation. For example, consider the system of nonlinear equations:

$$\begin{aligned} f_1(x_1, x_2, \dots, x_n) &= 0 \\ f_2(x_1, x_2, \dots, x_n) &= 0 \\ &\vdots \\ f_n(x_1, x_2, \dots, x_n) &= 0 \end{aligned}$$

The Jacobian for this system consists of partial derivatives of the form $\frac{\partial f_i}{\partial x_j}$, which can now be approximated using one of the approximation methods introduced. For example, using the center difference

$$\frac{\partial f_i}{\partial x_j} = \frac{f_i(x_j + \Delta x_j) - f_i(x_j - \Delta x_j)}{2\Delta x_j}$$

where Δx_j is usually chosen to be a small incremental change (of about 1%).

3.3.5 Newton–GMRES

The discussions of the previous sections all assumed that the series of linear equations obtained from the Newton–Raphson series expansion were solved using the LU factorization direct method. There is no underlying principle that requires that direct methods be used. The linear equations could just as easily be solved with one of the iterative methods described in Chapter 2. Note that using an iterative method to solve linear system equations leads to a set of nested iteration loops: an outer loop for the nonlinear equations and an inner loop for the linear equations.

Consider the solution of the linear equations. At each iteration of the outer loop $x^{k+1} = x^k - M_k^{-1}f(x^k)$, the linear iterates will be driven to convergence to solve for x^k , even though, for low values of k , these updates can be significantly inaccurate. Therefore, it is plausible to reduce the number of linear system iterations in the inner loop without having an effect on the overall convergence of the outer iteration loop. These mixed linear–nonlinear iterative methods are named to reflect the linear method they use, such as Newton–SOR, Newton–CG, and Newton–GMRES. In this section, the approach to the Newton–GMRES method will be discussed, although the basic procedure is applicable to the Newton–SOR and Newton–CG methods as well.

If the linear iterative method is one of the Krylov subspace methods, then each inner loop iteration requires the evaluation of the action of the Jacobian $\frac{\partial f}{\partial x} = J(x^k)$ on a vector. In many implementations, the action of $\frac{\partial f}{\partial x}$ on a vector w can be approximated by a forward difference $D_h f(x : w)$ for some h . Note that this is entirely different from forming the Jacobian matrix using a forward difference approximation and then multiplying that matrix by w . This approach of directly finding the vector product is often referred to as a Jacobian-free Newton–GMRES method, since the Jacobian matrix is not explicitly calculated.

Jacobian-Free Newton–GMRES Algorithm for Solving $F(x) = 0$

1. Outer Loop (Newton):

Initialize outer loop:

$$\begin{aligned} x &= x^0 \\ F^0 &= F(x^0) \\ q &= 0 \end{aligned}$$

While $\|F(x^k)\| \geq \text{tol}$ and $q \leq q_{\max}$,

2. Inner Loop (GMRES):

$$\begin{aligned} k &= 1 \\ r_0 &= -F(x^q) \end{aligned}$$

$$\begin{aligned} v_1 &= r_0 / \|r_0\| \\ y &= [0 \ 0 \ 0 \ \dots 0]^T \\ \text{error} &= \|r_0\| \end{aligned}$$

While $\text{error} \geq \varepsilon$

- (a) Calculate $D_h f(x : v_k)$
- (b) $H(j, k) = (D_h f(x : v_k))^T v_j, \ j = 1, \dots, k$
- (c) $v_{k+1} = D_h f(x : v_k) - \sum_{j=1}^k H(j, k)v_j$
- (d) $H(k+1, k) = \|v_{k+1}\|$
- (e) $v_{k+1} = v_{k+1} / \|v_{k+1}\|$
- (f) Givens rotation:

i.

$$\begin{bmatrix} H(j, k) \\ H(j+1, k) \end{bmatrix} = \begin{bmatrix} \text{cs}(j) & \text{sn}(j) \\ -\text{sn}(j) & \text{cs}(j) \end{bmatrix} \begin{bmatrix} H(j, k) \\ H(j+1, k) \end{bmatrix} \quad j = 1, \dots, k-1$$

ii.

$$\begin{aligned} \text{cs}(k) &= \frac{H(k, k)}{\sqrt{H(k+1, k)^2 + H(k, k)^2}} \\ \text{sn}(k) &= \frac{H(k+1, k)}{\sqrt{H(k+1, k)^2 + H(k, k)^2}} \end{aligned}$$

iii. Approximate residual norm

$$\begin{aligned} \alpha &= \text{cs}(k)s(k) \\ s(k+1) &= -\text{sn}(k)s(k) \\ s(k) &= \alpha \\ \text{error} &= |s(k+1)| \end{aligned}$$

iv. Set

$$\begin{aligned} H(k, k) &= \text{cs}(k)H(k, k) + \text{sn}(k)H(k+1, k) \\ H(k+1, k) &= 0 \end{aligned}$$

3. (outer loop)

Solve $Hz = s$ for z

$$x^q = x^q + v_k z$$

$$q = q + 1$$

Example 3.9

Repeat Example 3.5 using the Jacobian-Free Newton–GMRES method using $h = 0.001$ for the directional derivative calculation.

Solution 3.9

1. Initialize the outer loop:

$$\begin{aligned} q &= 1 \\ x^{(0)} &= \begin{bmatrix} 3 \\ 3 \end{bmatrix} \end{aligned}$$

and

$$F(x^0) = \begin{bmatrix} 4 \\ -12 \end{bmatrix}$$

2. Starting the inner loop (GMRES iteration):

$$\begin{aligned} k &= 1 \\ r_0 &= -F(x^0) = \begin{bmatrix} -4 \\ 12 \end{bmatrix} \\ v_1 &= r_0 / \|r_0\| = \begin{bmatrix} -0.3162 \\ 0.9487 \end{bmatrix} \end{aligned}$$

Note that the norm used in this example is the 2-norm.

Inner loop $k = 1$:

- (a) Calculate $D_h f(x : v_1)$

$$\begin{aligned} D_h f(x : v_1) &= \frac{F(x^0 + hv_1) - F(x^0)}{h} \\ x^0 + hv_1 &= \begin{bmatrix} 2.9997 \\ 3.0009 \end{bmatrix} \\ F(x^0 + hv_1) &= \begin{bmatrix} 4.0054 \\ -12.0104 \end{bmatrix} \\ D_h f(x : v_1) &= \begin{bmatrix} 5.3769 \\ -10.4363 \end{bmatrix} \end{aligned}$$

$$(b) H(1, 1) = (D_h f(x : v_1))^T v_1 = -11.6011$$

$$(c) v_2 = D_h f(x : v_1) = \begin{bmatrix} 5.3769 \\ -10.4363 \end{bmatrix}$$

$$(d) H(2, 1) = \|v_2\| = 1.8007$$

$$(e) v_2 = v_2 / \|v_2\| = \begin{bmatrix} 0.9487 \\ 0.3162 \end{bmatrix}$$

(f) Givens rotation:

$$H = 11.7400$$

error $= |s(2)| = 1.9401 > \text{tol}$, therefore $k = 2$, repeat

Inner loop $k = 2$:

(a) Calculate $D_h f(x : v_2)$

$$D_h f(x : v_2) = \frac{F(x^0 + hv_2) - F(x^0)}{h}$$

$$x^0 + hv_2 = \begin{bmatrix} 3.0009 \\ 3.0003 \end{bmatrix}$$

$$F(x^0 + hv_2) = \begin{bmatrix} 4.0028 \\ -11.9972 \end{bmatrix}$$

$$D_h f(x : v_2) = \begin{bmatrix} 2.8470 \\ 2.8468 \end{bmatrix}$$

(b)

$$H = \begin{bmatrix} 11.7400 & 1.8004 \\ 0 & 3.6012 \end{bmatrix}$$

(c)

$$v_3 = D_h f(x : v_2) - \sum_{j=1}^2 H(j, 2)v_j = \begin{bmatrix} 0.1104 \\ 0.9939 \end{bmatrix}$$

(d) $H(3, 2) = 0$

$$(e) v_3 = v_3 / \|v_3\| = \begin{bmatrix} 0.1104 \\ 0.9939 \end{bmatrix}$$

(f) Givens rotation:

$$H = \begin{bmatrix} 11.7400 & -1.2268 \\ 0 & 3.8347 \end{bmatrix}$$

error $= |s(3)| = 2.0346 \times 10^{-15} < \text{tol}$, therefore proceed

3. Solve $Hz = s$ for z

$$\begin{bmatrix} 11.7400 & -1.2268 \\ 0 & 3.8347 \end{bmatrix} z = \begin{bmatrix} -12.4994 \\ 1.9401 \end{bmatrix}$$

$$z = \begin{bmatrix} -1.0118 \\ 0.5059 \end{bmatrix}$$

$$x^1 = x^0 + v_3 z$$

$$\begin{aligned}x^1 &= \begin{bmatrix} 3 \\ 3 \end{bmatrix} + \begin{bmatrix} -0.3163 & 0.9487 \\ 0.9487 & 0.3162 \end{bmatrix} \begin{bmatrix} -1.0118 \\ 0.5059 \end{bmatrix} = \begin{bmatrix} 3.7999 \\ 2.2001 \end{bmatrix} \\F(x^1) &= \begin{bmatrix} 1.2803 \\ -0.0012 \end{bmatrix} \\q &= 2\end{aligned}$$

Since $\|F(x^1)\| > \text{tol}$, the outer loop continues and the inner loop is started again with $k = 1$. The process repeats until the convergence tolerance criterion is met, and x converges to

$$x = \begin{bmatrix} 3.6100 \\ 2.0045 \end{bmatrix}$$

which is the same solution as before. ■

In the solution to Example 3.9, the error at each iteration is

iteration	error
0	12.0000
1	1.2803
2	0.0661
3	0.0002

which are nearly identical to the errors obtained in the true Newton–Raphson solution in Example 3.5. The potential savings from using a Jacobian-free Newton–GMRES method arise from not having to calculate the system Jacobian at every step. Furthermore, it is possible that the GMRES iteration may not require a full number of iterations ($k < n$) as x^q approaches convergence.

3.4 Continuation Methods

Many of the iteration methods described so far will, in general, converge to a solution x^* of $f(x) = 0$ only if the initial condition is sufficiently close to x^* . The continuation method approach may be considered to be an attempt to widen the region of convergence of a given method. In many physical systems, the problem defined by the mathematical equation $f(x) = 0$ may in some way depend in a natural way on a parameter λ of the system. When this parameter is set equal to 0, the system $f_0(x) = 0$ has a known solution x^0 . However, for varying λ , an entire family of functions $H(x, \lambda)$ exists such that

$$H(x, 0) = f_0(x), \quad H(x, 1) = f(x) \tag{3.75}$$

where a solution x^0 of $H(x, 0) = 0$ is known, and the equation $H(x, 1) = 0$ is the desired problem to be solved.

Even if $f(x)$ does not depend naturally on a suitable parameter λ , a family of problems satisfying Equation (3.75) can be defined by

$$H(x, \lambda) = \lambda f(x) + (1 - \lambda) f_0(x), \quad \lambda \in [0, 1] \quad (3.76)$$

when the solution x^0 of $f_0(x) = 0$ is known. As λ varies from 0 to 1, the family of mappings varies from $f_0(x) = 0$ to $f_1(x) = 0$ where the solution of $f_1(x) = f(x) = 0$ is the desired value $x^1 = x^*$.

As a first approach to obtaining $x^* = x^1$, the interval $[0, 1]$ can be partitioned as

$$0 = \lambda_0 < \lambda_1 < \lambda_2 < \dots < \lambda_N = 1$$

and consider solving the problems

$$H(x, \lambda_i) = 0, \quad i = 1, \dots, N \quad (3.77)$$

Assuming that a Newton–Raphson iteration is used to solve each problem i in Equation (3.77), then the initial condition for the i th problem is the solution from $H(x, \lambda_{i-1}) = 0$. For small enough intervals between i and $i + 1$, this solves the problem of identifying a good initial condition.

The relationship given in Equation (3.76) is an example of a *homotopy* in which two functions $f(x)$ and $f_0(x)$ are embedded in a single continuous function. Formally, a homotopy between any two functions is a continuous mapping $f, f_0 : X \rightarrow Y$:

$$H : [0, 1] \times X \rightarrow Y \quad (3.78)$$

such that Equation (3.75) holds. If such a mapping exists, then it is said that f is *homotopic* to f_0 .

Homotopy functions are used to define a path from the known solution to a relatively simple problem ($f_0(x) = 0$) to the solution of a more complex problem to which the solution is desired ($f(x) = 0$). This path comprises the solutions to a family of problems which represent the continuous deformation from the simple problem to the desired problem. Continuation methods are a numerical approach to following the deformation path.

Homotopy continuation methods can be constructed to be exhaustive and globally convergent, meaning that all solutions to a given system of nonlinear equations can be found and will converge regardless of choice of initial condition [63]. Since a homotopy problem is equal to zero at every point $\lambda \in [0, 1]$ along the path, it is therefore equal to zero at $\lambda = \lambda^k$ and $\lambda = \lambda^{k+1}$, which are two successive points along the path. This gives

$$0 = H(x^k, \lambda^k) = \lambda^k f(x^k) + (1 - \lambda^k) f_0(x^k) \quad (3.79)$$

$$0 = H(x^{k+1}, \lambda^{k+1}) = \lambda^{k+1} f(x^{k+1}) + (1 - \lambda^{k+1}) f_0(x^{k+1}) \quad (3.80)$$

For paths along the path, the homotopy parameter λ^{k+1} is associated with the parameter set

$$x^{k+1} = x^k + \Delta x$$

If the changes in the parameters are small, the functions $f_0(x^{k+1})$ and $f(x^{k+1})$ can be linearly approximated by using a Taylor series expansion about x^k and neglecting all terms higher than second order. Applying this technique to Equation (3.80) yields

$$(\lambda^k + \Delta\lambda) [F_x(x^k) \Delta x] + (1 - \lambda^k - \Delta\lambda) [f_0(x^k) + F_{0x}(x^k) \Delta x] = 0 \quad (3.81)$$

where F_x and F_{0x} are the Jacobians of $f(x)$ and $f_0(x)$ with respect to x , respectively. Subtracting Equation (3.79) from Equation (3.81) and canceling like terms yields

$$0 = [\lambda^{k+1} F_x(x^k) + (1 - \lambda^{k+1}) F_{0x}(x^k)] \Delta x + [f(x^k) - f_0(x^k)] \Delta\lambda \quad (3.82)$$

Using $x^{k+1} = x^k + \Delta x$, Equation (3.82) can be rewritten in terms of the homotopy function to yield the update equation

$$x^{k+1} = x^k - \Delta\lambda H_x(x^k, \lambda^{k+1})^{-1} \frac{\partial}{\partial \lambda} H(x^k, \lambda^{k+1}) \quad (3.83)$$

where

$$\lambda^{k+1} = \lambda^k + \Delta\lambda$$

and $H_x(x, \lambda)$ is the Jacobian of $H(x, \lambda)$ with respect to x .

Example 3.10

Solve

$$0 = f_1(x_1, x_2) = x_1^2 - 3x_2^2 + 3 \quad (3.84)$$

$$0 = f_2(x_1, x_2) = x_1 x_2 + 6 \quad (3.85)$$

using the homotopic mapping with

$$0 = f_{01}(x_1, x_2) = x_1^2 - 4 \quad (3.86)$$

$$0 = f_{02}(x_1, x_2) = x_2^2 - 9 \quad (3.87)$$

Solution 3.10 Set up a homotopy such that

$$H(x, \lambda) = \lambda f(x) + (1 - \lambda) f_0(x), \quad \lambda \in [0, 1] \quad (3.88)$$

$$0 = \lambda(x_1^2 - 3x_2^2 + 3) + (1 - \lambda)(x_1^2 - 4) \quad (3.89)$$

$$0 = \lambda(x_1 x_2 + 6) + (1 - \lambda)(x_2^2 - 9) \quad (3.90)$$

The continuation method advances the solution via Equation (3.83):

$$\lambda^{k+1} = \lambda^k + \Delta\lambda \quad (3.91)$$

$$\begin{bmatrix} x_1^{k+1} \\ x_2^{k+1} \end{bmatrix} = \begin{bmatrix} x_1^k \\ x_2^k \end{bmatrix} - \Delta\lambda \left[\lambda^{k+1} \begin{bmatrix} 2x_1^k & -6x_2^k \\ x_2^k & x_1^k \end{bmatrix} + (1 - \lambda^{k+1}) \begin{bmatrix} 2x_1^k & 0 \\ 0 & 2x_2^k \end{bmatrix} \right]^{-1} \times \\ \begin{bmatrix} (x_1^k)^2 - 3(x_2^k)^2 + 3 - ((x_1^k)^2 - 4) \\ x_1^k x_2^k + 6 - ((x_2^k)^2 - 9) \end{bmatrix} \quad (3.92)$$

The solution is then refined through the Newton–Raphson solution for x_1^{k+1} and x_2^{k+1} from

$$0 = \lambda^{k+1} ((x_1^{k+1})^2 - 3(x_2^{k+1})^2 + 3) + (1 - \lambda^{k+1}) ((x_1^{k+1})^2 - 4) \quad (3.93)$$

$$0 = \lambda^{k+1} (x_1^{k+1} x_2^{k+1} + 6) + (1 - \lambda^{k+1}) ((x_2^{k+1})^2 - 9) \quad (3.94)$$

Starting with $\lambda^0 = 0$ and $\Delta\lambda = 0.1$ yields the easily obtained initial solution of the system:

$$\begin{aligned} x_1^0 &= 2 \\ x_2^0 &= 3 \end{aligned}$$

Predicting the values for $k = 1$ from Equations (3.91) and (3.92) yields

$$\begin{aligned} x_1^1 &= 2.3941 \\ x_2^1 &= 2.7646 \end{aligned}$$

Refining the solution via the Newton–Raphson solution to Equations (3.93) and (3.94) yields

$$\begin{aligned} x_1^1 &= 2.3628 \\ x_2^1 &= 2.7585 \end{aligned}$$

This process is repeated until $\lambda = 1$ and $x_1 = -3$ and $x_2 = 2$, which are the correct solutions to the desired problem.

The same process will work if the initial solutions are chosen as $x_1^0 = -2$ and $x_2^0 = -3$. In this case, the values obtained are the alternate solution $x_1 = 3$ and $x_2 = -2$ to the desired problem. ■

3.5 Power System Applications

The solution and analysis procedures outlined in this chapter form the basis of a set of powerful tools that can be used for a myriad of power system

applications. One of the most outstanding features of power systems is that they are modeled as an extremely large set of nonlinear equations. The North American transmission grid is one of the largest nonlinear engineering systems. Most types of power system analysis require the solution in one form or another of this system of nonlinear equations. The applications described below are a handful of the more common applications, but are certainly not a complete coverage of all possible nonlinear problems that arise in power system analysis.

3.5.1 Power Flow

Many power system problems give rise to systems of nonlinear equations that must be solved. Probably the most common nonlinear power system problem is the *power flow* or *power flow* problem. The underlying principle of a power flow problem is that, given the system loads, generation, and network configuration, the system bus voltages and line flows can be found by solving the nonlinear power flow equations. This is typically accomplished by applying Kirchoff's law at each power system bus throughout the system. In this context, Kirchoff's law can be interpreted as *the sum of the powers entering a bus must be zero*, or that the power at each bus must be conserved. Since power is comprised of two components, active power and reactive power, each bus gives rise to two equations – one for active power and one for reactive power. These equations are known as the *power flow equations*:

$$0 = \Delta P_i = P_i^{inj} - V_i \sum_{j=1}^{N_{bus}} V_j Y_{ij} \cos(\theta_i - \theta_j - \phi_{ij}) \quad (3.95)$$

$$0 = \Delta Q_i = Q_i^{inj} - V_i \sum_{j=1}^{N_{bus}} V_j Y_{ij} \sin(\theta_i - \theta_j - \phi_{ij}) \quad (3.96)$$

$$i = 1, \dots, N_{bus}$$

where P_i^{inj} , Q_i^{inj} are the active and reactive power injected at the bus i , respectively. Loads are modeled by negative power injection. The values V_i and V_j are the voltage magnitudes at bus i and bus j , respectively. The values θ_i and θ_j are the corresponding phase angles. The value $Y_{ij}\angle\phi_{ij}$ is the (ij) th element of the network admittance matrix Y . The constant N_{bus} is the number of buses in the system. The updates ΔP_i^k and ΔQ_i^k of Equations (3.95) and (3.96) are called the *mismatch* equations because they give a measure of the power difference, or mismatch, between the calculated power values as functions of voltage and phase angle, and the actual injected powers. As the Newton–Raphson iteration continues, this mismatch is driven to zero until the power leaving a bus, calculated from the voltages and phase angles, equals the injected power. At this point, the converged values of voltages and phase angles are used to calculate line flows, slack bus powers, and the injected reactive powers at the generator buses.

The formulation in Equations (3.95) and (3.96) is called the *polar* formulation of the power flow equations. If $Y_{ij}\angle\phi_{ij}$ is instead given by the complex sum $g_{ij} + jb_{ij}$, then the power flow equations may be written in *rectangular form* as

$$0 = P_i^{inj} - V_i \sum_{j=1}^{N_{bus}} V_j (g_{ij} \cos(\theta_i - \theta_j) + b_{ij} \sin(\theta_i - \theta_j)) \quad (3.97)$$

$$0 = Q_i^{inj} - V_i \sum_{j=1}^{N_{bus}} V_j (g_{ij} \sin(\theta_i - \theta_j) - b_{ij} \cos(\theta_i - \theta_j)) \quad (3.98)$$

$$i = 1, \dots, N_{bus}$$

In either case, the power flow equations are a system of nonlinear equations. They are nonlinear in both the voltage and phase angle. There are, at most, $2N_{bus}$ equations to solve. This number is then further reduced by removing one power flow equation for each known voltage (at voltage controlled buses) and the slack bus angle. This reduction is necessary since the number of equations must equal the number of unknowns in a fully determined system. Once the nonlinear power flow equations have been determined, the Newton–Raphson method may be directly applied.

The most common approach to solving the power flow equations by the Newton–Raphson method is to arrange the equations by phase angle followed by the voltage magnitudes as

$$\begin{bmatrix} J_1 & J_2 \\ J_3 & J_4 \end{bmatrix} \begin{bmatrix} \Delta\theta_1 \\ \Delta\theta_2 \\ \Delta\theta_3 \\ \vdots \\ \Delta\theta_{N_{bus}} \\ \Delta V_1 \\ \Delta V_2 \\ \Delta V_3 \\ \vdots \\ \Delta V_{N_{bus}} \end{bmatrix} = - \begin{bmatrix} \Delta P_1 \\ \Delta P_2 \\ \Delta P_3 \\ \vdots \\ \Delta P_{N_{bus}} \\ \Delta Q_1 \\ \Delta Q_2 \\ \Delta Q_3 \\ \vdots \\ \Delta Q_{N_{bus}} \end{bmatrix} \quad (3.99)$$

where

$$\begin{aligned} \Delta\theta_i &= \theta_i^{k+1} - \theta_i^k \\ \Delta V_i &= V_i^{k+1} - V_i^k \end{aligned}$$

These equations are then solved using LU factorization with forward/backward substitution. The Jacobian is typically divided into four submatrices, where

$$\begin{bmatrix} J_1 & J_2 \\ J_3 & J_4 \end{bmatrix} = \begin{bmatrix} \frac{\partial \Delta P}{\partial \theta} & \frac{\partial \Delta P}{\partial V} \\ \frac{\partial \Delta Q}{\partial \theta} & \frac{\partial \Delta Q}{\partial V} \end{bmatrix} \quad (3.100)$$

Each submatrix represents the partial derivatives of each of the mismatch equations with respect to each of the unknowns. These partial derivatives yield eight types – two for each mismatch equation, where one is for the diagonal element and the other is for off-diagonal elements. The derivatives are summarized as

$$\frac{\partial \Delta P_i}{\partial \theta_i} = V_i \sum_{j=1}^{N_{bus}} V_j Y_{ij} \sin(\theta_i - \theta_j - \phi_{ij}) + V_i^2 Y_{ii} \sin \phi_{ii} \quad (3.101)$$

$$\frac{\partial \Delta P_i}{\partial \theta_j} = -V_i V_j Y_{ij} \sin(\theta_i - \theta_j - \phi_{ij}) \quad (3.102)$$

$$\frac{\partial \Delta P_i}{\partial V_i} = - \sum_{i=1}^{N_{bus}} V_j Y_{ij} \cos(\theta_i - \theta_j - \phi_{ij}) - V_i Y_{ii} \cos \phi_{ii} \quad (3.103)$$

$$\frac{\partial \Delta P_i}{\partial V_j} = -V_i Y_{ij} \cos(\theta_i - \theta_j - \phi_{ij}) \quad (3.104)$$

$$\frac{\partial \Delta Q_i}{\partial \theta_i} = -V_i \sum_{j=1}^{N_{bus}} V_j Y_{ij} \cos(\theta_i - \theta_j - \phi_{ij}) + V_i^2 Y_{ii} \cos \phi_{ii} \quad (3.105)$$

$$\frac{\partial \Delta Q_i}{\partial \theta_j} = V_i V_j Y_{ij} \cos(\theta_i - \theta_j - \phi_{ij}) \quad (3.106)$$

$$\frac{\partial \Delta Q_i}{\partial V_i} = - \sum_{j=1}^{N_{bus}} V_j Y_{ij} \sin(\theta_i - \theta_j - \phi_{ij}) + V_i Y_{ii} \sin \phi_{ii} \quad (3.107)$$

$$\frac{\partial \Delta Q_i}{\partial V_j} = -V_i Y_{ij} \sin(\theta_i - \theta_j - \phi_{ij}) \quad (3.108)$$

A common modification to the power flow solution is to replace the unknown update ΔV_i by the normalized value $\frac{\Delta V_i}{V_i}$. This formulation yields a more symmetric Jacobian, as the Jacobian submatrices J_2 and J_4 are now multiplied by V_i to compensate for the scaling of ΔV_i by V_i . All partial derivatives of each submatrix then become quadratic in voltage magnitude.

The Newton–Raphson method for the solution of the power flow equations is relatively straightforward to program, since both the function evaluations and the partial derivatives use the same expressions. Thus it takes little extra computational effort to compute the Jacobian once the mismatch equations have been calculated.

Example 3.11

Find the voltage magnitudes, phase angles, and line flows for the small power system shown in Figure 3.12 with the following per unit system parameters:

bus	type	V	P_{gen}	Q_{gen}	P_{load}	Q_{load}
1	slack	1.02	—	—	0.0	0.0
2	PV	1.00	0.5	—	0.0	0.0
3	PQ	—	0.0	0.0	1.2	0.5

i	j	R_{ij}	X_{ij}	B_{ij}
1	2	0.02	0.3	0.15
1	3	0.01	0.1	0.1
2	3	0.01	0.1	0.1

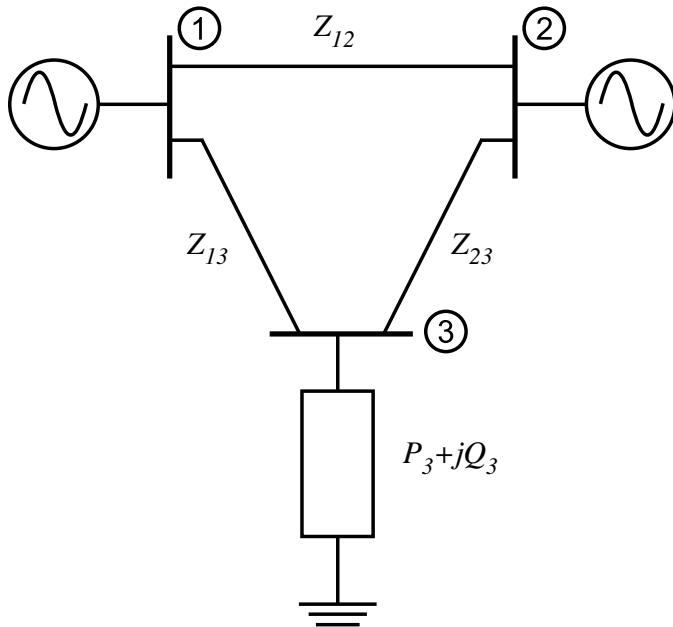


FIGURE 3.12
Example power system

Solution 3.11 The first step in any power flow solution is to calculate the admittance matrix Y for the power system. A simple procedure for calculating the elements of the admittance matrix is

$Y(i, j)$	negative of the admittance between buses i and j
$Y(i, i)$	sum of all admittances connected to bus i

Calculating the admittance matrix for this system yields

$$Y = \begin{bmatrix} 13.1505\angle -84.7148^\circ & 3.3260\angle 93.8141^\circ & 9.9504\angle 95.7106^\circ \\ 3.3260\angle 93.8141^\circ & 13.1505\angle -84.7148^\circ & 9.9504\angle 95.7106^\circ \\ 9.9504\angle 95.7106^\circ & 9.9504\angle 95.7106^\circ & 19.8012\angle -84.2606^\circ \end{bmatrix} \quad (3.109)$$

By inspection, this system has three unknowns: θ_2 , θ_3 , and V_3 ; thus three power flow equations are required. These power flow equations are

$$0 = \Delta P_2 = 0.5 - V_2 \sum_{j=1}^3 V_j Y_{ij} \cos(\theta_2 - \theta_j - \theta_{ij}) \quad (3.110)$$

$$0 = \Delta P_3 = -1.2 - V_3 \sum_{j=1}^3 V_j Y_{ij} \cos(\theta_3 - \theta_j - \theta_{ij}) \quad (3.111)$$

$$0 = \Delta Q_3 = -0.5 - V_3 \sum_{j=1}^3 V_j Y_{ij} \sin(\theta_3 - \theta_j - \theta_{ij}) \quad (3.112)$$

Substituting in the known quantities for $V_1 = 1.02$, $V_2 = 1.00$, and $\theta_1 = 0$ and the admittance matrix quantities yields

$$\begin{aligned} \Delta P_2 &= 0.5 - (1.00)((1.02)(3.3260) \cos(\theta_2 - 0 - 93.8141^\circ) \\ &\quad + (1.00)(13.1505) \cos(\theta_2 - \theta_2 + 84.7148^\circ) \\ &\quad + (V_3)(9.9504) \cos(\theta_2 - \theta_3 - 95.7106^\circ)) \end{aligned} \quad (3.113)$$

$$\begin{aligned} \Delta P_3 &= -1.2 - (V_3)((1.02)(9.9504) \cos(\theta_3 - 0 - 95.7106^\circ) \\ &\quad + (1.00)(9.9504) \cos(\theta_3 - \theta_2 - 95.7106^\circ) \\ &\quad + (V_3)(19.8012) \cos(\theta_3 - \theta_3 + 84.2606^\circ)) \end{aligned} \quad (3.114)$$

$$\begin{aligned} \Delta Q_3 &= -0.5 - (V_3)((1.02)(9.9504) \sin(\theta_3 - 0 - 95.7106^\circ) \\ &\quad + (1.00)(9.9504) \sin(\theta_3 - \theta_2 - 95.7106^\circ) \\ &\quad + ((V_3)(19.8012) \sin(\theta_3 - \theta_3 + 84.2606^\circ))) \end{aligned} \quad (3.115)$$

The Newton–Raphson iteration for this system is then given by

$$\begin{bmatrix} \frac{\partial \Delta P_2}{\partial \theta_2} & \frac{\partial \Delta P_2}{\partial \theta_3} & \frac{\partial \Delta P_2}{\partial V_3} \\ \frac{\partial \Delta P_3}{\partial \theta_2} & \frac{\partial \Delta P_3}{\partial \theta_3} & \frac{\partial \Delta P_3}{\partial V_3} \\ \frac{\partial \Delta Q_3}{\partial \theta_2} & \frac{\partial \Delta Q_3}{\partial \theta_3} & \frac{\partial \Delta Q_3}{\partial V_3} \end{bmatrix} \begin{bmatrix} \Delta \theta_2 \\ \Delta \theta_3 \\ \Delta V_3 \end{bmatrix} = - \begin{bmatrix} \Delta P_2 \\ \Delta P_3 \\ \Delta Q_3 \end{bmatrix} \quad (3.116)$$

where

$$\begin{aligned} \frac{\partial \Delta P_2}{\partial \theta_2} &= 3.3925 \sin(\theta_2 - 93.8141^\circ) \\ &\quad + 9.9504 V_3 \sin(\theta_2 - \theta_3 - 95.7106^\circ) \end{aligned}$$

$$\begin{aligned}
\frac{\partial \Delta P_2}{\partial \theta_3} &= -9.9504 V_3 \sin(\theta_2 - \theta_3 - 95.7106^\circ) \\
\frac{\partial \Delta P_2}{\partial V_3} &= -9.9504 \cos(\theta_2 - \theta_3 - 95.7106^\circ) \\
\frac{\partial \Delta P_3}{\partial \theta_2} &= -9.9504 V_3 \sin(\theta_3 - \theta_2 - 95.7106^\circ) \\
\frac{\partial \Delta P_3}{\partial \theta_3} &= 10.1494 V_3 \sin(\theta_3 - 95.7106^\circ) \\
&\quad + 9.9504 V_3 \sin(\theta_3 - \theta_2 - 95.7106^\circ) \\
\frac{\partial \Delta P_3}{\partial V_3} &= -10.1494 \cos(\theta_3 - 95.7106^\circ) \\
&\quad - 9.9504 \cos(\theta_3 - \theta_2 - 95.7106^\circ) \\
&\quad - 39.6024 V_3 \cos(84.2606^\circ) \\
\frac{\partial \Delta Q_3}{\partial \theta_2} &= 9.9504 V_3 \cos(\theta_3 - \theta_2 - 95.7106^\circ) \\
\frac{\partial \Delta Q_3}{\partial \theta_3} &= -10.1494 V_3 \cos(\theta_3 - 95.7106^\circ) \\
&\quad - 9.9504 V_3 \cos(\theta_3 - \theta_2 - 95.7106^\circ) \\
\frac{\partial \Delta Q_3}{\partial V_3} &= -10.1494 \sin(\theta_3 - 95.7106^\circ) \\
&\quad - 9.9504 \sin(\theta_3 - \theta_2 - 95.7106^\circ) \\
&\quad - 39.6024 V_3 \sin(84.2606^\circ)
\end{aligned}$$

Recall that one of the underlying assumptions of the Newton–Raphson iteration is that the higher-order terms of the Taylor series expansion are negligible only if the initial guess is sufficiently close to the actual solution to the nonlinear equations. Under most operating conditions, the voltages throughout the power system are within $\pm 10\%$ of the nominal voltage and therefore fall in the range $0.9 \leq V_i \leq 1.1$ per unit. Similarly, under most operating conditions, the phase angle differences between adjacent buses are typically small. Thus, if the slack bus angle is taken to be zero, then all phase angles throughout the system will also be close to zero. Therefore, in initializing a power flow, it is common to choose a “flat start” initial condition. That is, all voltage magnitudes are set to 1.0 per unit and all angles are set to zero.

Iteration 1

Evaluating the Jacobian and the mismatch equations at the flat start initial conditions yields

$$[J^0] = \begin{bmatrix} -13.2859 & 9.9010 & 0.9901 \\ 9.9010 & -20.0000 & -1.9604 \\ -0.9901 & 2.0000 & -19.4040 \end{bmatrix}$$

$$\begin{bmatrix} \Delta P_2^0 \\ \Delta P_3^0 \\ \Delta Q_3^0 \end{bmatrix} = \begin{bmatrix} 0.5044 \\ -1.1802 \\ -0.2020 \end{bmatrix}$$

Solving

$$[J^0] \begin{bmatrix} \Delta \theta_2^1 \\ \Delta \theta_3^1 \\ \Delta V_3^1 \end{bmatrix} = - \begin{bmatrix} \Delta P_2^0 \\ \Delta P_3^0 \\ \Delta Q_3^0 \end{bmatrix}$$

by LU factorization yields

$$\begin{bmatrix} \Delta \theta_2^1 \\ \Delta \theta_3^1 \\ \Delta V_3^1 \end{bmatrix} = \begin{bmatrix} -0.0096 \\ -0.0621 \\ -0.0163 \end{bmatrix}$$

Therefore,

$$\begin{aligned} \theta_2^1 &= \theta_2^0 + \Delta \theta_2^1 = 0 - 0.0096 = -0.0096 \\ \theta_3^1 &= \theta_3^0 + \Delta \theta_3^1 = 0 - 0.0621 = -0.0621 \\ V_3^1 &= V_3^0 + \Delta V_3^1 = 1 - 0.0163 = 0.9837 \end{aligned}$$

Note that the angles are given in *radians* and not degrees. The error at the first iteration is the largest absolute value of the mismatch equations, which is

$$\varepsilon^1 = 1.1802$$

One quick check of this process is to note that the voltage update V_3^1 is slightly less than 1.0 per unit, which would be expected given the system configuration. Note also that the diagonals of the Jacobian are all equal or greater in magnitude than the off-diagonal elements. This is because the diagonals are summations of terms, whereas the off-diagonal elements are single terms.

Iteration 2

Evaluating the Jacobian and the mismatch equations at the updated values θ_2^1 , θ_3^1 , and V_3^1 yields

$$\begin{aligned} [J^1] &= \begin{bmatrix} -13.1597 & 9.7771 & 0.4684 \\ 9.6747 & -19.5280 & -0.7515 \\ -1.4845 & 3.0929 & -18.9086 \end{bmatrix} \\ \begin{bmatrix} \Delta P_2^1 \\ \Delta P_3^1 \\ \Delta Q_3^1 \end{bmatrix} &= \begin{bmatrix} 0.0074 \\ -0.0232 \\ -0.0359 \end{bmatrix} \end{aligned}$$

Solving for the update yields

$$\begin{bmatrix} \Delta \theta_2^2 \\ \Delta \theta_3^2 \\ \Delta V_3^2 \end{bmatrix} = \begin{bmatrix} -0.0005 \\ -0.0014 \\ -0.0021 \end{bmatrix}$$

and

$$\begin{bmatrix} \theta_2^2 \\ \theta_3^2 \\ V_3^2 \end{bmatrix} = \begin{bmatrix} -0.0101 \\ -0.0635 \\ 0.9816 \end{bmatrix}$$

where

$$\varepsilon^2 = 0.0359$$

Iteration 3

Evaluating the Jacobian and the mismatch equations at the updated values θ_2^2 , θ_3^2 , and V_3^2 yields

$$\begin{aligned} [J^2] &= \begin{bmatrix} -13.1392 & 9.7567 & 0.4600 \\ 9.6530 & -19.4831 & -0.7213 \\ -1.4894 & 3.1079 & -18.8300 \end{bmatrix} \\ \begin{bmatrix} \Delta P_2^0 \\ \Delta P_3^0 \\ \Delta Q_3^0 \end{bmatrix} &= \begin{bmatrix} 0.1717 \\ -0.5639 \\ -0.9084 \end{bmatrix} \times 10^{-4} \end{aligned}$$

Solving for the update yields

$$\begin{bmatrix} \Delta \theta_2^2 \\ \Delta \theta_3^2 \\ \Delta V_3^2 \end{bmatrix} = \begin{bmatrix} -0.1396 \\ -0.3390 \\ -0.5273 \end{bmatrix} \times 10^{-5}$$

and

$$\begin{bmatrix} \theta_2^3 \\ \theta_3^3 \\ V_3^3 \end{bmatrix} = \begin{bmatrix} -0.0101 \\ -0.0635 \\ 0.9816 \end{bmatrix}$$

where

$$\varepsilon^3 = 0.9084 \times 10^{-4}$$

At this point, the iterations have converged, since the mismatch is sufficiently small and the values are no longer changing significantly.

The last task in power flow is to calculate the generated reactive powers, the slack bus active power output, and the line flows. The generated powers can be calculated directly from the power flow equations

$$\begin{aligned} P_i^{inj} &= V_i \sum_{j=1}^{N_{bus}} V_j Y_{ij} \cos(\theta_i - \theta_j - \phi_{ij}) \\ Q_i^{inj} &= V_i \sum_{j=1}^{N_{bus}} V_j Y_{ij} \sin(\theta_i - \theta_j - \phi_{ij}) \end{aligned}$$

Therefore,

$$P_{gen,1} = P_1^{inj} = 0.7087$$

$$Q_{gen,1} = Q_1^{inj} = 0.2806$$

$$Q_{gen,2} = Q_2^{inj} = -0.0446$$

The active power losses in the system are the difference between the sum of the generation and the sum of the loads; in this case

$$P_{loss} = \sum P_{gen} - \sum P_{load} = 0.7087 + 0.5 - 1.2 = 0.0087 \text{ pu} \quad (3.117)$$

The line losses for line $i-j$ are calculated at both the sending and receiving ends of the line. Therefore, the power sent from bus i to bus j is

$$S_{ij} = V_i \angle \theta_i I_{ij}^* \quad (3.118)$$

and the power received at bus j from bus i is

$$S_{ji} = V_j \angle \theta_j I_{ji}^* \quad (3.119)$$

Thus

$$P_{ij} = V_i V_j Y_{ij} \cos(\theta_i - \theta_j - \phi_{ij}) - V_i^2 Y_{ij} \cos(\phi_{ij}) \quad (3.120)$$

$$Q_{ij} = V_i V_j Y_{ij} \sin(\theta_i - \theta_j - \phi_{ij}) + V_i^2 Y_{ij} \sin(\phi_{ij}) - V_i^2 B_{ij}/2 \quad (3.121)$$

Similarly, the powers P_{ji} and Q_{ji} can be calculated. The active power loss on any given line is the difference between the active power sent from bus i and the active power received at bus j . Calculating the reactive power losses is more complex, since the reactive power generated by the line-charging (shunt capacitances) must also be included. ■

3.5.2 Regulating Transformers

One of the most common controllers found in the power system network is the *regulating transformer*. This is a transformer that is able to change the winding ratios (tap settings) in response to changes in load-side voltage. If the voltage on the secondary side (or load side) is lower than a desired voltage (such as during heavy loading), the tap will change so as to increase the secondary voltage while maintaining the primary side voltage. A regulating transformer is also frequently referred to as an *under-load-tap-changing* or ULTC transformer. The tap setting t may be real or complex, and per unit, the tap ratio is defined as $1 : t$ where t is typically within 10% of 1.0. A phase-shifting transformer is achieved by allowing the tap t to be complex with both magnitude and angle.

The effect of the regulating transformer is incorporated into the power flow algorithm through the admittance matrix. To incorporate a regulating transformer into the admittance matrix, consider the regulating transformer as a two-port network relating the input currents I_i and I_j to the input voltages V_i and V_j , as shown in Figure 3.13. The receiving end current is given by

$$I_j = (V_j - tV_i) Y \quad (3.122)$$

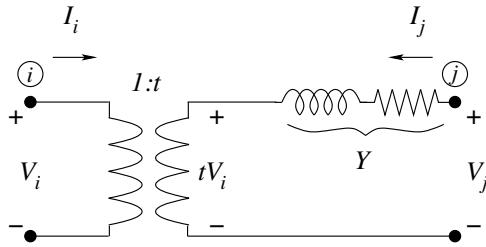


FIGURE 3.13
A regulating transformer

Note that the currents can be found from the power transfer equation

$$S_i = V_i I_i^* = -t V_i I_j^* \quad (3.123)$$

Therefore,

$$I_i = -t^* I_j \quad (3.124)$$

$$= -t^* (V_j - t V_i) Y \quad (3.125)$$

$$= t t^* Y V_i - t^* Y V_j \quad (3.126)$$

$$= |t|^2 Y V_i - t^* Y V_j \quad (3.127)$$

Therefore, the off-diagonal entries in the admittance matrix become

$$Y(i, j) = -t^* Y$$

$$Y(j, i) = -t Y$$

and $|t|^2 Y$ is added to $Y(i, i)$ and Y is added to $Y(j, j)$.

Since regulating transformers are used as voltage control devices, a common computational exercise is to find the tap setting t that will hold the secondary bus voltage magnitude V_j at a specified voltage \hat{V} . This may be interpreted as adding one additional variable to the system (t) and one additional constraint ($V_j = \hat{V}$). Since the additional constraint is counterbalanced by the additional degree of freedom, the dimension of the problem remains the same. There are two primary approaches for finding the tap setting t that results in $V_j = \hat{V}$. One approach is an iterative approach, while the second approach calculates t directly from the power flow equations.

The iterative approach may be summarized as

1. Set $t = t_0$
2. Run a power flow to calculate V_j

3. Is $V_j > \hat{V}$? If yes, then $t = t - \Delta t$, and go to step 2.
4. Is $V_j < \hat{V}$? If yes, then $t = t + \Delta t$, and go to step 2.
5. Done

This approach is conceptually simple and requires no changes to the power flow algorithm. However, it may require numerous runs of a power flow program if t_0 is far from the required tap setting.

The direct approach applies the Newton–Raphson method directly to the updated power flow equations as functions of the tap setting t .

1. Set $V_j = \hat{V}$ and let t be an unknown state
2. Modify the Newton–Raphson Jacobian such that the row of partial derivatives with respect to V_j is replaced by the row of partial derivatives with respect to t
3. Modify the state vector x such that

$$x = \begin{bmatrix} \theta_2 \\ \theta_3 \\ \vdots \\ \theta_n \\ V_2 \\ V_3 \\ \vdots \\ V_{j-1} \\ t \\ V_{j+1} \\ \vdots \\ V_n \end{bmatrix}$$

Note that the state V_j is replaced by t .

4. Perform the Newton–Raphson

In this case, the set of power flow equations is solved only once, but since the system Jacobian is modified, a standard power flow program cannot be used.

Since the tap cannot move continuously along the transformer windings, but must move vertically from one winding to the adjacent winding, the real tap setting is not a continuous state. Therefore, in both cases, the calculated tap setting must be rounded to the nearest possible physical tap setting.

Example 3.12

For the system shown in Figure 3.12, place a transformer with reactance X and real tap t between bus 3 and the load (introduce a new bus 4). Find the new admittance matrix and the corresponding Jacobian entries.

Solution 3.12 Let the admittance matrix of the subsystem containing buses 1 through 3 be given by

$$Y_{bus} = \begin{bmatrix} Y_{11}\angle\phi_{11} & Y_{12}\angle\phi_{12} & Y_{13}\angle\phi_{13} \\ Y_{21}\angle\phi_{21} & Y_{22}\angle\phi_{22} & Y_{23}\angle\phi_{23} \\ Y_{31}\angle\phi_{31} & Y_{32}\angle\phi_{32} & Y_{33}\angle\phi_{33} \end{bmatrix} \quad (3.128)$$

Adding the transformer between buses 3 and 4 yields the new admittance matrix

$$Y_{bus} = \begin{bmatrix} Y_{11}\angle\phi_{11} & Y_{12}\angle\phi_{12} & Y_{13}\angle\phi_{13} & 0 \\ Y_{21}\angle\phi_{21} & Y_{22}\angle\phi_{22} & Y_{23}\angle\phi_{23} & 0 \\ Y_{31}\angle\phi_{31} & Y_{32}\angle\phi_{32} & Y_{33}\angle\phi_{33} + \frac{t^2}{jX} & \frac{-t}{jX} \\ 0 & 0 & \frac{-t}{jX} & \frac{1}{jX} \end{bmatrix} \quad (3.129)$$

The power flow equations at bus 3 become

$$\begin{aligned} 0 &= P_3 - V_3 V_1 Y_{31} \cos(\theta_3 - \theta_1 - \phi_{31}) - V_3 V_2 Y_{32} \cos(\theta_3 - \theta_2 - \phi_{32}) \\ &\quad - V_3 V_4 \left(\frac{t}{X} \right) \cos(\theta_3 - \theta_4 - 90^\circ) - V_3^2 Y_{33} \cos(-\phi_{33}) - V_3^2 \left(\frac{t^2}{X} \right) \cos(90^\circ) \\ 0 &= Q_3 - V_3 V_1 Y_{31} \sin(\theta_3 - \theta_1 - \phi_{31}) - V_3 V_2 Y_{32} \sin(\theta_3 - \theta_2 - \phi_{32}) \\ &\quad - V_3 V_4 \left(\frac{t}{X} \right) \sin(\theta_3 - \theta_4 - 90^\circ) - V_3^2 Y_{33} \sin(-\phi_{33}) - V_3^2 \left(\frac{t^2}{X} \right) \sin(90^\circ) \end{aligned}$$

Since V_4 is specified, there is no partial derivative $\frac{\partial \Delta P_3}{\partial V_4}$; instead there is a partial derivative with respect to t :

$$\frac{\partial \Delta P_3}{\partial t} = -\frac{V_3 V_4}{X} \cos(\theta_3 - \theta_4 - 90^\circ) \quad (3.130)$$

Similarly, the partial derivative of $\frac{\partial \Delta Q_3}{\partial t}$ becomes

$$\frac{\partial \Delta Q_3}{\partial t} = -\frac{V_3 V_4}{X} \sin(\theta_3 - \theta_4 - 90^\circ) + 2V_3^2 \frac{t}{X} \quad (3.131)$$

The partial derivatives with respect to θ_1, θ_2, V_1 , and V_2 do not change, but the partial derivatives with respect to θ_3, θ_4 , and V_3 become

$$\begin{aligned} \frac{\partial \Delta P_3}{\partial \theta_3} &= V_3 V_1 Y_{31} \sin(\theta_3 - \theta_1 - \phi_{31}) + V_3 V_2 Y_{32} \sin(\theta_3 - \theta_2 - \phi_{32}) \\ &\quad + V_3 V_4 \frac{t}{X} \sin(\theta_3 - \theta_4 - 90^\circ) \\ \frac{\partial \Delta P_3}{\partial \theta_4} &= -V_3 V_4 \frac{t}{X} \sin(\theta_3 - \theta_4 - 90^\circ) \\ \frac{\partial \Delta P_3}{\partial V_3} &= -V_1 Y_{31} \cos(\theta_3 - \theta_1 - \phi_{31}) - V_2 Y_{32} \cos(\theta_3 - \theta_2 - \phi_{32}) \\ &\quad - V_4 \frac{t}{X} \cos(\theta_3 - \theta_4 - 90^\circ) - 2V_3 Y_{33} \cos(-\phi_{33}) \end{aligned}$$

$$\begin{aligned}\frac{\partial \Delta Q_3}{\partial \theta_3} &= -V_3 V_1 Y_{31} \cos(\theta_3 - \theta_1 - \phi_{31}) - V_3 V_2 Y_{32} \cos(\theta_3 - \theta_2 - \phi_{32}) \\ &\quad - V_3 V_4 \frac{t}{X} \cos(\theta_3 - \theta_4 - 90^\circ) \\ \frac{\partial \Delta Q_3}{\partial \theta_4} &= V_3 V_4 \frac{t}{X} \cos(\theta_3 - \theta_4 - 90^\circ) \\ \frac{\partial \Delta Q_3}{\partial V_3} &= -V_1 Y_{31} \sin(\theta_3 - \theta_1 - \phi_{31}) - V_2 Y_{32} \sin(\theta_3 - \theta_2 - \phi_{32}) \\ &\quad - V_4 \frac{t}{X} \sin(\theta_3 - \theta_4 - 90^\circ) - 2V_3 Y_{33} \sin(-\phi_{33}) - 2V_3 \frac{t^2}{X}\end{aligned}$$

These partial derivatives are used in developing the Newton–Raphson Jacobian for the iterative power flow method. ■

3.5.3 Decoupled Power Flow

The power flow is one of the most widely used computational tools in power systems analysis. It can be successfully applied to problems ranging from a single machine system to a power system containing tens of thousands of buses. For very large systems, the full power flow may require significant computational resources to calculate, store, and factorize the Jacobian matrix. As discussed previously, however, it is possible to replace the Jacobian matrix with a matrix M that is easier to calculate and factor and still retain good convergence properties. The power flow equations naturally lend themselves to several alternate matrices for the power flow solution that can be derived from the formulation of the system Jacobian. Recall that the system Jacobian has the form

$$\begin{bmatrix} J_1 & J_2 \\ J_3 & J_4 \end{bmatrix} = \begin{bmatrix} \frac{\partial \Delta P}{\partial \theta} & \frac{\partial \Delta P}{\partial V} \\ \frac{\partial \Delta Q}{\partial \theta} & \frac{\partial \Delta Q}{\partial V} \end{bmatrix} \quad (3.132)$$

The general form of the P submatrices is

$$\frac{\partial \Delta P_i}{\partial \theta_j} = -V_i V_j Y_{ij} \sin(\theta_i - \theta_j - \phi_{ij}) \quad (3.133)$$

$$\frac{\partial \Delta P_i}{\partial V_j} = V_i Y_{ij} \cos(\theta_i - \theta_j - \phi_{ij}) \quad (3.134)$$

For most transmission lines, the line resistance contributes only nominally to the overall line impedance; thus, the phase angles ϕ_{ij} of the admittance matrix entries are near $\pm 90^\circ$. Additionally, under normal operating conditions the phase angle difference between adjacent buses is typically small; therefore,

$$\cos(\theta_i - \theta_j - \phi_{ij}) \approx 0 \quad (3.135)$$

leading to

$$\frac{\partial \Delta P_i}{\partial V_j} \approx 0 \quad (3.136)$$

Similar arguments can be made such that

$$\frac{\partial \Delta Q_i}{\partial \theta_j} \approx 0 \quad (3.137)$$

Using the approximations of Equations (3.136) and (3.137), a possible substitution for the Jacobian matrix is the matrix

$$M = \begin{bmatrix} \frac{\partial \Delta P}{\partial \theta} & 0 \\ 0 & \frac{\partial \Delta Q}{\partial V} \end{bmatrix} \quad (3.138)$$

Using this matrix M as a replacement for the system Jacobian leads to a set of *decoupled* iterates for the power flow solution:

$$\theta^{k+1} = \theta^k - \left[\frac{\partial \Delta P}{\partial \theta} \right]^{-1} \Delta P \quad (3.139)$$

$$V^{k+1} = V^k - \left[\frac{\partial \Delta Q}{\partial V} \right]^{-1} \Delta Q \quad (3.140)$$

where the ΔP and ΔQ iterations can be carried out independently. The primary advantage of this decoupled power flow is that the LU factorization computation is significantly reduced. The LU factorization of the full Jacobian requires $(2n)^3 = 8n^3$ floating point operations per iteration, whereas the decoupled power flow requires only $2n^3$ floating point operations per iteration.

Example 3.13

Repeat Example 3.11 using the decoupled power flow algorithm.

Solution 3.13 The Jacobian of Example 3.11 evaluated at the initial condition is

$$[J^0] = \begin{bmatrix} -13.2859 & 9.9010 & 0.9901 \\ 9.9010 & -20.0000 & -1.9604 \\ -0.9901 & 2.0000 & -19.4040 \end{bmatrix} \quad (3.141)$$

Note that the off-diagonal submatrices are much smaller in magnitude than the diagonal submatrices. For example,

$$\|[J_2]\| = \left\| \begin{bmatrix} 0.9901 \\ -1.9604 \end{bmatrix} \right\| << \|[J_1]\| = \left\| \begin{bmatrix} -13.2859 & 9.9010 \\ 9.9010 & -20.0000 \end{bmatrix} \right\|$$

and

$$\|[J_3]\| = \|[-0.9901 \ 2.0000]\| << \|[J_4]\| = \|[-19.4040]\|$$

Thus it is reasonable to neglect the off-diagonal matrices J_2 and J_3 . Therefore, the first iteration of the decoupled power flow becomes

$$\begin{bmatrix} \Delta\theta_2^1 \\ \Delta\theta_3^1 \end{bmatrix} = [J_1]^{-1} \begin{bmatrix} \Delta P_2 \\ \Delta P_3 \end{bmatrix} \quad (3.142)$$

$$= \begin{bmatrix} -13.2859 & 9.9010 \\ 9.9010 & -20.000 \end{bmatrix}^{-1} \begin{bmatrix} 0.5044 \\ -1.1802 \end{bmatrix} \quad (3.143)$$

$$[\Delta V_3^1] = [J_4]^{-1} \Delta Q_3 \quad (3.144)$$

$$= -19.4040^{-1} (-0.2020) \quad (3.145)$$

leading to the updates

$$\begin{bmatrix} \theta_2^1 \\ \theta_3^1 \\ V_3^1 \end{bmatrix} = \begin{bmatrix} -0.0095 \\ -0.0637 \\ 0.9896 \end{bmatrix}$$

The iterative process continues similar to the full Newton–Raphson method by continually updating the J_1 and J_4 Jacobian submatrices and the mismatch equations. The iteration converges when both the ΔP mismatch equations and the ΔQ mismatch equations are both less than the convergence tolerance. Note that it is possible for one set of mismatch equations to meet the convergence criteria before the other; thus the number of “P” iterations required for convergence may differ from the number of “Q” iterations required for convergence. ■

3.5.4 Fast Decoupled Power Flow

In Example 3.13, each of the decoupled Jacobian submatrices is updated at every iteration. As discussed previously, however, it is often desirable to have constant matrices to minimize the number of function evaluations and LU factorizations. This is often referred to as the *fast decoupled power flow* and can be represented as

$$\begin{bmatrix} \Delta P^k \\ V \end{bmatrix} = [B'] [\Delta\theta^{k+1}] \quad (3.146)$$

$$\begin{bmatrix} \Delta Q^k \\ V \end{bmatrix} = [B''] [\Delta V^{k+1}] \quad (3.147)$$

where the B' and B'' are constant [53]. There are a number of variations that are typically denoted as the BB, XB, BX, and XX methods [1].

To derive these matrices from the power flow Jacobian, consider the decoupled power flow relationships for the Newton–Raphson method:

$$[\Delta P] = -[J_1] [\Delta\theta] \quad (3.148)$$

$$\begin{bmatrix} \Delta Q \\ V \end{bmatrix} = -[J_4] [\Delta V] \quad (3.149)$$

where the Jacobian submatrices in rectangular form are

$$J_1(i, i) = V_i \sum_{j \neq i} V_j (g_{ij} \sin \theta_{ij} - b_{ij} \cos \theta_{ij}) \quad (3.150)$$

$$J_1(i, j) = -V_i V_j (g_{ij} \sin \theta_{ij} - b_{ij} \cos \theta_{ij}) \quad (3.151)$$

$$J_4(i, i) = 2V_i b_{ii} - \sum_{j \neq i} V_j (g_{ij} \sin \theta_{ij} - b_{ij} \cos \theta_{ij}) \quad (3.152)$$

$$J_4(i, j) = -V_i (g_{ij} \sin \theta_{ij} - b_{ij} \cos \theta_{ij}) \quad (3.153)$$

where $b_{ij} = |Y_{ij} \sin \phi_{ij}|$ are the imaginary elements of the admittance matrix and $g_{ij} = |Y_{ij} \cos \phi_{ij}|$ are the real elements of the admittance matrix.

By noting that $\phi_{ij} \approx 90^\circ$, $\cos \phi_{ij} \approx 0$, which implies that $g_{ij} \approx 0$. By further approximating all voltage magnitudes as 1.0 per unit,

$$J_1(i, i) = - \sum_{j \neq i} b_{ij} \quad (3.154)$$

$$J_1(i, j) = b_{ij} \quad (3.155)$$

$$J_4(i, i) = 2b_{ii} - \sum_{j \neq i} b_{ij} \quad (3.156)$$

$$J_4(i, j) = b_{ij} \quad (3.157)$$

Since the J_1 submatrix relates the changes in active power to changes in angle, elements that affect mainly reactive power flow can be omitted from this matrix with negligible impact on the convergence properties. Thus shunt capacitors (including line-charging) and external reactances as well as the shunts formed due to representation of off-nominal non-phase-shifting transformers (i.e., taps are set to 1.0) are neglected. Hence, the admittance matrix diagonal elements are devoid of these shunts. Similarly, the J_4 submatrix relates the changes in reactive power to changes in voltage magnitude; therefore, elements that primarily affect active power flow are omitted. Thus all phase-shifting transformers are neglected.

Note that these approximations do not require that the line resistances R_{ij} be neglected, since

$$b_{ij} = \frac{-X_{ij}}{R_{ij}^2 + X_{ij}^2} \quad (3.158)$$

The different variations of the fast decoupled power flow arise primarily from how the line resistances are incorporated into the Jacobian approximations.

When all of the resistances are included, the approximation is known as the BB method. The approximation is seldom used since it usually suffers from poor convergence properties. Similarly, the XX method completely ignores the impact of the resistances. It has slightly better convergence than the BB method.

The most common method is the XB version of the fast decoupled power flow. In this method, the following B' and B'' matrices are defined:

$$B'_{ij} = \frac{1}{X_{ij}} \quad (3.159)$$

$$B'_{ii} = -\sum_{j \neq i} B'_{ij} \quad (3.160)$$

where B' is an approximation to the J_1 matrix.

B'' is an approximation to the J_4 matrix, where

$$B''_{ij} = b_{ij} \quad (3.161)$$

$$B''_{ii} = 2b_i - \sum_{j \neq i} B''_{ij} \quad (3.162)$$

and b_i is the shunt susceptance at bus i (i.e., the sum of susceptances of all the shunt branches connected to bus i).

This method results in a set of constant matrices that can be used to approximate the power flow Jacobian in the Newton–Raphson iteration. Both B' and B'' are real, sparse, and contain only network or admittance matrix elements. In the Newton–Raphson method, these matrices are only factorized once for the LU factorization, and are then stored and held constant throughout the iterative solution process. These matrices were derived based on the application of certain assumptions. If these assumptions do not hold (i.e., the voltage magnitudes deviate substantially from 1.0 per unit, the network has high R/X ratios, or the angle differences between adjacent buses are not small), then convergence problems with the fast decoupled power flow iterations can arise. Work still continues on developing modifications to the XB method to improve convergence [36], [37], [43]. The BX method is similar to the XB method, except that the resistances are ignored in the B'' matrix and not in the B' matrix.

Example 3.14

Repeat Example 3.11 using the fast decoupled power flow algorithm.

Solution 3.14 The line data for the example system are repeated below for convenience:

i	j	R_{ij}	X_{ij}	B_{ij}
1	2	0.02	0.3	0.15
1	3	0.01	0.1	0.1
2	3	0.01	0.1	0.1

and lead to the following admittance matrix:

$$Y_{bus} = \begin{bmatrix} 13.1505\angle -84.7148^\circ & 3.3260\angle 93.8141^\circ & 9.9504\angle 95.7106^\circ \\ 3.3260\angle 93.8141^\circ & 13.1505\angle -84.7148^\circ & 9.9504\angle 95.7106^\circ \\ 9.9504\angle 95.7106^\circ & 9.9504\angle 95.7106^\circ & 19.8012\angle -84.2606^\circ \end{bmatrix} \quad (3.163)$$

Taking the imaginary part of this matrix yields the following B matrix:

$$B = \begin{bmatrix} -13.0946 & 3.3186 & 9.9010 \\ 3.3186 & -13.0946 & 9.9010 \\ 9.9010 & 9.9010 & -19.7020 \end{bmatrix} \quad (3.164)$$

From the line data and the associated B matrix, the following B' and B'' matrices result:

$$B' = \begin{bmatrix} -\frac{1}{x_{21}} - \frac{1}{x_{23}} & \frac{1}{x_{23}} \\ \frac{1}{x_{23}} & -\frac{1}{x_{31}} - \frac{1}{x_{32}} \end{bmatrix} = \begin{bmatrix} -13.3333 & 10 \\ 10 & -20 \end{bmatrix} \quad (3.165)$$

$$\begin{aligned} B'' &= [2b_3 - (B_{31} + B_{32})] \\ &= [2(0.05 + 0.05) - (9.9010 + 9.9010)] = -19.6020 \end{aligned} \quad (3.166)$$

Compare these matrices to the J_1 and J_4 submatrices of Example 3.11 evaluated at the initial condition

$$\begin{aligned} J_1 &= \begin{bmatrix} -13.2859 & 9.9010 \\ 9.9010 & -20.000 \end{bmatrix} \\ J_4 &= [-19.4040] \end{aligned}$$

The similarity between the matrices is to be expected as a result of the defining assumptions of the fast decoupled power flow method.

Iteration 1

The updates can be found by solving the linear set of equations

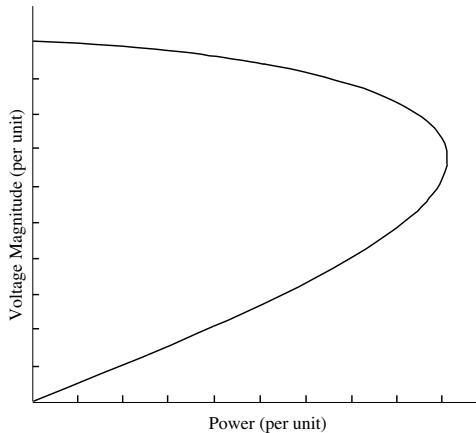
$$\begin{bmatrix} \Delta P_2^0 \\ \Delta P_3^0 \\ \Delta Q_3^0 \end{bmatrix} = \begin{bmatrix} 0.5044 \\ -1.1802 \\ -0.2020 \end{bmatrix} = - \begin{bmatrix} -13.3333 & 10 \\ 10 & -20 \end{bmatrix} \begin{bmatrix} \Delta \theta_2^1 \\ \Delta \theta_3^1 \end{bmatrix}$$

where $\Delta \theta_2^1 = \theta_2^{(1)} - \theta_2^{(0)}$, $\Delta \theta_3^1 = \theta_3^{(1)} - \theta_3^{(0)}$, and $\Delta V_3^1 = V_3^{(1)} - V_3^{(0)}$ and the initial conditions are a “flat start.” Solving for the updates yields

$$\begin{bmatrix} \theta_2^1 \\ \theta_3^1 \\ V_3^1 \end{bmatrix} = \begin{bmatrix} -0.0103 \\ -0.0642 \\ 0.9897 \end{bmatrix}$$

where the phase angles are in radians. This process is continued until convergence in both the “P” and “Q” iterations is achieved. ■

Note that, in both the decoupled power flow cases, the objective of the iterations is the same as for the full Newton–Raphson power flow algorithm.

**FIGURE 3.14**

A PV curve

The objective is to drive the mismatch equations ΔP and ΔQ to within some tolerance. Therefore, regardless of the number of iterations required to achieve convergence, the accuracy of the answer is the same as for the full Newton–Raphson method. In other words, the voltages and angles of the decoupled power flow methods will be the same as with the full Newton–Raphson method as long as the iterates converge.

3.5.5 PV Curves and Continuation Power Flow

The power flow is a useful tool for monitoring system voltages as a function of load change. One common application is to plot the voltage at a particular bus as the load is varied from the base case to a loadability limit (often known as the point of maximum loadability). If the load is increased to the loadability limit and then decreased back to the original loading, it is possible to trace the entire power-voltage or “PV” curve. This curve, shown in Figure 3.14, is sometimes called the *nose curve* for its shape.

At the loadability limit, or tip of the nose curve, the system Jacobian of the power flow equations will become singular as the slope of the nose curve becomes infinite. Thus the traditional Newton–Raphson method of obtaining the power flow solution will break down. In this case, a modification of the Newton–Raphson method known as the *continuation method* is employed. The continuation method introduces an additional equation and unknown into the basic power flow equations. The additional equation is chosen specifically to ensure that the augmented Jacobian is no longer singular at the loadability limit. The additional unknown is often called the continuation parameter.

Continuation methods usually depend on a predictor-corrector scheme and

the means to change the continuation parameter as necessary. The basic approach to tracing the PV curve is to choose a new value for the continuation parameter (either in power or voltage) and then predict the power flow solution for this value. This is frequently accomplished using a tangential (or linear) approximation. Using the predicted value as the initial condition for the nonlinear iteration, the augmented power flow equations are then solved (or corrected) to achieve the solution. So the solution is first predicted and then corrected. This prediction/correction step is shown in Figure 3.15.

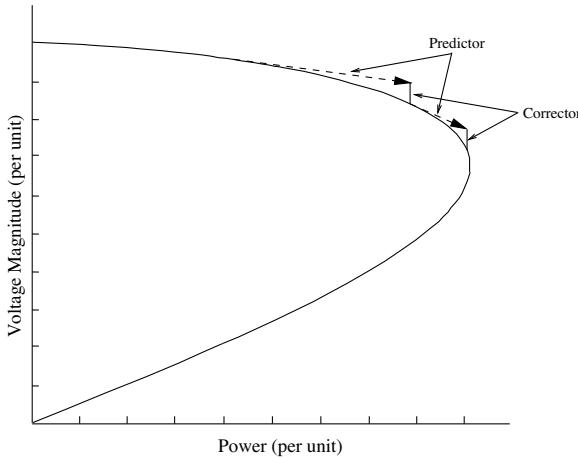


FIGURE 3.15

The predictor/corrector step

Let the set of power flow equations be given as

$$\lambda K - f(\theta, V) = 0 \quad (3.167)$$

or

$$F(\theta, V, \lambda) = 0 \quad (3.168)$$

where K is the loading profile (i.e., the base case relationship between P and Q) and λ is the loading parameter which will vary from unity (at the base case) to the point of maximum loadability. Equation (3.168) may be linearized to yield

$$\frac{\partial F}{\partial \theta} d\theta + \frac{\partial F}{\partial V} dV + \frac{\partial F}{\partial \lambda} d\lambda = 0 \quad (3.169)$$

Due to λ , the number of unknowns in Equation (3.169) is one larger than the number of equations, so one more equation is required:

$$e_k \begin{bmatrix} d\theta \\ dV \\ d\lambda \end{bmatrix} = 1 \quad (3.170)$$

where e_k is a row vector of zeros with a single ± 1 at the position of the unknown that is chosen to be the continuation parameter. The sign of the one in e_k is chosen based on whether the continuation parameter is increasing or decreasing. When the continuation parameter is λ (power), the sign is positive, indicating that the load is increasing. When voltage is the continuation parameter, the sign is negative, indicating that the voltage magnitude is expected to decrease toward the tip of the nose curve.

The unknowns are predicted such that

$$\begin{bmatrix} \theta \\ V \\ \lambda \end{bmatrix}^{\text{predicted}} = \begin{bmatrix} \theta_0 \\ V_0 \\ \lambda_0 \end{bmatrix} + \sigma \begin{bmatrix} d\theta \\ dV \\ d\lambda \end{bmatrix} \quad (3.171)$$

where

$$\begin{bmatrix} d\theta \\ dV \\ d\lambda \end{bmatrix} = \begin{bmatrix} & & \vdots \\ J_{LF} & & K \\ & & \vdots \\ \dots & \dots & \dots & \dots \\ [e_k] \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

and σ is the step size (or length) for the next prediction. Note that the continuation state $dx_k = 1$; thus

$$x_k^{\text{predicted}} = x_{k0} + \sigma$$

so σ should be chosen to represent a reasonable step size in terms of what the continuation parameter is (usually voltage or power).

The corrector step involves the solution of the set of equations

$$F(\theta, V, \lambda) = 0 \quad (3.172)$$

$$x_k - x_k^{\text{predicted}} = 0 \quad (3.173)$$

where x_k is the chosen continuation parameter. Typically, the continuation parameter is chosen as the state that exhibits the greatest rate of change.

Example 3.15

Plot the PV curve (P versus V) of the system shown in Figure 3.16 using the continuation power flow method as the load varies from zero to the point of maximum loadability.

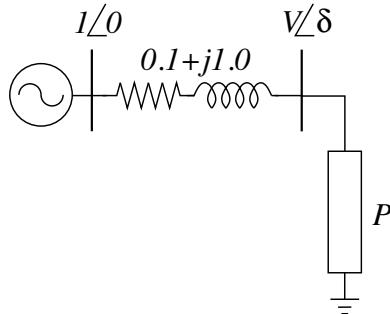


FIGURE 3.16
System for Example 3.15

Solution 3.15 The power flow equations for the system shown in Figure 3.16 are

$$0 = -P - 0.995V \cos(\theta - 95.7^\circ) - 0.995V^2 \cos(84.3^\circ) \quad (3.174)$$

$$0 = -0.995V \sin(\theta - 95.7^\circ) - 0.995V^2 \sin(84.3^\circ) \quad (3.175)$$

During the continuation power flow, the vector of injected active and reactive powers will be replaced by the vector λK . The loading vector λK is

$$\lambda K = \lambda \begin{bmatrix} -1 \\ 0 \end{bmatrix}$$

where λ will vary from zero to the maximum loading value. Typically, the vector K will contain the base case values for all injected active and reactive powers in the system. In this case, the entry for the load P is negative, indicating that the injected power is negative (i.e., a load).

The power flow Jacobian for this set of power flow equations is

$$J_{LF} = \begin{bmatrix} 0.995V \sin(\theta - 95.7^\circ) & -0.995 \cos(\theta - 95.7^\circ) - 1.99V \cos(84.3^\circ) \\ -0.995V \cos(\theta - 95.7^\circ) & -0.995 \sin(\theta - 95.7^\circ) - 1.99V \sin(84.3^\circ) \end{bmatrix}$$

Iteration 1

Initially, the continuation parameter is chosen to be λ since the load will change more rapidly than the voltage at points far from the tip of the nose curve. At $\lambda = 0$, the circuit is under no-load and the initial voltage magnitude

and angle are $1\angle 0^\circ$. With $\sigma = 0.1$ pu, the predictor step yields

$$\begin{bmatrix} \theta \\ V \\ \lambda \end{bmatrix}^{\text{predicted}} = \begin{bmatrix} \theta \\ V \\ \lambda \end{bmatrix}^i + \sigma \begin{bmatrix} & & & \vdots & \\ J_{LF} & & & K & \\ & & & \vdots & \\ \dots & \dots & \dots & \dots & \\ & & & [e_k] & \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \quad (3.176)$$

$$= \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} + \sigma \begin{bmatrix} -0.9901 & -0.0988 & -1 \\ 0.0988 & -0.9901 & 0 \\ 0 & 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad (3.177)$$

$$= \begin{bmatrix} -0.1000 \\ 0.9900 \\ 0.1000 \end{bmatrix} \quad (3.178)$$

where θ is in radians. Note that the predicted value for λ is 0.1 pu.

The corrector step solves the system of equations

$$0 = -\lambda - 0.995V \cos(\theta - 95.7^\circ) - 0.995V^2 \cos(84.3^\circ) \quad (3.179)$$

$$0 = -0.995V \sin(\theta - 95.7^\circ) - 0.995V^2 \sin(84.3^\circ) \quad (3.180)$$

with the load parameter λ set to 0.1 pu. Note that this is a regular power flow problem and can be solved without program modification.

The first corrector step yields

$$\begin{bmatrix} \theta \\ V \\ \lambda \end{bmatrix} = \begin{bmatrix} -0.1017 \\ 0.9847 \\ 0.1000 \end{bmatrix}$$

Note that this procedure is consistent with the illustration in Figure 3.15. The prediction step is of length σ taken tangentially to the PV at the current point. The corrector step will then occur along a vertical path because the power (λK) is held constant during the correction.

Iteration 2

The second iteration proceeds as the first. The predictor step yields the following guess:

$$\begin{bmatrix} \theta \\ V \\ \lambda \end{bmatrix} = \begin{bmatrix} -0.2060 \\ 0.9637 \\ 0.2000 \end{bmatrix}$$

where λ is increased by the step size $\sigma = 0.1$ pu.

Correcting the values yields the second update:

$$\begin{bmatrix} \theta \\ V \\ \lambda \end{bmatrix} = \begin{bmatrix} -0.2105 \\ 0.9570 \\ 0.2000 \end{bmatrix}$$

Iterations 3 and 4

The third and fourth iterations progress similarly. The values to this point are summarized:

λ	V	θ	σ
0.1000	0.9847	-0.1017	0.1000
0.2000	0.9570	-0.2105	0.1000
0.3000	0.9113	-0.3354	0.1000
0.4000	0.8268	-0.5050	0.1000

Beyond this point, the power flow fails to converge for a step size of $\sigma = 0.1$. The method is nearing the point of maximum power flow (the tip of the nose curve), as indicated by the rapid decline in voltage for relatively small changes in λ . At this point, the continuation parameter is switched from λ to V to ensure that the corrector step will converge. The predictor step is modified such that

$$\begin{bmatrix} d\theta \\ dV \\ d\lambda \end{bmatrix} = \begin{bmatrix} d\theta_0 \\ dV_0 \\ d\lambda_0 \end{bmatrix} + \sigma \begin{bmatrix} [J_{LF}] & K \\ 0 & 0 \\ 0 & -1 & 0 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

where the -1 in the last row (the e_k vector) now corresponds to V rather than λ . The step size σ is changed to 0.025 pu because changes in voltage are smaller than power (λ). Note that the minus sign will ensure that the voltage magnitude will *decrease* by 0.025 pu at every step.

The corrector step is also modified when the continuation parameter switches to voltage magnitude. The new augmented equations become:

$$0 = f_1(\theta, V, \lambda) = -\lambda - 0.995V(\cos(\theta - 95.7^\circ) + V \cos(84.3^\circ)) \quad (3.181)$$

$$0 = f_2(\theta, V, \lambda) = -0.995V \sin(\theta - 95.7^\circ) - 0.995V^2 \sin(84.3^\circ) \quad (3.182)$$

$$0 = f_3(\theta, V, \lambda) = V^{\text{predicted}} - V \quad (3.183)$$

which cannot be solved with a traditional power flow program due to the last equation. This equation is necessary to keep the Newton–Raphson iteration nonsingular. Fortunately, the Newton–Raphson iteration uses the same iteration matrix as the predictor matrix:

$$\begin{bmatrix} [J_{LF}] & -\lambda \\ 0 & 0 \\ 0 & -1 & 0 \end{bmatrix}^{-1} \left(\begin{bmatrix} \theta \\ V \\ \lambda \end{bmatrix}^{(k+1)} - \begin{bmatrix} \theta \\ V \\ \lambda \end{bmatrix}^{(k)} \right) = - \begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix} \quad (3.184)$$

thus minimizing the computational requirement.

Note that the corrector step is now a horizontal correction in voltage. The voltage magnitude is held constant while λ and θ are corrected. These iterates proceed as

λ	Predicted				Corrected		
	V	θ	σ		λ	V	θ
0.4196	0.8019	-0.5487	0.0250		0.4174	0.8019	-0.5474
0.4326	0.7769	-0.5887	0.0250		0.4307	0.7769	-0.5876
0.4422	0.7519	-0.6268	0.0250		0.4405	0.7519	-0.6260
0.4487	0.7269	-0.6635	0.0250		0.4472	0.7269	-0.6627
0.4525	0.7019	-0.6987	0.0250		0.4511	0.7019	-0.6981
0.4538	0.6769	-0.7328	0.0250		0.4525	0.6769	-0.7323
0.4528	0.6519	-0.7659	0.0250		0.4517	0.6519	-0.7654

Note that, at the last updates, the load parameter λ has started to decrease with decreasing voltage. This indicates that the continuation power flow is now starting to map out the lower half of the nose curve. However, while the iterations are still close to the tip of the nose curve, the Jacobian will still be ill conditioned, so it is a good idea to take several more steps before switching the continuation parameter from voltage magnitude back to λ .

λ	Predicted				Corrected		
	V	θ	σ		λ	V	θ
0.4497	0.6269	-0.7981	0.0250		0.4487	0.6269	-0.7977
0.4447	0.6019	-0.8295	0.0250		0.4438	0.6019	-0.8291
0.4380	0.5769	-0.8602	0.0250		0.4371	0.5769	-0.8598
0.4296	0.5519	-0.8902	0.0250		0.4288	0.5519	-0.8899
0.4197	0.5269	-0.9197	0.0250		0.4190	0.5269	-0.9194
0.4084	0.5018	-0.9486	0.0250		0.4077	0.5018	-0.9483
0.3958	0.4768	-0.9770	0.0250		0.3951	0.4768	-0.9768
0.3820	0.4518	-1.0051	0.0250		0.3814	0.4518	-1.0049
0.3670	0.4268	-1.0327	0.0250		0.3665	0.4268	-1.0325
0.3510	0.4018	-1.0600	0.0250		0.3505	0.4018	-1.0598
0.3340	0.3768	-1.0870	0.0250		0.3335	0.3768	-1.0868

After switching the continuation parameter back to λ , the e_k vector becomes

$$e_k = [0 \quad 0 \quad -1]$$

and σ is set to 0.1 pu which indicates that λ will decrease by 0.1 pu at every step (i.e., the power is decreasing back to the base case). The predictor/corrector steps proceed as above, yielding

λ	Predicted				Corrected		
	V	θ	σ		λ	V	θ
0.2335	0.2333	-1.2408	0.1000		0.2335	0.2486	-1.2212
0.1335	0.1312	-1.3417	0.1000		0.1335	0.1374	-1.3340
0.0335	0.0309	-1.4409	0.1000		0.0335	0.0338	-1.4375

These values are combined in the PV curve shown in Figure 3.17. Note the change in step size when the continuation parameter switches from λ to voltage near the tip of the PV curve. The appropriate choice of step size is problem dependent and can be adaptively changed to improve computational efficiency. ■

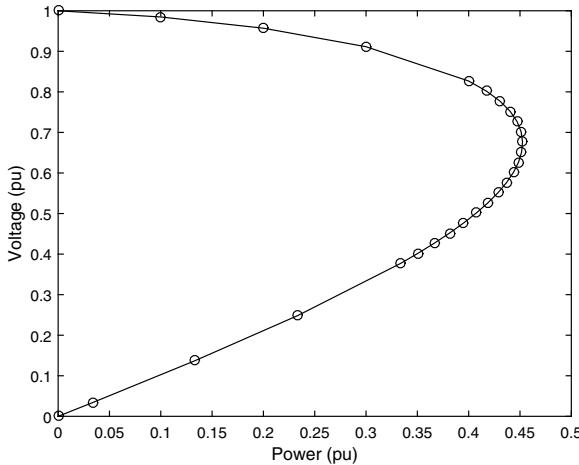


FIGURE 3.17

PV curve for system of Example 3.15

3.5.6 Three-Phase Power Flow

Another special-purpose power flow application is for three-phase power flow analysis. Although much of the power system analysis is performed on balanced three-phase systems using one line equivalent diagrams, certain situations call for a three-phase analysis. In particular, in the situation where the transmission system is not balanced due to nontransposed lines or when the loads are considerably unbalanced, it may be desirable to perform a full three-phase power flow to ascertain the effect on the individual line flows and bus voltages. The development of the three-phase power flow is similar to that of the single-phase equivalent, except that the mutual coupling between lines must be incorporated into the admittance matrix, yielding a $3n \times 3n$ matrix with elements Y_{ij}^{pq} where the subscript ij indicates bus number ($1 \leq (i, j) \leq n$) and the superscript pq indicates phase ($p, q \in [a, b, c]$).

The incorporation of each phase individually leads to the similar, but slightly more complex, three-phase power flow equations

$$0 = \Delta P_i^p = P_i^{inj,p} - V_i^p \sum_{q \in (a,b,c)} \sum_{j=1}^{N_{bus}} V_j^q Y_{ij}^{pq} \cos(\theta_i^p - \theta_j^q - \phi_{ij}^{pq}) \quad (3.185)$$

$$0 = \Delta Q_i^p = Q_i^{inj,p} - V_i^p \sum_{q \in (a,b,c)} \sum_{j=1}^{N_{bus}} V_j^q Y_{ij}^{pq} \sin(\theta_i^p - \theta_j^q - \phi_{ij}^{pq}) \quad (3.186)$$

$i = 1, \dots, N_{bus}$ and $p \in (a, b, c)$

There are three times as many power flow equations as in the single-phase equivalent power flow equations. Generator (PV) buses are handled similarly with the following exceptions:

1. $\theta^a = 0^\circ, \theta^b = -120^\circ, \theta^c = 120^\circ$ for the slack bus
2. All generator voltage magnitudes and active powers in each phase must be equal since generators are designed to produce balanced output

A three-phase power flow “flat start” is to set each bus voltage magnitude to

$$\begin{aligned} V_i^a &= 1.0\angle 0^\circ \\ V_i^b &= 1.0\angle -120^\circ \\ V_i^c &= 1.0\angle 120^\circ \end{aligned}$$

The system Jacobian used in the Newton–Raphson solution of the power flow equations will have a possible $(3(2n) \times 3(2n))$ or $36n^2$ entries. The Jacobian partial derivatives are found in the same manner as with the single-phase equivalent system except that the derivatives must also be taken with respect to phase differences. For example,

$$\frac{\partial \Delta P_i^a}{\partial \theta_j^b} = V_i^a V_j^b Y_{ij}^{ab} \sin(\theta_i^a - \theta_j^b - \phi_{ij}^{ab}) \quad (3.187)$$

which is similar to the single-phase equivalent system. Similarly,

$$\frac{\partial \Delta P_i^a}{\partial \theta_i^a} = -V_i^a \sum_{q \in (a,b,c)} \sum_{j=1}^{N_{bus}} V_j^q Y_{ij}^{pq} \sin(\theta_i^p - \theta_j^q - \phi_{ij}^{pq}) + (V_i^a)^2 Y_{ii}^{pp} \cos(\phi_{ii}^{pp}) \quad (3.188)$$

The remaining partial derivatives can be calculated in a similar manner and the solution process of the three-phase power flow follows the method outlined in Section 3.5.1.

3.6 Problems

1. Prove that the Newton–Raphson iteration will diverge for the following functions regardless of choice of initial condition

(a) $f(x) = x^2 + 1$

- (b) $f(x) = 7x^4 + 3x^2 + \pi$
2. Devise an iteration formula for computing the fifth root of any positive real number.
 3. Using the Newton–Raphson method, solve

$$0 = 4y^2 + 4y + 52x - 19$$

$$0 = 169x^2 + 3y^2 + 111x - 10y - 10$$

with $[x^0 \ y^0]^T = [1 \ 1]^T$.

4. Using the Newton–Raphson method, solve

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

$$0 = -xy + 2y^2 - 1$$

with $[x^0 \ y^0]^T = [1 \ 1]^T$.

5. Repeat Problems 3 and 4 using numerical differentiation to compute the Jacobian. Use a perturbation of 1% in each variable.
6. Repeat Problems 3 and 4 using Broyden’s method.
7. Repeat Problems 3 and 4 using the Jacobian-Free Newton–GMRES method.
8. Repeat Problems 3 and 4 using homotopic mapping with $0 = f_{01} = x_1 - 2$ and $0 = f_{02} = x_2 - 4$. Use $(2, 4)$ as the initial guess and a $\Delta\lambda$ of 0.05.
9. Write a *generalized* (for any system) power flow program. Your program should
 - (a) Read in load, voltage, and generation data. You may assume that bus #1 corresponds to the slack bus.
 - (b) Read in line and transformer data and create the Y_{bus} matrix.
 - (c) Solve the power flow equations using the Newton–Raphson algorithm, for a stopping criterion of $\|f(x^k)\| \leq \epsilon = 0.0005$.
 - (d) Calculate all dependent unknowns, line flows, and line losses.

The Newton–Raphson portion of the program should call the lufact and permute subroutines. Your program should give you the option of using either a “flat start” or “previous values” as an initial guess. The easiest way to accomplish this is to read and write to the same data file. Note that the first run must be a “flat start” case.

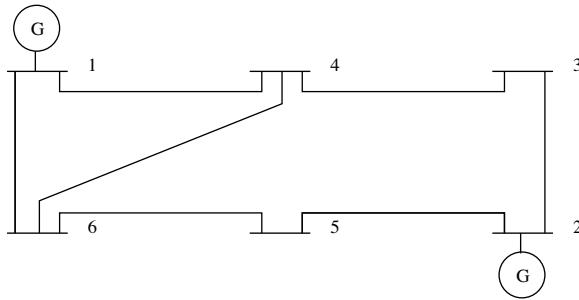


FIGURE 3.18
Ward–Hale six-bus system

10. The data for the system shown in Figure 3.18 are given below:

No.	Type	$ V $	θ	P_{gen}	Q_{gen}	P_{load}	Q_{load}
1	0	1.05	0	0	0	0.25	0.1
2	1	1.05	0	0.5	0	0.15	0.05
3	2	1.00	0	0	0	0.275	0.11
4	2	1.00	0	0	0	0	0
5	2	1.00	0	0	0	0.15	0.09
6	2	1.00	0	0	0	0.25	0.15
No.	To	From	R	X	B		
1	1	4	0.020	0.185	0.009		
2	1	6	0.031	0.259	0.010		
3	2	3	0.006	0.025	0.000		
4	2	5	0.071	0.320	0.015		
5	4	6	0.024	0.204	0.010		
6	3	4	0.075	0.067	0.000		
7	5	6	0.025	0.150	0.017		

Calculate the power flow solution for the system data given above. Remember to calculate all dependent unknowns, line flows, and line losses.

11. Modify your power flow program so that you are using a *decoupled power flow* (i.e., assume $\left[\frac{\partial \Delta P}{\partial V}\right] = 0$ and $\left[\frac{\partial \Delta Q}{\partial \theta}\right] = 0$). Repeat Problem 10. Discuss the difference in convergence between the decoupled and the full Newton–Raphson power flows.
12. Increase the line resistances by 75% (i.e., multiply all resistances by 1.75) and repeat Problem 10 and Problem 11. Discuss your findings.

13. Using a continuation power flow, map out the “PV” curve for the original system data by increasing/decreasing the load on bus 6, holding a constant P/Q ratio from $P = 0$ to the point of maximum power transfer.
14. Making the following assumptions, find a **constant, decoupled** Jacobian that could be used in a fast, decoupled three-phase power flow.
 - $V_i^p \approx 1.0$ pu for all i and p
 - $\theta_{ij}^{pp} \approx 0$
 - $\theta_{ij}^{pm} \approx \pm 120^\circ$ $p \neq m$
 - $g_{ij}^{pm} \ll b_{ij}^{pm}$

4

Sparse Matrix Solution Techniques

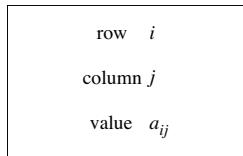
A sparse matrix is one that has “very few” nonzero elements. A sparse system is one in which its mathematical description gives rise to sparse matrices. Any large system that can be described by coupled nodes may give rise to a sparse matrix if the majority of nodes in the system have very few connections. Many systems in engineering and science result in sparse matrix descriptions. Large systems in which each node is connected to only a handful of other nodes include the mesh points in a finite-element analysis, nodes in an electronic circuit, and the busbars in an electric power network. For example, power networks may contain thousands of nodes (busbars), but the average connectivity of electric power network nodes is three; each node is connected, on average, to three other nodes. This means that, in a system composed of a thousand nodes, the nonzero percentage of the descriptive system matrix is

$$\frac{\frac{4 \text{ nonzero elements}}{\text{row}} \times 1000 \text{ rows}}{1000 \times 1000 \text{ elements}} \times 100\% = 0.4\% \text{ nonzero elements}$$

Thus, if only the nonzero elements were stored in memory, they would require only 0.4% of the memory requirements of the full 1000×1000 matrix. Full storage of an $n \times n$ system matrix grows as n^2 , whereas the sparse storage of the same system matrix increases only linearly as $\sim n$. Thus significant storage and computational savings can be realized by exploiting sparse storage and solution techniques. Another motivating factor in exploiting sparse matrix solution techniques is the computational effort involved in solving matrices with large percentages of zero elements. Consider the solution of the linear problem

$$Ax = b$$

where A is sparse. The factorization of L and U from A requires a significant number of multiplications where one or both of the factors may be zero. If it is known ahead of time where the zero elements reside in the matrix, these multiplications can be avoided (since their product will be zero) and significant computational effort can be saved. The salient point here is that these computations are skipped altogether. A person performing an LU factorization by hand can note which values are zero and skip those particular multiplications. A computer, however, does not have the ability to “see” the zero elements. Therefore, the sparse solution technique must be formulated in such a way as to avoid zero computations altogether and operate only upon nonzero elements.

**FIGURE 4.1**

Basic storage element for a_{ij}

In this chapter, both the storage and computational aspects of sparse matrix solutions will be explored. Several storage techniques will be discussed and ordering techniques to minimize computational effort will be developed.

4.1 Storage Methods

In sparse storage methods, only the nonzero elements of the $n \times n$ matrix A are stored in memory, along with the indexing information needed to traverse the matrix from element to element. Thus each element must be stored with its real value (a_{ij}) and its position indices (row and column) in the matrix. The basic storage unit may be visualized as the object shown in Figure 4.1.

In addition to the basic information, indexing information must also be included in the object, such as a link to the next element in the row, or the next element in the column, as shown in Figure 4.2.

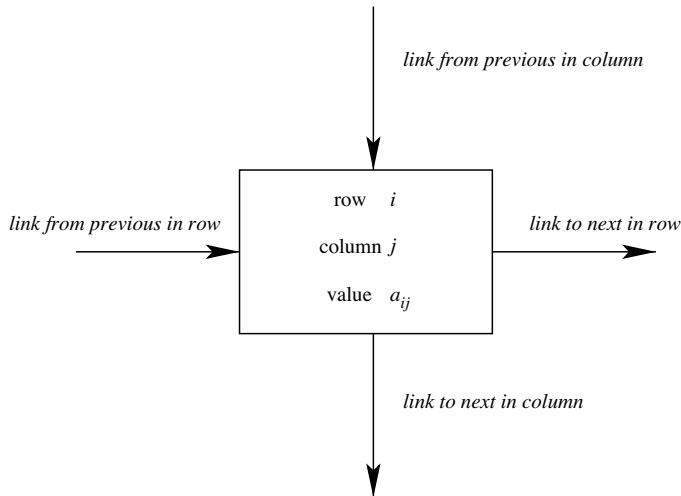
The only additional information necessary to fully traverse the matrix either by row or column is an indication of the first element in each row or column. This is a stand-alone set of links that point to the first element in each row or column.

Example 4.1

Determine a linked list representation for the sparse matrix

$$A = \begin{bmatrix} -1 & 0 & -2 & 0 & 0 \\ 2 & 8 & 0 & 1 & 0 \\ 0 & 0 & 3 & 0 & -2 \\ 0 & -3 & 2 & 0 & 0 \\ 1 & 2 & 0 & 0 & -4 \end{bmatrix}$$

Solution 4.1 A linked list representation of this matrix is shown in Figure 4.3. The last element of each column and row are linked to a null point. Note that each object is linked to its adjacent neighbors in the matrix, both

**FIGURE 4.2**

Storage element for element a_{ij} with links

by column and by row. In this way, the entire matrix can be traversed in any direction by starting at the first element and following the links until the desired element is reached. ■

If a command requires a particular matrix element, then by choosing either the column or row, the element can be located by progressing along the links. If, during the search, the null point is reached, then the desired element does not exist and a value of zero is returned. Additionally, if the matrix elements are linked by increasing index and an element is reached that has a greater index than the desired element, then the progression terminates and a value of zero is returned.

A linked list representation for a sparse matrix is not unique and the elements do not necessarily have to be linked in increasing order by index. However, ordering the elements by increasing index leads to a simplified search since the progression can be terminated before reaching the null point if the index of the linked object exceeds the desired element. If the elements are not linked in order, the entire row or column must be traversed to determine whether or not the desired element is nonzero. The drawback to an ordered list is that the addition of new nonzero elements to the matrix requires the update of both row and column links.

Example 4.2

Insert the matrix element $A(4, 5) = 10$ to the linked list of Example 4.1.

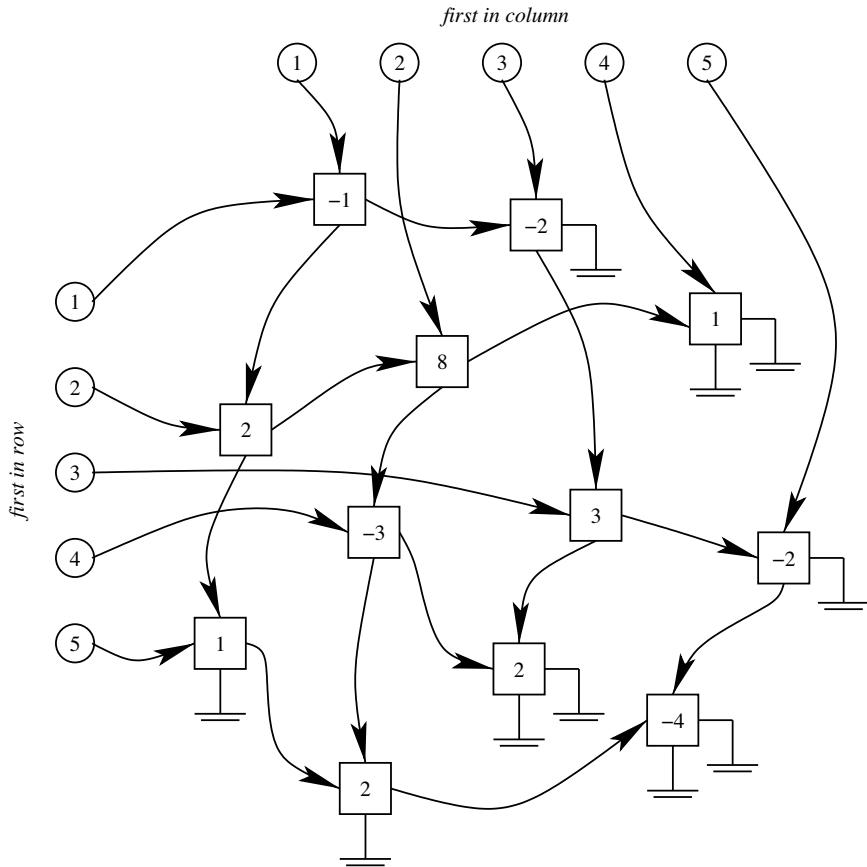
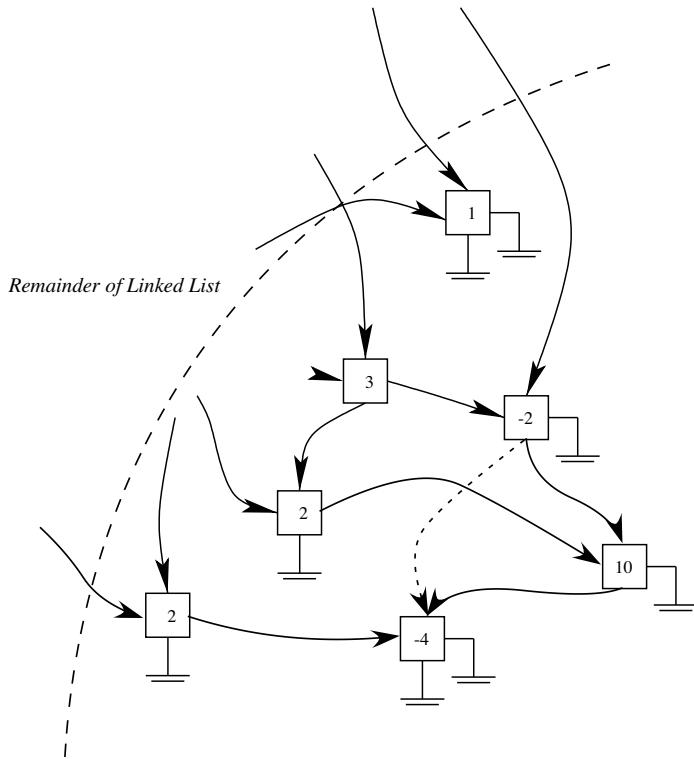


FIGURE 4.3

Linked list for Example 4.1

**FIGURE 4.4**Insertion of matrix element $A(4, 5) = 10$

Solution 4.2 The insertion of the new element is shown in Figure 4.4. The addition of this element requires two traversals of the matrix to insert the new element and to update the links: one traversal by row and one by column. Starting at the *first in row* link for row 4 (value=-3), the elements are progressed by link. Following the links and monitoring the column indices, it is discovered that the element with column index 3 (value=2) is the last element in the row since it points to null. Since the new element has column index 5, it is inserted between the (value=2) element and the null point in the row linked list. Similarly, starting at the *first in column* link for column 5 (value=-2), the column is traversed and inserted between the elements with row indices 3 (value=-2) and 5 (value=-4). The column links are updated to reflect the insertion. ■

If the linked lists of the matrix are not ordered by index, then new elements can be added without transversing the rows or columns. A new element can be inserted in each row or column by inserting them before the first element and updating the *first in row* and *first in column* pointers.

Many software languages, however, do not support the use of objects, pointers, and linked lists. In this case it is necessary to develop a procedure to mimic a linked list format by the use of vectors. Three vectors are required to represent each nonzero element object: one vector containing the row number (NROW), one vector containing the column number (NCOL), and one vector containing the value of the element (VALUE). These vectors are of length nnz where nnz is the number of nonzero elements. Two vectors, also of length nnz , are required to represent the next-in-row links (NIR) and the next-in-column (NIC) links. If an element is the last in the row or column, then the NIR or NIC value for that element is 0. Finally, two vectors of length n contain the *first in row* (FIR) and *first in column* (FIC) links.

The elements of the matrix are assigned a (possibly arbitrary) numbering scheme that corresponds to their order in the NROW, NCOL, VALUE, NIR, and NIC vectors. This order is the same for each of these five vectors. The FIR and FIC vectors will also refer to this number scheme.

Example 4.3

Find the vectors NROW, NCOL, VALUE, NIR, NIC, FIR, and FIC for the sparse matrix of Example 4.1.

Solution 4.3 The matrix of Example 4.1 is reproduced below with the numbering scheme given in parentheses to the left of each nonzero element. The numbering scheme is sequential by row and goes from 1 to $nnz = 12$.

$$A = \begin{bmatrix} (1) -1 & 0 & (2) -2 & 0 & 0 \\ (3) 2 & (4) 8 & 0 & (5) 1 & 0 \\ 0 & 0 & (6) 3 & 0 & (7) -2 \\ 0 & (8) -3 & (9) 2 & 0 & 0 \\ (10) 1 & (11) 2 & 0 & 0 & (12) -4 \end{bmatrix}$$

The ordering scheme indicated yields the following nnz vectors:

k	VALUE	NROW	NCOL	NIR	NIC
1	-1	1	1	2	3
2	-2	1	3	0	6
3	2	2	1	4	10
4	8	2	2	5	8
5	1	2	4	0	0
6	3	3	3	7	9
7	-2	3	5	0	12
8	-3	4	2	9	11
9	2	4	3	0	0
10	1	5	1	11	0
11	2	5	2	12	0
12	-4	5	5	0	0

and the following n vectors:

	FIR	FIC
1	1	1
2	3	4
3	6	2
4	8	5
5	10	7

Consider the matrix element $A(2, 2) = 8$. It is the fourth element in the numbering scheme, so its information is stored in the fourth place in vectors VALUE, NROW, NCOL, NIR, and NIC. Thus VALUE(4)=8, NROW(4)=2, and NCOL(4)=2. The next element in row 2 is $A(2, 4) = 1$ and it is element 5 in the numbering scheme. Therefore NROW(4)=5, signifying that element 5 follows element 4 in its row (note, however, that it does not indicate which row they are in). Similarly, the next element in column 2 is $A(4, 2) = -3$ and it is element 8 in the numbering scheme. Therefore NCOL(4)=8. ■

Example 4.4

Find the sparse storage vectors for the matrix data given below.

i	j	$A(i, j)$
8	8	-28
7	2	5
10	5	7
5	10	7
5	7	3
6	6	-33
6	5	10
1	8	8
5	5	-44
1	4	19
4	3	6
8	3	1
3	4	6
10	3	9
2	1	2
9	7	13
10	2	10
3	8	1
3	10	9
4	1	19
7	7	-68
8	1	8
2	10	10

i	j	$A(i, j)$
3	3	-40
6	7	19
7	8	15
4	4	-38
5	6	10
3	5	11
7	5	3
5	4	9
1	2	2
3	7	9
2	7	5
7	3	9
2	2	-21
1	1	-33
7	9	13
4	5	9
5	3	11
10	10	-30
7	6	19
9	9	-17
8	7	15

Solution 4.4 Since data is read in from a file, there is no guarantee that the data will be given in any particular order. The elements are numbered in the order in which they are read, and not in the order of the matrix A . The full matrix (including the zero elements) should never be explicitly created.

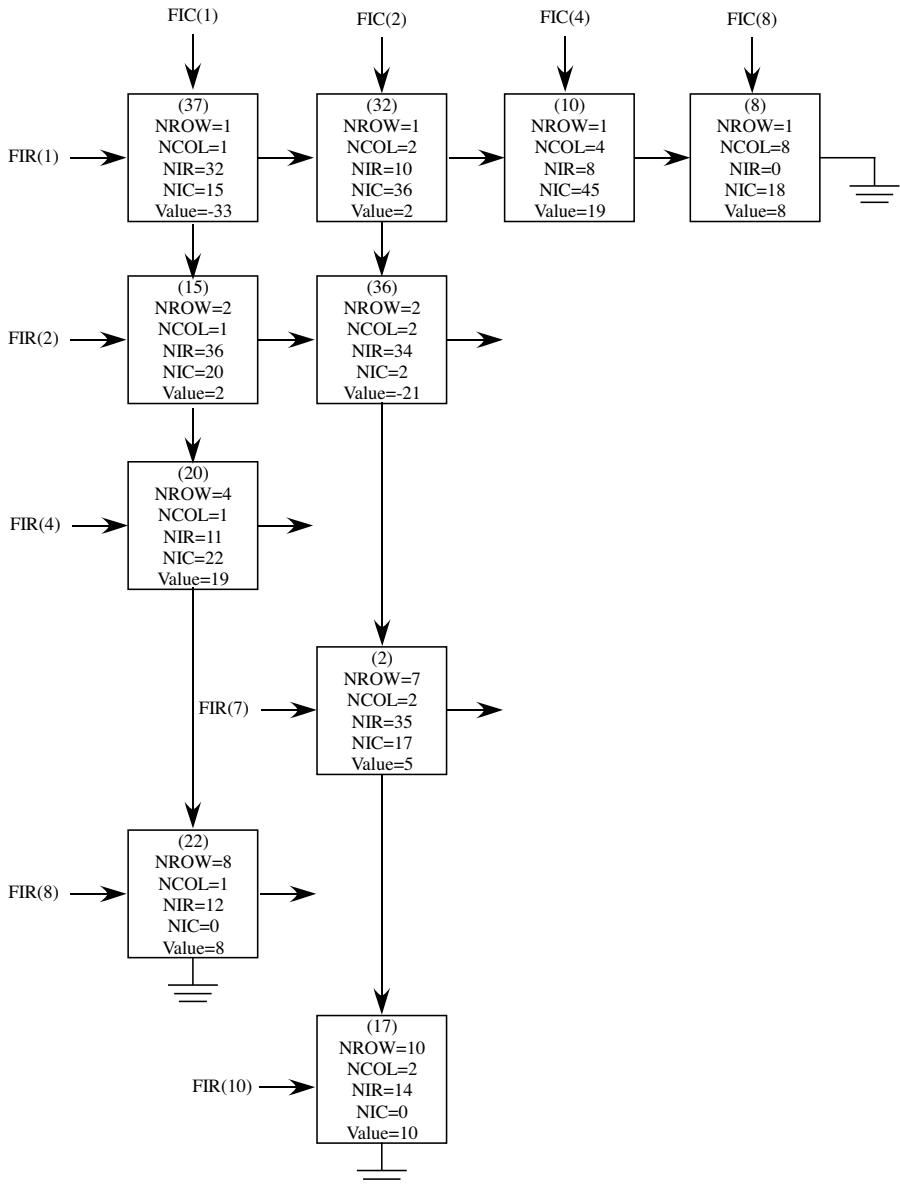
i	NROW	NCOL	NIR	NIC	Value
1	8	8	0	0	-28
2	7	2	35	17	5
3	10	5	41	0	7
4	5	10	0	41	7
5	5	7	4	25	3
6	6	6	25	42	-33
7	6	5	6	30	10
8	1	8	0	18	8
9	5	5	28	7	-44
10	1	4	8	13	19
11	4	3	27	40	6
12	8	3	44	14	1
13	3	4	29	27	6
14	10	3	3	0	9
15	2	1	36	20	2
16	9	7	43	0	13
17	10	2	14	0	10
18	3	8	19	26	1

<i>i</i>	NROW	NCOL	NIR	NIC	Value
19	3	10	0	4	9
20	4	1	11	22	19
21	7	7	26	44	-68
22	8	1	12	0	8
23	2	10	0	19	10
24	3	3	13	11	-40
25	6	7	0	21	19
26	7	8	38	1	15
27	4	4	39	31	-38
28	5	6	5	6	10
29	3	5	33	39	11
30	7	5	42	3	3
31	5	4	9	0	9
32	1	2	10	36	2
33	3	7	18	5	9
34	2	7	23	33	5
35	7	3	30	12	9
36	2	2	34	2	-21
37	1	1	32	15	-33
38	7	9	0	43	13
39	4	5	0	9	9
40	5	3	31	35	11
41	10	10	0	0	-30
42	7	6	21	0	19
43	9	9	0	0	-17
44	8	7	1	16	15

and

<i>i</i>	FIR	FIC
1	37	37
2	15	32
3	24	24
4	20	10
5	40	29
6	7	28
7	2	34
8	22	8
9	16	38
10	17	23

Figure 4.5 shows a visualization of the linked list produced by the sparse vectors. The first-in-row and first-in-column pointers point to the first element in each row/column. Note that not all of the first elements are shown. The next-in-row and next-in-column links are shown. Each row or column can be traversed by starting at the first element and subsequently moving through the links. ■

**FIGURE 4.5**Elements in columns 1 and 2 of A in Example 4.4

4.2 Sparse Matrix Representation

Sparse matrices arise as the result of the mathematical modeling of a sparse system. In many cases, the system has a naturally occurring physical network representation or lends itself to a physically intuitive representation. In these cases, it is often informative to visualize the connectivity of the system by graphical means. In the graphical representation, each node of the graph corresponds to a node in the system. Each edge of the graph corresponds to a branch of the network. As with a network, the graph, consisting of vertices and edges, is often represented by a set of points in the plane joined by a line representing each edge. Matrices that arise from the mathematical model of a graphically represented network are structurally symmetric. In other words, if the matrix element a_{ij} is nonzero, then the matrix element a_{ji} is also nonzero. This implies that, if node i is connected to node j , then node j is also connected to node i . Matrices that are not structurally symmetric can be made symmetric by adding an element of value zero in the appropriate position within the matrix.

In addition to a graphical representation, it is also common to visualize sparse matrices by a matrix that is clear except for an identifying symbol (such as a \times , \bullet , $*$, or other mark) to represent the position of the nonzero elements in the matrix. The finite element grid of the trapezoid shown in Figure 4.6(a) gives rise to the sparse matrix structure shown in Figure 4.6(b). Note that the ordering of the matrix is not unique; another numbering scheme for the nodes will result in an alternate matrix structure.

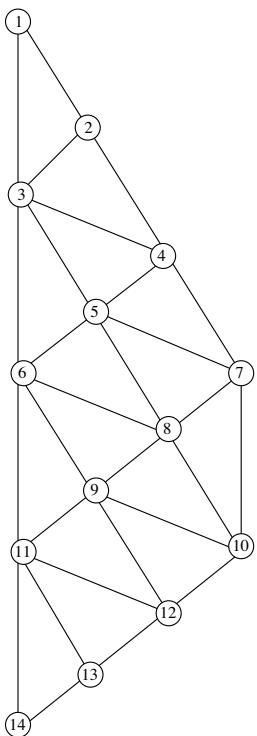
Example 4.5

Find the sparse storage matrix representation for the data given in Example 4.4.

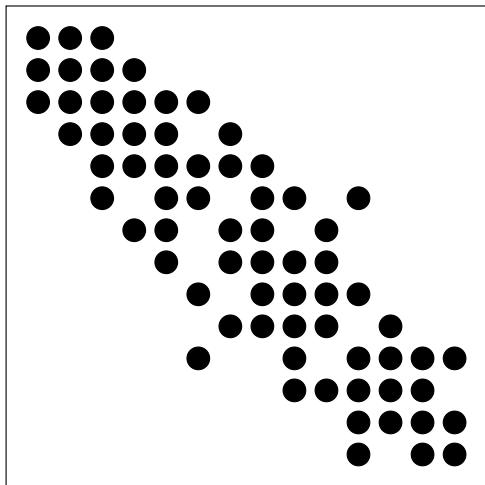
Solution 4.5 Figure 4.7 shows the sparse matrix visual representation of the matrix data given in Example 4.5. While it is not necessary to visualize the matrix explicitly to perform an LU factorization, it is often informative. ■

4.3 Ordering Schemes

Node ordering schemes are important in minimizing the number of multiplications and divisions required for both L and U triangularization and forward/backward substitution. A good ordering will result in the addition of few *fills* to the triangular factors during the LU factorization process. A *fill*



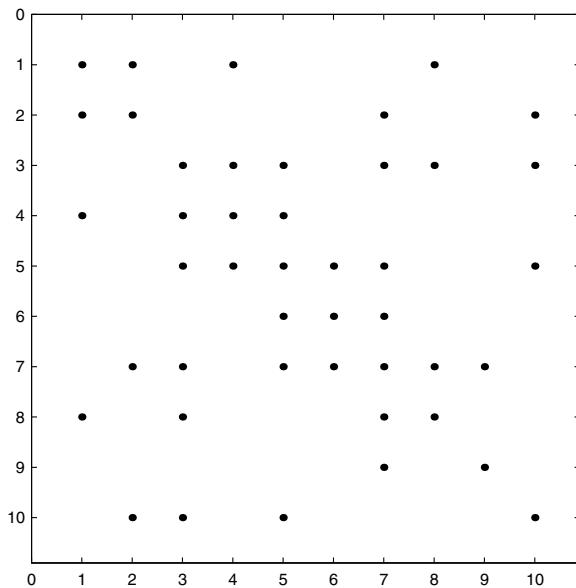
(a)



(b)

FIGURE 4.6

(a) A finite element grid model (b) The corresponding matrix

**FIGURE 4.7**

Matrix visualization for Example 4.5

is a nonzero element in the L or U matrix that was zero in the original A matrix. If A is a full matrix, $\alpha = \frac{n^3-n}{3}$ multiplications and divisions are required for the LU factorization process and $\beta = n^2$ multiplications and divisions are required for the forward/backward substitution process. The number of multiplications and divisions required can be substantially reduced in sparse matrix solutions if a proper node ordering is used.

Example 4.6

Determine the number of multiplications, divisions, and fills required for the solution of the system shown in Figure 4.8.

Solution 4.6 The LU factorization steps yield

$$q_{11} = a_{11}$$

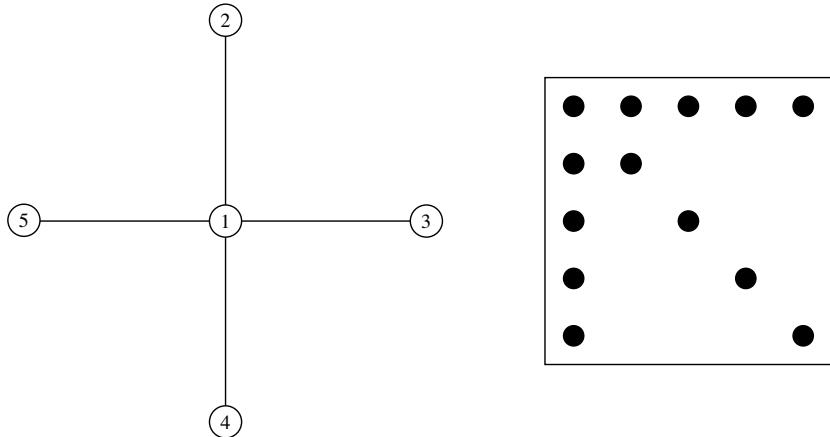
$$q_{21} = a_{21}$$

$$q_{31} = a_{31}$$

$$q_{41} = a_{41}$$

$$q_{51} = a_{51}$$

$$q_{12} = a_{12}/q_{11}$$

**FIGURE 4.8**

Graph and matrix for Example 4.6

$$q_{13} = a_{13}/q_{11}$$

$$q_{14} = a_{14}/q_{11}$$

$$q_{15} = a_{15}/q_{11}$$

$$q_{22} = a_{22} - q_{21}q_{12}$$

$$q_{32} = a_{32} - q_{31}q_{12}$$

$$q_{42} = a_{42} - q_{41}q_{12}$$

$$q_{52} = a_{52} - q_{51}q_{12}$$

$$q_{23} = (a_{23} - q_{21}q_{13})/q_{22}$$

$$q_{24} = (a_{24} - q_{21}q_{14})/q_{22}$$

$$q_{25} = (a_{25} - q_{21}q_{15})/q_{22}$$

$$q_{33} = a_{33} - q_{31}q_{13} - q_{32}q_{23}$$

$$q_{43} = a_{43} - q_{41}q_{13} - q_{42}q_{23}$$

$$q_{53} = a_{53} - q_{51}q_{13} - q_{52}q_{23}$$

$$q_{34} = (a_{34} - q_{31}q_{14} - q_{32}q_{24})/q_{33}$$

$$q_{35} = (a_{35} - q_{31}q_{15} - q_{32}q_{25})/q_{33}$$

$$q_{44} = a_{44} - q_{41}q_{14} - q_{42}q_{24} - q_{43}q_{34}$$

$$q_{54} = a_{54} - q_{51}q_{14} - q_{52}q_{24} - q_{53}q_{34}$$

$$q_{45} = (a_{45} - q_{41}q_{15} - q_{42}q_{25} - q_{43}q_{35}) / q_{44}$$

$$q_{55} = a_{55} - q_{51}q_{15} - q_{52}q_{25} - q_{53}q_{35} - q_{54}q_{45}$$

The multiplications and divisions required for the LU factorization are summarized by row and column.

row	column	multiplications	divisions	fills
1		0	0	
	1	0	4	
2		4	0	a_{32}, a_{42}, a_{52}
	2	3	3	a_{23}, a_{24}, a_{25}
3		6	0	a_{43}, a_{53}
	3	4	2	a_{34}, a_{35}
4		6	0	a_{54}
	4	3	1	a_{45}
5		4	0	

Therefore $\alpha = 40$ is the total number of multiplications and divisions in the LU factorization. The forward ($Ly = b$) and backward ($Ux = y$) substitution steps yield

$$y_1 = b_1 / q_{11}$$

$$y_2 = (b_2 - q_{21}y_1) / q_{22}$$

$$y_3 = (b_3 - q_{31}y_1 - q_{32}y_2) / q_{33}$$

$$y_4 = (b_4 - q_{41}y_1 - q_{42}y_2 - q_{43}y_3) / q_{44}$$

$$y_5 = (b_5 - q_{51}y_1 - q_{52}y_2 - q_{53}y_3 - q_{54}y_4) / q_{55}$$

$$x_5 = y_5$$

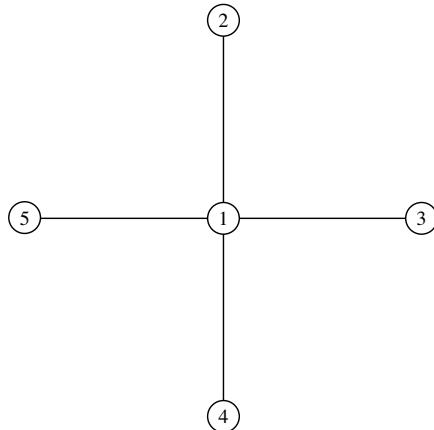
$$x_4 = y_4 - q_{45}x_5$$

$$x_3 = y_3 - q_{35}x_5 - q_{34}x_4$$

$$x_2 = y_2 - q_{25}x_5 - q_{24}x_4 - q_{23}x_3$$

$$x_1 = y_1 - q_{15}x_5 - q_{14}x_4 - q_{13}x_3 - q_{12}x_2$$

row	forward		backward	
	multiplications	divisions	multiplications	divisions
1	0	1	4	0
2	1	1	3	0
3	2	1	2	0
4	3	1	1	0
5	4	1	0	0

**FIGURE 4.9**

Graph for Example 4.6

Thus $\beta = 25$ is the total number of multiplications and divisions in the forward and backward substitution steps. The total number of multiplications and divisions for the solution of $Ax = b$ is $\alpha + \beta = 65$. ■

A fill occurs when a matrix element that was originally zero becomes nonzero during the factorization process. This can be visually simulated using a graphical approach. Consider the graph of Example 4.6 shown again in Figure 4.9.

In this numbering scheme, the row and column corresponding to node 1 is factorized first. This corresponds to the removal of node 1 from the graph. When node 1 is removed, all of the vertices to which it was connected must then be joined. Each edge added represents two fills in the Q matrix (q_{ij} and q_{ji}) since Q is symmetric. The graph after the removal of node 1 is shown in Figure 4.10. The dashed lines indicate that six fills will occur as a result: q_{23} , q_{24} , q_{25} , q_{34} , q_{35} , and q_{45} . These are the six fills that are also listed in the solution of the example.

Example 4.7

Determine the number of multiplications, divisions, and fills required for the solution of the system shown in Figure 4.11.

Solution 4.6 The LU factorization steps yield

$$q_{11} = a_{11}$$

$$q_{51} = a_{51}$$

$$q_{15} = a_{15}/q_{11}$$

$$q_{22} = a_{22}$$

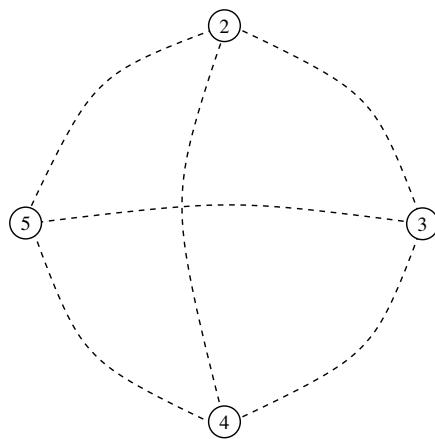


FIGURE 4.10
Resulting fills after removing node 1

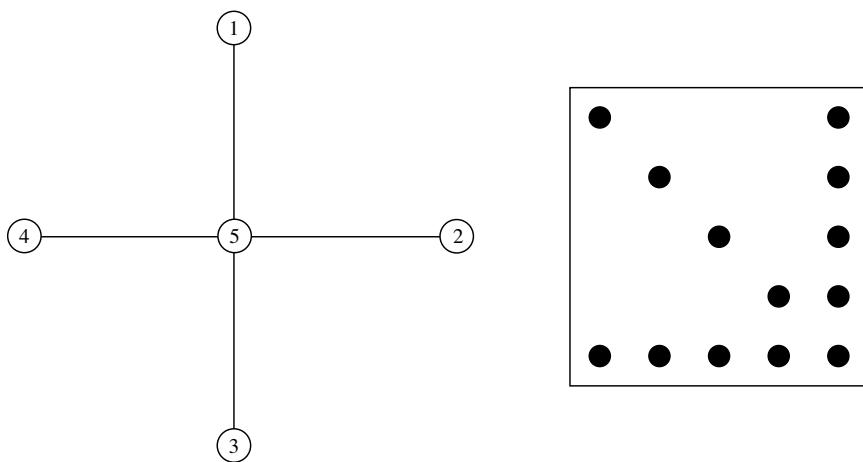


FIGURE 4.11
Graph and matrix for Example 4.7

$$q_{25} = a_{25}/q_{22}$$

$$q_{52} = a_{52}$$

$$q_{33} = a_{33}$$

$$q_{53} = a_{53}$$

$$q_{35} = a_{35}/q_{33}$$

$$q_{44} = a_{44}$$

$$q_{54} = a_{54}$$

$$q_{45} = a_{45}/q_{44}$$

$$q_{55} = a_{55} - q_{51}q_{15} - q_{52}q_{25} - q_{53}q_{35} - q_{54}q_{45}$$

The multiplications and divisions required for the LU factorization are summarized by row and column.

row	column	multiplications	divisions	fills
1		0	0	
	1	0	1	
2		0	0	
	2	0	1	
3		0	0	
	3	0	1	
4		0	0	
	4	0	1	
5		4	0	

Therefore $\alpha = 8$ is the total number of multiplications and divisions in the LU factorization. The forward ($Ly = b$) and backward ($Ux = y$) substitution steps yield

$$y_1 = b_1/q_{11}$$

$$y_2 = b_2/q_{22}$$

$$y_3 = b_3/q_{33}$$

$$y_4 = b_4/q_{44}$$

$$y_5 = (b_5 - q_{51}y_1 - q_{52}y_2 - q_{53}y_3 - q_{54}y_4) / q_{55}$$

$$x_5 = y_5$$

$$x_4 = y_4 - q_{45}x_5$$

$$x_3 = y_3 - q_{35}x_5$$

$$x_2 = y_2 - q_{25}x_5$$

$$x_1 = y_1 - q_{15}x_5$$

row	forward		backward	
	multiplications	divisions	multiplications	divisions
1	0	1	1	0
2	0	1	1	0
3	0	1	1	0
4	0	1	1	0
5	4	1	0	0

Thus $\beta = 13$ is the total number of multiplications and divisions in the forward and backward substitution steps. The total number of multiplications and divisions for the solution of $Ax = b$ is $\alpha + \beta = 21$. ■

Even though both original matrices had the same number of nonzero elements, there is a significant reduction in the number of multiplications and divisions by simply renumbering the vertices of the matrix graph. This is due, in part, to the number of fills that occurred during the LU factorization of the matrix. The Q matrix of Example 4.6 became full, whereas the Q matrix of Example 4.7 retained the same sparse structure as the original A matrix. From these two examples, it can be concluded that, although various node orders do not affect the accuracy of the linear solution, the ordering scheme greatly affects the time in which the solution is achieved. A good ordering scheme is one in which the resulting Q matrix has a structure similar to the original A matrix. This means that the number of fills is minimized. This objective forms the basis for a variety of ordering schemes. The problem of optimal ordering is an NP-complete problem [59], but several schemes have been developed that provide near-optimal results.

Example 4.8

Determine number of fills, α , and β for the matrix shown in Figure 4.12 as currently ordered. This is the same matrix as in Example 4.5.

Solution 4.8 The first step is to determine where the fills from LU factorization will occur. By observation, the fills will occur in the places designated by the \triangle in the matrix shown in Figure 4.13. From the figure, the number of fills is 24.

Rather than calculating the number of multiplications and divisions required for LU factorization and forward/backward substitution, there is a handy way of calculating α and β directly from the filled matrix.

$$\alpha = \sum_{i=1}^n (\text{nnz in column } i \text{ below } q_{ii} + 1) \times (\text{nnz in row } i \text{ to right of } q_{ii}) \quad (4.1)$$

$$\beta = \text{nnz of matrix } Q \quad (4.2)$$

Using Equations (4.1) and (4.2),

$$\begin{aligned} \alpha &= (3 \times 4) + (4 \times 5) + (5 \times 6) + (4 \times 5) + (4 \times 5) + (3 \times 4) \\ &\quad + (3 \times 4) + (2 \times 3) + (1 \times 2) + (0 \times 1) = 134 \end{aligned}$$

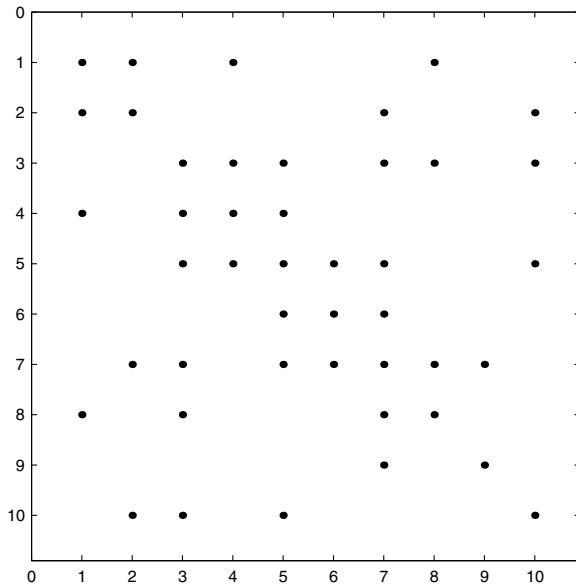


FIGURE 4.12
Matrix for Example 4.8

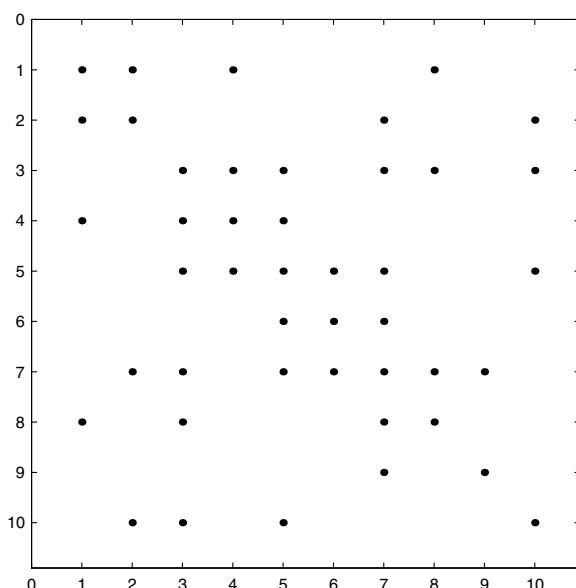


FIGURE 4.13
Matrix with fills for Example 4.8

and $\beta = nnz = 68$ for the Q matrix shown in Figure 4.13; thus $\alpha + \beta = 202$. ■

Even without an ordering scheme, the sparse matrix solution process yields over a 50% reduction in computation. One goal of an ordering scheme is to introduce the least number of fills in the factored matrix Q to minimize the number of multiplications and divisions α . A second goal is also to minimize β , which is the number of multiplications and divisions in the forward/backward substitution step. These dual objectives lead to several approaches to ordering.

Example 4.9

Update the sparse vectors of Example 4.4 to include the new fills.

Solution 4.9 The locations of the fills were calculated in Example 4.8. These fills are added as zero elements in the sparse storage vectors.

Recall that Example 4.8 indicated that β would be 68. Note that elements 45 to 68 correspond to the new fills, but since the LU factorization has not yet been conducted, the Value element for each of the fills is still zero. Note that the NIR and NIC columns have been updated throughout the vectors to accommodate the insertion of the new elements.

The FIR and FIC vectors have not changed. ■

<i>i</i>	NROW	NCOL	NIR	NIC	Value	<i>i</i>	NROW	NCOL	NIR	NIC	Value
1	8	8	65	66	-28	51	4	7	49	5	0
2	7	2	35	48	5	52	7	4	30	50	0
3	10	5	64	0	7	53	4	10	0	4	0
4	5	10	0	63	7	54	10	4	3	0	0
5	5	7	59	25	3	55	7	10	0	57	0
6	6	6	25	42	-33	56	10	7	58	0	0
7	6	5	6	30	10	57	8	10	0	67	0
8	1	8	0	47	8	58	10	8	68	0	0
9	5	5	28	7	-44	59	5	8	4	61	0
10	1	4	8	45	19	60	8	5	62	3	0
11	4	3	27	40	6	61	6	8	63	26	0
12	8	3	50	14	1	62	8	6	44	64	0
13	3	4	29	27	6	63	6	10	0	55	0
14	10	3	54	0	9	64	10	6	56	0	0
15	2	1	36	20	2	65	8	9	57	43	0
16	9	7	66	56	13	66	9	8	43	58	0
17	10	2	14	0	10	67	9	10	0	41	0
18	3	8	19	49	1	68	10	9	41	0	0
19	3	10	0	53	9						
20	4	1	46	22	19						
21	7	7	26	44	-68						
22	8	1	48	0	8						
23	2	10	0	19	10						
24	3	3	13	11	-40						
25	6	7	61	21	19						
26	7	8	38	1	15						
27	4	4	39	31	-38						
28	5	6	5	6	10						
29	3	5	33	39	11						
30	7	5	42	60	3						
31	5	4	9	52	9						
32	1	2	10	36	2						
33	3	7	18	51	9						
34	2	7	47	33	5						
35	7	3	52	12	9						
36	2	2	45	46	-21						
37	1	1	32	15	-33						
38	7	9	55	65	13						
39	4	5	51	9	9						
40	5	3	31	35	11						
41	10	10	0	0	-30						
42	7	6	21	62	19						
43	9	9	67	68	-17						
44	8	7	1	16	15						
45	2	4	34	13	0						
46	4	2	11	2	0						
47	2	8	23	18	0						
48	8	2	12	17	0						
49	4	8	53	59	0						
50	8	4	60	54	0						

Example 4.10

In the following LU factors, find the missing element (26) (which corresponds to $q(7,8)$).

<i>i</i>	NROW	NCOL	NIR	NIC	Value	<i>i</i>	NROW	NCOL	NIR	NIC	Value
1	8	8	65	66	-19.5512	35	7	3	52	12	9.0000
2	7	2	35	48	5.0000	36	2	2	45	46	-20.8788
3	10	5	64	0	10.2510	37	1	1	32	15	-33.0000
4	5	10	0	63	-0.2799	38	7	9	55	65	-0.2625
5	5	7	59	25	-0.1676	39	4	5	51	9	-0.4081
6	6	6	25	42	-30.2699	40	5	3	31	35	11.0000
7	6	5	6	30	10.0000	41	10	10	0	0	-16.9270
8	1	8	0	47	-0.2424	42	7	6	21	62	20.6759
9	5	5	28	7	-36.6288	43	9	9	67	68	-12.6365
10	1	4	8	45	-0.5758	44	8	7	1	16	16.4275
11	4	3	27	40	6.0000	45	2	4	34	13	-0.0552
12	8	3	50	14	1.0000	46	4	2	11	2	1.1515
13	3	4	29	27	-0.1500	47	2	8	23	18	-0.0232
14	10	3	54	0	9.0000	48	8	2	12	17	0.4848
15	2	1	36	20	2.0000	49	4	8	53	59	-0.1833
16	9	7	66	56	13.0000	50	8	4	60	54	4.7828
17	10	2	14	0	10.0000	51	4	7	49	5	-0.0623
18	3	8	19	49	-0.0250	52	7	4	30	50	1.6258
19	3	10	0	53	-0.2250	53	4	10	0	4	-0.0729
20	4	1	46	22	19.0000	54	10	4	3	0	1.9015
21	7	7	26	44	-49.5250	55	7	10	0	57	-0.1649
22	8	1	48	0	8.0000	56	10	7	58	0	8.1677
23	2	10	0	19	-0.4790	57	8	10	0	67	-0.2145
24	3	3	13	11	-40.0000	58	10	8	68	0	4.1944
25	6	7	61	21	-0.6830	59	5	8	4	61	-0.0608
26	7	8	38	1	*****	60	8	5	62	3	2.2268
27	4	4	39	31	-26.0971	61	6	8	63	26	-0.0201
28	5	6	5	6	-0.2730	62	8	6	44	64	0.6079
29	3	5	33	39	-0.2750	63	6	10	0	55	-0.0925
30	7	5	42	60	6.1385	64	10	6	56	0	2.7986
31	5	4	9	52	10.6500	65	8	9	57	43	-0.2206
32	1	2	10	36	-0.0606	66	9	8	43	58	4.3121
33	3	7	18	51	-0.2250	67	9	10	0	41	-0.2429
34	2	7	47	33	-0.2395	68	10	9	41	0	3.0691

Solution 4.10 The locations of the fills were calculated in Examples 4.8 and 4.9.

To begin the calculation of element (26), recall that row calculations are given by

$$q_{jk} = \frac{1}{q_{jj}} \left(a_{jk} - \sum_{i=1}^{j-1} q_{ji}q_{ik} \right) \text{ for } k = j+1, \dots, n \quad (4.3)$$

where in this case $j = 7$ and $k = 8$. Note that the product $q_{ji}q_{ik}$ will only be nonzero if both q_{ji} and q_{ik} are nonzero.

To determine which q_{ji} s are nonzero, the first-in-row pointer will begin at FIR(7)=2. Traversing the row, the nonzero elements in the row to the left of the diagonal are

$$\begin{array}{cccccc} (2) & \rightarrow & (35) & \rightarrow & (52) & \rightarrow (30) \rightarrow (42) \\ [2] & & [3] & & [4] & & [5] & & [6] \end{array}$$

The corresponding NCOL values are shown in the square brackets below each element. The diagonal element q_{77} (NROW=7, NCOL=7) is the next-in-row NIR(42)=21.

Similarly, to determine which q_{ik} s are nonzero, the first-in-column pointer will begin at FIC(8)=8. Traversing the column, the nonzero elements in the column above the diagonal are

$$\begin{array}{c} (8) \quad [1] \\ \downarrow \\ (47) \quad [2] \\ \downarrow \\ (18) \quad [3] \\ \downarrow \\ (49) \quad [4] \\ \downarrow \\ (59) \quad [5] \\ \downarrow \\ (61) \quad [6] \end{array}$$

The corresponding NROW values are shown in the square brackets to the right of each element.

Matching up the corresponding elements:

$$q_{78} = q_{78} - (q_{72}q_{28} + q_{73}q_{38} + q_{74}q_{48} + q_{75}q_{58} + q_{76}q_{68}) / q_{77}$$

(elements)

$$(26) = ((26) - ((2)(47) + (35)(18) + (52)(49) + (30)(59) + (42)(61))) / (21)$$

$$\begin{aligned}
 & (\text{Values (elements)}) \\
 & = (15 - ((5.0000)(-0.0232) + (9.0000)(-0.0250) + (1.6258)(-0.1833) + \\
 & \quad (6.1385)(-0.0608) + (20.6759)(-0.0201))) / (-49.5250) \\
 & = -0.3317
 \end{aligned}$$

Note that the initial value of element (26) is a_{78} . ■

4.3.1 Scheme 0

From Examples 4.6 and 4.7, it can be generalized that a better ordering is achieved if the nodes are ordered into a lower-right pointing “arrow” matrix. One rapid method of achieving this effect is to number the nodes according to their degree, where the degree of a node is defined as the number of edges connected to it. In this scheme, the nodes are ordered from lowest degree to highest degree.

Scheme 0

1. Calculate the degree of all vertices.
2. Choose the node with the lowest degree. Place in the ordering scheme.
3. In case of a tie, choose the node with the lowest natural ordering.
4. Return to step 2.

Example 4.11

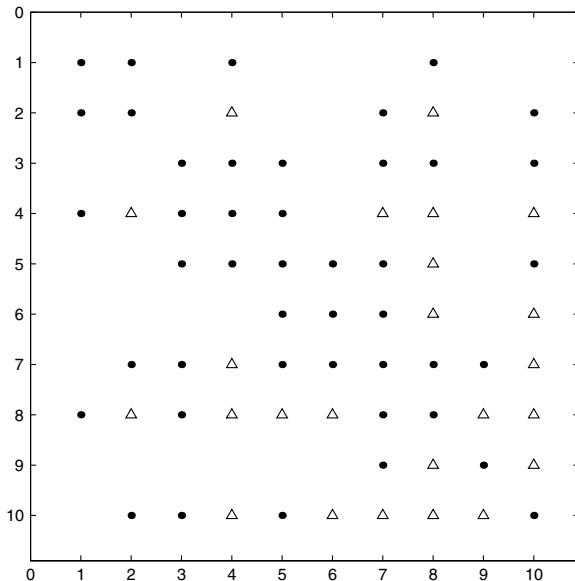
Using Scheme 0, reorder the matrix of Example 4.8. Calculate α , β , and the number of fills for this ordering.

Solution 4.11 The degrees of each of the nodes are given below:

node	degree
1	3
2	3
3	5
4	3
5	5
6	2
7	6
8	3
9	1
10	3

Applying Scheme 0, the new ordering is

$$\text{Ordering } 0 = [9 \ 6 \ 1 \ 2 \ 4 \ 8 \ 10 \ 3 \ 5 \ 7]$$

**FIGURE 4.14**

Matrix with fills for Example 4.11

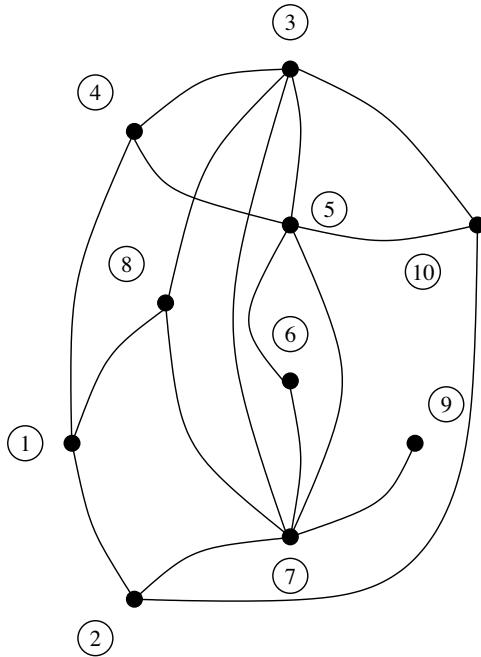
Reordering the matrix of Example 4.8 to reflect this ordering yields the matrix (with fills) shown in Figure 4.14. Note how the nonzero elements begin to resemble the desired lower-right pointing arrow. The Scheme 0 ordering results in 16 fills as compared to 24 with the original ordering. From the matrix and Equations (4.1) and (4.2), $\alpha = 110$ and $\beta = 60$; thus $\alpha + \beta = 170$, which is a considerable reduction over the original $\alpha + \beta = 202$. ■

4.3.2 Scheme I

Scheme 0 offers simplicity and speed of generation, but does not directly take into account the effect of fills on the ordering procedure. To do this, the effect of eliminating the nodes as they are ordered must be taken into account. This modification is given in Scheme I.

Scheme I

1. Calculate the degree of all vertices.
2. Choose the node with the lowest degree. Place in the ordering scheme. Eliminate it and update degrees accordingly.
3. In case of a tie, choose the node with the lowest natural ordering.
4. Return to step 1.

**FIGURE 4.15**

Graph of the matrix in Figure 4.12

Scheme I is also known by many names, including the Markowitz algorithm [35], the Tinney I algorithm [55], or most generally as the minimum degree algorithm.

Example 4.12

Using Scheme I, reorder the matrix of Example 4.8. Calculate α , β , and the number of fills for this ordering.

Solution 4.12 The ordering for Scheme I takes into account the effect of fills on the ordering as nodes are placed in the ordering scheme and eliminated. This algorithm is best visualized using the graphical representation of the matrix. The graph of the original unordered matrix of Figure 4.12 is shown in Figure 4.15.

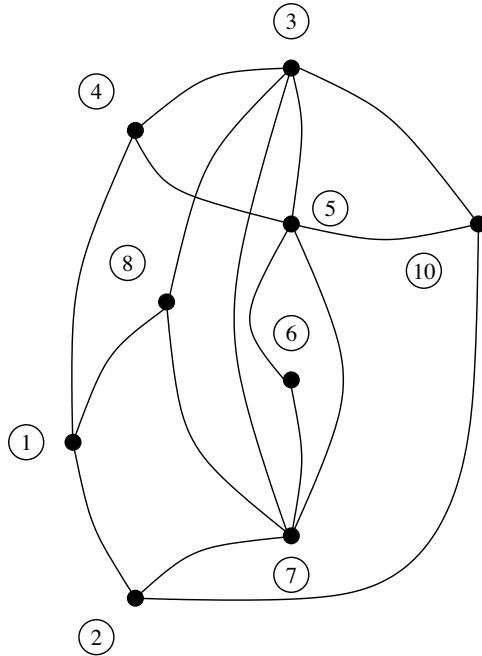


FIGURE 4.16
Updated graph with the removal of node 9

The degrees of each of the nodes are given below:

node	degree
1	3
2	3
3	5
4	3
5	5
6	2
7	6
8	3
9	1
10	3

From the degrees, the node with the lowest degree is ordered first. Node 9 has the lowest degree, with only one connection. Its elimination does not cause any fills. The updated graph is shown in Figure 4.16.

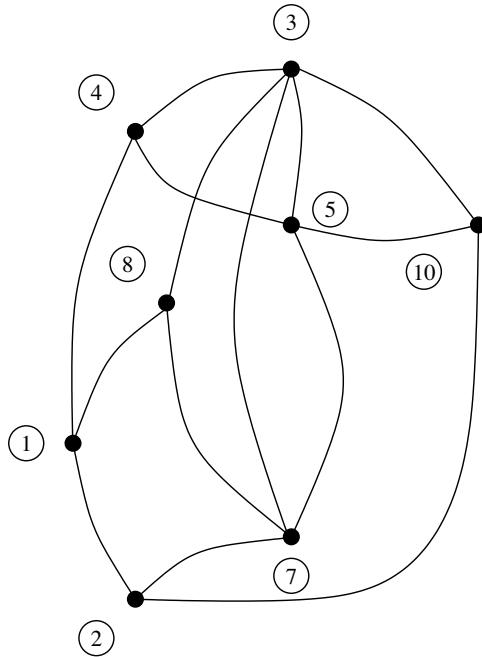


FIGURE 4.17
Updated graph with the removal of node 6

The updated degree of each of the nodes is given below:

node	degree
1	3
2	3
3	5
4	3
5	5
6	2
7	5
8	3
10	3

Node 7 now has one less degree. Applying the Scheme I algorithm again indicates that the next node to be chosen is node 6, with a degree of 2. Node 6 is connected to both node 5 and node 7. Since there is a preexisting connection between these nodes, the elimination of node 6 does not create a fill between nodes 5 and 6. The elimination of node 6 is shown in Figure 4.17.

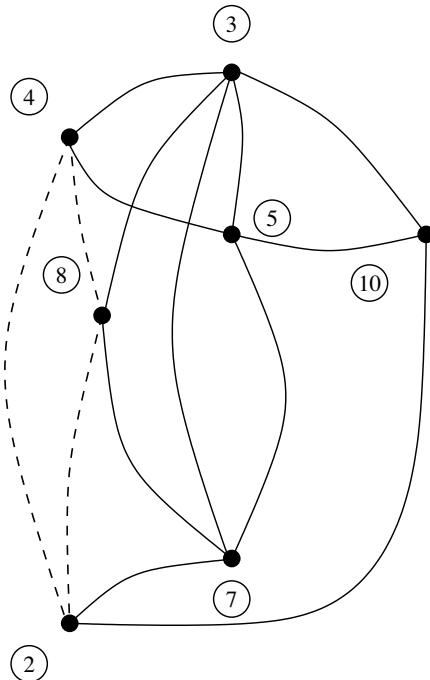
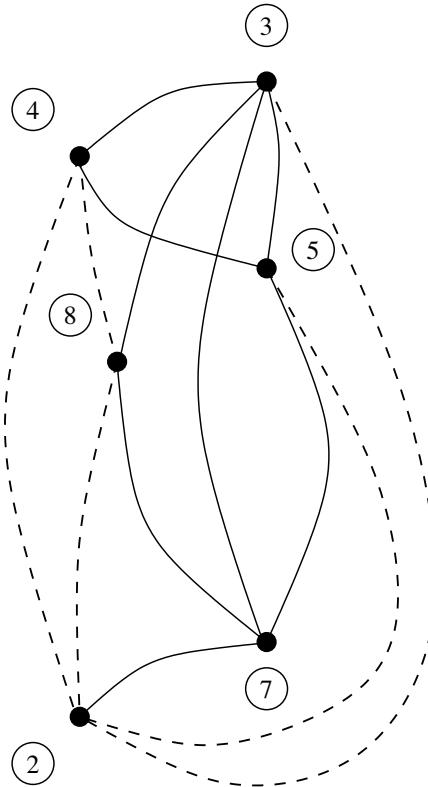


FIGURE 4.18
Updated graph with the removal of node 1

The new node degrees are

node	degree
1	3
2	3
3	5
4	3
5	4
7	4
8	3
10	3

As a result of the elimination of node 6, the degrees of nodes 5 and 7 decrease by one. Applying the Scheme I algorithm again indicates that the nodes with the lowest degrees are nodes [1 2 4 8 10]. Since there is a tie between these nodes, the node with the lowest natural ordering, node 1, is chosen and eliminated. Node 1 is connected to nodes 2, 4, and 8. None of these nodes is connected; therefore, the elimination of node 1 creates three fills: 4–8, 4–2, and 2–8. These fills are shown with the dashed edges in Figure 4.18.

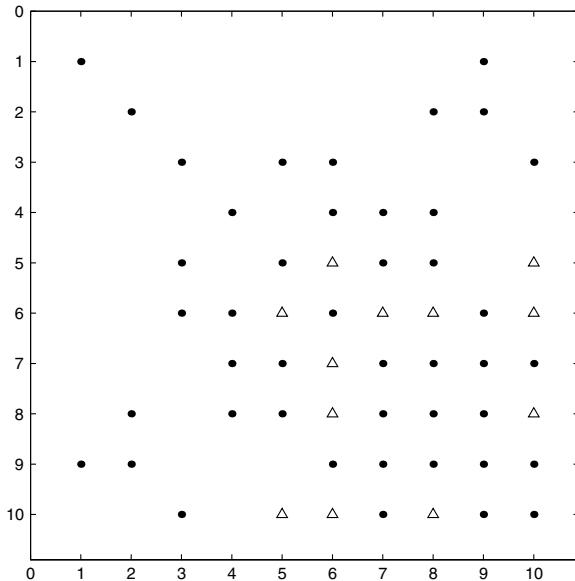
**FIGURE 4.19**

Updated graph with the removal of node 10

The new node degrees after the removal of node 1 are

node	degree
2	4
3	5
4	4
5	5
7	5
8	4
10	3

The addition of the three fills increased the degrees of nodes 2, 4, and 8. Applying the Scheme I algorithm again indicates that the node with the lowest degree is node 10. There is no tie in degree this time. Node 10 is chosen and eliminated. The elimination of node 10 creates two fills between nodes 2–5 and 2–3. These fills are shown with the dashed edges in Figure 4.19.

**FIGURE 4.20**

Matrix with fills for Example 4.12

Continuing to apply the Scheme I algorithm successively until all nodes have been chosen and eliminated yields the following final ordering:

$$\text{Ordering I} = [9 \ 6 \ 1 \ 10 \ 4 \ 2 \ 3 \ 5 \ 7 \ 8]$$

Reordering the matrix of Example 4.12 to reflect this ordering yields the matrix (with fills) shown in Figure 4.20. Note how the nonzero elements continue to resemble the desired lower-right pointing arrow. The Scheme I ordering results in 12 fills as compared to 24 with the original ordering, and 16 with Scheme 0. From the matrix and Equations (4.1) and (4.2), $\alpha = 92$ and $\beta = 56$; thus $\alpha + \beta = 148$, which is a considerable reduction over the original $\alpha + \beta = 202$ and the $\alpha + \beta = 170$ of Scheme 0. ■

4.3.3 Scheme II

Scheme 0 offers a rapid way to order the nodes to give a quick “once-over” and obtain a reasonable ordering. It requires little computation beyond calculating the degrees of each node of the matrix. Scheme I takes this approach one step further. It still relies on the minimum-degree approach, but it includes a simulation of the LU factorization process to update the node degrees at each step of the factorization. One further improvement to this approach is

to develop a scheme that endeavors to minimize the number of fills at each step of the factorization. Thus, at each step, each elimination alternative is considered and the number of resulting fills is calculated. This scheme is also known as the Berry algorithm and the Tinney II algorithm and is summarized below:

Scheme II

1. For each node, calculate the number of fills that would result from its elimination.
2. Choose the node with the lowest number of fills.
3. In case of a tie, choose the node with the lowest degree.
4. In case of a tie, choose the node with the lowest natural ordering.
5. Place the node in the ordering scheme. Eliminate it and update fills and degrees accordingly.
6. Return to step 1.

Example 4.13

Using Scheme II, reorder the matrix of Example 4.8. Calculate α , β , and the number of fills for this ordering.

Solution 4.13 The ordering for Scheme II takes into account the effect of fills on the ordering as nodes are placed in the ordering scheme and eliminated. The degrees and resulting fills are given below:

node	degree	fills if eliminated	edges created
1	3	3	2–4, 2–8, 4–8
2	3	3	1–7, 1–10, 7–10
3	5	6	4–7, 4–8, 4–10, 5–8, 7–10, 8–10
4	3	2	1–3, 1–5
5	5	6	3–6, 4–6, 4–7, 4–10, 6–10, 7–10
6	2	0	none
7	6	12	2–3, 2–5, 2–6, 2–8, 2–9, 3–6, 3–9, 5–8, 5–9, 6–8, 6–9, 8–9
8	3	2	1–3, 1–7
9	1	0	none
10	3	2	2–3, 2–5

From this list, the elimination of nodes 6 or 9 will not result in any additional edges or fills. Since there is a tie, the node with the lowest degree is chosen.

Thus node 9 is chosen and eliminated. The number of fills and degrees is updated to apply the Scheme II algorithm again.

node	degree	fills if eliminated	edges created
1	3	3	2–4, 2–8, 4–8
2	3	3	1–7, 1–10, 7–10
3	5	6	4–7, 4–8, 4–10, 5–8, 7–10, 8–10
4	3	2	1–3, 1–5
5	5	6	3–6, 4–6, 4–7, 4–10, 6–10, 7–10
6	2	0	none
7	5	7	2–3, 2–5, 2–6, 2–8, 3–6, 5–8, 6–8
8	3	2	1–3, 1–7
10	3	2	2–3, 2–5

The next node to be eliminated is node 6 because it creates the fewest fills if eliminated. This node is therefore chosen and eliminated. The number of fills and degrees is again updated.

node	degree	fills if eliminated	edges created
1	3	3	2–4, 2–8, 4–8
2	3	3	1–7, 1–10, 7–10
3	5	6	4–7, 4–8, 4–10, 5–8, 7–10, 8–10
4	3	2	1–3, 1–5
5	5	6	3–6, 4–6, 4–7, 4–10, 6–10, 7–10
7	5	7	2–3, 2–5, 2–6, 2–8, 3–6, 5–8, 6–8
8	3	2	1–3, 1–7
10	3	2	2–3, 2–5

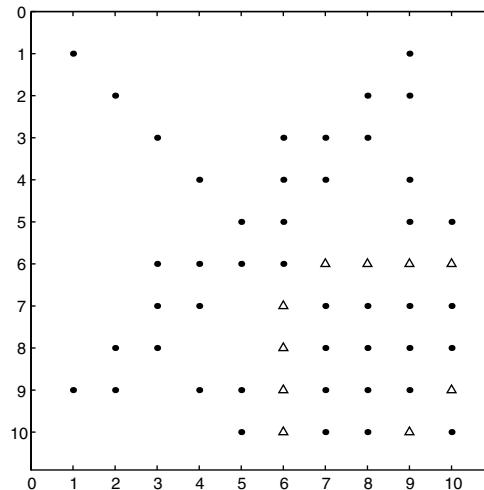
The two nodes that create the fewest fills are nodes 4 and 8. Both nodes have the same number of degrees; therefore, the node with the lowest natural ordering, node 4, is chosen and eliminated.

The Scheme II algorithm continues until all nodes have been added to the ordering scheme and subsequently eliminated. Scheme II results in the following ordering:

$$\text{Ordering II} = [9 \ 6 \ 4 \ 8 \ 2 \ 1 \ 3 \ 5 \ 7 \ 10]$$

Reordering the matrix of Example 4.8 to reflect the ordering of Scheme II yields the ordering with fills shown in Figure 4.21. This ordering yields only 10 fills, leading to $\alpha = 84$, $\beta = 54$, and $\alpha + \beta = 138$. This represents a computational effort of only 68% of the original unordered system. ■

Scheme I endeavors to reduce the number of multiplications and divisions in the LU factorization process, whereas Scheme II focuses on reducing the multiplications and divisions in the forward/backward substitution process. Scheme 0 offers simplicity and speed of generation, but the performance improvement of Scheme I offsets the additional algorithm complexity [55].

**FIGURE 4.21**

Matrix with fills for Example 4.13

Scheme II, however, frequently does not offer enough of an improvement to merit implementation. The decision of which scheme to implement is problem dependent and is best left up to the user.

4.3.4 Other Schemes

Modifications to these algorithms have been introduced to reduce computational requirements. These modifications are summarized below [19]. The first modification to the minimum-degree algorithm is the use of mass elimination, inspired by the concept of indistinguishable nodes [18]. This modification allows a subset of nodes to be eliminated at one time. If nodes x and y satisfy

$$\text{Adj}(y) \cup \{y\} = \text{Adj}(x) \cup \{x\} \quad (4.4)$$

where $\text{Adj}(y)$ indicates the set of nodes adjacent to y , then nodes x and y are said to be indistinguishable and can be numbered consecutively in the ordering. This also reduces the number of nodes to be considered in an ordering, since only a representative node from each set of indistinguishable nodes needs to be considered. This accelerates the degree update step of the minimum-degree algorithm, which is typically the most computationally intensive step. Using mass elimination, the degree update is required only for the representative nodes.

The idea of incomplete degree update allows avoiding degree update for nodes that are known not to be minimum degree. Between two nodes u and

v , node v is said to be outmatched by u if [12]

$$\text{Adj}(u) \cup \{u\} \subseteq \text{Adj}(v) \cup \{v\} \quad (4.5)$$

Thus, if a node v becomes outmatched by u in the elimination process, the node u can be eliminated before v in the minimum-degree ordering algorithm. From this, it follows that it is not necessary to update the degree of v until node u has been eliminated. This further reduces the time-consuming degree update steps.

Another modification to the minimum-degree algorithm is one in which all possible nodes of minimum degree are eliminated before the degree update step. At a given step in the elimination process, the elimination of node y does not change the structure of the remaining nodes not in $\text{Adj}(y)$. The multiple-minimum-degree (MMD) algorithm delays degree update of the nodes in $\text{Adj}(y)$ and chooses another node with the same degree as y to eliminate. This process continues until there are no more nodes left with the same degree as y . This algorithm was found to perform as well as the minimum-degree algorithm regarding the number of fills introduced [33]. In addition, it was found that the MMD algorithm performed faster. This was attributed to the identification of indistinguishable and outmatched nodes earlier in the algorithm, as well as the reduced number of degree updates.

Ties often occur for a given criteria (degrees or fills) in an ordering algorithm. The tie breaker often falls back on the natural ordering of the original matrix. It has been recognized that the natural ordering greatly affects the factorization in terms of number of fills and computation time. Thus it is often preferable to use a rapid “preconditioning” ordering before applying the ordering algorithm. Scheme 0 offers one such preordering, but to date no consistent optimum method for preordering has been developed that works well for all types of problems.

4.4 Power System Applications

Large sparse matrices occur frequently in power system applications, including state estimation, power flow analysis, and transient and dynamic stability simulations. Computational efficiency of these applications depends heavily on their formulation and the use of sparse matrix techniques. To better understand the impact of sparsity on power system problems, consider the power flow Jacobian of the IEEE 118 bus system shown in Figure 4.22.

The Jacobian of this system has 1051 nonzero elements and has the structure shown in Figure 4.23(a). Note the dominance of the main diagonal and then the two subdiagonals which result from the $\frac{\partial \Delta Q}{\partial \delta}$ and $\frac{\partial \Delta P}{\partial V}$ sub-Jacobians. The LU factorization of this Jacobian yields the structure shown in Figure

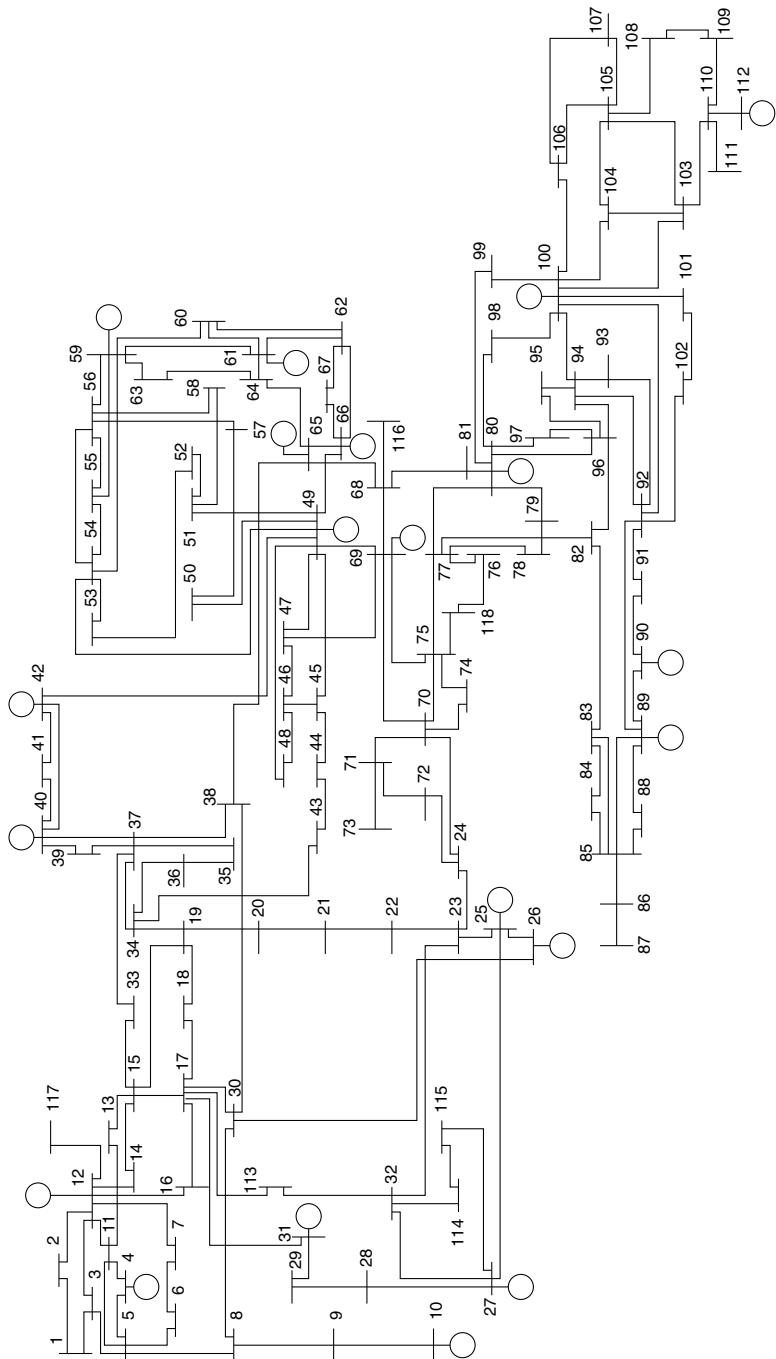


FIGURE 4.22
IEEE 118 bus system

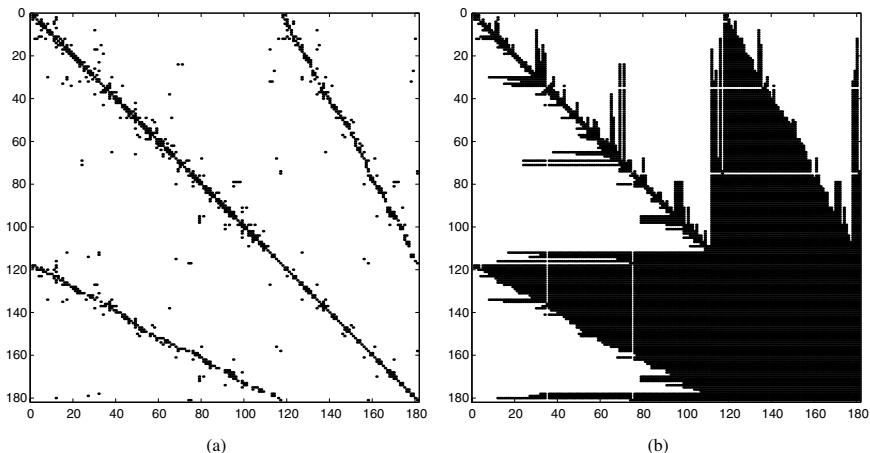


FIGURE 4.23
IEEE 118 bus system (a) Jacobian and (b) LU factors

4.23(b). This matrix has 14849 nonzero elements. Notice that the two subdiagonals have created a large number of fills extending between them and the main diagonal.

Figure 4.24(a) shows the structure of the Jacobian reordered according to Scheme 0. In this reordering, the presence of the subdiagonals is gone. The LU factorization of the Scheme 0 reordered Jacobian yields the structure shown in Figure 4.24(b). This matrix has only 1869 nonzero elements, which is almost an order of magnitude reduction from the nonordered Jacobian.

Figure 4.25(a) shows the structure of the power flow Jacobian reordered according to Scheme I. Note how the elements are gradually pulling into the main diagonal, which leads to a decrease in the number of fills. The LU factorization of the Scheme I ordering is shown in Figure 4.25(b), which has 1455 nonzero elements.

Finally, Figure 4.26(a) shows the structure of the Scheme II reordered Jacobian which yields the LU factorization in Figure 4.26(b). This ordering yields only 1421 nonzero elements, which is more than a full order of magnitude reduction. The LU factorization solution time for a sparse matrix is on the order of n^2 multiplications and divisions. The nonreordered power flow solution would require on the order of 220.5×10^6 multiplications and divisions per iteration, whereas the Scheme II reordered power flow solution would require only 2.02×10^6 multiplications and divisions. Thus the solution of the reordered system is over 100 times faster than the original system! When the solution time is multiplied by the number of iterations in a Newton–Raphson power flow or by the number of time steps in a time-domain integration, it would be computationally foolhardy not to use a reordering scheme.

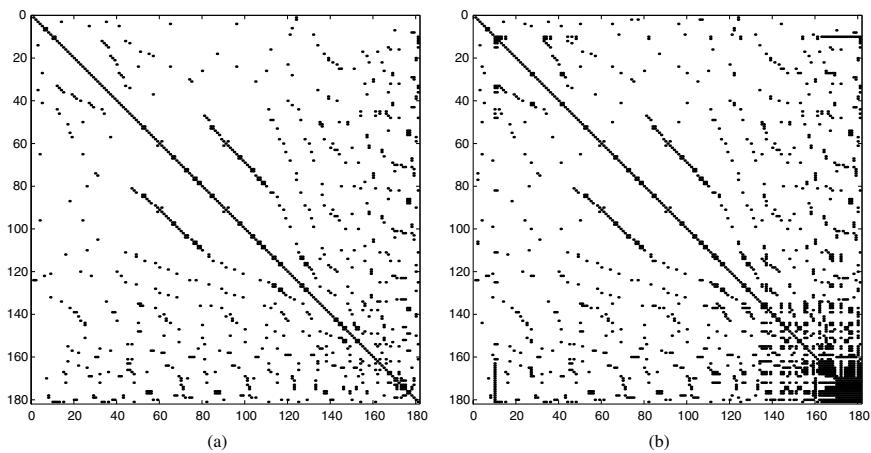


FIGURE 4.24
IEEE 118 bus system Scheme 0 (a) Jacobian and (b) LU factors

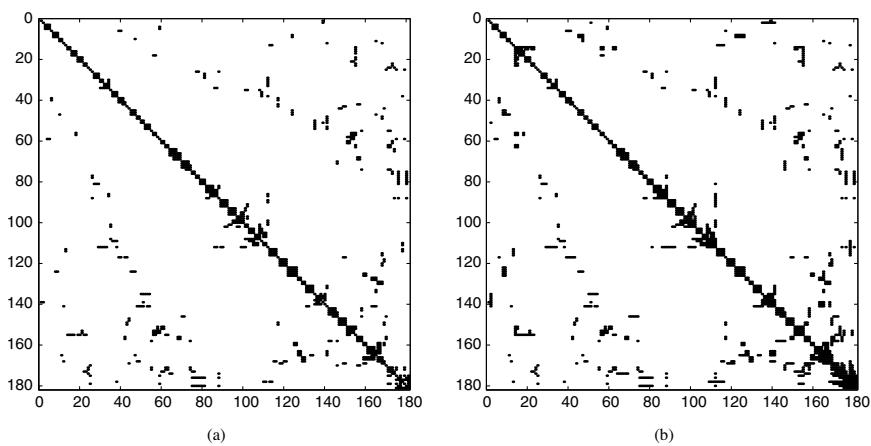


FIGURE 4.25
IEEE 118 bus system Scheme I (a) Jacobian and (b) LU factors

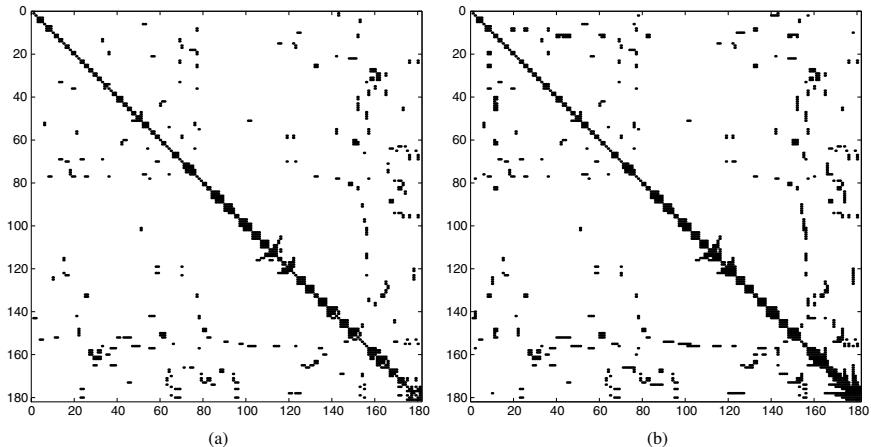


FIGURE 4.26
IEEE 118 bus system Scheme II (a) Jacobian and (b) LU factors

4.5 Problems

1. Verify Equations (4.1) and (4.2) for calculating α and β .
 2. Let A and B be two sparse (square) matrices of the same dimension. How can the graph of $C = A + B$ be characterized with respect to the graphs of A and B ?
 3. Consider the matrix

$$A = \begin{bmatrix} * & * & & * \\ * & * & * & \\ & * & * & \\ & & * & * \\ * & & * & * \\ * & & * & * \end{bmatrix}$$

- (a) Draw the graphical representation of A . How many multiplications and divisions will the LU factorization of A require?

(b) Reorder the matrix using the new node numbering $\phi = [1, 3, 4, 2, 5, 6]$. Draw the graphical representation of the reordered matrix. How many multiplications and divisions will the LU factorization of the reordered matrix require?

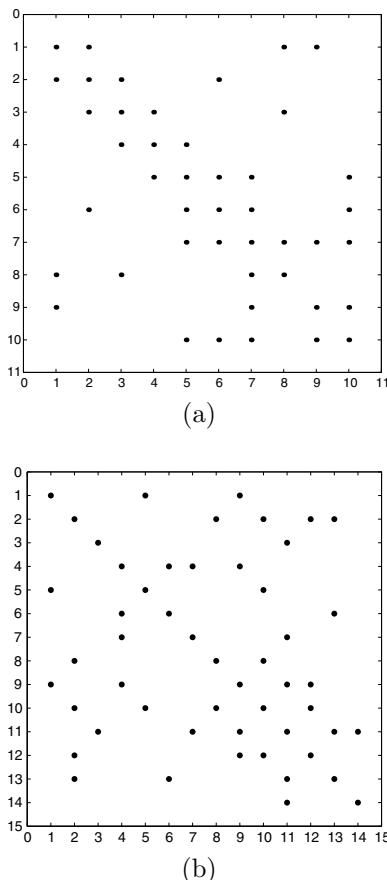


FIGURE 4.27
Sparse Test Systems

4. For the matrices shown in Figure 4.27
- Using the given ordering, compute $\alpha + \beta$.
 - Reorder the nodes in the network using Scheme 0 discussed in class. Compute $\alpha + \beta$ for this ordering.
 - Reorder the nodes in the network using Scheme I. Compute $\alpha + \beta$ for this ordering.
 - Reorder the nodes in the network using Scheme II. Compute $\alpha + \beta$ for this ordering.

5. For the sparse A given below:

NROW	NCOL	NIR	NIC	Value	FIR	FIC
1	8	4	16	0	5.8824	27
2	9	9	0	0	-30.2732	34
3	10	7	9	0	21.2766	17
4	7	10	0	5	21.2766	33
5	8	10	0	21	10.0000	25
6	7	9	4	7	11.7647	26
7	8	9	5	2	18.5185	19
8	5	4	10	1	7.1429	13
9	10	8	21	0	10.0000	24
10	5	5	22	26	-24.5829	3
11	4	5	32	10	7.1429	
12	3	5	18	11	5.0000	
13	8	3	1	0	11.1111	
14	1	4	0	29	13.3333	
15	6	6	20	16	-23.8376	
16	8	6	28	0	11.3636	
17	3	2	31	0	8.3333	
18	3	8	0	32	11.1111	
19	7	7	6	24	-33.0413	
20	6	8	0	28	11.3636	
21	10	10	0	0	-31.2566	
22	5	6	0	15	12.5000	
23	9	8	2	9	18.5185	
24	9	7	23	3	11.7647	
25	5	3	8	13	5.0000	
26	6	5	15	0	12.5000	
27	1	1	14	33	-13.3258	
28	8	8	7	23	-56.8046	
29	4	4	11	8	-26.3260	
30	2	3	0	31	8.3333	
31	3	3	12	25	-24.3794	
32	4	8	0	20	5.8824	
33	4	1	29	0	13.3333	
34	2	2	30	17	-8.3183	

Complete the LU factors in the sparse matrix given below. Do *not* explicitly create the matrix A . Show fully how you arrived at your answer.

	NROW	NCOL	NIR	NIC	Value	FIR	FIC
1	8	4	36	0	5.8824	27	27
2	9	9	37	38	-13.9405	34	34
3	10	7	9	0	21.2766	17	30
4	7	10	0	5	-0.6439	33	14
5	8	10	0	37	-0.3541	25	12
6	7	9	4	7	-0.3561	26	22
7	8	9	5	2	-0.6558	19	19
8	5	4	10	1	7.1429	13	18
9	10	8	38	0	10.0000	24	6
10	5	5	22	26	-19.0943	3	4
11	4	5	32	10	-0.5501		
12	3	5	18	11	-0.3119		
13	8	3	1	0	11.1111		
14	1	4	0	29	-1.0006		
15	6	6	20	16	-15.6545		
16	8	6	28	0	15.7506		
17	3	2	31	0	8.3333		
18	3	8	0	32	-0.6931		
19	7	7	6	24	-33.0413		
20	6	8	0	28	-1.0061		
21	10	10	0	0	0.3143		
22	5	6	[]	[]	[]		
23	9	8	2	9	18.5185		
24	9	7	23	3	11.7647		
25	5	3	8	13	5.0000		
26	6	5	[]	[]	[]		
27	1	1	14	33	-13.3258		
28	8	8	7	23	-28.2397		
29	4	4	11	8	-12.9852		
30	2	3	0	31	-1.0018		
31	3	3	12	25	-16.0311		
32	4	8	0	35	-0.4530		
33	4	1	29	0	13.3333		
34	2	2	30	17	-8.3183		
35	5	8	[]	[]	[]		
36	8	5	[]	[]	[]		
37	9	10	[]	[]	[]		
38	10	9	[]	[]	[]		

6. Write a subroutine *sparmat* for sparse matrix storage that will

- Read in data line by line in the format

$$i \ j \ a_{ij}$$

where the end of the data is signified by a 0 in the first column.

- Sequentially build the vectors FIR, FIC, NIR, NIC, NROW, NCOL, and Value as defined in class. **Do not explicitly create the matrix A .**

7. Write a subroutine *sparvec* for sparse vector storage that will

- Read in data line by line in the format

$$i \ b_i$$

where the end of the data is signified by a 0 in the first column.

- Sequentially build the vectors index, next, and Value. **Do not explicitly create the vector b .**

8. For the data given below, use *sparmat* and *sparvec* to create the sparse storage vectors.

<i>A</i> matrix			<i>b</i> vector	
<i>i</i>	<i>j</i>	<i>a_{ij}</i>	<i>i</i>	<i>b_i</i>
7	10	2.0	2	5
2	6	1.5	9	2
9	1	4.7	3	-1
5	5	-18.5		
8	7	2.8		
1	1	-15.0		
4	3	3.8		
6	7	6.1		
8	3	3.3		
5	7	4.4		
10	6	2.5		
6	5	1.1		
3	2	5.2		
7	8	2.9		
9	9	-12.1		
3	4	3.0		
7	6	5.6		
10	9	4.7		
8	8	-10.8		
1	9	4.5		
7	5	3.9		
5	6	7.2		
9	10	4.9		
5	4	0.8		
8	1	3.4		
5	10	4.5		
2	3	5.0		
6	6	-9.8		
7	9	1.8		
4	5	0.7		
7	7	-21.2		
1	2	4.4		
10	5	5.4		
3	8	3.1		
9	7	1.6		
4	4	-5.1		
6	10	2.7		
10	10	-16.9		
2	1	4.7		
3	3	-17.7		
1	8	3.5		
10	7	2.1		
2	2	-13.0		
6	2	1.2		

9. Write a subroutine *sparLU* that modifies your LU factorization routine to incorporate the sparse storage vectors of Problem 5 and apply it to the data of Problem 7 to compute the sparse LU factors (in sparse vector form).
10. Write a subroutine *sparsub* that modifies your forward/backward substitution routine *sub* to incorporate the sparse storage vectors of Problem 2 and apply it to the data of Problem 7 to solve the sparse linear system

$$Ax = b$$

11. Write a subroutine *scheme0* that will input the sparse vectors FIR, FIC, NIR, NIC, NROW, NCOL, and Value and will output the same vectors reordered according to Scheme 0, and calculate $\alpha + \beta$.
12. Write a subroutine *scheme1* that will input the sparse vectors FIR, FIC, NIR, NIC, NROW, NCOL, and Value and will output the same vectors reordered according to Scheme I, and calculate $\alpha + \beta$.
13. Write a subroutine *scheme2* that will input the sparse vectors FIR, FIC, NIR, NIC, NROW, NCOL, and Value and will output the same vectors reordered according to Scheme II, and calculate $\alpha + \beta$.

5

Numerical Integration

Dynamic systems may frequently be modeled by systems of *ordinary differential equations* (or ODEs) of the form

$$\dot{x}(t) = f(x, t) \quad x(t_0) = x_0 \quad (5.1)$$

where $x(t) \in R^n$ is a time-varying function that depends on the initial condition x_0 . Such problems are often referred to as “initial value problems.” A system of nonlinear differential equations cannot typically be solved analytically. In other words, a closed form expression for $x(t)$ cannot be found directly from Equation (5.1), but rather must be solved numerically.

In the numerical solution of Equation (5.1), a sequence of points x_0, x_1, x_2, \dots , is computed that approximates the true solution at a set of time points t_0, t_1, t_2, \dots . The time interval between adjacent time points is called the *time step* and an integration algorithm advances the numerical solution by one time step with each application. The time step $h_{n+1} = t_{n+1} - t_n$ may be constant for all time intervals over the entire integration interval $t \in [t_0, t_N]$ or may vary at each step.

The basic integration algorithm advances the solution from t_n to t_{n+1} with integration step size h_{n+1} based on a calculation that involves previously computed values x_n, x_{n-1}, \dots and functions $f(x_n, t_n), f(x_{n-1}, t_{n-1}), \dots$. Each practical integration algorithm must satisfy certain criteria concerning

1. numerical accuracy,
2. numerical stability, and
3. numerical efficiency.

Numerical accuracy ensures that the numerical error incurred at each step of the integration remains bounded. The global error of an integration error is the total error accrued over a given time interval. The global error at time t_n is given by

$$\text{global error} = \|x(t_n) - x_n\|$$

where $x(t_n)$ is the exact solution to Equation (5.1) at time t_n and x_n is the approximated solution. Of course, it is impossible to determine the global error exactly if the solution $x(t)$ is not known analytically, but it is possible to establish bounds on the error incurred at each step of the integration method.

The numerical stability of an integration algorithm implies that errors incurred at each step do not propagate to future times. Numerical efficiency is a function of the amount of computation required at each time step and the size of the steps between adjacent time intervals. Each of these criteria will be discussed in greater detail later in this chapter after an introduction to several different forms of integration algorithms.

5.1 One-Step Methods

The basic form of an integration algorithm is one that advances the solution from x_n to x_{n+1} using only the information currently available. This type of solution is called a *one-step* method, in that only information from one step of the integration algorithm is used. The family of one-step methods has the advantage of conserving memory, since only the previous solution must be retained. Several well-known methods fall into this category.

5.1.1 Taylor Series-Based Methods

One important class of integration methods is derived from using the Taylor series expansion of Equation (5.1). Let $\hat{x}(t)$ denote the exact solution to Equation (5.1). Expanding $\hat{x}(t)$ in a Taylor series about $t = t_n$ and evaluating the series at $t = t_{n+1}$ yields the following series expansion for $\hat{x}(t_{n+1})$:

$$\begin{aligned}\hat{x}(t_{n+1}) &= \hat{x}(t_n) + \dot{\hat{x}}(t_n)(t_{n+1} - t_n) \\ &\quad + \frac{1}{2!}\ddot{\hat{x}}(t_n)(t_{n+1} - t_n)^2 + \dots + \frac{1}{p!}x^{(p)}(t_n)(t_{n+1} - t_n)^p + h.o.t.\end{aligned}$$

where *h.o.t.* stands for *higher-order terms* of the expansion. If the time step $h = t_{n+1} - t_n$, then

$$\hat{x}(t_{n+1}) = \hat{x}(t_n) + h\dot{\hat{x}}(t_n) + \frac{h^2}{2!}\ddot{\hat{x}}(t_n) + \dots + \frac{h^p}{p!}x^{(p)}(t_n) + h.o.t.$$

From Equation (5.1), $\dot{x}(t) = f(x, t)$; therefore,

$$\begin{aligned}\hat{x}(t_{n+1}) - h.o.t. &= \hat{x}(t_n) + hf(x(t_n), t_n) \\ &\quad + \frac{h^2}{2!}f'(x(t_n), t_n) + \dots + \frac{h^p}{p!}f^{(p-1)}(x(t_n), t_n)\end{aligned}\quad (5.2)$$

If the higher-order terms are small, then a good approximation x_{n+1} to $\hat{x}(t_{n+1})$ is given by the right-hand side of Equation (5.2).

In general, the Taylor series-based integration methods can be expressed as

$$x_{n+1} = x_n + hT_p(x_n) \quad (5.3)$$

where

$$T_p(x_n) = f(x(t_n), t_n) + \frac{h^2}{2!} f'(x(t_n), t_n) + \dots + \frac{h^p}{p!} f^{(p-1)}(x(t_n), t_n)$$

and the integer p is called the *order* of the integration method. This method is very accurate for large p , but is not computationally efficient for large p since it requires a large number of function derivatives and evaluations.

5.1.2 Forward Euler Method

For $p = 1$, the Taylor series-based integration algorithm is given by:

$$x_{n+1} = x_n + hf(x_n, t_n) \quad (5.4)$$

which is also the well-known *Euler* or *forward Euler* method.

5.1.3 Runge–Kutta Methods

A second order Taylor method can be derived for $p = 2$.

$$\begin{aligned} x_{n+1} &= x_n + hT_2(x_n, t_n) \\ &= x_n + hf(x_n, t_n) + \frac{h^2}{2} f'(x_n, t_n) \end{aligned}$$

As the order of the Taylor method increases, so does the number of derivatives and partial derivatives. In many cases, the analytic derivation of the derivatives can be replaced by a numerical approximation. One of the most commonly known higher-order Taylor series-based integration methods is the Runge–Kutta method, where the derivatives are replaced by approximations. The fourth-order Runge–Kutta method is given by

$$x_{n+1} = x_n + hK_4(x_n, t_n) \quad (5.5)$$

where K_4 is an approximation to T_4 :

$$\begin{aligned} K_4 &= \frac{1}{6} [k_1 + 2k_2 + 2k_3 + k_4] \\ k_1 &= f(x_n, t_n) \\ k_2 &= f\left(x_n + \frac{h}{2}k_1, t_n + \frac{h}{2}\right) \\ k_3 &= f\left(x_n + \frac{h}{2}k_2, t_n + \frac{h}{2}\right) \\ k_4 &= f(x_n + hk_3, t_n + h) \end{aligned}$$

where each k_i represents the slope (derivative) of the function at four different points. The slopes are then weighted $\left[\frac{1}{6} \frac{2}{6} \frac{2}{6} \frac{1}{6}\right]$ to approximate the T_4 function.

The advantages of Taylor series-based methods is that the method is straightforward to program and only depends on the previous time step. These methods (especially the Runge–Kutta methods) suffer from difficult error analysis, however, since the derivatives are approximated and not found analytically. Therefore, the integration step size is typically chosen conservatively (small), and computational efficiency may be lost.

5.2 Multistep Methods

Another approach to approximating the solution $x(t)$ of Equation (5.1) is to approximate the nonlinear function as a polynomial of degree k such that

$$\hat{x}(t) = \alpha_0 + \alpha_1 t + \alpha_2 t^2 + \dots + \alpha_k t^k \quad (5.6)$$

where the coefficients $\alpha_0, \alpha_1, \dots, \alpha_k$ are constant. It can be proven that any function can be approximated arbitrarily closely (within a predetermined ε) with a polynomial of sufficiently high degree on a finite interval $[t_0, t_N]$. The polynomial approximation can be related to the solution of Equation (5.1) through the introduction of *multistep* methods. A multistep method is one in which the approximation x_{n+1} can be a function of any number of previous numerical approximations x_n, x_{n-1}, \dots and corresponding functions $f(x_n, t_n), f(x_{n-1}, t_{n-1}), \dots$ unlike one-step methods (such as the Runge–Kutta), which depend only on the information from the immediately previous step. In general,

$$x_{n+1} = a_0 x_n + a_1 x_{n-1} + \dots + a_p x_{n-p} + h [b_{-1} f(x_{n+1}, t_{n+1}) + b_0 f(x_n, t_n) + b_1 f(x_{n-1}, t_{n-1}) + \dots + b_p f(x_{n-p}, t_{n-p})] \quad (5.7)$$

$$= \sum_{i=0}^p a_i x_{n-i} + h \sum_{i=-1}^p b_i f(x_{n-i}, t_{n-i}) \quad (5.8)$$

To relate the integration method to the polynomial approximation, a relationship between the coefficients must be determined. A k -degree polynomial is uniquely determined by $k + 1$ coefficients $(\alpha_0, \dots, \alpha_k)$. The numerical integration method has $2p + 3$ coefficients; therefore, the coefficients must be chosen such that

$$2p + 3 \geq k + 1 \quad (5.9)$$

The order of the numerical integration method is the highest degree k of a polynomial in t for which the numerical solution coincides with the exact

solution. The coefficients may be determined by selecting a set of linear basis functions $[\phi_1(t) \ \phi_2(t), \dots, \ \phi_k(t)]$ such that

$$\phi_j(t) = t^j \quad j = 0, 1, \dots, k$$

and solving the set of multistep equations

$$\phi_j(t_{n+1}) = \sum_{i=0}^p a_i \phi_j(t_{n-i}) + h_{n+1} \left[\sum_{i=-1}^p b_i \dot{\phi}(t_{n-i}) \right]$$

for all $j = 0, 1, \dots, k$.

This method can be applied to derive several first-order numerical integration methods. Consider the case where $p = 0$ and $k = 1$. This satisfies the constraint of Equation (5.9); thus it is possible to determine multistep coefficients that will result in an exact polynomial of degree 1. The set of basis functions for $k = 1$ is

$$\phi_0(t) = 1 \tag{5.10}$$

$$\phi_1(t) = t \tag{5.11}$$

which lead to the derivatives

$$\dot{\phi}_0(t) = 0 \tag{5.12}$$

$$\dot{\phi}_1(t) = 1 \tag{5.13}$$

and the multistep equation

$$x_{n+1} = a_0 x_n + b_{-1} h_{n+1} f(x_{n+1}, t_{n+1}) + b_0 h_{n+1} f(x_n, t_n) \tag{5.14}$$

Representing the multistep method of Equation (5.14) in terms of basis functions yields the following two equations:

$$\phi_0(t_{n+1}) = a_0 \phi_0(t_n) + b_{-1} h_{n+1} \dot{\phi}_0(t_{n+1}) + b_0 h_{n+1} \dot{\phi}_0(t_n) \tag{5.15}$$

$$\phi_1(t_{n+1}) = a_0 \phi_1(t_n) + b_{-1} h_{n+1} \dot{\phi}_1(t_{n+1}) + b_0 h_{n+1} \dot{\phi}_1(t_n) \tag{5.16}$$

Substituting the choice of basis functions of Equations (5.10) and (5.11) into Equations (5.15) and (5.16) results in

$$1 = a_0(1) + b_{-1} h_{n+1}(0) + h_{n+1} b_0(0) \tag{5.17}$$

$$t_{n+1} = a_0 t_n + b_{-1} h_{n+1}(1) + b_0 h_{n+1}(1) \tag{5.18}$$

From Equation (5.17), the coefficient $a_0 = 1$. Recalling that $t_{n+1} - t_n = h_{n+1}$, Equation (5.18) yields

$$b_{-1} + b_0 = 1 \tag{5.19}$$

This choice of order and degree leads to two equations in three unknowns; therefore, one of them may be chosen arbitrarily. By choosing $a_0 = 1, b_{-1} = 0$, and $b_0 = 1$, Euler's method is once again obtained:

$$x_{n+1} = x_n + h_{n+1} f(x_n, t_n)$$

However, if $a_0 = 1$, $b_{-1} = 1$, and $b_0 = 0$, a different integration method is obtained:

$$x_{n+1} = x_n + h_{n+1} f(x_{n+1}, t_{n+1}) \quad (5.20)$$

This particular integration method is frequently called the *backward Euler* method. Note that, in this method, the coefficient b_{-1} is not identically zero; thus the expression for x_{n+1} depends implicitly on the function $f(x_{n+1}, t_{n+1})$. Methods in which $b_{-1} \neq 0$ are referred to as *implicit* methods; otherwise they are *explicit*. Since there is an implicit (and often nonlinear) dependence on x_{n+1} , implicit integration methods must usually be solved iteratively at each time interval.

Consider now the case where $p = 0$, and $k = 2$. In this case, $2p + 3 = k + 1$ and the coefficients can be uniquely determined. Choosing the basis functions as previously with $\phi_2(t) = t^2$ and $\dot{\phi}_2(t) = 2t$ yields the following three equations:

$$1 = a_0(1) + b_{-1}h_{n+1}(0) + h_{n+1}b_0(0) \quad (5.21)$$

$$t_{n+1} = a_0t_n + b_{-1}h_{n+1}(1) + b_0h_{n+1}(1) \quad (5.22)$$

$$t_{n+1}^2 = a_0t_n^2 + h_{n+1}(b_{-1}(2t_{n+1}) + b_0(2t_n)) \quad (5.23)$$

If $t_n = 0$, then $t_{n+1} = h_{n+1}$, and Equations (5.21) through (5.23) yield $a_0 = 1$, $b_{-1} = \frac{1}{2}$, and $b_0 = \frac{1}{2}$; thus

$$x_{n+1} = x_n + \frac{1}{2}h_{n+1}[f(x_{n+1}, t_{n+1}) + f(x_n, t_n)] \quad (5.24)$$

This second-order integration method is called the *trapezoidal* method and it is also implicit. This formula is called the trapezoidal method since the second term of Equation (5.24) can be interpreted as being the area under a trapezoid.

Example 5.1

Numerically solve

$$\ddot{x}(t) = -x(t) \quad x(0) = 1 \quad (5.25)$$

using the Euler, backward Euler, trapezoidal, and Runge–Kutta methods for different fixed step sizes.

Solution 5.1 This second-order differential equation must first be converted to ODE format by defining $x_1 = x$ and $x_2 = \dot{x}$. Then

$$\dot{x}_1 = x_2 = f_1(x_1, x_2) \quad x_1(0) = 1 \quad (5.26)$$

$$\dot{x}_2 = -x_1 = f_2(x_1, x_2) \quad (5.27)$$

By inspection, the analytic solution to this set of equations is

$$x_1(t) = \cos t \quad (5.28)$$

$$x_2(t) = -\sin t \quad (5.29)$$

Typically, it is not possible to find the exact solution, but in this example, the exact solution will be used to compare the numerical solutions against.

Forward Euler

Applying the forward Euler method to the ODEs yields

$$x_{1,n+1} = x_{1,n} + h f_1(x_{1,n}, x_{2,n}) \quad (5.30)$$

$$= x_{1,n} + h x_{2,n} \quad (5.31)$$

$$x_{2,n+1} = x_{2,n} + h f_2(x_{1,n}, x_{2,n}) \quad (5.32)$$

$$= x_{2,n} - h x_{1,n} \quad (5.33)$$

or in matrix form

$$\begin{bmatrix} x_{1,n+1} \\ x_{2,n+1} \end{bmatrix} = \begin{bmatrix} 1 & h \\ -h & 1 \end{bmatrix} \begin{bmatrix} x_{1,n} \\ x_{2,n} \end{bmatrix} \quad (5.34)$$

Backward Euler

Applying the backward Euler method to the ODEs yields

$$x_{1,n+1} = x_{1,n} + h f_1(x_{1,n+1}, x_{2,n+1}) \quad (5.35)$$

$$= x_{1,n} + h x_{2,n+1} \quad (5.36)$$

$$x_{2,n+1} = x_{2,n} + h f_2(x_{1,n+1}, x_{2,n+1}) \quad (5.37)$$

$$= x_{2,n} - h x_{1,n+1} \quad (5.38)$$

or in matrix form

$$\begin{bmatrix} x_{1,n+1} \\ x_{2,n+1} \end{bmatrix} = \begin{bmatrix} 1 & -h \\ h & 1 \end{bmatrix}^{-1} \begin{bmatrix} x_{1,n} \\ x_{2,n} \end{bmatrix} \quad (5.39)$$

In the solution of Equation (5.39), the inverse of the matrix is not found explicitly, but rather the equations would be solved using LU factorization.

Trapezoidal

Applying the trapezoidal method to the ODEs yields

$$x_{1,n+1} = x_{1,n} + \frac{1}{2}h [f_1(x_{1,n}, x_{2,n}) + f_1(x_{1,n+1}, x_{2,n+1})] \quad (5.40)$$

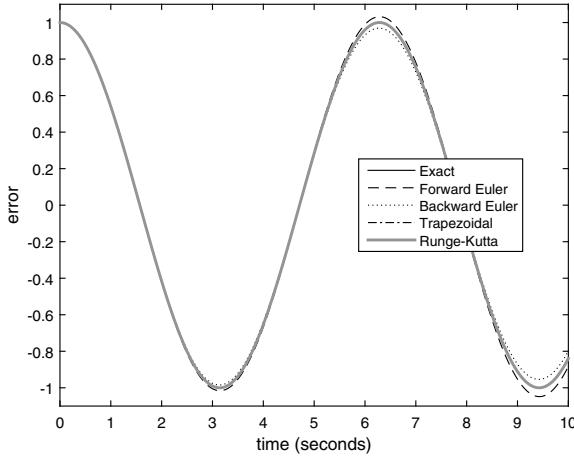
$$= x_{1,n} + \frac{1}{2}h [x_{2,n} + x_{2,n+1}] \quad (5.41)$$

$$x_{2,n+1} = x_{2,n} + \frac{1}{2}h [f_2(x_{1,n}, x_{2,n}) + f_2(x_{1,n+1}, x_{2,n+1})] \quad (5.42)$$

$$= x_{2,n} - \frac{1}{2}h [x_{1,n} + x_{1,n+1}] \quad (5.43)$$

or in matrix form

$$\begin{bmatrix} x_{1,n+1} \\ x_{2,n+1} \end{bmatrix} = \begin{bmatrix} 1 & -\frac{1}{2}h \\ \frac{1}{2}h & 1 \end{bmatrix}^{-1} \begin{bmatrix} 1 & \frac{1}{2}h \\ -\frac{1}{2}h & 1 \end{bmatrix} \begin{bmatrix} x_{1,n} \\ x_{2,n} \end{bmatrix} \quad (5.44)$$

**FIGURE 5.1**

Numerical solutions for Example 5.1

Runge–Kutta

Applying the Runge–Kutta method to the ODEs yields

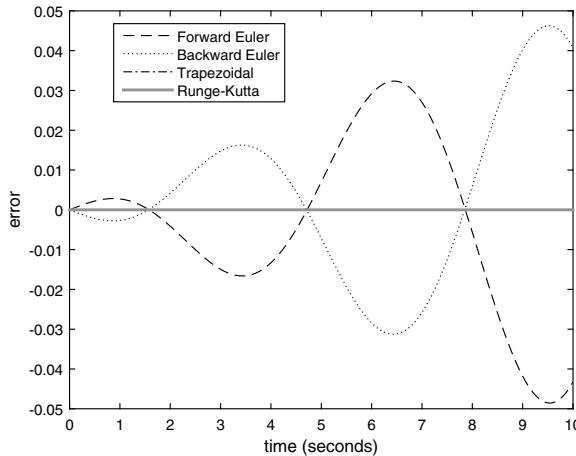
$$\begin{aligned} k_{11} &= x_{2,n} & k_{21} &= -x_{1,n} \\ k_{12} &= x_{2,n} + \frac{h}{2}k_{11} & k_{22} &= -x_{1,n} - \frac{h}{2}k_{21} \\ k_{13} &= x_{2,n} + \frac{h}{2}k_{12} & k_{23} &= -x_{1,n} - \frac{h}{2}k_{22} \\ k_{14} &= x_{2,n} + hk_{13} & k_{24} &= -x_{1,n} - hk_{23} \end{aligned}$$

and

$$x_{1,n+1} = x_{1,n} + \frac{h}{6} (k_{11} + 2k_{12} + 2k_{13} + k_{14}) \quad (5.45)$$

$$x_{2,n+1} = x_{2,n} + \frac{h}{6} (k_{21} + 2k_{22} + 2k_{23} + k_{24}) \quad (5.46)$$

The numerical solution of Equation (5.25) for each of the methods is shown in Figure 5.1 including the exact solution $\cos t$. Note that the trapezoidal and Runge–Kutta methods are nearly indistinguishable from the exact solution. Since the forward and backward Euler methods are first-order methods, they are not as accurate as the higher order trapezoidal and Runge–Kutta methods. Note that the forward Euler method generates a numerical solution whose magnitude is slightly larger than the exact solution and is increasing with time. Conversely, the backward Euler method generates a numerical solution whose magnitude is slightly less than the exact solution and is decreasing with time. Both properties are due to the *local truncation error* of the methods. The forward Euler method has a tendency to generate numerical solutions that increase with time (underdamped), whereas the backward Euler method

**FIGURE 5.2**

Error in numerical solutions for Example 5.1

tends to add damping to the numerical solution. Therefore, caution must be used when using either of these first-order methods for numerical integration.

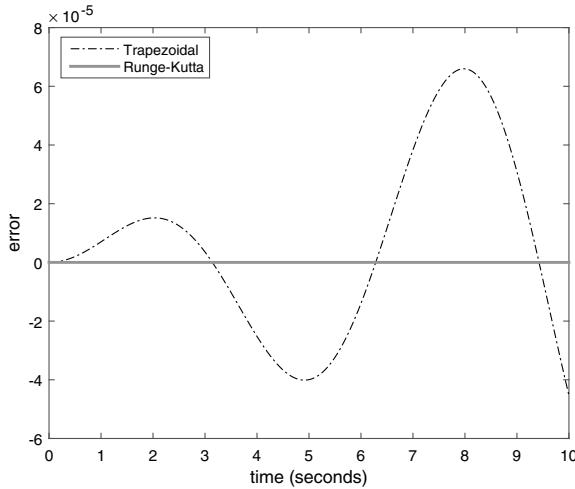
Figure 5.2 shows the global error with time for each of the numerical methods. Note that the errors for the forward and backward Euler methods are equal in magnitude, but opposite in sign. This relationship will be further discussed later in this chapter. The numerical errors for the trapezoidal and the Runge–Kutta methods are reproduced in Figure 5.3 on a magnified scale. Note that even in this case, the Runge–Kutta error is much smaller than the trapezoidal method error. This is because the Runge–Kutta method is a fourth-order Taylor method whereas the trapezoidal method is a second-order polynomial approximation method. Section 5.3 will further explore the development of expressions for estimating the error of various integration methods. ■

When implicit methods, such as the trapezoidal method, are used to solve nonlinear systems of differential equations, the system of equations must be solved iteratively at each time step. For example, consider the following nonlinear system of equations:

$$\dot{x} = f(x(t), t) \quad x_0 = x(t_0) \quad (5.47)$$

Applying the trapezoidal method to numerically integrate this system results in the following discretized system:

$$x_{n+1} = x_n + \frac{h}{2} [f(x_n, t_n) + f(x_{n+1}, t_{n+1})] \quad (5.48)$$

**FIGURE 5.3**

Error in trapezoidal and Runge–Kutta numerical solutions for Example 5.1

Since this nonlinear expression is implicit in x_{n+1} , it must be solved numerically:

$$x_{n+1}^{k+1} = x_{n+1}^k - \left[I - \frac{h}{2} \frac{\partial f}{\partial x} \right]^{-1} \Bigg|_{x_{n+1}^k} \left(x_{n+1}^k - x_n - \frac{h}{2} [f(x_n) + f(x_{n+1}^k)] \right) \quad (5.49)$$

where k is the Newton–Raphson iteration index, I is the identity matrix, and x_n is the converged value from the previous time step.

5.2.1 Adams Methods

Recall that the general class of multistep methods may be represented by

$$x_{n+1} = \sum_{i=0}^p a_i x_{n-i} + h \sum_{i=-1}^p b_i f(x_{n-i}, t_{n-i}) \quad (5.50)$$

A numerical multistep algorithm will give the exact value for x_{n+1} if $x(t)$ is a polynomial of degree less than or equal to k if the following *exactness constraints* are satisfied:

$$\sum_{i=0}^p a_i = 1 \quad (5.51)$$

$$\sum_{i=1}^p (-i)^p a_i + j \sum_{i=-1}^p (-i)^{(j-1)} b_i = 1 \text{ for } j = 1, 2, \dots, k \quad (5.52)$$

The exactness constraint of Equation (5.51) is frequently referred to as the *consistency* constraint. Numerical multistep integration algorithms that satisfy Equation (5.51) are said to be “consistent.” For a desired polynomial of degree k , these constraints can be satisfied by a wide variety of possibilities. Several families of methods have been developed by predefining some of the relationships between the coefficients. The family of *Adams* methods is defined by setting the coefficients $a_1 = a_2 = \dots = a_p = 0$. By the consistency constraint, the coefficient a_0 must therefore equal 1.0. Thus, the Adams methods are reduced to

$$x_{n+1} = x_n + h \sum_{i=-1}^p b_i f(x_{n-i}, t_{n-i}) \quad (5.53)$$

where $p = k - 1$. The Adams methods can be further classified by the choice of implicit or explicit integration. The explicit class, frequently referred to as the “Adams–Bashforth” methods, is specified by setting $b_{-1} = 0$ and applying the second exactness constraint as

$$\sum_{i=0}^{k-1} (-i)^{(j-1)} b_i = \frac{1}{j} \quad j = 1, \dots, k \quad (5.54)$$

In matrix form, Equation (5.54) becomes

$$\begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 0 & -1 & -2 & \dots & -(k-1) \\ 0 & 1 & 4 & \dots & -(k-1)^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & (-1)^{(k-1)} & (-2)^{(k-1)} & \dots & (-(k-1))^{(k-1)} \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \\ b_{k-1} \end{bmatrix} = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \frac{1}{3} \\ \vdots \\ \frac{1}{k} \end{bmatrix} \quad (5.55)$$

By choosing the desired degree k (and subsequently the order p), the remaining b_i coefficients may be found from solving Equation (5.55).

Example 5.2

Find the third-order Adams–Bashforth integration method.

Solution 5.2 Setting $k = 3$ yields the following linear system:

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & -1 & -2 \\ 0 & 1 & 4 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \frac{1}{3} \end{bmatrix}$$

Solving this system yields

$$\begin{aligned} b_0 &= \frac{23}{12} \\ b_1 &= -\frac{16}{12} \\ b_2 &= \frac{5}{12} \end{aligned}$$

Thus the third-order Adams–Bashforth method is given by

$$x_{n+1} = x_n + \frac{1}{12}h [23f(x_n, t_n) - 16f(x_{n-1}, t_{n-1}) + 5f(x_{n-2}, t_{n-2})] \quad (5.56)$$

When implementing this algorithm, the values of x_n , x_{n-1} , and x_{n-2} must be saved in memory. ■

The implicit versions of the Adams methods have $b_{-1} \neq 0$, $p = (k - 2)$, are called the “Adams–Moulton” methods, and are given by

$$x_{n+1} = x_n + h \sum_{i=-1}^{k-2} b_i f(x_{n-i}, t_{n-i}) \quad (5.57)$$

The second exactness constraint yields

$$\sum_{i=-1}^{k-2} (-i)^{(j-1)} b_i = \frac{1}{j} \quad j = 1, \dots, k \quad (5.58)$$

or in matrix form

$$\begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & 0 & -1 & \dots & -(k-2) \\ 1 & 0 & 1 & \dots & (-k-2)^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & (-1)^{(k-2)} & (-2)^{(k-2)} & \dots & (-k-2)^{(k-2)} \end{bmatrix} \begin{bmatrix} b_{-1} \\ b_0 \\ b_1 \\ \vdots \\ b_{k-2} \end{bmatrix} = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \frac{1}{3} \\ \vdots \\ \frac{1}{k} \end{bmatrix} \quad (5.59)$$

Example 5.3

Find the third-order Adams–Moulton integration method.

Solution 5.3 Setting $k = 3$ yields the following linear system:

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & -1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} b_{-1} \\ b_0 \\ b_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \frac{1}{3} \end{bmatrix}$$

Solving this system yields

$$\begin{aligned} b_{-1} &= \frac{5}{12} \\ b_0 &= \frac{8}{12} \\ b_1 &= -\frac{1}{12} \end{aligned}$$

Thus the third-order Adams–Moulton method is given by

$$x_{n+1} = x_n + \frac{1}{12}h [5f(x_{n+1}, t_{n+1}) + 8f(x_n, t_n) - f(x_{n-1}, t_{n-1})] \quad (5.60)$$

When implementing this algorithm, the values of x_n and x_{n-1} must be saved in memory and the equations must be solved iteratively if the function $f(x)$ is nonlinear. ■

The Adams–Moulton method is implicit and must be solved iteratively using the Newton–Raphson method (or other similar method), as shown in Equation (5.49). Iterative methods require an initial value for the iterative process to reduce the number of required iterations. The explicit Adams–Bashforth method is frequently used to estimate the initial value for the implicit Adams–Moulton method. If sufficiently high-order predictor methods are used, the Adams–Moulton method iteration will typically converge in only one iteration. This process is often called a *predictor corrector* approach; the Adams–Bashforth method predicts the solution and the implicit Adams–Moulton corrects the solution.

Another implementation issue for multistep methods is how to start up the integration at the beginning of the simulation since a high-order method requires several previous values. The usual procedure is to use a high order one-step method or to increase the number of steps of the method with each time step to generate the required number of values for the desired multistep method.

5.2.2 Gear's Methods

Another well-known family of multistep methods is Gear's methods [15]. This family of methods is particularly well suited for the numerical solution of stiff systems. As opposed to the Adams family of methods, where all the a_i coefficients except a_0 are zero, Gear's methods are identified by having all of the b_i coefficients equal to zero except b_{-1} . Obviously, since $b_{-1} \neq 0$, all Gear's methods are implicit methods. The k th order Gear's algorithm is obtained by setting $p = k - 1$ and $b_0 = b_1 = \dots = 0$, yielding

$$x_{n+1} = a_0 x_n + a_1 x_{n-1} + \dots + a_{k-1} x_{n-k+1} + h b_{-1} f(x_{n+1}, t_{n+1}) \quad (5.61)$$

The $k + 1$ coefficients can be calculated explicitly by applying the exactness constraints, as illustrated with the Adams methods

$$\begin{bmatrix} 1 & 1 & 1 \dots & 1 & 0 \\ 0 & -1 & -2 \dots & -(k-1) & 1 \\ 0 & 1 & 4 \dots & [-(k-1)]^2 & 2 \\ \vdots & \vdots & \vdots \ddots & \vdots & \vdots \\ 0 & (-1)^k & (-2)^k \dots & [-(k-1)]^k & k \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ b_{-1} \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \quad (5.62)$$

The solution of Equation (5.62) uniquely determines the $k + 1$ coefficients of the k th order Gear's method.

Example 5.4

Find the third-order Gear's integration method.

Solution 5.4 Setting $k = 3$ yields the following linear system:

$$\begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & -1 & -2 & 1 \\ 0 & 1 & 4 & 2 \\ 0 & -1 & -8 & 3 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ b_{-1} \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

Solving this system yields

$$\begin{aligned} b_{-1} &= \frac{6}{11} \\ a_0 &= \frac{18}{11} \\ a_1 &= -\frac{9}{11} \\ a_2 &= \frac{2}{11} \end{aligned}$$

Thus the third-order Gear's method is given by

$$x_{n+1} = \frac{18}{11}x_n - \frac{9}{11}x_{n-1} + \frac{2}{11}x_{n-2} + \frac{6}{11}hf(x_{n+1}, t_{n+1}) \quad (5.63)$$

When implementing this algorithm, the values of x_n through x_{n-2} must be saved in memory and the equations must be solved iteratively if the function $f(x)$ is nonlinear. ■

5.3 Accuracy and Error Analysis

The accuracy of numerical integration methods is impacted by two primary causes: computer round-off error and truncation error. Computer round-off error occurs as a result of the finite precision of the computer upon which the algorithm is implemented, and little can be done to reduce this error short of using a computer with greater precision. A double precision word length is normally used for scientific computation. The difference between the exact solution and the calculated solution is dominated by truncation error, which arises from the truncation of the Taylor series or polynomial being used to approximate the solution.

In the implementation of numerical integration algorithms, the most effective methods are those methods which require the least amount of calculation to yield the most accurate results. In general, higher-order methods produce

the most accurate results, but also require the greatest amount of computation. Therefore, it is desirable to compute the solution as infrequently as possible by taking the largest time step possible between intervals. Several factors impact the size of the time step, including the error introduced at each step by the numerical integration process itself. This error is the *local truncation error* (LTE) and arises from the truncation of the polynomial approximation and/or the truncation of the Taylor series expansion, depending on the method used. The term *local* emphasizes that the error is introduced locally and is not residual global error from earlier time steps. The error introduced at a single step of an integration method is given by

$$\varepsilon_T \triangleq x(t_{n+1}) - x_{n+1} \quad (5.64)$$

where $x(t_{n+1})$ is the exact solution at time t_{n+1} and x_{n+1} is the numerical approximation. This definition assumes that this is the error *introduced in one step*; therefore, $x(t_n) = x_n$. The local truncation error is shown graphically in Figure 5.4. To compute the error, the solution $x(t_{n-i})$ is expanded about t_{n+1} :

$$x_{n-i} = x(t_{n-i}) = \sum_{j=0}^{\infty} \frac{(t_{n-i} - t_{n+1})^j}{j!} \frac{d^{(j)}}{dt^j} x(t_{n+1}) \quad (5.65)$$

Recall that

$$\begin{aligned} f(x_{n-i}, t_{n-i}) &= \dot{x}(t_{n-i}) \\ &= \sum_{j=0}^{\infty} \frac{(t_{n-i} - t_{n+1})^j}{j!} \frac{d^{(j+1)}}{dt^{j+1}} x(t_{n+1}) \end{aligned} \quad (5.66)$$

Solving for $x(t_{n+1}) - x_{n+1}$ yields

$$\begin{aligned} \varepsilon_T &= C_0 x(t_n) + C_1 x(t_{n-1}) \\ &\quad + C_2 x(t_{n-2}) + \dots + C_k x(t_{n-k}) + C_{k+1} x(t_{n-k-1}) + \dots \end{aligned} \quad (5.67)$$

If the order of this method is k , then the first k coefficients are equal to zero and the local truncation error is given by

$$\varepsilon_T = C_{k+1} h^{k+1} x^{(k+1)}(t_{n+1}) + O(h^{k+2}) \quad (5.68)$$

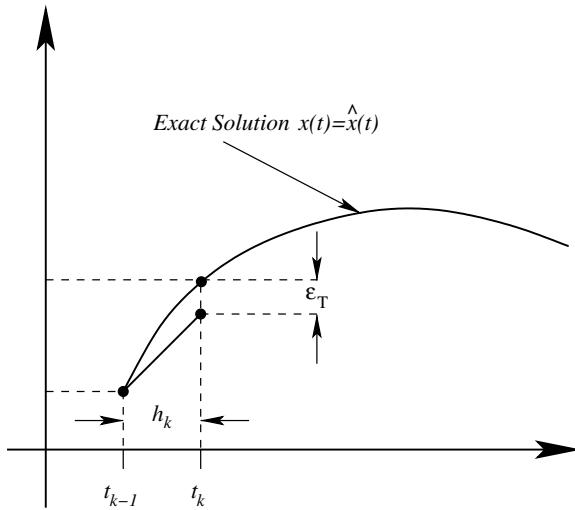
where $O(h^{k+2})$ indicates an error on the order of h^{k+2} .

Example 5.5

Find expressions for the local truncation error of the forward Euler, backward Euler, and trapezoidal integration methods.

Solution 5.5

Forward Euler

**FIGURE 5.4**

Graphical depiction of local truncation error

Recall that the expression for the forward Euler integration algorithm is

$$x_{n+1} = x_n + h f(x_n, t_n)$$

If $x_n = x(t_n)$ (by the definition of the local truncation error), then

$$x_n = x(t_{n+1}) - h \dot{x}(t_{n+1}) + \frac{1}{2!} h^2 \ddot{x}(t_{n+1}) + \dots \quad (5.69)$$

and

$$f(x_n, t_n) = \dot{x}(t_n) = \dot{x}(t_{n+1}) - h \ddot{x}(t_{n+1}) + \dots \quad (5.70)$$

Thus

$$\varepsilon_T = x(t_{n+1}) - x_{n+1} \quad (5.71)$$

$$= x(t_{n+1}) - x_n - h f(x_n, t_n) \quad (5.72)$$

$$= x(t_{n+1}) - \left[x(t_{n+1}) - h \dot{x}(t_{n+1}) + \frac{1}{2!} h^2 \ddot{x}(t_{n+1}) + \dots \right] \\ - h [\dot{x}(t_{n+1}) - h \ddot{x}(t_{n+1}) + \dots] \quad (5.73)$$

$$= \frac{h^2}{2} \ddot{x}(t_{n+1}) + O(h^3) \quad (5.74)$$

Backward Euler

The expression for the backward Euler integration algorithm is

$$x_{n+1} = x_n + h f(x_{n+1}, t_{n+1})$$

Following the same approach as the forward Euler method, but using

$$f(x_{n+1}, t_{n+1}) = \dot{x}(t_{n+1}) \quad (5.75)$$

then

$$\varepsilon_T = x(t_{n+1}) - x_{n+1} \quad (5.76)$$

$$= x(t_{n+1}) - x_n - hf(x_{n+1}, t_{n+1}) \quad (5.77)$$

$$= x(t_{n+1}) - \left[x(t_{n+1}) - h\dot{x}(t_{n+1}) + \frac{1}{2!}h^2\ddot{x}(t_{n+1}) + \dots \right] \\ - h\dot{x}(t_{n+1}) \quad (5.78)$$

$$= -\frac{h^2}{2}\ddot{x}(t_{n+1}) - O(h^3) \quad (5.79)$$

Note that the local truncation errors for the forward and backward Euler methods are equal, but opposite in sign. This property is consistent with the results of Example 5.1, shown in Figure 5.2, where the respective errors were identical except in sign.

Trapezoidal

The expression for the second-order trapezoidal integration algorithm is

$$x_{n+1} = x_n + \frac{1}{2}h[f(x_{n+1}, t_{n+1}) + f(x_n, t_n)]$$

Following the same approach as the previous methods using similar substitutions,

$$\varepsilon_T = x(t_{n+1}) - x_{n+1} \quad (5.80)$$

$$= x(t_{n+1}) - x_n - \frac{1}{2}hf(x_n, t_n) - \frac{1}{2}hf(x_{n+1}, t_{n+1}) \quad (5.81)$$

$$= x(t_{n+1}) - \left[x(t_{n+1}) - h\dot{x}(t_{n+1}) + \frac{h^2}{2}\ddot{x}(t_{n+1}) - \frac{h^3}{3!}x^{(3)}(t_{n+1}) + \dots \right]$$

$$- \frac{h}{2} \left[\dot{x}(t_{n+1}) - h\ddot{x}(t_{n+1}) + \frac{h^2}{2}x^{(3)}(t_{n+1}) + \dots \right] - \frac{h}{2}\dot{x}(t_{n+1}) \quad (5.82)$$

$$= \frac{h^3}{6}x^{(3)}(t_{n+1}) - \frac{h^3}{4}x^{(3)}(t_{n+1}) + O(h^4) \quad (5.83)$$

$$= -\frac{1}{12}h^3x^{(3)}(t_{n+1}) + O(h^4) \quad (5.84)$$

■

Both the first-order Euler methods had errors on the order of h^2 , whereas the second-order method (trapezoidal) had an error on the order of h^3 . Both methods are implicit and must be solved iteratively at each time step. Consider the iterative solution of the trapezoidal method repeated here from Equation (5.49):

$$x_{n+1}^{k+1} = x_{n+1}^k - \left[I - \frac{h}{2} \frac{\partial f}{\partial x} \right]^{-1} \Bigg|_{x_{n+1}^k} \left(x_{n+1}^k - x_n - \frac{h}{2} [f(x_n) + f(x_{n+1}^k)] \right) \quad (5.85)$$

Similarly, the iterative solution of the backward Euler method is given by

$$x_{n+1}^{k+1} = x_{n+1}^k - \left[I - h \frac{\partial f}{\partial x} \right]^{-1} \Bigg|_{x_{n+1}^k} (x_{n+1}^k - x_n - h [f(x_{n+1}^k)]) \quad (5.86)$$

Note that both methods require the same function evaluations and comparable computational effort, yet the trapezoidal method yields a much smaller local truncation error for the same time step size h . For this reason, the trapezoidal method is a more widely used general purpose implicit numerical integration algorithm than the backward Euler method.

For multistep methods, a generalized expression for the local truncation error has been developed [8]. For a multistep method

$$x_{n+1} = \sum_{i=0}^p a_i x_{n-i} + h \sum_{i=-1}^p b_i f(x_{n-i}, t_{n-i}) \quad (5.87)$$

which is exact for a polynomial solution of degree less than or equal to k , the local truncation error is given by

$$\varepsilon_T = C_k x^{(k+1)}(\tau) h^{k+1} = O(h^{k+1}) \quad (5.88)$$

where $-ph \leq \tau \leq h$ and

$$C_k \triangleq \frac{1}{(k+1)!} \left\{ (p+1)^{k+1} - \left[\sum_{i=0}^{p-1} a_i (p-i)^{k+1} + (k+1) \sum_{i=-1}^{p-1} b_i (p-i)^k \right] \right\} \quad (5.89)$$

This expression provides an approach for approximating the local truncation error at each time step as a function of h and x .

5.4 Numerical Stability Analysis

From the previous discussion, it was shown that the choice of integration step size directly impacts the numerical accuracy of the solution. Less obvious is how the choice of step size impacts the numerical stability of the integration method. Numerical stability guarantees that the global truncation error remains bounded. This guarantees that the error introduced at each time step

does not accrue with time, but rather dissipates such that the choice of step size can be made by considering the local truncation error only. To analyze the effect of step size on the numerical stability of integration methods, consider the simple, scalar ODE

$$\dot{x} = f(x) = \lambda x(t) \quad x_0 = x(t_0) \quad (5.90)$$

By inspection, the solution to this equation is

$$x(t) = x_0 e^{(\lambda t)} \quad (5.91)$$

If $\lambda < 0$, then $x(t)$ approaches zero as t goes to infinity. Conversely, if $\lambda > 0$, then $x(t)$ approaches infinity as t goes to infinity. Numerical stability ensures that the global behavior of the estimated system matches that of the actual system. Consider the forward Euler method applied to the scalar system of Equation (5.90):

$$\begin{aligned} x_{n+1} &= x_n + h\lambda x_n \\ &= (1 + h\lambda)x_n \end{aligned}$$

thus

$$\begin{aligned} x_1 &= (1 + h\lambda)x_0 \\ x_2 &= (1 + h\lambda)x_1 = (1 + h\lambda)^2 x_0 \\ &\vdots \\ x_n &= (1 + h\lambda)^n x_0 \end{aligned}$$

If $\lambda < 0$, then $x(t)$ should approach zero as t goes to infinity. This will be achieved only if

$$|1 + h\lambda| < 1 \quad (5.92)$$

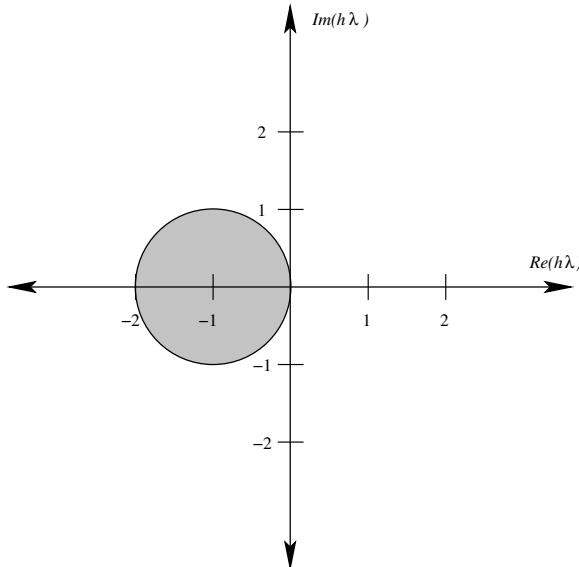
Therefore, this system is stable for $\lambda < 0$ only if $h\lambda$ lies within the unit circle centered at $(-1, 0)$, shown in Figure 5.5. Thus the larger the value of λ , the smaller the integration step size must be.

Similarly, consider the backward Euler integration method applied to the same scalar ODE system:

$$\begin{aligned} x_{n+1} &= x_n + h\lambda x_{n+1} \\ &= \frac{x_n}{(1 - h\lambda)} \end{aligned}$$

thus

$$\begin{aligned} x_1 &= \frac{x_0}{(1 - h\lambda)} \\ x_2 &= \frac{x_1}{(1 - h\lambda)} = \frac{x_0}{(1 - h\lambda)^2} \\ &\vdots \\ x_n &= \frac{x_0}{(1 - h\lambda)^n} \end{aligned}$$

**FIGURE 5.5**

Region of absolute stability of the forward Euler method

If $\lambda < 0$, then $x(t)$ should approach zero as t goes to infinity. This will be achieved only if

$$|1 - h\lambda| > 1 \quad (5.93)$$

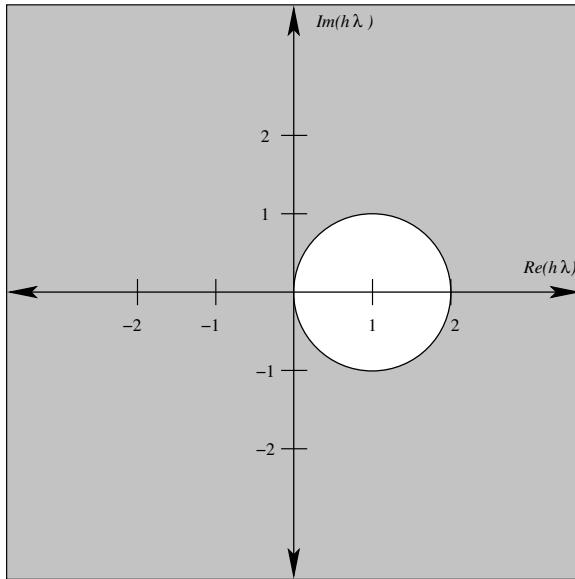
Therefore, this system is stable for $\lambda < 0$ only if $h\lambda$ does not lie within the unit circle centered at $(1, 0)$, shown in Figure 5.6. This implies that, for all $\lambda < 0$, the backward Euler method is numerically stable. Thus, if the ODE system is stable, the integration step size may be chosen arbitrarily large without affecting the numerical stability of the solution. Thus the selection of integration step size will be dependent only on the local truncation error. Note that, if $h\lambda$ is large, then x_n will rapidly approach zero. This characteristic manifests itself as a tendency to overdamp the numerical solution. This characteristic was illustrated in Figure 5.1.

Extending this approach to the general family of multistep methods yields

$$x_{n+1} = \sum_{i=0}^p a_i x_{n-i} + h\lambda \sum_{i=-1}^p b_i x_{n-i} \quad (5.94)$$

Rearranging the terms of the multistep method gives

$$x_{n+1} = \frac{(a_0 + h\lambda b_0)}{(1 - h\lambda b_{-1})} x_n + \frac{(a_1 + h\lambda b_1)}{(1 - h\lambda b_{-1})} x_{n-1} + \dots$$

**FIGURE 5.6**

Region of absolute stability of the backward Euler method

$$+ \frac{(a_p + h\lambda b_p)}{(1 - h\lambda b_{-1})} x_{n-p} \quad (5.95)$$

$$= \gamma_0 x_n + \gamma_1 x_{n-1} + \dots + \gamma_p x_{n-p} \quad (5.96)$$

This relationship specifies the characteristic equation

$$P(z, h\lambda) = z^{p+1} + \gamma_0 z^p + \dots + \gamma_p = 0 \quad (5.97)$$

where z_1, z_2, \dots, z_{p+1} are the (complex) roots of Equation (5.97). Therefore,

$$x_{n+1} = \sum_{i=1}^{p+1} C_i z_i^{n+1} \quad (5.98)$$

If $\lambda < 0$, the solution x_{n+1} will go to zero as n goes to infinity only if $|z_j| < 1$ for all $j = 1, 2, \dots, p+1$. Thus a multistep method is said to be *absolutely stable* for a given value of $h\lambda$ if the roots of $P(z, h\lambda) = 0$ satisfy $|z_i| < 1$ for $i = 1, \dots, k$. Absolute stability implies that the global error decreases with increasing n . The region of absolute stability is defined to be the region in the complex $h\lambda$ plane where the roots of $P(z, h\lambda) = 0$ satisfy $|z_i| < 1$ for $i = 1, \dots, k$. Let

$$P(z, h\lambda) = P_a(z) + h\lambda P_b(z) = 0$$

where

$$\begin{aligned} P_a(z) &\stackrel{\Delta}{=} z^{p+1} - a_0 z^p - a_1 z^{p-1} - \dots - a_p \\ P_b(z) &\stackrel{\Delta}{=} b_{-1} z^{p+1} + b_0 z^p + b_1 z^{p-1} + \dots + b_p \end{aligned}$$

then

$$h\lambda = -\frac{P_a(z)}{P_b(z)} \quad (5.99)$$

Since z is a complex number, it can also be represented as

$$z = e^{(j\theta)}$$

The boundary of the region can be mapped by plotting $h\lambda$ in the complex plane as θ varies from 0 through 2π where

$$h\lambda(\theta) = -\frac{e^{j(p+1)\theta} - a_0 e^{jp\theta} - a_1 e^{j(p-1)\theta} - \dots - a_{p-1} e^{j\theta} - a_p}{b_{-1} e^{j(p+1)\theta} + b_0 e^{jp\theta} + b_1 e^{j(p-1)\theta} + \dots + b_{p-1} e^{j\theta} + b_p} \quad (5.100)$$

Example 5.6

Plot the regions of absolute stability of Gear's third-order and the Adams third-order (both implicit and explicit) methods.

Solution 5.6

Gear's

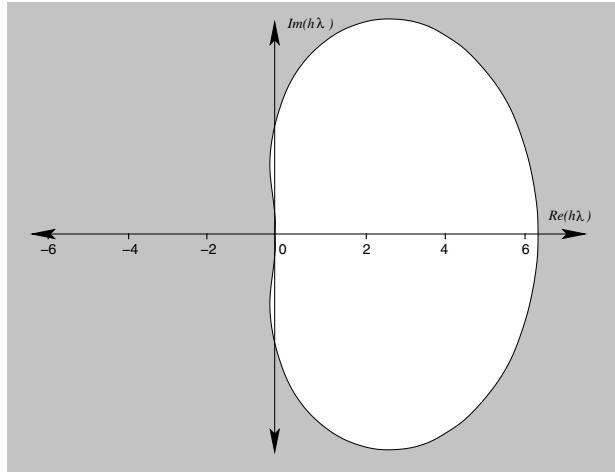
The region of stability of a Gear's method can be found by plotting $h\lambda$ in the complex plane from Equation (5.100) by setting $p = k - 1$ and $b_0, b_1, \dots = 0$.

$$h\lambda(\theta) = \frac{e^{jk\theta} - a_0 e^{j(k-1)\theta} - \dots - a_{k-1}}{b_{-1} e^{jk\theta}} \quad (5.101)$$

Substituting in the coefficients for the third-order Gear's method yields the following expression in θ :

$$h\lambda(\theta) = \frac{e^{j3\theta} - \frac{18}{11}e^{j2\theta} + \frac{9}{11}e^{j\theta} - \frac{2}{11}}{\frac{6}{11}e^{j3\theta}} \quad (5.102)$$

By varying θ from zero to 2π , the region of absolute stability of Gear's third-order method is shown as the shaded region of Figure 5.7.

**FIGURE 5.7**

Region of absolute stability for Gear's third-order method

Adams–Moulton

The region of absolute stability of the Adams–Moulton methods can be developed from Equation (5.100) by setting $p = k - 1$, and $a_1, a_2, \dots = 0$:

$$h\lambda(\theta) = \frac{e^{jk\theta} - a_0 e^{j(k-1)\theta}}{b_{-1} e^{jk\theta} + b_0 e^{j(k-1)\theta} + b_1 e^{j(k-2)\theta} + \dots + b_{k-2} e^{j\theta}} \quad (5.103)$$

After substituting in the third-order coefficients, the expression for the region of absolute stability as a function of θ is given by

$$h\lambda(\theta) = \frac{e^{j3\theta} - e^{j2\theta}}{\frac{5}{12}e^{j3\theta} + \frac{8}{12}e^{j2\theta} - \frac{1}{12}e^{j\theta}} \quad (5.104)$$

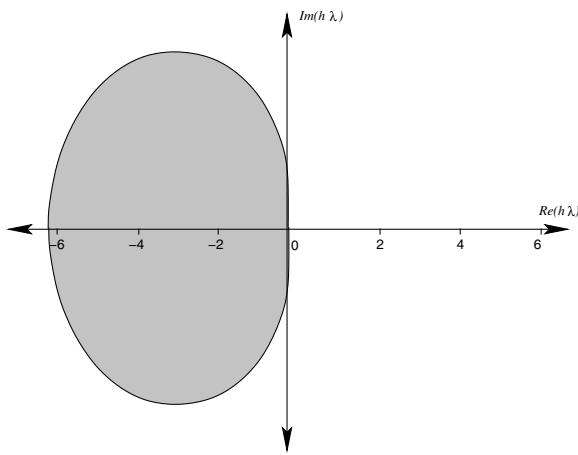
The region of absolute stability of the Adams–Moulton third-order method is shown as the shaded region of Figure 5.8.

Adams–Bashforth

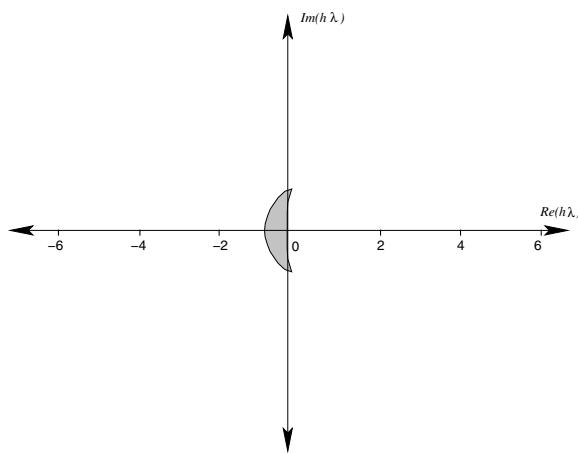
The stability of the family of Adams–Bashforth methods can be derived from Equation (5.100) by setting $p = k - 1$, $b_{-1} = 0$ and $a_1, a_2, \dots = 0$:

$$h\lambda(\theta) = \frac{e^{jk\theta} - a_0 e^{j(k-1)\theta}}{b_0 e^{j(k-1)\theta} + b_1 e^{j(k-2)\theta} + \dots + b_{k-1}} \quad (5.105)$$

$$h\lambda(\theta) = \frac{e^{j3\theta} - e^{j2\theta}}{\frac{23}{12}e^{j2\theta} - \frac{16}{12}e^{j\theta} + \frac{5}{12}} \quad (5.106)$$

**FIGURE 5.8**

Region of absolute stability for the Adams–Moulton third-order method

**FIGURE 5.9**

Region of absolute stability for the Adams–Bashforth third-order method

The region of absolute stability of the Adams–Bashforth method is shown as the shaded region in Figure 5.9.

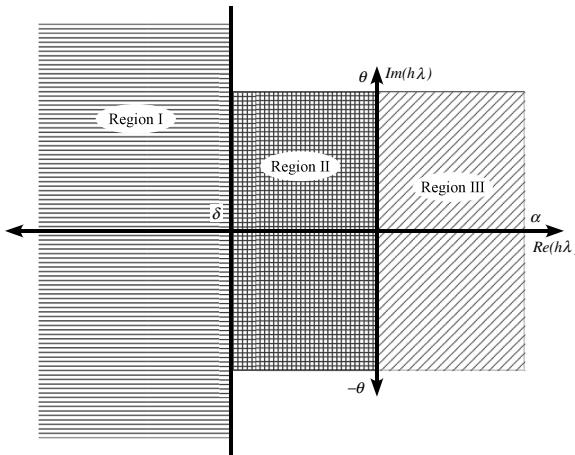
This example illustrates one of the primary differences between implicit and explicit methods. For the same order, the two implicit methods (Gear’s and Adams–Moulton) exhibit much larger regions of absolute stability than does the explicit method (Adams–Bashforth). Gear’s region of absolute stability contains nearly the entire left half $h\lambda$ plane; thus, for any stable dynamic system, the step size can be chosen as large as desired without any consideration for numerical stability. Even the Adams–Moulton region of absolute stability is quite large compared to the Adams–Bashforth. Typically, the region of stability of explicit methods is much smaller than the region of absolute stability of corresponding order implicit methods. For this reason, implicit methods are frequently used in commercial integration packages so that the integration step size can be chosen based purely on the local truncation error criteria. The region of absolute stability shrinks as the order of the method increases, whereas the accuracy increases. This is a trade-off between accuracy, stability, and numerical efficiency in integration algorithms. ■

5.5 Stiff Systems

Gear’s methods were originally developed for the solution of systems of *stiff* ordinary differential equations. Stiff systems are systems that exhibit a wide range of time varying dynamics from “very fast” to “very slow.” A stiff linear system exhibits eigenvalues that span several orders of magnitude. A nonlinear system is stiff if its associated Jacobian matrix exhibits widely separated eigenvalues when evaluated at the operating points of interest. For efficient and accurate solutions of stiff differential equations, it is desirable for a multi-step method to be “stiffly stable.” An appropriate integration algorithm will allow the step size to be varied over a wide range of values and yet will remain numerically stable. A stiffly stable method exhibits the three stability regions shown in Figure 5.10 such that:

1. Region I is a region of absolute stability
2. Region II is a region of accuracy and stability
3. Region III is a region of accuracy and relative stability

Only during the initial period of the solution of the ODE do the large negative eigenvalues significantly impact the solution, yet they must be accounted for throughout the whole solution. Large negative eigenvalues ($\lambda < 0$) will decay rapidly by a factor of $1/e$ in time $1/\lambda$. If $h\lambda = \gamma + j\beta$, then the change in magnitude in one step is e^γ . If $\gamma \leq \delta \leq 0$, where δ defines the interface

**FIGURE 5.10**

Regions required for stiff stability

between Regions I and II, then the component is reduced by at least e^δ in one step. After a finite number of steps, the impact of the fast components is negligible and their numerical accuracy is unimportant. Therefore, the integration method is required to be absolutely stable in Region I.

Around the origin, numerical accuracy becomes more significant and relative or absolute stability is required. A region of *relative stability* consists of those values of $h\lambda$ for which the extraneous eigenvalues of the characteristic polynomial of Equation (5.97) are less in magnitude than the principal eigenvalue. The principal eigenvalue is the eigenvalue which governs the system response most closely. If the method is relatively stable in Region III, then the system response will be dominated by the principal eigenvalue in that region. If $\gamma > \alpha > 0$, one component of the system response is increasing by at least e^α one step. This increase must be limited by choosing the step sizes small enough to track this change.

If $\|\beta\| > \theta$, there are at least $\theta/2\pi$ complete cycles of oscillation in one step. Except in Region I, where the response is rapidly decaying and where $\gamma > \alpha$ is not used, the oscillatory responses must be captured. In practice, it is customary to have eight or more time points per cycle (to accurately capture the magnitude and frequency of the oscillation); thus θ is chosen to be bounded by $\pi/4$ in Region II.

Examination of the family of Adams–Bashforth methods shows that they all fail to satisfy the criteria to be stiffly stable and are not suitable for integrating stiff systems. Only the first- and second-order Adams–Moulton (backward Euler and trapezoidal, respectively) satisfy the stiffly stable criteria. Gear's algorithms, on the other hand, were developed specifically to address stiff

system integration [15]. Gear's algorithms up to order six satisfy the stiff properties with the following choice of δ [8]:

Order	δ
1	0
2	0
3	0.1
4	0.7
5	2.4
6	6.1

Example 5.7

Compare the application of the third-order Adams–Bashforth, Adams–Moulton, and Gear's method, to the integration of the following system:

$$\dot{x}_1 = 48x_1 + 98x_2 \quad x_1(0) = 1 \quad (5.107)$$

$$\dot{x}_2 = -49x_1 - 99x_2 \quad x_2(0) = 0 \quad (5.108)$$

Solution 5.7 The exact solution to Example 5.7 is

$$x_1(t) = 2e^{-t} - e^{-50t} \quad (5.109)$$

$$x_2(t) = -e^{-t} + e^{-50t} \quad (5.110)$$

This solution is shown in Figure 5.11. Both states contain both fast and slow components, with the fast component dominating the initial response and the slow component dominating the longer-term dynamics. Since Gear's, Adams–Bashforth, and Adams–Moulton methods are multistep methods, each method is initialized using the absolutely stable trapezoidal method for the first two to three steps using a small step size.

The Adams–Bashforth algorithm with a time step of 0.0111 seconds is shown in Figure 5.12. Note that, even with a small step size of 0.0111 seconds, the inherent error in the integration algorithm eventually causes the system response to exhibit numerical instabilities. The step size can be decreased to increase the stability properties, but this requires more time steps in the integration window ($t \in [0, 2]$) than is computationally necessary.

The Adams–Moulton response to the stiff system for an integration step size of 0.15 seconds is shown in Figure 5.13. Although a much larger time step can be used for integration as compared to the Adams–Bashforth algorithm, the Adams–Moulton algorithm does not exhibit numerical absolute stability. For an integration step size of $h = 0.15$ seconds, the Adams–Moulton algorithm exhibits numerical instability. Note that the solution is oscillating with growing magnitude around the exact solution.

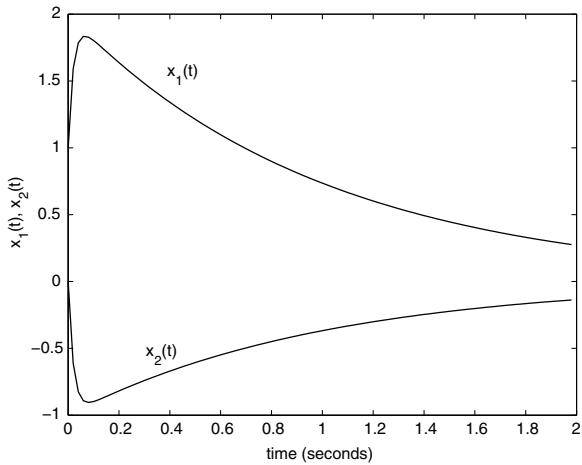


FIGURE 5.11
Stiff system response

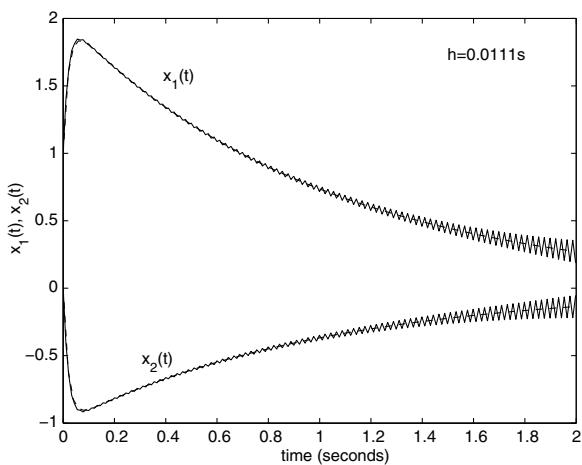
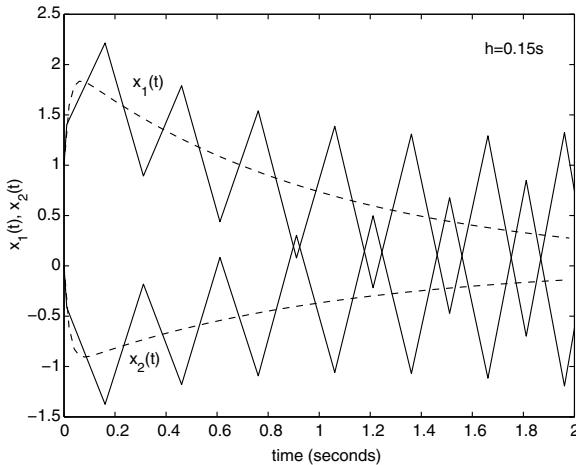


FIGURE 5.12
Adams-Basforth stiff system response with $h = 0.0111$ s step size

**FIGURE 5.13**

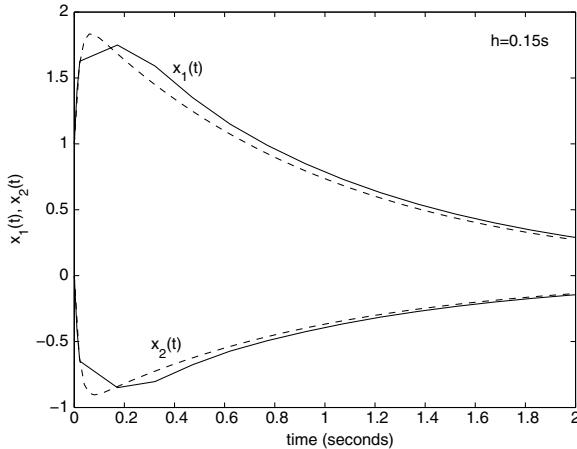
Adams–Moulton stiff system response with $h = 0.15$ s step size

Figure 5.14 shows the Gear's method numerical solution to the stiff system using an integration step size of 0.15 s, which is the same step size as used in the Adams–Moulton response of Figure 5.13. Gear's method is numerically stable since the global error decreases with increasing time.

The comparison of the three integration methods supports the necessity of using integration algorithms developed specifically for stiff systems. Even though both the third-order Adams–Bashforth and Adams–Moulton have regions of absolute stability, these regions are not sufficient to ensure accurate stiff system integration. ■

5.6 Step Size Selection

For computational efficiency, it is desirable to choose the largest integration step size possible while satisfying a predetermined level of accuracy. The level of accuracy can be maintained by constant vigilance of the local truncation error of the method if the chosen method is numerically stable. If the function $x(t)$ is varying rapidly, the step size should be chosen small enough to capture the significant dynamics. Conversely, if the function $x(t)$ is not varying significantly (nearly linear) over a finite time interval, then the integration time steps can be chosen quite large over that interval while still maintaining numerical accuracy. The true challenge to a numerical integration algorithm is when the dynamic response of $x(t)$ has intervals of both rapid variance and latency. In this case, it is desirable for the integration step size to have

**FIGURE 5.14**

Gear's stiff system response with $h = 0.15$ s step size

the ability to increase or decrease throughout the simulation interval. This can be accomplished by choosing the integration step size based on the local truncation error bounds.

Consider the trapezoidal method absolute local truncation error:

$$\varepsilon_T = \frac{1}{12} h^3 x^{(3)}(\tau) \quad (5.111)$$

The local truncation error is dependent on the integration step size h and the third derivative of the function $x^{(3)}(\tau)$. If the local truncation error is chosen to be in the interval

$$B_L \leq \varepsilon \leq B_U \quad (5.112)$$

where B_L and B_U represent the prespecified lower and upper bounds, respectively, then the integration step size can be bounded by

$$h \leq \sqrt[3]{\frac{12B_U}{x^{(3)}(\tau)}} \quad (5.113)$$

If $x(t)$ is rapidly varying, the $x^{(3)}(\tau)$ will be large and h must be chosen small to satisfy $\varepsilon \leq B_U$, whereas, if $x(t)$ is not varying rapidly, then $x^{(3)}(\tau)$ will be small and h can be chosen relatively large while still satisfying $\varepsilon \leq B_U$. This leads to the following procedure for calculating the integration step size:

Integration Step Size Selection

Attempt an integration step size h_{n+1} to calculate x_{n+1} from x_n, x_{n-1}, \dots

1. Using x_{n+1} , calculate the local truncation error ε_T . If the size of x is greater than 1, then
- $$\varepsilon_T = \max_i (|\varepsilon_{T,i}|)$$
2. If $B_L \leq \varepsilon_T \leq B_U$, then PASS, accept $h_{n+1}, h_{next} = h_{n+1}$, and continue.
 3. If $\varepsilon > B_U$, then FAIL (h_{n+1} is too large), set $h_{n+1} = \alpha h_{n+1}$, repeat integration for x_{n+1} .
 4. If $\varepsilon \leq B_L$, then PASS, accept h_{n+1} , set $h_{next} = \alpha h_{n+1}$, and continue.

where

$$\alpha = \left[\frac{B_{avg}}{\varepsilon_T} \right]^{\frac{1}{k+1}} \quad (5.114)$$

where $B_L \leq B_{avg} \leq B_U$ and k is the degree of the method.

Commercial integration packages may implement an algorithm that is slightly different. In these packages, if the local truncation error is smaller than the lower bound, then the attempted integration step size also FAILS, and the integration step is reattempted with a larger integration step size. Once again, there is a trade off between the time spent in recalculating x_{n+1} with a larger step size and the additional computational effort acquired by simply accepting the current value and continuing on.

The difficulty in implementing this step size selection approach is the calculation of the higher-order derivatives of $x(t)$. Since $x(t)$ is not known analytically, the derivatives must be calculated numerically. One common approach is to use *difference methods* to approximate the derivatives. The $(k+1)$ st derivative to $x(\tau)$ is approximated by

$$x^{(k+1)}(\tau) \approx (k+1)! \nabla_{k+1} x_{n+1} \quad (5.115)$$

where $\nabla_{k+1} x_{n+1}$ is found recursively:

$$\nabla_1 x_{n+1} = \frac{x_{n+1} - x_n}{t_{n+1} - t_n} \quad (5.116)$$

$$\nabla_1 x_n = \frac{x_n - x_{n-1}}{t_n - t_{n-1}} \quad (5.117)$$

$$\vdots \quad (5.118)$$

$$\nabla_2 x_{n+1} = \frac{\nabla_1 x_{n+1} - \nabla_1 x_n}{t_{n+1} - t_{n-1}} \quad (5.119)$$

$$\nabla_2 x_n = \frac{\nabla_1 x_n - \nabla_1 x_{n-1}}{t_n - t_{n-2}} \quad (5.120)$$

$$\vdots \quad (5.121)$$

$$\nabla_{k+1}x_{n+1} = \frac{\nabla_k x_{n+1} - \nabla_k x_n}{t_{n+1} - t_{n-k}} \quad (5.122)$$

Example 5.8

Find an expression for step size selection for the trapezoidal integration method.

Solution 5.8

The LTE for the trapezoidal method is

$$|\varepsilon_T| = \frac{1}{12}h^3x^{(3)}(\tau) \quad (5.123)$$

$$= \frac{1}{2}h^3\nabla_3x_{n+1} \quad (5.124)$$

where

$$\nabla_3x_{n+1} = \frac{\nabla_2x_{n+1} - \nabla_2x_n}{t_{n+1} - t_{n-2}} \quad (5.125)$$

$$= \frac{\nabla_2x_{n+1} - \nabla_2x_n}{h_{n+1} + h_n + h_{n-1}} \quad (5.126)$$

$$= \frac{1}{h_{n+1} + h_n + h_{n-1}} \left\{ \frac{1}{h_{n+1} + h_n} \left[\frac{x_{n+1} - x_n}{h_{n+1}} - \frac{x_n - x_{n-1}}{h_n} \right] - \frac{1}{h_n + h_{n-1}} \left[\frac{x_n - x_{n-1}}{h_n} - \frac{x_{n-1} - x_{n-2}}{h_{n-1}} \right] \right\} \quad (5.127)$$

■

Example 5.9

Numerically solve

$$\dot{x}(t) = \begin{bmatrix} -1 & 1 & 1 \\ 1 & -25 & 98 \\ 1 & -49 & -60 \end{bmatrix} x, \quad \text{with } x(0) = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix}$$

using a variable integration step size with the following parameters:

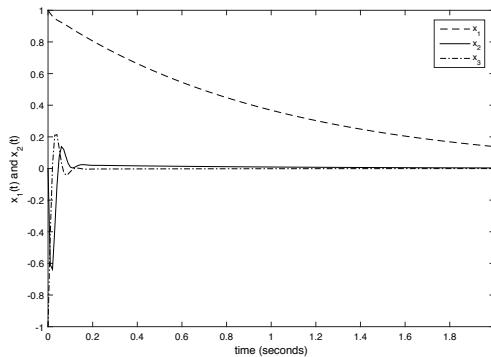
$$\begin{aligned} h_0 &= 0.01 & B_L &= 10^{-6} \\ B_U &= 10^{-4} & B_{avg} &= 10^{-5} \end{aligned}$$

Solution 5.9 This is a stiff system with the following eigenvalues:

$$-0.9788$$

$$-42.5106 + j67.0420$$

$$-42.5106 - j67.0420$$

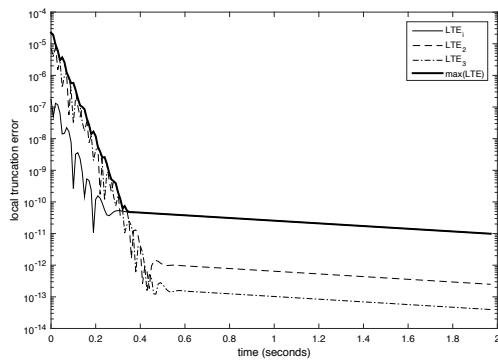
**FIGURE 5.15**

Numerical solution for Example 5.9 with constant step size $h = 0.01$ s

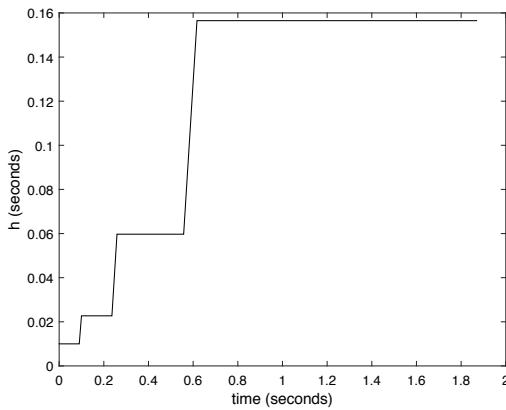
A constant step size integration yields the numerical solution in Figure 5.15 with local truncation error shown in Figure 5.16.

Note that the local truncation error continually decreases throughout the simulation. The LTE is initially dominated by the LTE of x_2 and x_3 due to the influence of the oscillatory response resulting from the complex eigenvalues. As these states decay as a result of the large negative real part of the complex eigenvalues, the LTE of the real eigenvalue x_1 begins to dominate. However, the LTE decreases throughout the simulation interval as the rate of change of all states decreases.

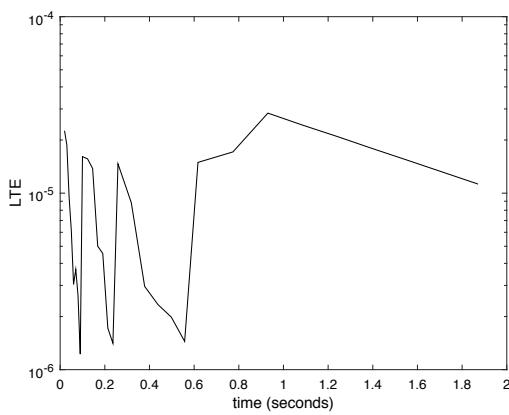
After implementing the step size selection process based on local truncation error, the variable step size and corresponding LTE are shown in Figures 5.17 and 5.18. Note that using the variable step size, the LTE is held between the specified upper and lower bounds. This is accomplished by increasing the step size when the LTE becomes too small in response to the decay of the response rates of the underlying states. Figure 5.19 shows the numerical results of the states using the variable step size. Note that the simulation results differ very little from the original results of Figure 5.15 because the error is not allowed to exceed the upper bound. ■

**FIGURE 5.16**

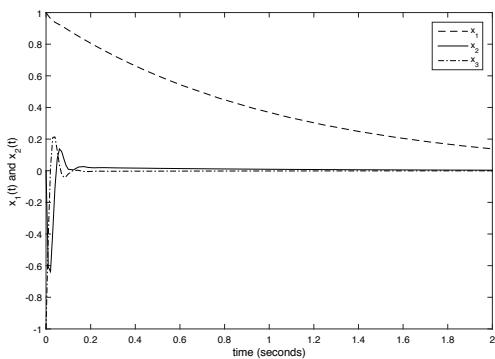
LTE solution for Example 5.9 with constant step size $h = 0.01$ s

**FIGURE 5.17**

Step sizes as varied through the step size selection process

**FIGURE 5.18**

LTE solution for Example 5.9 with variable step size

**FIGURE 5.19**

Numerical solution for Example 5.9 with variable step sizes

5.7 Differential-Algebraic Equations

Many classes of systems can be generically written in the form

$$F(t, y(t), y'(t)) = 0 \quad (5.128)$$

where F and $y \in R^m$. In some cases, Equation (5.128) can be rewritten as

$$F(t, x(t), \dot{x}(t), y(t)) = 0 \quad (5.129)$$

$$g(t, x(t), y(t)) = 0 \quad (5.130)$$

In this form, the system of equations is typically known as a system of *differential-algebraic equations*, or DAEs [7]. This form of DAEs may be considered to be a system of differential equations (Equation (5.129)) constrained to an algebraic manifold (Equation (5.130)). Often the set of equations (5.130) are invertible (i.e., $y(t)$ can be obtained from $y(t) = g^{-1}(t, x(t))$) and Equation (5.129) can be rewritten as an ordinary differential equation:

$$F(t, x(t), \dot{x}(t), g^{-1}(t, x(t))) = f(t, x(t), \dot{x}(t)) = 0 \quad (5.131)$$

Although the set of ODEs may be conceptually simpler to solve, there are usually several compelling reasons to leave the system in its original DAE form. Many DAE models are derived from physical problems in which each variable has individual characteristics and physical significance. Converting the DAE into an ODE may result in a loss of physical information in the solution. Furthermore, it may be more computationally expensive to obtain the solution of the system of ODEs since inherent sparsity may have been lost. Solving the original DAE directly provides greater information regarding the behavior of the system.

A special case of DAEs is the semi-explicit DAE:

$$\dot{x} = f(x, y, t) \quad x \in R^n \quad (5.132)$$

$$0 = g(x, y, t) \quad y \in R^m \quad (5.133)$$

where y has the same dimension as g . DAE systems that can be written in semi-explicit form are often referred to as *index 1* DAE systems [7]. The first concerted effort to solve semi-explicit DAEs was proposed in [16] and later refined in [17], and consisted of replacing $\dot{x}(t)$ by a k -step backwards difference formula (BDF) approximation

$$\dot{x}(t) \approx \frac{\rho_n x_n}{h_n} = \frac{1}{h_n} \sum_{i=0}^k \alpha_i x_{n-i} \quad (5.134)$$

and then solving the resulting equations

$$\rho_n x_n = h_n f(x_n, y_n, t_n) \quad (5.135)$$

$$0 = g(x_n, y_n, t_n) \quad (5.136)$$

for approximations to x_n and y_n .

Various other numerical integration techniques have been studied for application to DAE systems. Variable step size/fixed formula code has been proposed to solve the system of DAEs [48]. Specifically, a classic fourth-order Runge–Kutta method was used to solve for x and a third-order BDF to solve for y .

In general, however, many DAE systems may be solved by any multistep numerical integration method which is convergent when applied to ODEs [22]. The application of multistep integration methods to a DAE system is straightforward. A general multistep method of the form of Equation (5.8) applied to Equations (5.132) and (5.133) yields

$$x_{n+1} = \sum_{i=0}^p a_i x_{n-i} + h \sum_{i=-1}^p b_i f(x_{n-i}, y_{n-i}, t_{n-i}) \quad (5.137)$$

$$0 = g(x_{n+1}, y_{n+1}, t_{n+1}) \quad (5.138)$$

Multistep methods applied to semi-explicit index 1 DAEs are stable and convergent to the same order of accuracy for the DAE as for standard nonstiff ODEs [7].

There are two basic approaches to solving Equations (5.137) and (5.138). One is the *iterative* approach in which Equation (5.138) is solved for y_{n+1} , which is then substituted into Equation (5.137). This equation is then solved for x_{n+1} . This process is repeated at t_{n+1} until the values for x_{n+1} and y_{n+1} converge, and then the solution is advanced to the next time step.

Another approach is the *simultaneous* approach, in which both sets of equations are solved simultaneously for x_{n+1} and y_{n+1} using a nonlinear solver such as the Newton–Raphson method. In this case, the system of equations is recast as

$$\begin{aligned} 0 &= F(x_{n+1}, y_{n+1}, t_{n+1}) \\ &= x_{n+1} - \sum_{i=0}^p a_i x_{n-i} - h \sum_{i=-1}^p b_i f(x_{n-i}, y_{n-i}, t_{n-i}) \end{aligned} \quad (5.139)$$

$$0 = g(x_{n+1}, y_{n+1}, t_{n+1}) \quad (5.140)$$

and the Newton–Raphson Jacobian becomes

$$J_{xy} = \begin{bmatrix} I_n - hb_{-1} \frac{\partial f}{\partial x} & -hb_{-1} \frac{\partial f}{\partial y} \\ \frac{\partial g}{\partial x} & \frac{\partial g}{\partial y} \end{bmatrix} \quad (5.141)$$

Solving the full $n + m$ system of equations yields both x_{n+1} and y_{n+1} simultaneously.

The solvability of a DAE system implies that there exists a unique solution for sufficiently different inputs and *consistent* initial conditions, unlike the ODE system, which has a unique solution for *arbitrary* initial conditions. For

a DAE system, there exists only one set of initial conditions to the nonstate (algebraic) variables. The initial values that are consistent with the system input are called *admissible* initial values. The admissible initial values are those which satisfy

$$y_0 = g^{-1}(x_0, t_0) \quad (5.142)$$

5.8 Power System Applications

Power systems in general are considered to be large scale, often involving several hundred equations to describe the behavior of an interconnected power system during and following a fault on the system. As power system operations become increasingly complex, it becomes necessary to be able to perform analyses of voltage conditions and system stability. A medium-sized power company serving a mixed urban and rural population of two to three million people operates a network that may typically contain hundreds of buses and thousands of transmission lines, excluding the distribution system [13]. Under certain assumptions, such as instantaneous transmission lines, interconnected power systems are often modeled in DAE form with over 1000 differential equations and 10,000 algebraic constraint equations. One traditional approach to solving large-scale systems of this type has been to replace the full system model with a reduced-order state space model.

5.8.1 Transient Stability Analysis

The “classical model” of a synchronous machine is often used to study the transient stability of a power system during the period of time in which the system dynamics depend largely on the stored kinetic energy in the rotating masses. This is usually on the order of a second or two. The classical model is derived under several simplifying assumptions [2]:

1. Mechanical power input, P_m , is constant.
2. Damping is negligible.
3. The constant voltage behind the transient reactance model for the synchronous machines is valid.
4. The rotor angle of the machine coincides with the voltage behind the transient reactance angle.
5. Loads are represented as constant impedances.

The equations of motion are given by

$$\dot{\omega}_i = \frac{1}{M_i} \left(P_{m_i} - E_i^2 G_{ii} - E_i \sum_{j \neq i}^n E_j (B_{ij} \sin \delta_{ij} + G_{ij} \cos \delta_{ij}) \right) \quad (5.143)$$

$$\dot{\delta}_i = \omega_i - \omega_s \quad i = 1, \dots, n \quad (5.144)$$

where n is the number of machines, ω_s is the synchronous angular frequency, $\delta_{ij} = \delta_i - \delta_j$, $M_i = \frac{2H_i}{\omega_s}$, and H_i is the inertia constant in seconds. B_{ij} and G_{ij} are elements of the reduced admittance matrix Y at the internal nodes of the machine. The loads are modeled as constant impedances, which are then absorbed into the admittance matrix. The classical model is appropriate for frequency studies that result from faults on the transmission system for the first or second swing of the rotor angle. The procedure for setting up a transient stability analysis is given below.

Transient Stability Analysis

1. Perform a load flow analysis to obtain system voltages, angles, active and reactive power generation.
2. For each generator i, \dots, n , in the system, calculate the internal voltage and initial rotor angle $E \angle \delta_0$:

$$I_{gen}^* = \frac{(P_{gen} + jQ_{gen})}{V_T \angle \theta_T} \quad (5.145)$$

$$E \angle \delta_0 = jx'_d I_{gen} + V_T \angle \theta_T \quad (5.146)$$

where $P_{gen} + jQ_{gen}$ are the generated active and reactive power obtained from the power flow solution and $V_T \angle \theta_T$ is the generator terminal voltage.

3. For each load $1, \dots, m$ in the system, convert the active and reactive power loads to admittances:

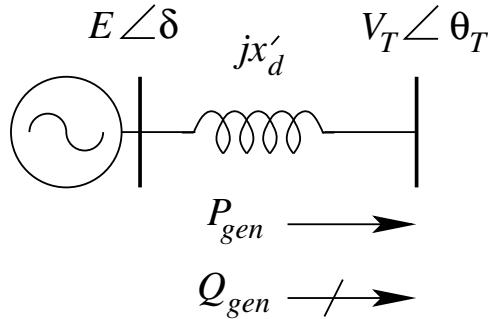
$$Y_L = G_L + jB_L = \frac{I_L}{V_L \angle \theta_L} \quad (5.147)$$

$$= \frac{S_L^*}{V_L^2} \quad (5.148)$$

$$= \frac{P_L - jQ_L}{V_L^2} \quad (5.149)$$

Add the shunt admittance Y_L to the corresponding diagonal of the admittance matrix.

4. For each generator in the system, augment the admittance matrix by adding an internal bus connected to the terminal bus with the transient reactance x'_d , as shown in Figure 5.20.

**FIGURE 5.20**

Voltage behind transient reactance

Then let

$$Y_{nn} = \begin{bmatrix} jx'_{d_1} & 0 & 0 & \dots & 0 \\ 0 & jx'_{d_2} & 0 & \dots & 0 \\ 0 & 0 & jx'_{d_3} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & jx'_{d_n} \end{bmatrix}$$

and

$$\begin{aligned} Y_{nm} &= [[-Y_{nn}] [0_{n \times m}]] \\ Y_{mn} &= Y_{nm}^T \end{aligned}$$

where $[0_{n \times m}]$ is an $(n \times m)$ matrix of zeros.

Further, let

$$Y_{mm} = \left[Y_{original} + \begin{bmatrix} Y_{nn} & 0_{n \times (m-n)} \\ 0_{(m-n) \times n} & 0_{(m-n) \times (m-n)} \end{bmatrix} \right]$$

This matrix layout assumes that the systems buses have been ordered such that the generators are numbered $1, \dots, n$, with the remaining load buses numbered $n+1, \dots, m$.

5. The reduced admittance matrix Y_{red} is found by

$$Y_{red} = [Y_{nn} - Y_{nm} Y_{mm}^{-1} Y_{mn}] \quad (5.150)$$

$$= G_{red} + jB_{red} \quad (5.151)$$

The reduced admittance matrix is now $n \times n$ where n is the number of generators, whereas the original admittance matrix was $m \times m$.

6. Repeat steps 4 and 5 to calculate the fault-on reduced admittance matrix and the postfault reduced admittance matrix (if different from the prefault reduced admittance matrix).
7. For $0 < t \leq t_{apply}$, integrate the transient stability equations (5.143) and (5.144) with the integration method of choice using the prefault reduced admittance matrix, where t_{apply} is the time the system fault is applied. In many applications, $t_{apply} = 0$.
8. For $t_{apply} < t \leq t_{clear}$, integrate the transient stability equations (5.143) and (5.144) with the integration method of choice using the fault-on reduced admittance matrix, where t_{clear} is the time the system fault is cleared.
9. For $t_{clear} < t \leq t_{max}$, integrate the transient stability equations (5.143) and (5.144) with the integration method of choice using the postfault reduced admittance matrix, where t_{max} is end of the simulation interval.

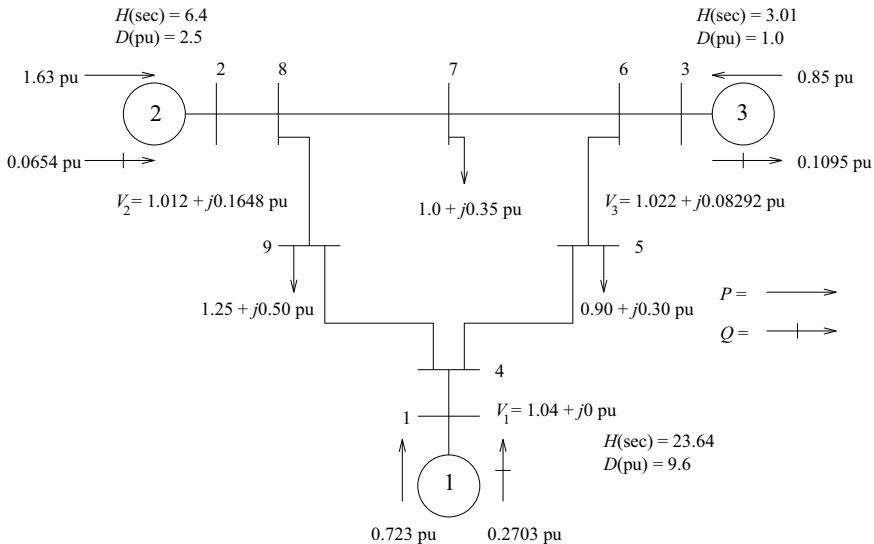
At the end of the simulation, the state variables (δ_i, ω_i) for each of the generators may be plotted against time. The rotor angle responses are in radians and may be converted to degrees if desired. The rotor angular frequency is in radians per second and may be converted to hertz (cycles per second) if desired. These waveforms may then be analyzed to determine whether or not system stability is maintained. If the system responses diverge from one another or exhibit growing oscillations, then the system is most probably unstable.

Example 5.10

For the three-machine, nine-bus system shown in Figure 5.21, a solid three-phase fault occurs at bus 8 at 0.1 seconds. The fault is cleared by opening line 8-9. Determine whether or not the system will remain stable for a clearing time of 0.12 seconds after the application of the fault.

Solution 5.10 Following the procedure for transient stability analysis outlined previously, the first step is to perform a power flow analysis of the system. The line and bus data are given in Figure 5.21. The load flow results are

i	V	θ	P_{gen}	Q_{gen}
1	1.0400	0	0.7164	0.2685
2	1.0253	9.2715	1.6300	0.0669
3	1.0254	4.6587	0.8500	-0.1080
4	1.0259	-2.2165		
5	1.0128	-3.6873		
6	1.0327	1.9625		
7	1.0162	0.7242		
8	1.0261	3.7147		
9	0.9958	-3.9885		

**FIGURE 5.21**

Three-machine, nine-bus system of Example 5.10

where the bus angles are given in degrees and all other data are in per unit. The admittance matrix for this system is given in Figure 5.22.

The generator data for this system are

i	x'_d	H
1	0.0608	23.64
2	0.1198	6.40
3	0.1813	3.01

The internal voltages and rotor angles for each generator are computed using the generated active and reactive powers, voltage magnitudes, and angles as:

$$I_1^* = \frac{(0.7164 + j0.2685)}{1.0400\angle 0^\circ} = 0.6888 + j0.2582$$

$$E_1\angle\delta_1 = (j0.0608)(0.6888 - j0.2582) + 1.0400\angle 0^\circ = 1.0565\angle 2.2718^\circ$$

$$I_2^* = \frac{(1.6300 + j0.0669)}{1.0253\angle 9.2715^\circ} = 1.5795 - j0.1918$$

$$E_2\angle\delta_2 = (j0.1198)(1.5795 + j0.1918) + 1.0253\angle 9.2715^\circ = 1.0505\angle 19.7162^\circ$$

$$I_3^* = \frac{(0.8500 - j0.1080)}{1.0254\angle 4.6587^\circ} = 0.8177 - j0.1723$$

$$[Y_{ms}] = \begin{bmatrix} 17.3611 \angle -90.00^\circ & 0 & 0 & 0 & 0 & 0 \\ 0 & 16.0000 \angle -90.00^\circ & 0 & 0 & 0 & 0 \\ 0 & 0 & 17.0648 \angle -90.00^\circ & 0 & 0 & 0 \\ 0 & 0 & 0 & 17.0648 \angle 90.00^\circ & 0 & 0 \\ 17.3611 \angle 90.00^\circ & 0 & 0 & 39.4478 \angle -85.19^\circ & 10.6886 \angle 100.47^\circ & 0 \\ 0 & 0 & 0 & 33.8085 \angle -90.00^\circ & 16.1657 \angle -78.50^\circ & 5.7334 \angle 102.92^\circ \\ 0 & 0 & 0 & 17.0648 \angle 90.00^\circ & 0 & 5.7334 \angle 102.92^\circ \\ 0 & 0 & 0 & 0 & 32.2461 \angle -85.67^\circ & 9.8522 \angle 96.73^\circ \\ 0 & 0 & 0 & 0 & 0 & 9.8522 \angle 96.73^\circ \\ 0 & 16.0000 \angle 90.00^\circ & 0 & 0 & 0 & 13.7931 \angle -83.22^\circ \\ 0 & 0 & 0 & 0 & 0 & 13.7931 \angle 96.73^\circ \\ 0 & 0 & 0 & 0 & 0 & 35.5564 \angle -85.48^\circ \\ 0 & 0 & 0 & 0 & 0 & 6.0920 \angle 101.24^\circ \\ 0 & 0 & 0 & 0 & 0 & 17.5252 \angle -81.62^\circ \end{bmatrix}$$

FIGURE 5.22

Admittance matrix for Example 5.10

$$E_3 \angle \delta_3 = (j0.1813)(0.8177 + j0.1723) + 1.0254 \angle 4.6587^\circ = 1.0174 \angle 13.1535^\circ$$

The next step is to convert the loads to equivalent impedances:

$$\begin{aligned} G_5 + jB_5 &= \frac{(0.90 - j0.30)}{1.0128^2} = 0.8773 - j0.2924 \\ G_7 + jB_7 &= \frac{(1.00 - j0.35)}{1.0162^2} = 0.9684 - j0.3389 \\ G_9 + jB_9 &= \frac{(1.25 - j0.50)}{0.9958^2} = 1.2605 - j0.5042 \end{aligned}$$

These values are added to the diagonal of the original admittance matrix.

The reduced admittance matrices can now be computed as outlined in steps 4 and 5 above. The prefault admittance matrix is

$$Y_{red}^{prefault} = \begin{bmatrix} 0.8453 - j2.9881 & 0.2870 + j1.5131 & 0.2095 + j1.2257 \\ 0.2870 + j1.5131 & 0.4199 - j2.7238 & 0.2132 + j1.0880 \\ 0.2095 + j1.2257 & 0.2132 + j1.0880 & 0.2769 - j2.3681 \end{bmatrix}$$

The fault-on matrix is found similarly, except that the Y_{mm} is altered to reflect the fault on bus 8. The solid three-phase fault is modeled by shorting the bus to ground. In the admittance matrix, the row and column corresponding to bus 8 are removed. The lines between bus 8 and adjacent buses are now connected to ground; thus they will still appear in the original admittance diagonals. The column of Y_{nm} and the row of Y_{mn} corresponding to bus 8 must also be removed. The matrix Y_{nn} remains unchanged. The fault-on reduced admittance matrix is

$$Y_{red}^{fault-on} = \begin{bmatrix} 0.6567 - j3.8159 & 0 & 0.0701 + j0.6306 \\ 0 & 0 - j5.4855 & 0 \\ 0.0701 + j0.6306 & 0 & 0.1740 - j2.7959 \end{bmatrix}$$

The postfault reduced admittance matrix is computed in much the same way, except that line 8-9 is removed from Y_{mm} . The elements of Y_{mm} are updated to reflect the removal of the line:

$$Y_{mm}(8, 8) = Y_{mm}(8, 8) + Y_{mm}(8, 9)$$

$$Y_{mm}(8, 9) = Y_{mm}(9, 9) + Y_{mm}(8, 9)$$

$$Y_{mm}(8, 9) = 0$$

$$Y_{mm}(9, 8) = 0$$

Note that the diagonals must be updated before the off-diagonals are zeroed out. The postfault reduced admittance is then computed:

$$Y_{red}^{postfault} = \begin{bmatrix} 1.1811 - j2.2285 & 0.1375 + j0.7265 & 0.1909 + j1.0795 \\ 0.1375 + j0.7265 & 0.3885 - j1.9525 & 0.1987 + j1.2294 \\ 0.1909 + j1.0795 & 0.1987 + j1.2294 & 0.2727 - j2.3423 \end{bmatrix}$$

These admittance matrices are then ready to be substituted into the transient stability equations at the appropriate time in the simulation.

Applying the trapezoidal algorithm to the transient stability equations yields the following system of equations:

$$\delta_1(n+1) = \delta_1(n) + \frac{h}{2} [\omega_1(n+1) - \omega_s + \omega_1(n) - \omega_s] \quad (5.152)$$

$$\omega_1(n+1) = \omega_1(n) + \frac{h}{2} [f_1(n+1) + f_1(n)] \quad (5.153)$$

$$\delta_2(n+1) = \delta_2(n) + \frac{h}{2} [\omega_2(n+1) - \omega_s + \omega_2(n) - \omega_s] \quad (5.154)$$

$$\omega_2(n+1) = \omega_2(n) + \frac{h}{2} [f_2(n+1) + f_2(n)] \quad (5.155)$$

$$\delta_3(n+1) = \delta_3(n) + \frac{h}{2} [\omega_3(n+1) - \omega_s + \omega_3(n) - \omega_s] \quad (5.156)$$

$$\omega_3(n+1) = \omega_3(n) + \frac{h}{2} [f_3(n+1) + f_3(n)] \quad (5.157)$$

where

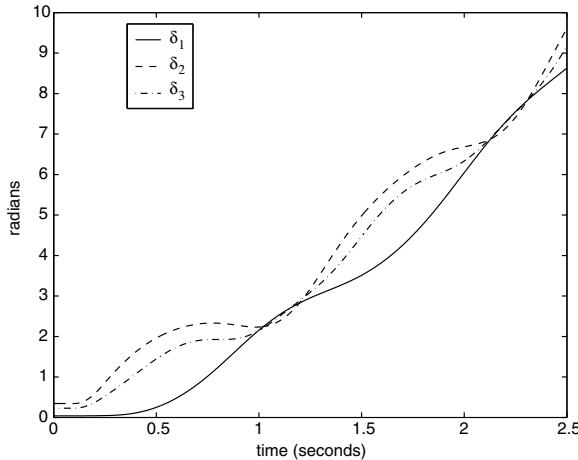
$$f_i(n+1) = \frac{1}{M_i} \left(P_{m_i} - E_i^2 G_{ii} - E_i \sum_{j \neq i}^n E_j (B_{ij} \sin \delta_{ij}(n+1) + G_{ij} \cos \delta_{ij}(n+1)) \right) \quad (5.158)$$

Since the transient stability equations are nonlinear and the trapezoidal method is an implicit method, they must be solved iteratively using the Newton–Raphson method at each time point. The iterative equations are

$$\begin{aligned} & \left[I - \frac{h}{2} [J(n+1)^k] \right] \begin{bmatrix} \delta_1(n+1)^{k+1} - \delta_1(n+1)^k \\ \omega_1(n+1)^{k+1} - \omega_1(n+1)^k \\ \delta_2(n+1)^{k+1} - \delta_2(n+1)^k \\ \omega_2(n+1)^{k+1} - \omega_2(n+1)^k \\ \delta_3(n+1)^{k+1} - \delta_3(n+1)^k \\ \omega_3(n+1)^{k+1} - \omega_3(n+1)^k \end{bmatrix} = \\ & - \left[\begin{bmatrix} \delta_1^k(n+1) \\ \omega_1^k(n+1) \\ \delta_2^k(n+1) \\ \omega_2^k(n+1) \\ \delta_3^k(n+1) \\ \omega_3^k(n+1) \end{bmatrix} - \begin{bmatrix} \delta_1(n) \\ \omega_1(n) \\ \delta_2(n) \\ \omega_2(n) \\ \delta_3(n) \\ \omega_3(n) \end{bmatrix} \right] - \frac{h}{2} \begin{bmatrix} \omega_1^k(n+1) + \omega_1(n) - 2\omega_s \\ f_1^k(n+1) + f_1(n) \\ \omega_2^k(n+1) + \omega_2(n) - 2\omega_s \\ f_2^k(n+1) + f_2(n) \\ \omega_3^k(n+1) + \omega_3(n) - 2\omega_s \\ f_3^k(n+1) + f_3(n) \end{bmatrix} \end{aligned} \quad (5.159)$$

where

$$[J] = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ \frac{\partial f_1}{\partial \delta_1} & 0 & \frac{\partial f_1}{\partial \delta_2} & 0 & \frac{\partial f_1}{\partial \delta_3} & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ \frac{\partial f_2}{\partial \delta_1} & 0 & \frac{\partial f_2}{\partial \delta_2} & 0 & \frac{\partial f_2}{\partial \delta_3} & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ \frac{\partial f_3}{\partial \delta_1} & 0 & \frac{\partial f_3}{\partial \delta_2} & 0 & \frac{\partial f_3}{\partial \delta_3} & 0 \end{bmatrix} \quad (5.160)$$

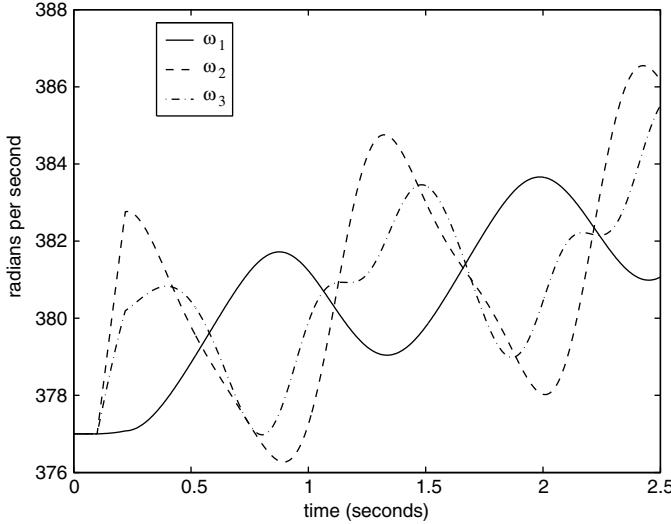
**FIGURE 5.23**

Rotor angle response for Example 5.10

Note that LU factorization must be employed to solve the discretized equations. These equations are iterated at each time point until convergence of the Newton–Raphson algorithm. The fault-on and postfault matrices are substituted in at the appropriate times in the integration. The simulation results are shown in Figures 5.23 and 5.24 for the rotor angles and angular frequencies, respectively. From the waveforms shown in these figures, it can be concluded that the system remains stable since the waveforms do not diverge during the simulation interval. ■

5.8.2 Midterm Stability Analysis

Reduction processes frequently destroy the natural physical structure and sparsity of the full-order system. Numerical solution algorithms which make use of structure and sparsity for efficiency perform poorly on the reduced-order system even though the reduced-order system is still quite large. After the first few seconds of a power system disturbance, the classical model representation may no longer be valid due to the dynamic behavior of the automatic voltage regulator, the turbine/governor system, under-load-tap-changing transformers, and the dynamic nature of some system loads. For midterm stability analyses, a more detailed model is required to capture a wider range of system behavior. Since the behavior of loads may significantly impact the stability of the system, it is desirable to be able to retain individual load buses during the simulation. This type of model is often referred to as a “structure-preserving” model, since the physical structure of the power system is retained. The inclusion of the load buses requires the solution of the set of power flow equations

**FIGURE 5.24**

Angular frequency response for Example 5.10

governing the system network. This constraint leads to the inclusion of the algebraic power flow equations in conjunction with the differential equations describing the states. One example of a structure-preserving DAE model is given below:

$$T_{d0_i} \dot{E}'_{q_i} = -E'_{q_i} - (x_{d_i} - x'_{d_i}) I_{d_i} + E_{fd_i} \quad (5.161)$$

$$T_{q0_i} \dot{E}'_{d_i} = -E'_{d_i} + (x_{q_i} - x'_{q_i}) I_{q_i} \quad (5.162)$$

$$\dot{\delta}_i = \omega_i - \omega_s \quad (5.163)$$

$$\frac{2H_i}{\omega_s} \dot{\omega}_i = T_{m_i} - E'_{d_i} I_{d_i} - E'_{q_i} I_{q_i} - (x'_{q_i} - x'_{d_i}) I_{d_i} I_{q_i} \quad (5.164)$$

$$T_{E_i} \dot{E}_{fd_i} = -(K_{E_i} + S_{E_i}(E_{fd_i})) E_{fd_i} + V_{R_i} \quad (5.165)$$

$$T_{F_i} \dot{R}_{F_i} = -R_{F_i} + \frac{K_{F_i}}{T_{F_i}} E_{fd_i} \quad (5.166)$$

$$T_{A_i} \dot{V}_{R_i} = -V_{R_i} + K_{A_i} R_{F_i} - \frac{K_{A_i} K_{F_i}}{T_{F_i}} E_{fd_i} + K_{A_i} (V_{ref_i} - V_{T_i}) \quad (5.167)$$

$$T_{RH_i} \dot{T}_{M_i} = -T_{M_i} + \left(1 - \frac{K_{HP_i} T_{RH_i}}{T_{CH_i}}\right) P_{CH_i} + \frac{K_{HP_i} T_{RH_i}}{T_{CH_i}} P_{SV_i} \quad (5.168)$$

$$T_{CH_i} \dot{P}_{CH_i} = -P_{CH_i} + P_{SV_i} \quad (5.169)$$

$$T_{SV_i} \dot{P}_{SV_i} = -P_{SV_i} + P_{C_i} - \frac{1}{R} \frac{\omega_i}{\omega_s} \quad (5.170)$$

and

$$0 = V_i e^{j\theta_i} + (r_s + jx'_{d_i}) (I_{d_i} + jI_{q_i}) e^{j(\delta_i - \frac{\pi}{2})} \\ - [E'_{d_i} + (x'_{q_i} - x'_{d_i}) I_{q_i} + jE'_{q_i}] e^{j(\delta_i - \frac{\pi}{2})} \quad (5.171)$$

$$0 = V_i e^{j\theta_i} (I_{d_i} - jI_{q_i}) e^{-j(\delta_i - \frac{\pi}{2})} - \sum_{k=1}^N V_i V_k Y_{ik} e^{j(\theta_i - \theta_k - \phi_{ik})} \quad (5.172)$$

$$0 = P_i + jQ_i - \sum_{k=1}^N V_i V_k Y_{ik} e^{j(\theta_i - \theta_k - \phi_{ik})} \quad (5.173)$$

These equations describe the behavior of a two-axis generator model, a simple automatic voltage regulator and exciter, a simple turbine/governor, and constant power loads. This set of dynamic equations is listed here for illustration purposes and is not intended to be inclusive of all possible representations. A detailed development of these equations may be found in [47].

These equations may be modeled in a more general form as

$$\dot{x} = f(x, y) \quad (5.174)$$

$$0 = g(x, y) \quad (5.175)$$

where the state vector x contains the dynamic state variables of the generators. The vector y is typically much larger and contains all of the network variables, including bus voltage magnitude and angle. It may also contain generator states such as currents that are not inherently dynamic. There are two typical approaches to solving this set of differential/algebraic equations. The first is the method suggested by Gear whereby all equations are solved simultaneously. The second approach is to solve the differential and algebraic sets of equations separately and iterate between each.

Consider the first approach applied to the DAE system using the trapezoidal integration method:

$$x(n+1) = x(n) + \frac{h}{2} [f(x(n+1), y(n+1)) + f(x(n), y(n))] \quad (5.176)$$

$$0 = g(x(n+1), y(n+1)) \quad (5.177)$$

This set of nonlinear equations must then be solved using the Newton–Raphson method for the combined state vector $[x(n+1) \ y(n+1)]^T$:

$$= \begin{bmatrix} I - \frac{h}{2} \frac{\partial f}{\partial x} & -\frac{h}{2} \frac{\partial f}{\partial y} \\ \frac{\partial g}{\partial x} & \frac{\partial g}{\partial y} \end{bmatrix} \begin{bmatrix} x(n+1)^{k+1} - x(n+1)^k \\ y(n+1)^{k+1} - y(n+1)^k \end{bmatrix} \\ \begin{bmatrix} x(n+1)^k - x(n) - \frac{h}{2} [f(x(n+1)^k, y(n+1)^k) + f(x(n), y(n))] \\ g(x(n+1)^k, y(n+1)^k) \end{bmatrix} \quad (5.178)$$

The advantages of this method are that, since the whole set of system equations is used, the system matrices are quite sparse, and sparse solution techniques can be used efficiently. Second, since the set of equations is solved

simultaneously, the iterations are more likely to converge since the only iteration involved is that of the Newton–Raphson algorithm. Note, however, that, in a system of ODEs, the left-hand side matrix (the matrix to be factored in LU factorization) can be made to be diagonally dominant (and therefore well conditioned) by decreasing the step size h . In DAE systems, however, the left-hand side matrix may be ill conditioned under certain operating points where $\frac{\partial g}{\partial y}$ is ill conditioned, causing difficulty in solving the system of equations. This situation may occur during the simulation of voltage collapse where the states encounter a bifurcation. The subject of bifurcation and voltage collapse is complex and outside the scope of this book; however, several excellent texts have been published that study these phenomena in great detail [25] [47] [60].

The second approach to solving the system of DAEs is to solve each of the subsystems independently and iteratively. The system of differential equations is solved first for $x(n+1)$ while holding $y(n+1)$ constant as an input. After $x(n+1)$ has been found, it is then used as an input to solve the algebraic system. The updated value of $y(n+1)$ is then substituted back into the differential equations, and $x(n+1)$ is recalculated. This back and forth process is repeated until the values $x(n+1)$ and $y(n+1)$ converge. The solution is then advanced to the next time point. The advantage of this method is simplicity in programming, since each subsystem is solved independently and the Jacobian elements $\frac{\partial f}{\partial y}$ and $\frac{\partial g}{\partial x}$ are not used. In some cases, this may speed up the computation, although more iterations may be required to reach convergence.

5.9 Problems

- Determine a Taylor series expansion for the solution of the equation

$$\dot{x} = x^2 \quad x(0) = 1$$

about the point $\hat{x} = 0$ (a McLaurin series expansion). Use this approximation to compute x for $\hat{x} = 0.2$ and $\hat{x} = 1.2$. Compare with the exact solution and explain the results.

- Use the following algorithms to solve the initial value problem.

$$\begin{aligned}\dot{x}_1 &= -2x_2 + 2t^2 & x_1(0) &= -4 \\ \dot{x}_2 &= \frac{1}{2}x_1 + 2t & x_2(0) &= 0\end{aligned}$$

on the interval $0 \leq t \leq 5$ with a fixed integration step of 0.25 seconds.

- Backward Euler

- (b) Forward Euler
- (c) Trapezoidal rule
- (d) Fourth-order Runge–Kutta

Compare the answers to the exact solution

$$\begin{aligned}x_1(t) &= -4 \cos t \\x_2(t) &= -2 \sin t + t^2\end{aligned}$$

3. Consider the following linear multistep formula.

$$x_{n+1} = a_0 x_n + a_1 x_{n-1} + a_2 x_{n-2} + a_3 x_{n-3} + h b_{-1} f(x_{n+1}, t_{n+1})$$

- (a) What are the number of steps in the formula?
- (b) What is the maximum order that will enable you to determine all the coefficients of the formula?
- (c) Assuming uniform step size h , find all the coefficients of the formula such that its order is the answer of part (b). (This is a Gear's method)
- (d) Is the formula implicit or explicit?
- (e) Assuming uniform step size h , find the expression for the local truncation error of the formula.
- (f) Is the formula absolutely stable?
- (g) Is the formula stiffly-stable?

4. Consider the following initial value problem.

$$\dot{x} = 100(\sin(t) - x), \quad x(0) = 0$$

The exact solution is

$$x(t) = \frac{\sin(t) - 0.01 \cos(t) + 0.01 e^{-100t}}{1.0001}$$

Solve the initial value problem with $h = 0.02$ s using the following integration methods. Plot each of the numerical solutions against the exact solution over the range $t \in [0, 3.0]$ seconds. Plot the global error for each method over the range $t \in [0, 3.0]$ seconds. Discuss your results.

- (a) Backward Euler
- (b) Forward Euler

- (c) Trapezoidal rule
 (d) Fourth-order Runge–Kutta
 (e) Repeat (a)–(d) with step size $h = 0.03$ s.
5. Consider a simple ecosystem consisting of rabbits that have an infinite food supply and foxes that prey upon the rabbits for their food. A classical mathematical “predator-prey” model due to Volterra describes this system by a pair of nonlinear, first-order differential equations:
- $$\begin{aligned}\dot{r} &= \alpha r + \beta r f & r(0) &= r_0 \\ \dot{f} &= \gamma f + \delta r f & f(0) &= f_0\end{aligned}$$
- where $r = r(t)$ is the number of rabbits, $f = f(t)$ is the number of foxes. When $\beta = 0$, the two populations do not interact, and so the rabbits multiply and the foxes die off from starvation.
- Investigate the behavior of this system for $\alpha = -1$, $\beta = 0.01$, $\gamma = 0.25$, and $\delta = -0.01$. Use the trapezoidal integration method with $h = 0.1$ and $T = 50$, and the set of initial conditions $r_0 = 40$ and $f_0 = 70$. Plot (1) r and f vs. t , and (2) r vs. f for each case.
6. The following system of nonlinear equations is known as Duffing’s differential equation.

$$\ddot{x} + \delta \dot{x} + \alpha x + \beta x^3 = \gamma \cos \omega t$$

For the following parameters, plot x versus y in two dimensions and x and y versus t for $0 \leq t \leq 200$ seconds.

$$\begin{aligned}\alpha &= 0.1 \\ \beta &= 0.5 \\ \gamma &= -0.8 \\ \delta &= 0.1 \\ \omega &= \frac{\pi}{2}\end{aligned}$$

with $x(0) = 1$ and $\dot{x} = -1$ and $h = 0.01$.

7. The following system of equations is known as the Lorenz equations.

$$\begin{aligned}\dot{x} &= \sigma(y - x) \\ \dot{y} &= \rho x - y - xz \\ \dot{z} &= xy - \beta z\end{aligned}$$

Let $\sigma = 10$, $\rho = 28$, and $\beta = 4$. Use the trapezoidal method to plot the response of the system for $0 \leq t \leq 10$ seconds using a fixed integration time step of $h = 0.005$ seconds and a Newton–Raphson convergence error of 10^{-5} . Plot x versus y in two dimensions and x vs. y vs. z in three dimensions.

- (a) Use $[x(0) \ y(0) \ z(0)]^T = [20 \ 20 \ 20]^T$.
- (b) Use $[x(0) \ y(0) \ z(0)]^T = [21 \ 20 \ 20]^T$. Explain the difference.
- 8. Consider the following system of “stiff” equations.

$$\begin{aligned}\dot{x}_1 &= -2x_1 + x_2 + 100 & x_1(0) &= 0 \\ \dot{x}_2 &= 10^4x_1 - 10^4x_2 + 50 & x_2(0) &= 0\end{aligned}$$
 - (a) Determine the *maximum step size* h_{\max} for the forward Euler algorithm to remain numerically stable.
 - (b) Use the forward Euler algorithm with step size $h = \frac{1}{2}h_{\max}$ to solve for $x_1(t)$ and $x_2(t)$ for $t > 0$.
 - (c) Repeat using a step size $h = 2h_{\max}$.
 - (d) Use the backward Euler algorithm to solve the stiff equations. Choose the following step sizes in terms of the maximum step size h_{\max} .
 - i. $h = 10h_{\max}$
 - ii. $h = 100h_{\max}$
 - iii. $h = 1000h_{\max}$
 - iv. $h = 10,000h_{\max}$
 - (e) Repeat using the multistep method (Gear’s method) of Problem 4.
- 9. Consider a multistep method of the form

$$x_{n+2} - x_{n-2} + \alpha(x_{n+1} - x_{n-1}) = h[\beta(f_{n+1} + f_{n-1}) + \gamma f_n]$$
 - (a) Show that the parameters α , β , and γ can be chosen uniquely so that the method has order $p = 6$.
 - (b) Discuss the stability properties of this method. For what region is it absolutely stable? Stiffly stable?

10. A power system can be described by the following ODE system.

$$\begin{aligned}\dot{\delta}_i &= \omega_i - \omega_s \\ M_i \dot{\omega}_i &= P_i - E_i \sum_{k=1}^n E_k Y_{ik} \sin(\delta_i - \delta_k - \phi_{ik})\end{aligned}$$

Using a step size of $h = 0.01$ s and the trapezoidal integration method, determine whether or not this system is stable for a clearing time of 4 cycles.

Let

$$\begin{aligned}E_1 &= 1.0566\angle 2.2717^\circ \\E_2 &= 1.0502\angle 19.7315^\circ \\E_3 &= 1.0170\angle 13.1752^\circ \\P_1 &= 0.716 \\P_2 &= 1.630 \\P_3 &= 0.850 \\H_1 &= 23.64 \\H_2 &= 6.40 \\H_3 &= 3.01\end{aligned}$$

where $M_i = \frac{2H_i}{\omega_s}$.

and the following admittance matrices:

Prefault:

$$\begin{aligned}0.846 - j2.988 &\quad 0.287 + j1.513 \quad 0.210 + j1.226 \\0.287 + j1.513 &\quad 0.420 - j2.724 \quad 0.213 + j1.088 \\0.210 + j1.226 &\quad 0.213 + j1.088 \quad 0.277 - j2.368\end{aligned}$$

Fault-on:

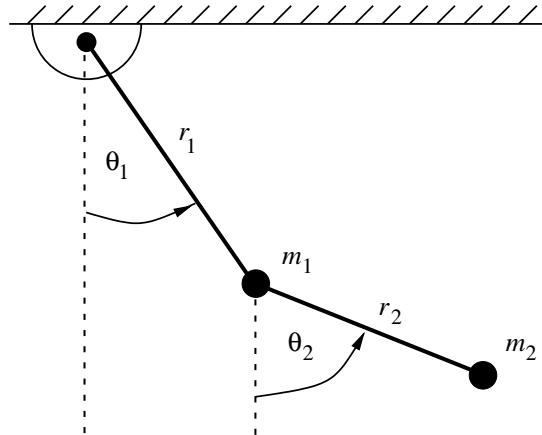
$$\begin{aligned}0.657 - j3.816 &\quad 0.000 + j0.000 \quad 0.070 + j0.631 \\0.000 + j0.000 &\quad 0.000 - j5.486 \quad 0.000 + j0.000 \\0.070 + j0.631 &\quad 0.000 + j0.000 \quad 0.174 - j2.796\end{aligned}$$

Postfault:

$$\begin{aligned}1.181 - j2.229 &\quad 0.138 + j0.726 \quad 0.191 + j1.079 \\0.138 + j0.726 &\quad 0.389 - j1.953 \quad 0.199 + j1.229 \\0.191 + j1.079 &\quad 0.199 + j1.229 \quad 0.273 - j2.342\end{aligned}$$

11. A double pendulum system is shown in Figure 5.25. Masses m_1 and m_2 are connected by massless rods of length r_1 and r_2 . The equations of motion of the two masses, expressed in terms of the angles θ_1 and θ_2 as indicated, are

$$\begin{aligned}-(m_1 + m_2)gr_1 \sin \theta_1 &= (m_1 + m_2)r_1^2 \ddot{\theta}_1 + m_2 r_1 r_2 \ddot{\theta}_2 \cos(\theta_1 - \theta_2) \\&\quad + m_2 r_1 r_2 \dot{\theta}_2^2 \sin(\theta_1 - \theta_2) \\-m_2 gr_2 \sin \theta_2 &= m_2 r_2^2 \ddot{\theta}_2 + m_2 r_1 r_2 \ddot{\theta}_1 \cos(\theta_1 - \theta_2) \\&\quad - m_2 r_1 r_2 \dot{\theta}_1^2 \sin(\theta_1 - \theta_2)\end{aligned}$$

**FIGURE 5.25**

Double pendulum system

- (a) Choose $x_1 = \theta_1, x_2 = \dot{\theta}_1, x_3 = \theta_2, x_4 = \dot{\theta}_2$, and show that $[0; 0; 0; 0]$ is an equilibrium of the system.
- (b) Show that the second-order Adams–Bashforth method is given by

$$x_{n+1} = x_n + h \left\{ \frac{3}{2} f(x_n, t_n) - \frac{1}{2} f(x_{n-1}, t_{n-1}) \right\}$$

with

$$\epsilon_T = \frac{5}{12} \hat{x}^{(3)}(\tau) h^3$$

- (c) Show that the third-order Adams–Bashforth method is given by

$$x_{n+1} = x_n + h \left\{ \frac{23}{12} f(x_n, t_n) - \frac{16}{12} f(x_{n-1}, t_{n-1}) + \frac{5}{12} f(x_{n-2}, t_{n-2}) \right\}$$

with

$$\epsilon_T = \frac{3}{8} \hat{x}^{(4)}(\tau) h^4$$

- (d) Let $r_1 = 1, r_2 = 0.5$, and $g = 10$. Using the second-order Adams–Bashforth method with $h = 0.005$, plot the behavior of the system for an initial displacement of $\theta_1 = 25^\circ$ and $\theta_2 = 10^\circ$, for $T \in [0, 10]$ with

i. $m_1 = 10$ and $m_2 = 10$.ii. $m_1 = 10$ and $m_2 = 5$.iii. $m_1 = 10$ and $m_2 = 1$.

- (e) Repeat (d) using a variable step method with an upper LTE bound of 0.001 and $B_{avg} = 0.1B_U$. Plot h versus t for each case. Discuss your solution.
- (f) Repeat (d) and (e) using a third-order Adams–Bashforth method.

6

Optimization

The basic objective of any optimization method is to find the values of the system state variables and/or parameters that minimize some cost function of the system. The types of cost functions are system dependent and can vary widely from application to application and are not necessarily strictly measured in terms of dollars. Examples of engineering optimizations can range from minimizing

- the error between a set of measured and calculated data,
- active power losses,
- the weight of a set of components that comprise the system,
- particulate output (emissions),
- system energy, or
- the distance between actual and desired operating points,

to name a few possibilities. The basic formulation of any optimization can be represented as minimizing a defined cost function subject to any physical or operational constraints of the system:

$$\begin{aligned} & \text{minimize } f(x, u) \quad x \in R^n \\ & \quad u \in R^m \end{aligned} \tag{6.1}$$

subject to

$$g(x, u) = 0 \text{ equality constraints} \tag{6.2}$$

$$h(x, u) = 0 \text{ inequality constraints} \tag{6.3}$$

where x is the vector of system states and u is the vector of system parameters. The basic approach is to find the vector of system parameters that, when substituted into the system model, will result in the state vector x that minimizes the cost function $f(x, u)$.

6.1 Least Squares State Estimation

In many physical systems, the system operating condition cannot be determined directly by an analytical solution of known equations using a given set of known, dependable quantities. More frequently, the system operating condition is determined by the measurement of system states at different points throughout the system. In many systems, more measurements are made than are necessary to uniquely determine the operating point. This redundancy is often purposely designed into the system to counteract the effect of inaccurate or missing data due to instrument failure. Conversely, not all of the states may be available for measurement. High temperatures, moving parts, or inhospitable conditions may make it difficult, dangerous, or expensive to measure certain system states. In this case, the missing states must be estimated from the rest of the measured information of the system. This process is often known as *state estimation* and is the process of estimating unknown states from measured quantities. State estimation gives the “best estimate” of the state of the system in spite of uncertain, redundant, and/or conflicting measurements. A good state estimation will smooth out small random errors in measurements, detect and identify large measurement errors, and compensate for missing data. This process strives to minimize the error between the (unknown) true operating state of the system and the measured states.

The set of measured quantities can be denoted by the vector z , which may include measurements of system states (such as voltage and current) or quantities that are functions of system states (such as power flows). Thus

$$z^{true} = Ax \quad (6.4)$$

where x is the set of system states and A is usually not square. The error vector is the difference between the measured quantities z and the true quantities:

$$e = z - z^{true} = z - Ax \quad (6.5)$$

Typically, the minimum of the square of the error is desired to negate any effects of sign differences between the measured and true values. Thus a state estimator endeavors to find the minimum of the squared error, or a *least squares minimization*:

$$\text{minimize } \|e\|^2 = e^T \cdot e = \sum_{i=1}^m \left[z_i - \sum_{j=1}^m a_{ij} x_j \right]^2 \quad (6.6)$$

The squared error function can be denoted by $U(x)$ and is given by

$$U(x) = e^T \cdot e = (z - Ax)^T (z - Ax) \quad (6.7)$$

$$= (z^T - x^T A^T) (z - Ax) \quad (6.8)$$

$$= z^T z - z^T A x - x^T A^T z + x^T A^T A x \quad (6.9)$$

Note that the product $z^T Ax$ is a scalar and so it can be equivalently written as

$$z^T Ax = (z^T Ax)^T = x^T A^T z$$

Therefore, the squared error function is given by

$$U(x) = z^T z - 2x^T A^T z + x^T A^T Ax \quad (6.10)$$

The minimum of the squared error function can be found by an unconstrained optimization where the derivative of the function with respect to the states x is set to zero:

$$\frac{\partial U(x)}{\partial x} = 0 = -2A^T z + 2A^T Ax \quad (6.11)$$

Thus

$$A^T Ax = A^T z \quad (6.12)$$

Thus, if $b = A^T z$ and $\hat{A} = A^T A$, then

$$\hat{A}x = b \quad (6.13)$$

which can be solved by LU factorization. This state vector x is the best estimate (in the squared error) to the system operating condition from which the measurements z were taken. The measurement error is given by

$$e = z^{meas} - Ax \quad (6.14)$$

Example 6.1

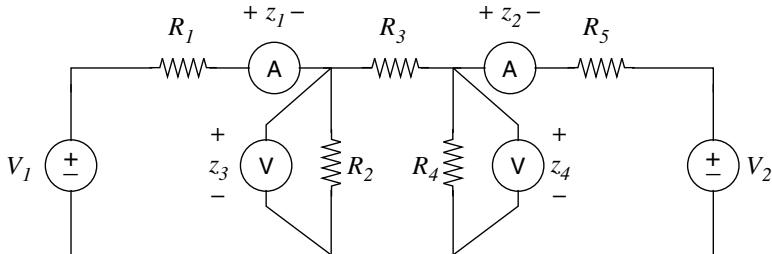
A set of measurements for the circuit shown in Figure 6.1 is given by

Ammeter 1	z_1	4.27 A
Ammeter 2	z_2	-1.71 A
Voltmeter 1	z_3	3.47 V
Voltmeter 2	z_4	2.50 V

where $R_1 = R_3 = R_5 = 1.5\Omega$ and $R_2 = R_4 = 1.0\Omega$. Find the node voltages V_1 and V_2 .

Solution 6.1 The Kirchoff voltage and current law equations for this system can be written as

$$\begin{aligned} -V_1 + R_1 z_1 + z_3 &= 0 \\ -V_2 - R_5 z_2 + z_4 &= 0 \\ z_3/R_2 - z_1 + (z_3 - z_4)/R_3 &= 0 \\ z_4/R_4 + z_2 + (z_4 - z_3)/R_3 &= 0 \end{aligned}$$

**FIGURE 6.1**

Circuit for Example 6.1

These equations can be rewritten in matrix form as

$$\begin{bmatrix} R_1 & 0 & 1 & 0 \\ 0 & -R_5 & 0 & 1 \\ 1 & 0 & -\left(\frac{1}{R_2} + \frac{1}{R_3}\right) & \frac{1}{R_3} \\ 0 & 1 & -\frac{1}{R_3} & \frac{1}{R_3} + \frac{1}{R_4} \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} \quad (6.15)$$

To find the relationship between the measurements z and x , this equation must be reformulated as $z = Ax$. Note that this equation can be solved easily by LU factorization by considering each column of A individually. Thus

$$\begin{bmatrix} R_1 & 0 & 1 & 0 \\ 0 & -R_5 & 0 & 1 \\ 1 & 0 & -\left(\frac{1}{R_2} + \frac{1}{R_3}\right) & \frac{1}{R_3} \\ 0 & 1 & -\frac{1}{R_3} & \frac{1}{R_3} + \frac{1}{R_4} \end{bmatrix} [A(:, 1)] = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad (6.16)$$

Similarly,

$$\begin{bmatrix} R_1 & 0 & 1 & 0 \\ 0 & -R_5 & 0 & 1 \\ 1 & 0 & -\left(\frac{1}{R_2} + \frac{1}{R_3}\right) & \frac{1}{R_3} \\ 0 & 1 & -\frac{1}{R_3} & \frac{1}{R_3} + \frac{1}{R_4} \end{bmatrix} [A(:, 2)] = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad (6.17)$$

yielding

$$\begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{bmatrix} = \begin{bmatrix} 0.4593 & -0.0593 \\ 0.0593 & -0.4593 \\ 0.3111 & 0.0889 \\ 0.0889 & 0.3111 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} \quad (6.18)$$

Thus

$$b = A^T z = \begin{bmatrix} 0.4593 & 0.0593 & 0.3111 & 0.0889 \\ -0.0593 & -0.4593 & 0.0889 & 0.3111 \end{bmatrix} \begin{bmatrix} 4.27 \\ -1.71 \\ 3.47 \\ 2.50 \end{bmatrix}$$

$$= \begin{bmatrix} 3.1615 \\ 1.6185 \end{bmatrix}$$

and

$$\hat{A} = A^T A = \begin{bmatrix} 0.3191 & 0.0009 \\ 0.0009 & 0.3191 \end{bmatrix}$$

leading to

$$\begin{bmatrix} 0.3191 & 0.0009 \\ 0.0009 & 0.3191 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = \begin{bmatrix} 3.1615 \\ 1.6185 \end{bmatrix} \quad (6.19)$$

Solving this equation yields

$$V_1 = 9.8929$$

$$V_2 = 5.0446$$

The error between the measured values and the estimated values of this system is given by

$$\begin{aligned} e &= z - Ax \\ &= \begin{bmatrix} 4.27 \\ -1.71 \\ 3.47 \\ 2.50 \end{bmatrix} - \begin{bmatrix} 0.4593 & -0.0593 \\ 0.0593 & -0.4593 \\ 0.3111 & 0.0889 \\ 0.0889 & 0.3111 \end{bmatrix} \begin{bmatrix} 9.8929 \\ 5.0446 \end{bmatrix} \\ &= \begin{bmatrix} 0.0255 \\ 0.0205 \\ -0.0562 \\ -0.0512 \end{bmatrix} \end{aligned} \quad (6.20)$$

■

6.1.1 Weighted Least Squares Estimation

If all measurements are treated equally in the least squares solution, then the less accurate measurements will affect the estimation as significantly as the more accurate measurements. As a result, the final estimation may contain large errors due to the influence of inaccurate measurements. By introducing a weighting matrix to emphasize the more accurate measurements more heavily than the less accurate measurements, the estimation procedure can then force the results to coincide more closely with the measurements of greater accuracy. This leads to the weighted least squares estimation

$$\text{minimize } \|e\|^2 = e^T \cdot e = \sum_{i=1}^m w_i \left[z_i - \sum_{j=1}^m a_{ij} x_j \right]^2 \quad (6.21)$$

where w_i is a weighting factor reflecting the level of confidence in the measurement z_i .

Example 6.2

Suppose that the ammeters are known to have been more recently calibrated than the voltmeters; thus the level of confidence in the current measurements is greater than the voltage measurements. Using the following weighting matrix, find the node voltages V_1 and V_2 .

$$W = \begin{bmatrix} 100 & 0 & 0 & 0 \\ 0 & 100 & 0 & 0 \\ 0 & 0 & 50 & 0 \\ 0 & 0 & 0 & 50 \end{bmatrix}$$

Solution 6.2 By introducing the weighting matrix, the new minimum is given by

$$A^T W A x = A^T W z \quad (6.22)$$

The matrix $A^T W A$ is also known as the *gain* matrix. Using the same procedure as before, the weighted node voltage values are given by

$$\begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = \begin{bmatrix} 9.9153 \\ 5.0263 \end{bmatrix} \quad (6.23)$$

and the error vector is given by

$$e = \begin{bmatrix} 0.0141 \\ 0.0108 \\ -0.0616 \\ 0.0549 \end{bmatrix} \quad (6.24)$$

■

Note that the added confidence in the current measurements has decreased the estimation error in the current, but the voltage measurement error is approximately the same.

Example 6.2 illustrates the impact of confidence weighting on the accuracy of the estimation. All instruments add some degree of error to the measured values, but the problem is how to quantify this error and account for it during the estimation process. In general, it can be assumed that the introduced errors have normal (Gaussian) distribution with zero mean and that each measurement is independent of all other measurements. This means that each measurement error is as likely to be greater than the true value as it is to be less than the true value. A zero mean Gaussian distribution has several attributes. The standard deviation of a zero mean Gaussian distribution is denoted by σ . This means that 68% of all measurements will fall within $\pm\sigma$ of the expected value, which is zero in a zero mean distribution. Further, 95% of all measurements will fall within $\pm 2\sigma$, and 99% of all measurements will fall within $\pm 3\sigma$. The variance of the measurement distribution is given by σ^2 .

This implies that, if the variance of the measurements is relatively small, then the majority of measurements are close to the mean. One interpretation of this is that accurate measurements lead to small variance in the distribution.

This relationship between accuracy and variance leads to a straightforward approach from which to develop a weighting matrix for the estimation. Consider the squared error matrix given by

$$e \cdot e^T = \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \vdots \\ e_m \end{bmatrix} [e_1 \ e_2 \ e_3 \ \dots \ e_m] \quad (6.25)$$

$$= \begin{bmatrix} e_1^2 & e_1e_2 & e_1e_3 & \dots & e_1e_m \\ e_2e_1 & e_2^2 & e_2e_3 & \dots & e_2e_m \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ e_me_1 & e_me_2 & e_me_3 & \dots & e_m^2 \end{bmatrix} \quad (6.26)$$

where each e_i is the error in the i th measurement. The expected, or mean, value of each error product is given by $E[\cdot]$. The expected value of each of the diagonal terms is the variance of the i th error distribution σ_i^2 . The expected value of each of the off-diagonal terms, or covariance, is zero because each measurement is assumed to be independent of every other measurement. Therefore, the expected value of the squared error matrix (also known as the covariance matrix) is

$$E[e \cdot e^T] = \begin{bmatrix} E[e_1^2] & E[e_1e_2] & \dots & E[e_1e_m] \\ E[e_2e_1] & E[e_2^2] & \dots & E[e_2e_m] \\ \vdots & \vdots & \vdots & \vdots \\ E[e_me_1] & E[e_me_2] & \dots & E[e_m^2] \end{bmatrix} \quad (6.27)$$

$$= \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & \sigma_m^2 \end{bmatrix} \quad (6.28)$$

$$= R \quad (6.29)$$

With measurements taken from a particular meter, the smaller the variance of the measurements (i.e., the more consistent they are), the greater the level of confidence in that set of measurements. A set of measurements that has a high level of confidence should have a higher weighting than a set of measurements that has a larger variance (and therefore less confidence). Therefore, a plausible weighting matrix that reflects the level of confidence in each measurement set is the inverse of the covariance matrix $W = R^{-1}$. Thus measurements that come from instruments with good consistency (small variance) will carry

greater weight than measurements that come from less accurate instruments (high variance). Thus one possible weighting matrix is given by

$$W = R^{-1} = \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2^2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{\sigma_m^2} \end{bmatrix} \quad (6.30)$$

6.1.2 Bad Data Detection

Frequently, a set of measurements will contain one or more data points from faulty or poorly calibrated instruments. Telemetered measurements are subject to noise or error in metering and communication. These “bad” data points typically fall outside of the standard deviation of the measurements and may affect the reliability of the state estimation process. In severe cases, the bad data may actually lead to grossly inaccurate results. Bad data may cause the accuracy of the estimate to deteriorate because of the “smearing” effect as the bad data will pull, or smear, the estimated values away from the true values. Therefore, it is desirable to develop a measure of the “goodness” of the data upon which the estimation is based. If the data lead to a good estimate of the states, then the error between the measured and calculated values will be small in some sense. If the error is large, then the data contain at least one bad data point. One error that is useful to consider is the estimated measurement error \hat{e} . This error is the difference between the actual measurements z and the estimated measurements \hat{z} . Recall the error vector from Equation (6.14) where $e = z - Ax$; then the estimated measurement error becomes

$$\hat{e} = z - \hat{z} \quad (6.31)$$

$$= z - A\hat{x} \quad (6.32)$$

$$= z - A (A^T W A)^{-1} A^T W z \quad (6.33)$$

$$= \left(I - A (A^T W A)^{-1} A^T W \right) z \quad (6.34)$$

$$= \left(I - A (A^T W A)^{-1} A^T W \right) (e + Ax) \quad (6.35)$$

$$= \left(I - A (A^T W A)^{-1} A^T W \right) e + A \left(I - (A^T W A)^{-1} A^T W A \right) x \quad (6.36)$$

$$= \left(I - A (A^T W A)^{-1} A^T W \right) e \quad (6.37)$$

Thus the variance of \hat{e} can be calculated from

$$\hat{e} \hat{e}^T = (z - \hat{z})(z - \hat{z})^T \quad (6.38)$$

$$= \left[I - A (A^T W A)^{-1} A^T W \right] e e^T \left[I - W A (A^T W A)^{-1} A^T \right] \quad (6.39)$$

The expected, or mean, value of $\hat{e}\hat{e}^T$ is given by

$$E [\hat{e}\hat{e}^T] = \left[I - A (A^T W A)^{-1} A^T W \right] E [ee^T] \left[I - W A (A^T W A)^{-1} A^T \right] \quad (6.40)$$

Recall that $E [ee^T]$ is just the covariance matrix $R = W^{-1}$, which is a diagonal matrix. Thus

$$E [\hat{e}\hat{e}^T] = \left[I - A (A^T W A)^{-1} A^T W \right] \left[I - A (A^T W A)^{-1} A^T W \right] R \quad (6.41)$$

The matrix

$$\left[I - A (A^T W A)^{-1} A^T W \right]$$

has the unusual property that it is an *idempotent* matrix. An idempotent matrix M has the property that $M^2 = M$; thus, no matter how many times M is multiplied by itself, it will still return the product M . Therefore,

$$E [\hat{e}\hat{e}^T] = \left[I - A (A^T W A)^{-1} A^T W \right] \left[I - A (A^T W A)^{-1} A^T W \right] R \quad (6.42)$$

$$= \left[I - A (A^T W A)^{-1} A^T W \right] R \quad (6.43)$$

$$= R - A (A^T W A)^{-1} A^T \quad (6.44)$$

$$= R' \quad (6.45)$$

To determine whether the estimated values differ significantly from the measured values, a useful statistical measure is the χ^2 (chi-squared) test of inequality. This measure is based on the χ^2 probability distribution, which differs in shape depending on its degrees of freedom k , which is the difference between the number of measurements and the number of states. By comparing the weighted sum of errors with the χ^2 value for a particular degree of freedom and significance level, it can be determined whether the errors exceed the bounds of what would be expected by chance alone. A significance level indicates the level of probability that the measurements are erroneous. A significance level of 0.05 indicates there is a 5% likelihood that bad data exist, or conversely, a 95% level of confidence in the goodness of the data. For example, for $k = 2$ and a significance level $\alpha = 0.05$, if the weighted sum of errors does not exceed a χ^2 of 5.99, then the set of measurements can be assured of being good with 95% confidence; otherwise the data must be rejected as containing at least one bad data point. Although the χ^2 test is effective in signifying the presence of bad data, it cannot identify locations. The identification of bad data locations continues to be an open research topic.

k	χ^2 Values			
	α			
	0.10	0.05	0.01	0.001
1	2.71	3.84	6.64	10.83
2	4.61	5.99	9.21	13.82
3	6.25	7.82	11.35	16.27
4	7.78	9.49	13.23	18.47
5	9.24	11.07	15.09	20.52
6	10.65	12.59	16.81	22.46
7	12.02	14.07	18.48	24.32
8	13.36	15.51	20.09	26.13
9	14.68	16.92	21.67	27.88
10	15.99	18.31	23.21	29.59
11	17.28	19.68	24.73	31.26
12	18.55	21.03	26.22	32.91
13	19.81	22.36	27.69	34.53
14	21.06	23.69	29.14	36.12
15	22.31	25.00	30.68	37.70
16	23.54	26.30	32.00	39.25
17	24.77	27.59	33.41	40.79
18	25.99	28.87	34.81	42.31
19	27.20	30.14	36.19	43.82
20	28.41	31.41	37.67	45.32
21	29.62	32.67	38.93	46.80
22	30.81	33.92	40.29	48.27
23	32.00	35.17	41.64	49.73
24	33.20	36.42	42.98	51.18
25	34.38	37.65	44.31	52.62
26	35.56	38.89	45.64	54.05
27	36.74	40.11	46.96	55.48
28	37.92	41.34	48.28	56.89
29	39.09	42.56	49.59	58.30
30	40.26	43.77	50.89	59.70

A test procedure to test for the existence of bad data is given by

Test Procedure for Bad Data

1. Use z to estimate x .

2. Calculate the error

$$e = z - Ax$$

3. Evaluate the weighted sum of squares

$$f = \sum_{i=1}^m \frac{1}{\sigma_i^2} e_i^2$$

4. For $k = m - n$ and a specified probability α , if $f < \chi_{k,\alpha}^2$, then the data are good; otherwise at least one bad data point exists.

Example 6.3

Using the chi-square test of inequality with $\alpha = 0.01$, check for the presence of bad data in the measurements of Example 6.1.

Solution 6.3 The number of states in Example 6.1 is 2 and the number of measurements is 4; therefore $k = 4 - 2 = 2$. The weighted sum of squares is given by

$$\begin{aligned} f &= \sum_{i=1}^{m=4} \frac{1}{\sigma_i^2} e_i^2 \\ &= 100(0.0141)^2 + 100(0.0108)^2 + 50(-0.0616)^2 + 50(0.0549) \\ &= 0.3720 \end{aligned}$$

From the table of chi-squared values, the chi-square value for this example is 9.21. The weighted least squares error is less than the chi-square value; this indicates that the estimated values are good to a confidence level of 99%. ■

6.1.3 Nonlinear Least Squares State Estimation

As in the linear least squares estimation, the nonlinear least squares estimation attempts to minimize the square of the errors between a known set of measurements and a set of weighted nonlinear functions:

$$\text{minimize } f = \|e\|^2 = e^T \cdot e = \sum_{i=1}^m \frac{1}{\sigma_i^2} [z_i - h_i(x)]^2 \quad (6.46)$$

where $x \in R^n$ is the vector of unknowns to be estimated, $z \in R^m$ is the vector of measurements, σ_i^2 is the variance of the i th measurement, and $h(x)$ is the

function vector relating x to z , where the measurement vector z can be a set of geographically distributed measurements, such as voltages and power flows.

In state estimation, the unknowns in the nonlinear equations are the state variables of the system. The state values that minimize the error are found by setting the derivatives of the error function to zero:

$$F(x) = H_x^T R^{-1} [z - h(x)] = 0 \quad (6.47)$$

where

$$H_x = \begin{bmatrix} \frac{\partial h_1}{\partial x_1} & \frac{\partial h_1}{\partial x_2} & \cdots & \frac{\partial h_1}{\partial x_n} \\ \frac{\partial h_2}{\partial x_1} & \frac{\partial h_2}{\partial x_2} & \cdots & \frac{\partial h_2}{\partial x_n} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial h_m}{\partial x_1} & \frac{\partial h_m}{\partial x_2} & \cdots & \frac{\partial h_m}{\partial x_n} \end{bmatrix} \quad (6.48)$$

and R is the matrix of measurement variances. Note that Equation (6.47) is a set of nonlinear equations that must be solved using the Newton–Raphson or another iterative numerical solver. In this case, the Jacobian of $F(x)$ is

$$J_F(x) = H_x^T(x) R^{-1} \frac{\partial}{\partial x} [z - h(x)] \quad (6.49)$$

$$= -H_x^T(x) R^{-1} H_x(x) \quad (6.50)$$

and the Newton–Raphson iteration becomes

$$[H_x^T(x^k) R^{-1} H_x(x^k)] [x^{k+1} - x^k] = H_x^T(x^k) R^{-1} [z - h(x^k)] \quad (6.51)$$

which is solved repeatedly using LU factorization. At convergence, x^{k+1} is equal to the set of states that minimize the error function f of Equation (6.46). The test procedure for bad data is the same as that for the linear state estimation.

6.2 Linear Programming

Linear programming is one of the most successful forms of optimization. Linear programming can be used when a problem can be expressed by a linear objective (or cost) function to be maximized (or minimized) subject to linear equality or inequality constraints. A general linear programming problem can be formulated as

$$\text{minimize } f(x) = c^T x \quad (6.52)$$

$$\text{subject to } Ax \leq b \quad (6.53)$$

$$x \geq 0 \quad (6.54)$$

Note that almost any linear optimization problem can be put into this form via one of the following transformations:

1. Maximizing $c^T x$ is the same as minimizing $-c^T x$,
2. Any constraint of the form $a^T x \geq \beta$ is equivalent to $-a^T x \leq -\beta$,
3. Any constraint of the form $a^T x = \beta$ is equivalent to $a^T x \leq \beta$ and $-a^T x \leq -\beta$,
4. If a problem does not require x_i to be nonnegative, then x_i can be replaced by the difference of two variables $x_i = u_i - v_i$ where u_i and v_i are nonnegative.

In any linear programming problem described by (A, b, c) , there exists another equivalent, or *dual* problem $(-A^T, -c, -b)$. If a linear programming problem and its dual both have feasible points (i.e., any point that satisfies $Ax \leq b$, $x \geq 0$ or $-A^T y \leq -c$, $y \geq 0$ for the dual problem), then both problems have solutions and their values are the negatives of each other.

6.2.1 Simplex Method

One of the most common methods of solving linear programming problems is the well-known *simplex* method. The simplex method is an iterative method that moves the x vector from one feasible basic vector to another in such a way that $f(x)$ always decreases. It gives the exact result after a number of steps, which is usually much less than $\binom{n+m}{m}$, generally taking $2m$ to $3m$ iterations at most (where m is the number of equality constraints) [9]. However, its worst-case complexity is exponential, as can be demonstrated with carefully constructed examples.

The simplex method is often accomplished by representing the problem in *tableau* form, which is then modified in successive steps according to given rules. Every step of the simplex method begins with a tableau. The top row contains the coefficients that pertain to the objective function $f(x)$. The current value of $f(x)$ is displayed in the top right corner of the tableau. The next m rows in the tableau represent the equality constraints. The last row of the tableau contains the current x vector. The rows in the tableau pertaining to the equality constraints can be transformed by elementary row operations without altering the solution.

The initial state of a simplex tableau has the form

$$\begin{array}{c|c|c|c} \hline & c^T & 0 & 0 \\ \hline & A & I & b \\ \hline 0 & \underbrace{b^T}_{x} & & \\ \hline \end{array}$$

The rules of the tableau that must be satisfied are

1. The x vector must satisfy the equality constraints $Ax = b$.
2. The x vector must satisfy the inequality $x \geq 0$.
3. There are n components of x (designated *nonbasic variables*) that are zero. The remaining m components are usually nonzero and are designated as *basic variables*.
4. In the matrix that defines the constraints, each basic variable occurs in only one row.
5. The objective function $f(x)$ must be expressed only in terms of nonbasic variables.
6. An artificial variable may be added to one or more constraints to obtain a starting solution.

The Simplex Algorithm is summarized:

Simplex Algorithm

1. If all coefficients in $f(x)$ (i.e., top row of the tableau) are greater than or equal to zero, then the current x vector is the solution.
2. Select the nonbasic variable whose coefficient in $f(x)$ is the largest negative entry. This variable becomes the new basic variable x_j .
3. Divide each b_i by the coefficient of the new basic variable in that row, a_{ij} . The value assigned to the new basic variable is the least of these ratios (i.e., $x_j = b_k/a_{kj}$).
4. Using pivot element a_{kj} , create zeros in column j of A with Gaussian elimination. Return to 1.

The series of inequalities in Equation (6.53), when taken together, form intersecting hyperplanes. The feasible region is the interior of this n -dimensional polytope, and the minimum $f(x)$ must occur on an edge or vertex of this polytope. The simplex method is an organized search of the vertices by moving along the steepest edge of the polytope until x^* is obtained, as shown in Figure 6.2.

Example 6.4

Minimize

$$f(x) : -6x_1 - 14x_2$$

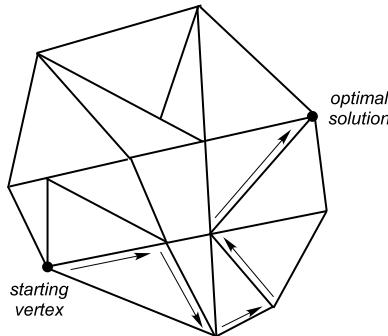
subject to the following constraints:

$$2x_1 + x_2 \leq 12$$

$$2x_1 + 3x_2 \leq 15$$

$$x_1 + 7x_2 \leq 21$$

$$x_1 \geq 0, x_2 \geq 0$$

**FIGURE 6.2**

Example of a simplex method search

Solution 6.4 Introduce slack variables x_3 , x_4 , and x_5 such that the problem becomes

Minimize

$$f(x) : -6x_1 - 14x_2 + 0x_3 + 0x_4 + 0x_5$$

subject to the following constraints:

$$\begin{aligned} 2x_1 + x_2 + x_3 &= 12 \\ 2x_1 + 3x_2 + x_4 &= 15 \\ x_1 + 7x_2 + x_5 &= 21 \end{aligned}$$

$$x_1 \geq 0, x_2 \geq 0, x_3 \geq 0, x_4 \geq 0, x_5 \geq 0$$

Form the tableau:

-6	-14	0	0	0	0
2	1	1	0	0	12
2	3	0	1	0	15
1	7	0	0	1	21
0	0	12	15	21	

The starting vector is $x = [0 \ 0 \ 12 \ 15 \ 21]^T$. The current value of $f(x)$ is displayed in the top right corner. In the next step, the function $f(x)$ is examined to determine which variable will cause the greatest decrease. Since -14 is more negative than -6 , a unit increase in x_2 will decrease $f(x)$ faster than a unit increase in x_1 . Therefore, hold x_1 constant at zero and allow x_2 to increase as much as possible (i.e., traverse one edge of the polytope to the next vertex). To determine the new value of x_2 consider the constraints

$$0 \leq x_3 = 12 - x_2$$

$$\begin{aligned} 0 &\leq x_4 = 15 - 3x_2 \\ 0 &\leq x_5 = 21 - 7x_2 \end{aligned}$$

From these constraints, the possible values of x_2 are $x_2 \leq 12$, $x_2 \leq 5$, $x_2 \leq 3$. The most stringent is $x_2 \leq 3$; therefore, x_2 can be increased to 3, and

$$\begin{aligned} x_3 &= 12 - x_2 = 9 \\ x_4 &= 15 - 3x_2 = 6 \\ x_5 &= 21 - 7x_2 = 0 \end{aligned}$$

yielding the new vector $x = [0 \ 3 \ 9 \ 6 \ 0]^T$ and $f(x) = -42$.

The new basic (nonzero) variables are x_2 , x_3 , and x_4 ; therefore, $f(x)$ must be expressed in terms of x_1 and x_5 . By substitution,

$$x_2 = \frac{(21 - x_5 - x_1)}{7}$$

and

$$f(x) = -6x_1 - 14x_2 = -6x_1 - 14\frac{(21 - x_5 - x_1)}{7} = -4x_1 + 2x_5 - 42$$

To satisfy the rule that a basic variable must occur in only one row, Gaussian elimination is used to eliminate x_2 from every row but one using the pivot as defined in Step 3 of the algorithm.

The new tableau after Gaussian elimination is

-4	0	0	0	2	-42
$\frac{13}{7}$	0	1	0	$-\frac{1}{7}$	9
$\frac{11}{7}$	0	0	1	$-\frac{3}{7}$	6
$\frac{1}{7}$	1	0	0	$\frac{1}{7}$	3
0	3	9	6	0	

The method is again repeated. Any increase in x_5 will increase $f(x)$, so x_1 is chosen to become a basic variable. Therefore, x_5 is held at zero and x_1 is allowed to increase as much as possible. The new constraints are

$$\begin{aligned} 0 &\leq x_3 = 9 - \frac{13}{7}x_1 \\ 0 &\leq x_4 = 6 - \frac{11}{7}x_1 \\ 0 &\leq 7x_2 = 21 - x_1 \end{aligned}$$

or $x_1 \leq \frac{63}{13}$, $x_1 \leq \frac{42}{11}$, and $x_1 \leq 21$. The most stringent constraint is $x_1 \leq \frac{42}{11}$, therefore, x_1 is set to $\frac{42}{11}$, and the new values of x_2 , x_3 , and x_4 are computed.

The new vector is $x = [\frac{42}{11} \ \frac{27}{11} \ \frac{21}{11} \ 0 \ 0]^T$ and $f(x)$ is rewritten in terms of x_4

and x_5 :

$$\begin{aligned}x_1 &= \frac{7}{11}(6 - x_4) \\f(x) &= -4x_1 + 2x_5 - 42 \\&= \frac{28}{11}x_4 + 2x_5 - \frac{630}{11}\end{aligned}$$

Since all coefficients in $f(x)$ are positive, the simplex method terminates because any increase in x_4 or x_5 will increase $f(x)$. This signifies that the current x vector is the solution and the final value of $f(x)$ is $-\frac{630}{11}$. ■

6.2.2 Interior Point Method

A different type of method for linear programming problems is interior point methods (also known as Karmarkar's methods), whose complexity is polynomial for both average and worst case. The simplex method has the potential to have a worst case scenario of exponential complexity that can occur in the situation in which the solution visits *every vertex* in the feasible region before reaching the optimum. For this reason, interior point methods have received considerable attention over the past few decades. Interior point methods construct a sequence of strictly feasible points (i.e., lying in the interior of the polytope but never on its boundary) that converges to the solution.

The interior point method constructs a series of feasible solutions x^0, x^1, \dots that must satisfy $Ax_i = b$. Since x^0 satisfies $Ax^0 = b$ and the next point must satisfy $Ax^1 = b$, the difference in solutions must satisfy $A\Delta x = 0$. In other words, each step must lie in the nullspace of A , which is parallel to the feasible set. Projecting $-c$ onto that nullspace gives the direction of most rapid change. However, if the iterate x^k is close to the boundary (as with \hat{x}^k in Figure 6.3), very little improvement will occur. If, however, the current iterate is near the center (as with \bar{x}^k), there could be significant improvement. One of the key aspects of the interior point method is that a transformation is applied such that the current feasible point is moved (through the transformation) to the center of the interior. The new direction is then computed and the interior point is transformed back to the original space.

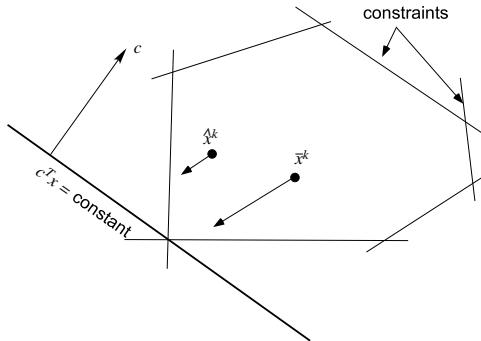
This direction of change is the *projected gradient direction*, or p^k , and the feasible points are updated through

$$x^{k+1} = x^k + \alpha p^k \quad (6.55)$$

where $\alpha > 0$ is the step length. Since the feasible points must lie in the null space of A , each p^k must be orthogonal to the rows of A . The projection matrix P

$$P = I - A^T(AA^T)^{-1}A \quad (6.56)$$

will transform any vector v into $Pv = p$, and p will be in the null space of A because AP is the zero matrix.

**FIGURE 6.3**

Comparison of two different points in the interior

Since projecting $-c$ onto the nullspace gives the direction of most rapid change, then to maintain feasibility a new iterate must satisfy $p^k = -Pc$. To remain in the interior of the space, the step length α is chosen at each step to ensure feasibility of the nonnegativity constraints. To ensure that the updates remain in the interior of the feasible space, the step length is chosen to be less than the full distance to the boundary, usually $0.5 \leq \alpha \leq 0.98$.

The last aspect is the transformation required to center the iterate in the feasible space. This is accomplished by a scaling, such that the iterate is equidistant from all constraint boundaries in the transformed feasible space. Therefore, after rescaling, $x^k = e$, where $e = [1 \ 1 \ \dots \ 1]^T$. Let $D = \text{diag}(x^k)$ be the diagonal matrix with each component of the current iterate x^k on the diagonals. This is accomplished by letting $x = D\hat{x}$ so that $\hat{x}^k = e$. The new problem to be solved then becomes

$$\text{minimize } \hat{c}^T \hat{x} = z \quad (6.57)$$

$$\text{subject to } \hat{A}\hat{x} \leq b \quad (6.58)$$

$$\hat{x} \geq 0 \quad (6.59)$$

where $\hat{c} = Dc$ and $\hat{A} = AD$. After scaling, the projection matrix becomes

$$\hat{P} = I - \hat{A}^T(\hat{A}\hat{A}^T)^{-1}\hat{A} \quad (6.60)$$

Hence, at each iteration k , the iterate x^k is rescaled to $\hat{x}^k = e$ and the update is given by

$$\hat{x}^{k+1} = e - \alpha \hat{P}\hat{c} \quad (6.61)$$

and then the updated iterate is transformed back to the original space:

$$x^{k+1} = D\hat{x}^{k+1} \quad (6.62)$$

This process repeats until $\|x^{k+1} - x^k\| < \varepsilon$. This process is often referred to as the *primal affine* interior point method [10]. The steps of the method are summarized:

Primal Affine Interior Point Method

1. Let $k = 0$.
 2. Let $D = \text{diag}(x^k)$.
 3. Compute $\hat{A} = AD$, $\hat{c} = Dc$.
 4. Compute \hat{P} from Equation (6.60).
 5. Set $p^k = \hat{P}\hat{c}$.
 6. Set $\theta = -\min_j p_j^k$. The factor θ is used to determine the maximum step length that can be taken before exiting the feasible region.
 7. Compute
- $$\hat{x}^{k+1} = e + \frac{\alpha}{\theta} p^k$$
8. Compute $x^{k+1} = D\hat{x}^{k+1}$.
 9. If $\|x^{k+1} - x^k\| < \varepsilon$, then done. Else set $k = k + 1$. Go to step 2.

Example 6.5

Repeat Example 6.4 using the primal affine interior point method.

Solution 6.5 The problem is restated (with slack variables included) for convenience:

Minimize

$$f(x) : -6x_1 - 14x_2 + 0x_3 + 0x_4 + 0x_5 = z$$

subject to the following constraints:

$$\begin{array}{rcl} 2x_1 + x_2 + x_3 & = 12 \\ 2x_1 + 3x_2 + x_4 & = 15 \\ x_1 + 7x_2 + x_5 & = 21 \end{array}$$

$$x_1 \geq 0, x_2 \geq 0, x_3 \geq 0, x_4 \geq 0, x_5 \geq 0$$

A feasible initial starting solution is

$$x^o = [1 \ 1 \ 9 \ 10 \ 13]$$

with $z^0 = c^T x^0 = -20$. The first scaling matrix is

$$D = \begin{bmatrix} 1 \\ & 1 \\ & & 9 \\ & & & 10 \\ & & & & 13 \end{bmatrix}$$

The rescaled matrix \hat{A} and objective function vector \hat{c} are computed as

$$\hat{A} = AD = \begin{bmatrix} 2 & 1 & 1 & 0 & 0 \\ 2 & 3 & 0 & 1 & 0 \\ 1 & 7 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ & 1 \\ & & 9 \\ & & & 10 \\ & & & & 13 \end{bmatrix} = \begin{bmatrix} 2 & 1 & 9 & 0 & 0 \\ 2 & 3 & 0 & 10 & 0 \\ 1 & 7 & 0 & 0 & 13 \end{bmatrix}$$

$$\hat{c} = Dc = \begin{bmatrix} 1 \\ & 1 \\ & & 9 \\ & & & 10 \\ & & & & 13 \end{bmatrix} \begin{bmatrix} -6 \\ -14 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} -6 \\ -14 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

The projection matrix \hat{P} is

$$\begin{aligned} \hat{P} &= I - \hat{A}^T (\hat{A} \hat{A}^T)^{-1} \hat{A} \\ &= \begin{bmatrix} 1 \\ & 1 \\ & & 1 \\ & & & 1 \\ & & & & 1 \end{bmatrix} - \begin{bmatrix} 2 & 2 & 1 \\ 1 & 3 & 7 \\ 9 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 13 \end{bmatrix} \begin{bmatrix} 2 & 1 & 9 & 0 & 0 \\ 2 & 3 & 0 & 10 & 0 \\ 1 & 7 & 0 & 0 & 13 \end{bmatrix} \begin{bmatrix} 2 & 2 & 1 \\ 1 & 3 & 7 \\ 9 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 13 \end{bmatrix}^{-1} \begin{bmatrix} 2 & 1 & 9 & 0 & 0 \\ 2 & 3 & 0 & 10 & 0 \\ 1 & 7 & 0 & 0 & 13 \end{bmatrix} \\ &= \begin{bmatrix} 0.9226 & -0.0836 & -0.1957 & -0.1595 & -0.0260 \\ -0.0836 & 0.7258 & -0.0621 & -0.2010 & -0.3844 \\ -0.1957 & -0.0621 & 0.0504 & 0.0578 & 0.0485 \\ -0.1595 & -0.2010 & 0.0578 & 0.0922 & 0.1205 \\ -0.0260 & -0.3844 & 0.0485 & 0.1205 & 0.2090 \end{bmatrix} \end{aligned}$$

The projected gradient is

$$\begin{aligned} p^0 &= -\hat{P}\hat{c} = - \begin{bmatrix} 0.9226 & -0.0836 & -0.1957 & -0.1595 & -0.0260 \\ -0.0836 & 0.7258 & -0.0621 & -0.2010 & -0.3844 \\ -0.1957 & -0.0621 & 0.0504 & 0.0578 & 0.0485 \\ -0.1595 & -0.2010 & 0.0578 & 0.0922 & 0.1205 \\ -0.0260 & -0.3844 & 0.0485 & 0.1205 & 0.2090 \end{bmatrix} \begin{bmatrix} -6 \\ -14 \\ 0 \\ 0 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 4.3657 \\ 9.6600 \\ -2.0435 \\ -3.7711 \\ -5.5373 \end{bmatrix} \end{aligned}$$

Calculate $\theta = -\min_j p_j^0 = 5.5373$. Rescale the current iterate to $\hat{x}^0 = D^{-1}x^0 = e$ and move to \hat{x}^1 in the transformed space with $\alpha = 0.9$:

$$\hat{x}^1 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + \alpha p^0 = \begin{bmatrix} 1.7096 \\ 2.5701 \\ 0.6679 \\ 0.3871 \\ 0.1000 \end{bmatrix}$$

Transforming this point back to the original space:

$$x^1 = D\hat{x}^1 = \begin{bmatrix} 1.7096 \\ 2.5701 \\ 6.0108 \\ 3.8707 \\ 1.3000 \end{bmatrix}$$

and the updated cost function is $c^T x^1 = -46.2383$.

Performing one more iteration (and omitting the detailed text) yields

$$\begin{aligned} \hat{A} &= \begin{bmatrix} 3.4191 & 2.5701 & 6.0108 & 0 & 0 \\ 3.4191 & 7.7102 & 0 & 3.8707 & 0 \\ 1.7096 & 17.9904 & 0 & 0 & 1.3000 \end{bmatrix} \\ \hat{c} &= \begin{bmatrix} -10.2574 \\ -35.9809 \\ 0 \\ 0 \\ 0 \end{bmatrix} \\ \hat{P} &= \begin{bmatrix} 0.5688 & -0.0584 & -0.2986 & -0.3861 & 0.0606 \\ -0.0584 & 0.0111 & 0.0285 & 0.0295 & -0.0766 \\ -0.2986 & 0.0285 & 0.1577 & 0.2070 & -0.0017 \\ -0.3861 & 0.0295 & 0.2070 & 0.2822 & 0.0991 \\ 0.0606 & -0.0766 & -0.0017 & 0.0991 & 0.9803 \end{bmatrix} \\ p^1 &= \begin{bmatrix} 3.7321 \\ -0.2004 \\ -2.0373 \\ -2.8975 \\ -2.1345 \end{bmatrix} \\ \hat{x}^1 &= \begin{bmatrix} 2.1592 \\ 0.9378 \\ 0.3672 \\ 0.1000 \\ 0.3370 \end{bmatrix} \end{aligned}$$

$$x^1 = \begin{bmatrix} 3.6914 \\ 2.4101 \\ 2.2072 \\ 0.3871 \\ 0.4381 \end{bmatrix}$$

and the updated cost function is $c^T x^1 = -55.8892$.

This continues until $\|x^{k+1} - x^k\| < \varepsilon$, at which time the solution is

$$x^* = \begin{bmatrix} 3.8182 \\ 2.4545 \\ 1.9091 \\ 0.0000 \\ 0.0000 \end{bmatrix}$$

and the updated cost function is $c^T x^1 = -57.2727$. Both the resulting solution and cost function are the same as the simplex method. ■

6.3 Nonlinear Programming

Continuous nonlinear optimization problems are typically of the following form:

$$\text{minimize } f(x) \quad x \in \mathbb{R}^n \quad (6.63)$$

$$\text{subject to } c_i(x) = 0, \quad i \in \xi \quad (6.64)$$

$$h_i(x) \geq 0, \quad i \in \Xi \quad (6.65)$$

where $[c(x) \ h(x)]$ is an m -vector of nonlinear constraint functions such that ξ and Ξ are nonintersecting index sets. The function $f(x)$ is sometimes referred to as a “cost” function. It is assumed throughout that f , c , and h are twice-continuously differentiable. Any point x satisfying the constraints of Equations (6.64) and (6.65) is called a *feasible* point, and the set of all such points is the *feasible region*. These types of problems are known generically as *nonlinear programming problems*, or NLP.

Often in optimization problems, it is convenient to refer to the “Karush–Kuhn–Tucker” (or KKT) conditions. The first-order KKT conditions for the inequality-constrained problem hold at the point x^* , if there exists an m -vector λ^* called a Lagrange-multiplier vector, such that [4]

$$c(x^*) \geq 0 \quad (\text{feasibility condition}) \quad (6.66)$$

$$g(x^*) = J(x^*)^T \lambda^* \quad (\text{stationarity condition}) \quad (6.67)$$

$$\lambda^* \geq 0 \quad (\text{nonnegativity of the multipliers}) \quad (6.68)$$

$$c(x^*) \cdot \lambda^* = 0 \quad (\text{complementarity}) \quad (6.69)$$

The stationarity condition (6.67) can be written as

$$\nabla_x L(x^*, \lambda^*) = 0, \text{ where } L(x, \lambda) \triangleq f(x) - \lambda^T c(x) \quad (6.70)$$

where λ is often generically known as a *Lagrangian multiplier* and Equation (6.70) as the Lagrangian equation. The Karush–Kuhn–Tucker conditions are necessary for a solution in nonlinear programming to be optimal.

6.3.1 Quadratic Programming

A special subset of nonlinear problems is *quadratic problems* that are characterized by the following formulation:

$$\text{minimize } f(x) = \frac{1}{2}x^T Qx + c^T x \quad x \in \mathbb{R}^n \quad (6.71)$$

$$\text{subject to } Ax \leq b \quad (6.72)$$

$$x \geq 0 \quad (6.73)$$

If Q is a positive semidefinite matrix, then $f(x)$ is a convex function. If Q is zero, then the problem becomes a linear program. The Lagrangian function for the quadratic program is given by

$$L(x, \lambda) = c^T x + \frac{1}{2}x^T Qx + \lambda(Ax - b) \quad (6.74)$$

where λ is an m dimensional row vector. The KKT conditions for a local minima are

$$c^T + x^T Q + \lambda A \geq 0 \quad (6.75)$$

$$Ax - b \leq 0 \quad (6.76)$$

$$x^T(c + Qx + A^T\lambda) = 0 \quad (6.77)$$

$$\lambda(Ax - b) = 0 \quad (6.78)$$

$$x \geq 0 \quad (6.79)$$

$$\lambda \geq 0 \quad (6.80)$$

To put these equations in a more manageable form, the nonnegative slack variable $y \in \mathbb{R}^n$ is introduced to the inequalities in Equation (6.75) and the slack variable $v \in \mathbb{R}^m$ in Equation (6.76) to obtain

$$c + Qx + A^T\lambda^T - y = 0 \quad (6.81)$$

$$Ax - b + v = 0 \quad (6.82)$$

and the KKT equations are

$$Qx + A^T\lambda^T - y = -c^T \quad (6.83)$$

$$Ax + v = b \quad (6.84)$$

$$x \geq 0 \quad (6.85)$$

$$y \geq 0 \quad (6.86)$$

$$v \geq 0 \quad (6.87)$$

$$\lambda \geq 0 \quad (6.88)$$

$$y^T x = 0 \quad (6.89)$$

$$\lambda v = 0 \quad (6.90)$$

Now any linear programming method can be used to solve this set of equations by treating the complementary slackness conditions (Equations (6.89 and 6.90)) implicitly with a restricted basis entry rule. The goal is to find the solution to the linear program problem with the additional requirement that the complementary slackness conditions be satisfied at each iteration. The objective function is satisfied by adding an artificial variable to each equation and minimizing the sum of the artificial variables.

Example 6.6

Minimize

$$f(x) : -10x_1 - 8x_2 + x_1^2 + 2x_2^2$$

subject to the following constraints:

$$x_1 + x_2 \leq 10$$

$$x_2 \leq 5$$

$$x_1 \geq 0, x_2 \geq 0$$

Solution 6.6 This can be written in matrix form as

Minimize

$$f(x) : [-10 \ -8] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \frac{1}{2} [x_1 \ x_2] \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

subject to the following constraints:

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \leq \begin{bmatrix} 10 \\ 5 \end{bmatrix}$$

The linear program problem becomes (where the artificial variables are given by $a_1 - a_4$)

$$\text{Minimize } a_1 + a_2 + a_3 + a_4$$

subject to

$$\begin{array}{rcl} 2x_1 & +\lambda_1 & -y_1 & +a_1 & = 10 \\ & 4x_2 & +\lambda_1 + \lambda_2 & -y_2 & +a_2 & = 8 \\ x_1 & +x_2 & & +v_1 & +a_3 & = 10 \\ & x_2 & & +v_2 & +a_4 & = 5 \end{array}$$

This can now be solved using either the simplex or interior point method. The solution is

$$\begin{aligned}x &= \begin{bmatrix} 5.0000 \\ 2.0000 \end{bmatrix} \\ \lambda &= \begin{bmatrix} 0.6882 \\ 0.7439 \end{bmatrix} \\ y &= \begin{bmatrix} 0.6736 \\ 1.3354 \end{bmatrix} \\ v &= \begin{bmatrix} 3.0315 \\ 3.0242 \end{bmatrix} \\ a &= [0 \ 0 \ 0 \ 0]\end{aligned}$$

with the cost function $f(x) = -33$. ■

6.3.2 Steepest Descent Algorithm

For engineering applications, general nonlinear programming problems are generally solved by two classes of approaches:

1. gradient methods such as steepest descent, or
2. iterative programming techniques such as successive quadratic programming

In an unconstrained system, the usual approach to minimizing the function $f(x)$ is to set the function derivatives to zero and then solve for the system states from the set of resulting equations. In the majority of applications, however, the system states evaluated at the unconstrained minimum will not satisfy the constraint equations. Thus an alternate approach is required to find the constrained minimum. One approach is to introduce an additional set of parameters λ , frequently known as *Lagrange multipliers*, to impose the constraints on the cost function. The augmented cost function then becomes

$$\text{minimize } f(x) - \lambda c(x) \quad (6.91)$$

The augmented function in Equation (6.91) can then be minimized by solving for the set of states that result from setting the derivatives of the augmented function to zero. Note that the derivative of Equation (6.91) with respect to λ effectively enforces the equality constraint of Equation (6.64).

Example 6.7

Minimize

$$C : \frac{1}{2} (x^2 + y^2) \quad (6.92)$$

subject to the following constraint:

$$2x - y = 5$$

Solution 6.7 Note that the function to be minimized is the equation for a circle. The unconstrained minimum of this function is the point at the origin with $x = 0$ and $y = 0$ which defines a circle with a radius of zero length. However, the circle must also intersect the line defined by the constraint equation; thus the constrained circle must have a nonzero radius. The augmented cost function becomes

$$C^* : \frac{1}{2} (x^2 + y^2) - \lambda (2x - y - 5) \quad (6.93)$$

where λ represents the Lagrange multiplier. Setting the derivatives of the augmented cost function to zero yields the following set of equations:

$$\begin{aligned} 0 &= \frac{\partial C^*}{\partial x} = x - 2\lambda \\ 0 &= \frac{\partial C^*}{\partial y} = y + \lambda \\ 0 &= \frac{\partial C^*}{\partial \lambda} = 2x - y - 5 \end{aligned}$$

Solving this set of equations yields $[x \ y \ \lambda]^T = [2 \ -1 \ 1]^T$. The cost function of Equation (6.92) evaluated at the minimum of the augmented cost function is

$$C : \frac{1}{2} ((2)^2 + (-1)^2) = \frac{5}{2}$$

■

Both the cost function f and the constraint equations c may be a function of an external input u , which may be varied to minimize the function $f(x)$. In this case, Equation (6.91) may be more generally written as

$$\text{minimize } f(x, u) - \lambda c(x, u) \quad (6.94)$$

If there is more than one equality constraint, then λ becomes a vector of multipliers and the augmented cost function becomes

$$C^* : f(x) - [\lambda]^T c(x) \quad (6.95)$$

where the derivatives of C^* become

$$\left[\frac{\partial C^*}{\partial \lambda} \right] = 0 = c(x) \quad (6.96)$$

$$\left[\frac{\partial C^*}{\partial x} \right] = 0 = \left[\frac{\partial f}{\partial x} \right] - \left[\frac{\partial c}{\partial x} \right]^T [\lambda] \quad (6.97)$$

$$\left[\frac{\partial C^*}{\partial u} \right] = 0 = \left[\frac{\partial f}{\partial u} \right] - \left[\frac{\partial c}{\partial u} \right]^T [\lambda] \quad (6.98)$$

Note that, for any *feasible* solution to the equality constraint, Equation (6.96) is satisfied, but the feasible solution may not be the optimal solution which minimizes the cost function. In this case, $[\lambda]$ can be obtained from Equation (6.97) and then only

$$\left[\frac{\partial C^*}{\partial u} \right] \neq 0$$

This vector can be used as a gradient vector $[\nabla C]$ which is orthogonal to the contour of constant values of the cost function C . Thus

$$[\lambda] = \left[\left[\frac{\partial c}{\partial x} \right]^T \right]^{-1} \left[\frac{\partial f}{\partial x} \right] \quad (6.99)$$

which gives

$$\nabla C = \left[\frac{\partial C^*}{\partial u} \right] = \left[\frac{\partial f}{\partial u} \right] - \left[\frac{\partial c}{\partial u} \right]^T [\lambda] \quad (6.100)$$

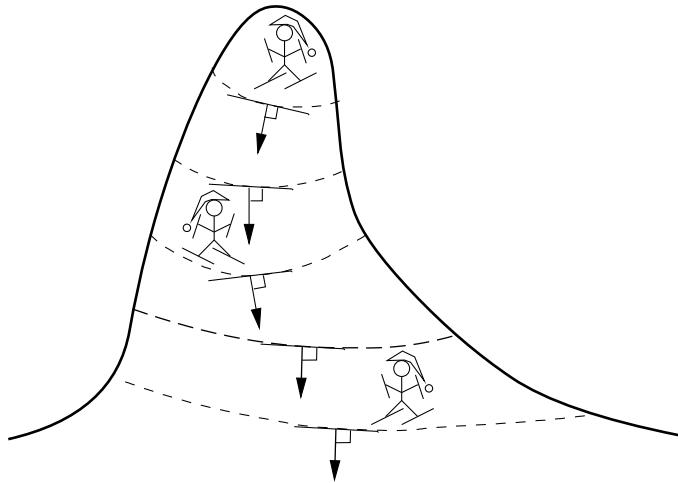
$$= \left[\frac{\partial f}{\partial u} \right] - \left[\frac{\partial c}{\partial u} \right]^T \left[\left[\frac{\partial c}{\partial x} \right]^T \right]^{-1} \left[\frac{\partial f}{\partial x} \right] \quad (6.101)$$

This relationship provides the foundation of the optimization method known as the *steepest descent* algorithm.

Steepest Descent Algorithm

1. Let $k = 0$. Guess an initial vector $u^k = u^0$.
2. Solve the (possibly nonlinear) system of Equations (6.96) for a feasible solution x .
3. Calculate C^{k+1} and ∇C^{k+1} from Equation (6.101). If $\|C^{k+1} - C^k\|$ is less than some predefined tolerance, stop.
4. Calculate the new vector $u^{k+1} = u^k - \gamma \nabla C^{k+1}$, where γ is a positive number which is the user-defined “step size” of the algorithm.
5. $k = k + 1$. Go to step 2.

In the steepest descent method, the u vector update direction is determined at each step of the algorithm by choosing the direction of the greatest change of the augmented cost function C^* . For example, consider a person skiing from the top of a mountain to the bottom, as illustrated in Figure 6.4. The skier will travel in a straight path for a certain distance. At that point, he may no longer be pointed directly down the mountain. Therefore, he will adjust his direction so that his skis point in the direction of steepest descent. The direction of steepest descent is perpendicular to the tangent of the curve of

**FIGURE 6.4**

Example of steepest descent

constant altitude (or cost). The distance the skier travels between adjustments is analogous to the step size γ of the algorithm. For small γ , the skier will frequently alter direction; thus his descent will be slow. For large γ , however, he may overshoot the foot of the mountain and will start ascending again. Thus the critical part of the steepest descent algorithm is the choice of γ . If γ is chosen small, then convergence to a minimum value is more likely, but may require many iterations, whereas a large value of γ may result in oscillations about the minimum.

Example 6.8

Minimize

$$C : x_1^2 + 2x_2^2 + u^2 = f(x_1, x_2, u) \quad (6.102)$$

subject to the following constraints:

$$0 = x_1^2 - 3x_2 + u - 3 \quad (6.103)$$

$$0 = x_1 + x_2 - 4u + 2 \quad (6.104)$$

Solution 6.8 To find ∇C of Equation (6.101), the following partial derivatives are required:

$$\left[\frac{\partial f}{\partial u} \right] = 2u$$

$$\begin{aligned}\left[\frac{\partial f}{\partial x} \right] &= \begin{bmatrix} 2x_1 \\ 4x_2 \end{bmatrix} \\ \left[\frac{\partial c}{\partial u} \right]^T &= [1 \ -4] \\ \left[\frac{\partial c}{\partial x} \right] &= \begin{bmatrix} 2x_1 & -3 \\ 1 & 1 \end{bmatrix}\end{aligned}$$

yielding

$$\begin{aligned}\nabla C &= \left[\frac{\partial f}{\partial u} \right] - \left[\frac{\partial c}{\partial u} \right]^T \left[\left[\frac{\partial c}{\partial x} \right]^T \right]^{-1} \left[\frac{\partial f}{\partial x} \right] \\ &= 2u - [1 \ -4] \left[\begin{bmatrix} 2x_1 & -3 \\ 1 & 1 \end{bmatrix}^T \right]^{-1} \begin{bmatrix} 2x_1 \\ 4x_2 \end{bmatrix}\end{aligned}$$

Iteration 1

Let $u = 1$, $\gamma = 0.05$, and choose a stopping criterion of $\epsilon = 0.0001$. Solving for x_1 and x_2 yields two values for each with a corresponding cost function:

$$\begin{aligned}x_1 &= 1.7016 & x_2 &= 0.2984 & f &= 4.0734 \\ x_1 &= -4.7016 & x_2 &= 6.7016 & f &= 23.2828\end{aligned}$$

The first set of values leads to the minimum cost function, so they are selected as the operating solution. Substituting $x_1 = 1.7016$ and $x_2 = 0.2984$ into the gradient function yields $\nabla C = 10.5705$ and the new value of u becomes

$$\begin{aligned}u^{(2)} &= u^{(1)} - \gamma \nabla C \\ &= 1 - (0.05)(10.5705) \\ &= 0.4715\end{aligned}$$

Iteration 2

With $u = 0.4715$, solving for x_1 and x_2 again yields two values for each with a corresponding cost function:

$$\begin{aligned}x_1 &= 0.6062 & x_2 &= -0.7203 & f &= 1.6276 \\ x_1 &= -3.6062 & x_2 &= 3.4921 & f &= 14.2650\end{aligned}$$

The first set of values again leads to the minimum cost function, so they are selected as the operating solution. The difference in cost functions is

$$|C^{(1)} - C^{(2)}| = |4.0734 - 1.6276| = 2.4458$$

which is greater than the stopping criterion. Substituting these values into the gradient function yields $\nabla C = 0.1077$ and the new value of u becomes

$$\begin{aligned}u^{(3)} &= u^{(2)} - \gamma \nabla C \\ &= 0.4715 - (0.05)(0.1077) \\ &= 0.4661\end{aligned}$$

Iteration 3

With $u = 0.4661$, solving for x_1 and x_2 again yields two values for each with a corresponding cost function:

$$\begin{aligned} x_1 &= 0.5921 \quad x_2 = -0.7278 \quad f = 1.6271 \\ x_1 &= -3.5921 \quad x_2 = 3.4565 \quad f = 14.1799 \end{aligned}$$

The first set of values again leads to the minimum cost function, so they are selected as the operating solution. The difference in cost functions is

$$|C^{(2)} - C^{(3)}| = |1.6276 - 1.6271| = 0.0005$$

which is greater than the stopping criterion. Substituting these values into the gradient function yields $\nabla C = 0.0541$ and the new value of u becomes

$$\begin{aligned} u^{(4)} &= u^{(3)} - \gamma \nabla C \\ &= 0.4661 - (0.05)(0.0541) \\ &= 0.4634 \end{aligned}$$

This process is continued until the stopping criterion is reached. The final values are $x_1 = 0.5774$, $x_2 = -0.7354$, and $u = 0.4605$ which yield the minimum cost function $f = 1.6270$. ■

6.3.3 Sequential Quadratic Programming Algorithm

Gradient descent techniques work well for small nonlinear systems, but become inefficient as the dimension of the search space grows. The nonlinear sequential quadratic programming (SQP) optimization method is computationally efficient and has been shown to exhibit super linear convergence for convex search spaces [5]. The SQP is an iterative procedure which models the nonlinear optimization problem for a given iterate x^k by a quadratic programming subproblem, whose solution is used to construct a new iterate x^{k+1} . Thus it sequentially solves the QP algorithm to produce the solution to the original nonlinear problem in much the same way that the Newton–Raphson sequentially solves a series of linear problems to find the solution to a nonlinear problem.

The SQP method also attempts to solve the system

$$\text{minimize } f(x) \quad x \in \mathbb{R}^n \tag{6.105}$$

$$\text{subject to } c_i(x) = 0, \quad i \in \xi \tag{6.106}$$

$$h_i(x) \geq 0, \quad i \in \Xi \tag{6.107}$$

As before, the usual approach to solving this optimization problem is to use Lagrangian multipliers and minimize the hybrid system:

$$L(x, \lambda) = f(x) + \lambda^T c(x) + \pi^T h(x) \tag{6.108}$$

The KKT conditions are

$$\begin{aligned}\nabla f(x) + C^T \lambda + H^T \pi &= 0 \\ c(x) &= 0 \\ h(x) + s &= 0 \\ \pi^T s &= 0 \\ \pi, s &\geq 0\end{aligned}$$

where λ is a vector of Lagrange multipliers for the equality constraints, π is a vector of Lagrange multipliers for the inequality constraints, s is a vector of slack variables, and

$$C = \frac{\partial c(x)}{\partial x} \quad (6.109)$$

$$H = \frac{\partial h(x)}{\partial x} \quad (6.110)$$

This nonlinear system can be solved for x, λ, π , and s using the Newton–Raphson method. Consider the case of only x and λ . Using the Newton–Raphson method to solve for $y = [x \ \lambda]^T$ with

$$F(y) = 0 = \begin{bmatrix} \nabla f(x) + C^T \lambda \\ c(x) \end{bmatrix}$$

and the Newton–Raphson update becomes

$$y^{k+1} = y^k - [\nabla F_k]^{-1} F(y^k)$$

or in original variables

$$\begin{bmatrix} x^{k+1} \\ \lambda^{k+1} \end{bmatrix} = \begin{bmatrix} x^k \\ \lambda^k \end{bmatrix} - \left[\begin{bmatrix} \nabla^2 L & \nabla c^T \\ \nabla c^T & 0 \end{bmatrix}_k \right]^{-1} \begin{bmatrix} \nabla L \\ c \end{bmatrix}_k \quad (6.111)$$

Example 6.9

Repeat Example 6.8 using the SQP method.

Solution 6.9 The problem is repeated here:

Minimize

$$x_1^2 + 2x_2^2 + u^2 = f(x_1, x_2, u) \quad (6.112)$$

subject to the following constraints:

$$0 = x_1^2 - 3x_2 + u - 3 \quad (6.113)$$

$$0 = x_1 + x_2 - 4u + 2 \quad (6.114)$$

Applying the KKT conditions yields the following nonlinear system of equations:

$$\begin{aligned} 0 &= 2x_1 + 2\lambda_1 x_1 + \lambda_2 \\ 0 &= 4x_2 - 3\lambda_1 + \lambda_2 \\ 0 &= 2u + \lambda_1 - 4\lambda_2 \\ 0 &= x_1^2 - 3x_2 + u - 3 \\ 0 &= x_1 + x_2 - 4u + 2 \end{aligned}$$

The Newton–Raphson iteration becomes

$$\begin{bmatrix} x_1 \\ x_2 \\ u \\ \lambda_1 \\ \lambda_2 \end{bmatrix}^{k+1} = \begin{bmatrix} x_1 \\ x_2 \\ u \\ \lambda_1 \\ \lambda_2 \end{bmatrix}^k - \begin{bmatrix} 2 + 2\lambda_1^k & 0 & 0 & 2x_1^k & 1 \\ 0 & 4 & 0 & -3 & 1 \\ 0 & 0 & 2 & 1 & -4 \\ 2x_1^k & -3 & 1 & 0 & 0 \\ 1 & 1 & -4 & 0 & 0 \end{bmatrix}^{-1} \begin{bmatrix} 2x_1^k + 2\lambda_1^k x_1^k + \lambda_2^k \\ 4x_2^k - 3\lambda_1^k + \lambda_2^k \\ 2u^k + \lambda_1^k - 4\lambda_2^k \\ (x_1^k)^2 - 3x_2^k + u^k - 3 \\ x_1^k + x_2^k - 4u^k + 2 \end{bmatrix} \quad (6.115)$$

Starting with an initial guess $[x_1 \ x_2 \ u \ \lambda_1 \ \lambda_2]^T = [1 \ 1 \ 1 \ 1 \ 1]^T$ yields the following updates:

k	x_1	x_2	u	λ_1	λ_2	$f(x)$
0	1.0000	1.0000	1.0000	1.0000	1.0000	4.0000
1	0.5774	-0.7354	0.4605	-0.9859	-0.0162	1.6270
2	0.8462	-0.5804	0.5664	-0.7413	0.0979	1.7106
3	0.6508	-0.7098	0.4853	-0.9442	0.0066	1.6666
4	0.5838	-0.7337	0.4625	-0.9831	-0.0145	1.6314
5	0.5774	-0.7354	0.4605	-0.9859	-0.0162	1.6270

which is the same solution as previously. ■

6.4 Power System Applications

6.4.1 Optimal Power Flow

Many power system applications, such as the power flow, offer only a snapshot of the system operation. Frequently, the system planner or operator is interested in the effect that making adjustments to the system parameters will have on the power flow through lines or system losses. Rather than making the adjustments in a random fashion, the system planner will attempt to optimize the adjustments according to some objective function. This objective function can be chosen to minimize generating costs, reservoir water levels, or system losses, among others. The optimal power flow problem is to formulate the power flow problem to find system voltages and generated powers within the framework of the objective function. In this application, the inputs to the power flow are systematically adjusted to maximize (or minimize) a scalar function of the power flow state variables. The two most common objective functions are minimization of generating costs and minimization of active power losses.

The time frame of optimal power flow is on the order of minutes to one hour; therefore, it is assumed that the optimization occurs using only those units that are currently on-line. The problem of determining whether or not to engage a unit, at what time, and for how long is part of the *unit commitment* problem and is not covered here. The minimization of active transmission losses saves both generating costs and creates a higher generating reserve margin.

Usually, generator cost curves (the curves that relate generated power to the cost of such generation) are given as piecewise linear incremental cost curves. This has its origin in the simplification of concave cost functions with the valve points as cost curve breakpoints [20]. Piecewise linear incremental cost curves correspond to piecewise quadratic cost curves by integrating the incremental cost curves. This type of objective function lends itself easily to the *economic dispatch*, or λ -dispatch problem where only generating units are considered in the optimization. In this process, system losses and constraints on voltages and line powers are neglected. This economic dispatch method is illustrated in the following example.

Example 6.10

Three generators with the following cost functions serve a load of 952 MW. Assuming a lossless system, calculate the optimal generation scheduling.

$$C_1 : P_1 + 0.0625P_1^2 \text{ \$/hr}$$

$$C_2 : P_2 + 0.0125P_2^2 \text{ \$/hr}$$

$$C_3 : P_3 + 0.0250P_3^2 \text{ \$/hr}$$

Solution 6.10 The first step in determining the optimal scheduling of the generators is to construct the problem in the general form. Thus the optimization statement is

$$\text{Minimize } C: P_1 + 0.0625P_1^2 + P_2 + 0.0125P_2^2 + P_3 + 0.0250P_3^2$$

$$\text{Subject to: } P_1 + P_2 + P_3 - 952 = 0$$

From this statement, the constrained cost function becomes

$$C^*: P_1 + 0.0625P_1^2 + P_2 + 0.0125P_2^2 + P_3 + 0.0250P_3^2 - \lambda(P_1 + P_2 + P_3 - 952) \quad (6.116)$$

Setting the derivatives of C^* to zero yields the following set of linear equations:

$$\begin{bmatrix} 0.125 & 0 & 0 & -1 \\ 0 & 0.025 & 0 & -1 \\ 0 & 0 & 0.050 & -1 \\ 1 & 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ \lambda \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ -1 \\ 952 \end{bmatrix} \quad (6.117)$$

Solving Equation (6.117) yields

$$P_1 = 112 \text{ MW}$$

$$P_2 = 560 \text{ MW}$$

$$P_3 = 280 \text{ MW}$$

$$\lambda = \$15/\text{MWhr}$$

for a constrained cost of \$7,616/hr. ■

This is the generation scheduling that minimizes the hourly cost of production. The value of λ is the *incremental* or *break-even* cost of production. This gives a company a price cut-off for buying or selling generation: if they can purchase generation for less than λ , then their overall costs will decrease. Likewise, if generation can be sold for greater than λ , their overall costs will decrease. Also note that at the optimal scheduling

$$\lambda = 1 + 0.125P_1 = 1 + 0.025P_2 = 1 + 0.050P_3 \quad (6.118)$$

Since λ is the incremental cost for the system, this point is also called the point of “equal incremental cost,” and the generation schedule is said to satisfy the “equal incremental cost criterion.” Any deviation in generation from the equal increment cost scheduling will result in an increase in the production cost C .

Example 6.11

If a buyer is willing to pay \$16/MWhr for generation, how much excess generation should be produced and sold, and what is the profit for this transaction?

Solution 6.11 From Example 6.10, the derivatives of the augmented cost function yield the following relationships between generation and λ :

$$P_1 = 8(\lambda - 1)$$

$$P_2 = 40(\lambda - 1)$$

$$P_3 = 20(\lambda - 1)$$

from which the equality constraint yields

$$8(\lambda - 1) + 40(\lambda - 1) + 20(\lambda - 1) - 952 = 0 \quad (6.119)$$

To determine the excess amount, the equality Equation (6.119) will be augmented and then evaluated at $\lambda = \$16/\text{MWhr}$:

$$8(16 - 1) + 40(16 - 1) + 20(16 - 1) - (952 + \text{excess}) = 0 \quad (6.120)$$

Solving Equation (6.120) yields a required excess of 68 MW, and $P_1 = 120$ MW, $P_2 = 600$ MW, and $P_3 = 300$ MW. The total cost of generation becomes

$$C : P_1 + 0.0625P_1^2 + P_2 + 0.0125P_2^2 + P_3 + 0.0250P_3^2 = \$8,670/\text{hr} \quad (6.121)$$

but the amount recovered by the sale of generation is the amount of excess times the incremental cost λ ,

$$68 \text{ MW} \times \$16/\text{MWhr} = \$1,088/\text{hr}$$

Therefore, the total cost is $\$8,670 - 1,088 = \$7,580/\text{hr}$. This amount is $\$34/\text{hr}$ less than the original cost of $\$7,616/\text{hr}$; thus $\$34/\text{hr}$ is the profit achieved from the sale of the excess generation at $\$16/\text{MWhr}$. ■

Figure 6.5 shows an incremental cost table for a medium size utility. The incremental cost of generation is listed vertically along the left-hand side of the table. The various generating units are listed across the top from least expensive to most expensive (left to right). Nuclear units are among the least expensive units to operate and the nuclear unit *Washington* at the far left can produce up to 1222 MW at an incremental cost of $7.00 \$/\text{MWhr}$. This incremental cost is half of the next least expensive unit, *Adams* at $14 \$/\text{MWhr}$, which is a coal unit. As the available units become increasingly more expensive to operate, the incremental cost also increases.

Example 6.12

What is the incremental cost for the utility to produce 12,500 MW?

Solution 6.12 To find the incremental cost that corresponds to 12,500 MW from the incremental cost table in Figure 6.5, the maximum generation available from each unit is summed until it equals 12,500 MW. This amount is

FIGURE 6.5

Incremental cost table

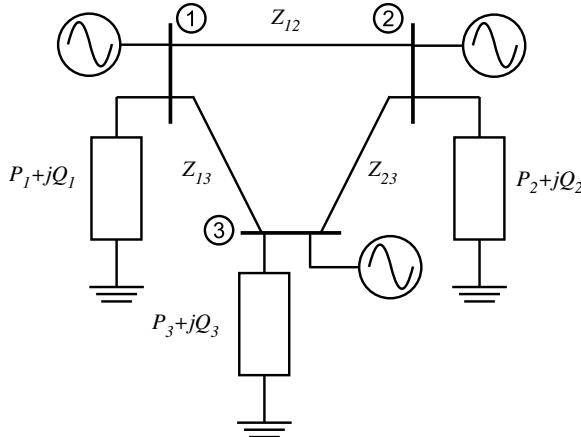
**FIGURE 6.6**

Figure for Example 6.13

reached by the gas unit of *Monroe 1-2*. This corresponds to an incremental cost of 28.00\$/MW hr. This is the breakeven point for 12,500 MW.

If power can be purchased for less than 28.00\$/MW hr, the utility should purchase generation. ■

The primary drawback with the equal incremental cost scheduling is that it neglects all losses in the system. The only enforced equality constraint is that the sum of the generation must equal the total load demand. In reality, however, the sum of the generation must equal the load demand plus any system losses. In the consideration of system losses, the equality constraints must include the set of power flow equations, and the optimization process must be extended to the steepest descent, or similar, approach [11].

Example 6.13

Consider the three machine system shown in Figure 6.6. This system has the same parameters as the three-bus system of Example 3.11 except that bus 3 has been converted to a generator bus with a voltage magnitude of 1.0 pu. The cost functions of the generators are the same as in Example 6.10. The per unit loads are

bus	P_{load}	Q_{load}
1	0.312	0.1
2	0.320	0.1
3	0.320	0.1

Using $P_{gen} = [0.112 \quad 0.560 \quad 0.280]'$ as a starting point (from Example 6.10), find the optimal scheduling of this system considering losses. Let $\gamma = 1$.

Solution 6.13 Following the steepest descent procedure detailed in Section 6.3.2, the first step is to develop an expression for the gradient ∇C , where

$$\nabla C = \left[\frac{\partial f}{\partial u} \right] - \left[\frac{\partial g}{\partial u} \right]^T \left[\left[\frac{\partial g}{\partial x} \right]^T \right]^{-1} \left[\frac{\partial f}{\partial x} \right] \quad (6.122)$$

where f is the sum of the generator costs

$$f : C_1 + C_2 + C_3 = P_1 + 0.0625P_1^2 + P_2 + 0.0125P_2^2 + P_3 + 0.0250P_3^2$$

g is the set of power flow equations

$$\begin{aligned} g_1 : 0 &= P_2 - P_{L2} - V_2 \sum_{i=1}^3 V_i Y_{2i} \cos(\theta_2 - \theta_i - \phi_{2i}) \\ g_2 : 0 &= P_3 - P_{L3} - V_3 \sum_{i=1}^3 V_i Y_{3i} \cos(\theta_3 - \theta_i - \phi_{3i}) \end{aligned}$$

where P_{Li} denotes the active power load at bus i , the set of inputs u is the set of independent generation settings

$$u = \begin{bmatrix} P_2 \\ P_3 \end{bmatrix}$$

and x is the set of unknown states

$$x = \begin{bmatrix} \theta_2 \\ \theta_3 \end{bmatrix}$$

The generator setting P_1 is not an input because it is the slack bus generation and cannot be independently set. From these designations, the various partial derivatives required for ∇C can be derived:

$$\left[\frac{\partial g}{\partial u} \right] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (6.123)$$

$$\left[\frac{\partial g}{\partial x} \right] = \begin{bmatrix} \frac{\partial g_1}{\partial \theta_2} & \frac{\partial g_1}{\partial \theta_3} \\ \frac{\partial g_2}{\partial \theta_2} & \frac{\partial g_2}{\partial \theta_3} \end{bmatrix} \quad (6.124)$$

where

$$\frac{\partial g_1}{\partial \theta_2} = V_2 (V_1 Y_{12} \sin(\theta_2 - \theta_1 - \phi_{21}) + V_3 Y_{13} \sin(\theta_2 - \theta_3 - \phi_{23})) \quad (6.125)$$

$$\frac{\partial g_1}{\partial \theta_3} = -V_2 V_3 Y_{32} \sin(\theta_2 - \theta_3 - \phi_{23}) \quad (6.126)$$

$$\frac{\partial g_2}{\partial \theta_2} = -V_3 V_2 Y_{23} \sin(\theta_3 - \theta_2 - \phi_{32}) \quad (6.127)$$

$$\frac{\partial g_2}{\partial \theta_3} = V_3 (V_1 Y_{13} \sin(\theta_3 - \theta_1 - \phi_{31}) + V_2 Y_{23} \sin(\theta_3 - \theta_2 - \phi_{32})) \quad (6.128)$$

and

$$\left[\frac{\partial f}{\partial u} \right] = \begin{bmatrix} 1 + 0.025P_2 \\ 1 + 0.050P_3 \end{bmatrix} \quad (6.129)$$

Finding the partial derivative $\left[\frac{\partial f}{\partial x} \right]$ is slightly more difficult since the cost function is not written as a direct function of x . Recall, however, that P_1 is not an input, but is actually a quantity that depends on x , i.e.,

$$P_1 = V_1 (V_1 Y_{11} \cos(\theta_1 - \theta_1 - \phi_{11}) + V_2 Y_{12} \cos(\theta_1 - \theta_2 - \phi_{12}) + V_3 Y_{13} \cos(\theta_1 - \theta_3 - \phi_{13})) \quad (6.130)$$

Thus, using the chain rule,

$$\left[\frac{\partial f}{\partial x} \right] = \left[\frac{\partial f}{\partial P_1} \right] \left[\frac{\partial P_1}{\partial x} \right] \quad (6.131)$$

$$= (1 + 0.125P_1) \begin{bmatrix} V_1 V_2 Y_{12} \sin(\theta_1 - \theta_2 - \phi_{12}) \\ V_1 V_3 Y_{13} \sin(\theta_1 - \theta_3 - \phi_{13}) \end{bmatrix} \quad (6.132)$$

The initial values of P_2 and P_3 are obtained from the equal incremental cost rule. Using $P_2 = 0.56$ pu and $P_3 = 0.28$ pu as inputs into the power flow yields the following states:

i	θ_i (rad)	V_i (pu)	$P_{gen,i}$ (pu)
1	0	1.02	0.1131
2	0.0279	1.00	0.5600
3	0.0128	1.00	0.2800

Substituting these values into the partial derivatives yields

$$\left[\frac{\partial g}{\partial u} \right] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (6.133)$$

$$\left[\frac{\partial g}{\partial x} \right] = \begin{bmatrix} -13.3048 & 9.9148 \\ 9.8849 & -19.9960 \end{bmatrix} \quad (6.134)$$

$$\left[\frac{\partial f}{\partial u} \right] = \begin{bmatrix} 1.0140 \\ 1.0140 \end{bmatrix} \quad (6.135)$$

$$\left[\frac{\partial f}{\partial x} \right] = 1.0141 \begin{bmatrix} -3.3773 \\ -10.0853 \end{bmatrix} = \begin{bmatrix} -3.4251 \\ -10.2278 \end{bmatrix} \quad (6.136)$$

which yields

$$\nabla C = \begin{bmatrix} 0.0048 \\ 0.0021 \end{bmatrix} \quad (6.137)$$

Thus the new values for the input generation are

$$\begin{bmatrix} P_2 \\ P_3 \end{bmatrix} = \begin{bmatrix} 0.560 \\ 0.280 \end{bmatrix} - \gamma \begin{bmatrix} 0.0048 \\ 0.0021 \end{bmatrix} = \begin{bmatrix} 0.5552 \\ 0.2779 \end{bmatrix} \quad (6.138)$$

Thus $P_2 = 555.2$ and $P_3 = 277.9$ MW.

Already the gradient ∇C is very small, indicating that the generation values from the equal incremental cost process were relatively close to the optimal values, even considering losses.

Proceeding until convergence is achieved yields the final generation values for all of the generators:

$$\begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix} = \begin{bmatrix} 136.6 \\ 517.9 \\ 298.4 \end{bmatrix} \text{ MW}$$

which yields a cost of \$7,698/MW hr. Note that this amount is greater than the calculated cost for the equal incremental cost function. This increase is due to the extra generation required to satisfy the losses in the system. ■

Often the steepest descent method may indicate either states or inputs lie outside of their physical constraints. For example, the algorithm may result in a power generation value that exceeds the physical maximum output of the generating unit. Similarly, the resulting bus voltages may lie outside of the desired range (usually $\pm 10\%$ of unity). These are violations of the *inequality constraints* of the problem. In these cases, the steepest descent algorithm must be modified to reflect these physical limitations. There are several approaches to account for limitations and these approaches depend on whether or not the limitation is on the input (independent) or on the state (dependent).

6.4.1.1 Limitations on Independent Variables

If the application of the steepest descent algorithm results in an updated value of input that exceeds the specified limit, then the most straightforward method of handling this violation is simply to set the input state equal to its limit and continue with the algorithm except with one less degree of freedom.

Example 6.14

Repeat Example 6.13 except that the generators must satisfy the following limitations:

$$80 \leq P_1 \leq 1200 \text{ MW}$$

$$450 \leq P_2 \leq 750 \text{ MW}$$

$$150 \leq P_3 \leq 250 \text{ MW}$$

Solution 6.14 From the solution of Example 6.13, the output of generator 3 exceeds the maximum limit of 0.25 pu. Therefore, after the first iteration in the previous example, P_3 is set to 0.25 pu. The new partial derivatives

become

$$\left[\frac{\partial g}{\partial u} \right] = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad (6.139)$$

$$\left[\frac{\partial g}{\partial x} \right] = \text{same} \quad (6.140)$$

$$\left[\frac{\partial f}{\partial u} \right] = [1 + 0.025P_2] \quad (6.141)$$

$$\left[\frac{\partial f}{\partial x} \right] = \text{same} \quad (6.142)$$

From the constrained steepest descent, the new values of generation become

$$\begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix} = \begin{bmatrix} 147 \\ 556 \\ 250 \end{bmatrix} \text{ MW}$$

with a cost of \$7,728/MWhr, which is higher than the unconstrained cost of generation of \$7,698/MWhr. As more constraints are added to the system, the system is moved away from the optimal operating point, increasing the cost of generation. ■

6.4.1.2 Limitations on Dependent Variables

In many cases, the physical limitations of the system are imposed upon states that are dependent variables in the system description. In this case, the inequality equations are functions of x and must be added to the cost function. Examples of limitations on dependent variables include maximum line flows or bus voltage levels. In these cases, the value of the states cannot be independently set, but must be enforced indirectly. One method of enforcing an inequality constraint is to introduce a *penalty function* into the cost function. A penalty function is a function that is small when the state is far away from its limit, but becomes increasingly larger the closer the state is to its limit. Typical penalty functions include

$$p(h) = e^{kh} \quad k > 0 \quad (6.143)$$

$$p(h) = x^{2n}e^{kh} \quad n, k > 0 \quad (6.144)$$

$$p(h) = ax^{2n}e^{kh} + be^{kh} \quad n, k, a, b > 0 \quad (6.145)$$

and the cost function becomes

$$C^* : \quad C(u, x) + \lambda^T g(u, x) + p(h(u, x) - h^{max}) \quad (6.146)$$

This cost equation is then minimized in the usual fashion by setting the appropriate derivatives to zero. This method has the advantage of simplicity of implementation, but also has several disadvantages. The first disadvantage

is that the choice of penalty function is often a heuristic choice and can vary by application. A second disadvantage is that this method cannot enforce *hard* limitations on states, i.e., the cost function becomes large if the maximum is exceeded, but the state is allowed to exceed its maximum. In many applications this is not a serious disadvantage. If the power flow on a transmission line slightly exceeds its maximum, it is reasonable to assume that the power system will continue to operate, at least for a finite length of time. If, however, the physical limit is the height above ground for an airplane, then even a slightly negative altitude will have dire consequences. Thus penalty functions to enforce limits must be used with caution and is not applicable for all systems.

Example 6.15

Repeat Example 6.13, except use penalty functions to limit the power flow across line 2-3 to 0.1 per unit.

Solution 6.15 The power flow across line 2-3 in Example 6.13 is given by

$$\begin{aligned} P_{23} &= V_2 V_3 Y_{23} \cos(\theta_2 - \theta_3 - \phi_{23}) - V_2^2 Y_{23} \cos \phi_{23} \\ &= 0.1211 \text{ per unit} \end{aligned} \quad (6.147)$$

If P_{23} exceeds 0.1 per unit, then the penalty function

$$p(h) = (P_{23} - 0.1)^2 \quad (6.148)$$

will be appended to the cost function. The partial derivatives remain the same, with the exception of $\left[\frac{\partial f}{\partial x} \right]$, which becomes:

$$\left[\frac{\partial f}{\partial x} \right] = \left[\frac{\partial f}{\partial P_1} \right] \left[\frac{\partial P_1}{\partial x} \right] + \left[\frac{\partial f}{\partial P_{23}} \right] \left[\frac{\partial P_{23}}{\partial x} \right] \quad (6.149)$$

$$\begin{aligned} &= (1 + 0.125 P_1) \begin{bmatrix} V_1 V_2 Y_{12} \sin(\theta_1 - \theta_2 - \phi_{1,2}) \\ V_1 V_3 Y_{13} \sin(\theta_1 - \theta_3 - \phi_{1,3}) \end{bmatrix} \\ &\quad + 2(P_{23} - 0.1) \begin{bmatrix} -V_2 V_3 Y_{23} \sin(\theta_2 - \theta_3 - \phi_{23}) \\ V_2 V_3 Y_{23} \sin(\theta_2 - \theta_3 - \phi_{23}) \end{bmatrix} \quad (6.150) \end{aligned}$$

Proceeding with the steepest gradient algorithm iterations yields the final constrained optimal generation scheduling:

$$\begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix} = \begin{bmatrix} 141.3 \\ 491.9 \\ 319.7 \end{bmatrix} \text{ MW}$$

and $P_{23} = 101.2$ MW. The cost for this constrained scheduling is \$7,781/MWhr, which is slightly greater than the nonconstrained cost. Also note that P_{23} is slightly greater than 100 MW since this was considered a soft limit. ■

In the case where hard limits must be imposed, an alternate approach to enforcing the inequality constraints must be employed. In this approach, the

inequality constraints are added as additional equality constraints with the inequality set equal to the limit (upper or lower) that is violated. This in essence introduces an additional set of Lagrangian multipliers. This is often referred to as the dual variable approach, because each inequality has the potential of resulting in two equalities: one for the upper limit and one for the lower limit. However, the upper and lower limits cannot be simultaneously violated; thus, out of the possible set of additional Lagrangian multipliers, only one of the two will be included at any given operating point and thus the dual limits are mutually exclusive.

Example 6.16

Repeat Example 6.15 using the dual variable approach.

Solution 6.16 Introducing the additional equation

$$P_{23} = V_2 V_3 Y_{23} \cos(\theta_2 - \theta_3 - \phi_{23}) - V_2^2 Y_{23} \cos \phi_{23} = 0.100 \text{ per unit} \quad (6.151)$$

to the equality constraints adds an additional equation to the set of $g(x)$. Therefore an additional unknown must be added to the state vector x to yield a solvable set of Equations (three equations in three unknowns). Either P_{G2} or P_{G3} can be chosen as the additional unknown. In this example, P_{G3} will be chosen. The new system Jacobian becomes

$$\left[\frac{\partial g}{\partial x} \right] = \begin{bmatrix} \frac{\partial g_1}{\partial x_1} & \frac{\partial g_1}{\partial x_2} & \frac{\partial g_1}{\partial x_3} \\ \frac{\partial g_2}{\partial x_1} & \frac{\partial g_2}{\partial x_2} & \frac{\partial g_2}{\partial x_3} \\ \frac{\partial g_3}{\partial x_1} & \frac{\partial g_3}{\partial x_2} & \frac{\partial g_3}{\partial x_3} \end{bmatrix} \quad (6.152)$$

where

$$\frac{\partial g_1}{\partial x_1} = V_2 (V_1 Y_{12} \sin(\theta_2 - \theta_1 - \phi_{21}) + V_3 Y_{13} \sin(\theta_2 - \theta_3 - \phi_{23}))$$

$$\frac{\partial g_1}{\partial x_2} = -V_2 V_3 Y_{32} \sin(\theta_2 - \theta_3 - \phi_{23})$$

$$\frac{\partial g_1}{\partial x_3} = 0$$

$$\frac{\partial g_2}{\partial x_1} = -V_3 V_2 Y_{23} \sin(\theta_3 - \theta_2 - \phi_{32})$$

$$\frac{\partial g_2}{\partial x_2} = V_3 V_1 Y_{13} \sin(\theta_3 - \theta_1 - \phi_{31}) + V_2 Y_{23} \sin(\theta_3 - \theta_2 - \phi_{32})$$

$$\frac{\partial g_2}{\partial x_3} = 1$$

$$\frac{\partial g_3}{\partial x_1} = -V_2 V_3 Y_{23} \sin(\theta_2 - \theta_3 - \phi_{23})$$

$$\frac{\partial g_3}{\partial x_2} = V_2 V_3 Y_{23} \sin(\theta_2 - \theta_3 - \phi_{23})$$

$$\frac{\partial g_3}{\partial x_3} = 0$$

and

$$\begin{bmatrix} \frac{\partial g}{\partial u} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}^T, \quad \begin{bmatrix} \frac{\partial f}{\partial u} \end{bmatrix} = [1 + 0.025P_{G2}]$$

Similar to Example 6.13, the chain rule is used to obtain $\left[\frac{\partial f}{\partial x} \right]$:

$$\begin{bmatrix} \frac{\partial f}{\partial x} \end{bmatrix} = \begin{bmatrix} \frac{\partial C}{\partial P_{G1}} \end{bmatrix} \begin{bmatrix} \frac{\partial P_{G1}}{\partial x} \end{bmatrix} + \begin{bmatrix} \frac{\partial C}{\partial P_{G3}} \end{bmatrix} \begin{bmatrix} \frac{\partial P_{G3}}{\partial x} \end{bmatrix} \quad (6.153)$$

$$= (1 + 0.125P_{G1}) \begin{bmatrix} V_1 V_2 Y_{12} \sin(\theta_1 - \theta_2 - \phi_{12}) \\ V_1 V_3 Y_{13} \sin(\theta_1 - \theta_3 - \phi_{13}) \\ 0 \end{bmatrix} + (1 + 0.050P_{G3}) \times$$

$$\begin{bmatrix} V_3 V_2 Y_{32} \sin(\theta_3 - \theta_2 - \phi_{32}) \\ -V_3 (V_1 Y_{13} \sin(\theta_3 - \theta_1 - \phi_{31}) + V_2 Y_{23} \sin(\theta_3 - \theta_2 - \phi_{32})) \\ 0 \end{bmatrix} \quad (6.154)$$

Substituting these partial derivatives into the expression for ∇C of Equation (6.122) yields

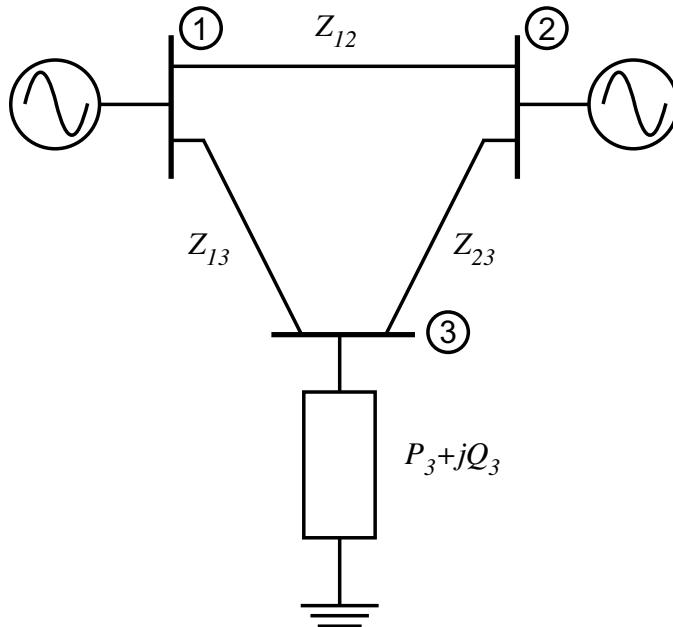
$$\begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix} = \begin{bmatrix} 141.3 \\ 490.3 \\ 321.2 \end{bmatrix} \text{ MW}$$

and $P_{23} = 100.0$ MW. ■

6.4.2 State Estimation

In power system state estimation, the estimated variables are the voltage magnitudes and the voltage phase angles at the system buses. The inputs to the state estimator are the active and reactive powers of the system, measured either at the injection sites or on the transmission lines. The state estimator is designed to give the best estimates of the voltages and phase angles, minimizing the effects of the measurement errors. Another consideration for the state estimator is to determine if a sufficient number of measurements are available to fully estimate the power system. This is the notion of system observability.

A set of specified measurements of a power system is said to be *observable* if the entire state vector of bus voltage magnitude and phase angles can be estimated from the set of available measurements. An unobservable system is one in which the set of measurements does not span the entire state space. The power system is observable if the matrix H_x in Equation (6.47) has rank n (full rank), where the number of measurements m is greater than or equal to the number of system states n . A *redundant* measurement is one whose addition to the measurement does not increase the rank of the matrix H_x .

**FIGURE 6.7**

Example power system

The observability of a power system can be determined by examining the measurement set and the structure of the power system. A *tree* is a set of measurements (either bus or line) that spans the entire set of power system buses. In other words, by graphically connecting the buses and lines that contribute to the set of measurements, the entire set of system buses can be connected by a single connected graph. A power system can be made observable by adding measurements at those lines which will connect disjoint trees.

Example 6.17

The SCADA system for the power network shown in Figure 6.7 reports the following measurements and variances:

z_i	state	measurement	variance (σ^2)
1	V_3	0.975	0.010
2	P_{13}	0.668	0.050
3	Q_{21}	-0.082	0.075
4	P_3	-1.181	0.050
5	Q_2	-0.086	0.075

This system has the same parameters as the three-bus system of Example 3.11. Note that P_3 and Q_3 are *injected* powers and are therefore negative. Estimate the power system states and use the chi-square test of inequality with $\alpha = 0.01$ to check for the presence of bad data in the measurements.

Solution 6.17 The first step in the estimation process is to identify and enumerate the unknown states. In this example, the unknowns are $[x_1 \ x_2 \ x_3]^T = [\theta_2 \ \theta_3 \ V_3]^T$. After the states are identified, the next step in the estimation process is to identify the appropriate functions $h(x)$ that correspond to each of the measurements. The nonlinear function that is being driven to zero to minimize the weighted error is

$$F(x) = H_x^T R^{-1} [z - h(x)] = 0 \quad (6.155)$$

where the set of $z - h(x)$ is given by

$$\begin{aligned} z_1 - h_1(x) &= V_3 - x_3 \\ z_2 - h_2(x) &= P_{13} - (V_1 x_3 Y_{13} \cos(-x_2 - \phi_{13}) - V_1^2 Y_{13} \cos \phi_{13}) \\ z_3 - h_3(x) &= Q_{21} - (V_2 V_1 Y_{21} \sin(x_1 - \phi_{21}) + V_2^2 Y_{21} \sin \phi_{21}) \\ z_4 - h_4(x) &= P_3 - (x_3 V_1 Y_{31} \cos(x_2 - \phi_{31}) + x_3 V_2 Y_{32} \cos(x_2 - x_1 - \phi_{32}) \\ &\quad + x_3^2 Y_{33} \cos \phi_{33}) \\ z_5 - h_5(x) &= Q_2 - (V_2 V_1 Y_{21} \sin(x_1 - \phi_{21}) - V_2^2 Y_{22} \sin \phi_{22} \\ &\quad + V_2 x_3 Y_{23} \sin(x_1 - x_2 - \phi_{23})) \end{aligned}$$

and the matrix of partial derivatives for the set of functions (6.155) is $H_x =:$

$$\left[\begin{array}{ccc} 0 & 0 & 1 \\ 0 & V_1 x_3 Y_{13} \sin(-x_2 - \phi_{13}) & V_1 Y_{13} \cos(-x_2 - \phi_{13}) \\ V_1 V_2 Y_{21} \cos(x_1 - \phi_{21}) & 0 & 0 \\ x_3 V_2 Y_{32} \sin(x_2 - x_1 - \phi_{32}) & -x_3 V_1 Y_{31} \sin(x_2 - \phi_{31}) & V_1 Y_{31} \cos(x_2 - \phi_{31}) \\ & -x_3 V_2 Y_{32} \sin(x_2 - x_1 - \phi_{32}) & +V_2 Y_{32} \cos(x_2 - x_1 - \phi_{32}) \\ & & +2x_3 Y_{33} \cos \phi_{33} \\ V_1 V_2 Y_{21} \cos(x_1 - \phi_{21}) & -V_2 x_3 Y_{23} \cos(x_1 - x_2 - \phi_{23}) & V_2 Y_{23} \sin(x_1 - x_2 - \phi_{23}) \\ +V_2 x_3 Y_{23} \cos(x_1 - x_2 - \phi_{23}) & & \end{array} \right] \quad (6.156)$$

This matrix has rank 3; therefore, this set of measurements spans the observable space of the power system.

The covariance matrix of the measurements is

$$R^{-1} = \begin{bmatrix} \frac{1}{0.010} & & & \\ & \frac{1}{0.050} & & \\ & & \frac{1}{0.075} & \\ & & & \frac{1}{0.050} \\ & & & & \frac{1}{0.075} \end{bmatrix} \quad (6.157)$$

The Newton–Raphson iteration to solve for the set of states x that minimize the weighted errors is

$$[H_x^T(x^k) R^{-1} H_x(x^k)] [x^{k-1} - x^k] = H_x^T(x^k) R^{-1} [z - h(x^k)] \quad (6.158)$$

Iteration 1

The initial condition for the state estimation solution is the same flat start as for the power flow equations; namely, all angles are set to zero and all unknown voltage magnitudes are set to unity. The measurement functions $h(x)$ evaluated at the initial conditions are

$$h(x^0) = \begin{bmatrix} 1.0000 \\ 0.0202 \\ -0.0664 \\ -0.0198 \\ -0.1914 \end{bmatrix}$$

The matrix of partials evaluated at the initial condition yields

$$H_x^0 = \begin{bmatrix} 0 & 0 & 1.0000 \\ 0 & -10.0990 & -1.0099 \\ -0.2257 & 0 & 0 \\ -9.9010 & 20.0000 & 1.9604 \\ -1.2158 & 0.9901 & -9.9010 \end{bmatrix}$$

The nonlinear functions (6.155) are

$$F(x^0) = \begin{bmatrix} 228.2791 \\ -593.9313 \\ -75.0229 \end{bmatrix}$$

The updated states are

$$\begin{bmatrix} \theta_2^1 \\ \theta_3^1 \\ V_3^1 \end{bmatrix} = \begin{bmatrix} -0.0119 \\ -0.0624 \\ 0.9839 \end{bmatrix}$$

where θ_2 and θ_3 are in radians. The error at iteration 1 is

$$\varepsilon^1 = 593.9313$$

Iteration 2

The updated values are used to recalculate the Newton–Raphson iterations:

$$h(x^1) = \begin{bmatrix} 0.9839 \\ 0.6582 \\ -0.0634 \\ -1.1593 \\ -0.0658 \end{bmatrix}$$

The matrix of partials is

$$H_x^1 = \begin{bmatrix} 0 & 0 & 1.0000 \\ 0 & -9.9791 & -0.3780 \\ -0.2659 & 0 & 0 \\ -9.6800 & 19.5351 & 0.7700 \\ -0.7470 & 0.4811 & -9.9384 \end{bmatrix}$$

The nonlinear function evaluated at the updated values yields

$$F(x^1) = \begin{bmatrix} 4.4629 \\ -10.5560 \\ 1.3749 \end{bmatrix}$$

The updated states are

$$\begin{bmatrix} \theta_2^2 \\ \theta_3^2 \\ V_3^2 \end{bmatrix} = \begin{bmatrix} -0.0113 \\ -0.0633 \\ 0.9851 \end{bmatrix}$$

The error at iteration 2 is

$$\varepsilon^2 = 10.5560$$

The iterations are obviously converging. At convergence, the states that minimize the weighted measurement errors are

$$x = \begin{bmatrix} -0.0113 \\ -0.0633 \\ 0.9851 \end{bmatrix}$$

To check for the presence of bad data, the weighted sum of squares of the measurement errors is compared to the chi-square distribution for $k = 2$ and $\alpha = 0.01$. The weighted sum of squares is

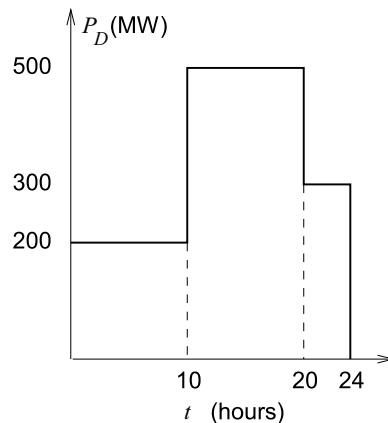
$$\begin{aligned} f &= \sum_{i=1}^5 \frac{1}{\sigma_i^2} (z_i - h_i(x))^2 \\ &= \frac{(-0.0101)^2}{0.010} + \frac{(0.0012)^2}{0.050} + \frac{(-0.0184)^2}{0.075} + \frac{(0.0007)^2}{0.050} + \frac{(-0.0076)^2}{0.075} \\ &= 0.0155 \end{aligned}$$

This value is less than $\chi_{2,0.01} = 9.21$; therefore, the data set is good and does not contain any spurious measurements. ■

6.5 Problems

1. The fuel costs for a three-unit plant are given by

$$\begin{aligned} F_1 &: 173.61 + 8.670P_1 + 0.00230P_1^2 \quad \text{\$/MWhr} \\ F_2 &: 180.68 + 9.039P_2 + 0.00238P_2^2 \quad \text{\$/MWhr} \\ F_3 &: 182.62 + 9.190P_3 + 0.00235P_3^2 \quad \text{\$/MWhr} \end{aligned}$$

**FIGURE 6.8**

Load curve for Problem 1

The daily load curve for the plant is given in Figure 6.8. Obtain and sketch the optimal power generated by each unit and the plant's incremental cost of power delivered (λ).

2. Use the method of least squares to find the “best fit” coefficients c_0 and c_1 in the function

$$f(x) = c_0 + c_1 x$$

for the following measured data:

x	$f(x)$
1	-2.1
3	-0.9
4	-0.6
6	0.6
7	0.9

3. Use the method of least squares to find the “best fit” coefficients a_0 , a_1 , and a_2 in the function

$$f(t) = a_0 + a_1 \sin \frac{2\pi t}{12} + a_2 \cos \frac{2\pi t}{12}$$

for the following measured data:

t	$f(t)$
0	1.0
2	1.6
4	1.4
6	0.6
8	0.2
10	0.8

This function describes the movement of the tide with a 12-hour period.

4. Minimize $-7x_1 - 3x_2 + x_3$ subject to the constraints

$$\begin{aligned} x_1 + x_2 + x_3 &\leq 15 \\ 2x_1 - 3x_2 + x_3 &\leq 10 \\ x_1 - 5x_2 - x_3 &\leq 0 \\ x_1, x_2, x_3 &\geq 0 \end{aligned}$$

- (a) using the simplex method
- (b) using the primal affine method with $\alpha = 0.9$.

Both methods should be augmented with slack variables x_4, x_5, x_6 and initiated with the feasible vector $x^0 = [1 \ 1 \ 1 \ 12 \ 10 \ 5]^T$.

5. In bulk transmission systems operation, it is useful to determine the total simultaneous interchange capacity (SIC) available via interties. This SIC may be limited by the transmission system or the available generation. The SIC may be found using the simplex method.

Missouri State Transmission (MST) has interchanges with three other utilities. The interchanges have the following generation limits:

$$\begin{aligned} \text{Interchange A: } x_A &\leq 1800\text{MW} \\ \text{Interchange B: } x_B &\leq 1500\text{MW} \\ \text{Interchange C: } x_C &\leq 2000\text{MW} \end{aligned}$$

Each interchange consists of several intertie lines which may limit the SIC. They can be described by their distribution factors:

$$\begin{bmatrix} 2.50 & 1.00 & 0.50 \\ 1.00 & 1.50 & 1.00 \\ 2.00 & 1.25 & 1.00 \end{bmatrix} \begin{bmatrix} x_A \\ x_B \\ x_C \end{bmatrix} = \Delta P$$

where the maximum incremental load ΔP is 5000 MW.

The simplex method to solve this SIC problem is to maximize $c(x) = x_A + x_B + x_C$, subject to

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 2.50 & 1.00 & 0.50 \\ 1.00 & 1.50 & 1.00 \\ 2.00 & 1.25 & 1.00 \end{bmatrix} \begin{bmatrix} x_A \\ x_B \\ x_C \end{bmatrix} \leq \begin{bmatrix} 1800 \\ 1500 \\ 2000 \\ 5000 \\ 5000 \\ 5000 \end{bmatrix}$$

Find the optimal SIC. Is the system generation or transmission limited? Support your answer.

6. Using the steepest descent method, minimize $x_1^2 + x_2^2$ subject to the constraint

$$x_1^2 + 2x_1x_2 + 3x_2^2 - 1 = 0$$

Start with $x_1=0.5$ and $x_2 = 0.5$.

7. Repeat Problem 6 using the SQP method. Start with $\lambda = 0$.

8. Find the minimum of

$$C : x_1^2 + x_2^2 + u_1x_1 + u_2x_2 + 1$$

(a) subject to

$$\begin{aligned} x_1 \cos(x_2) + x_2^2 - u_1 \cos(x_1) &= 1 \\ x_1 - x_2 + 3u_2 &= -3 \end{aligned}$$

(b) subject to

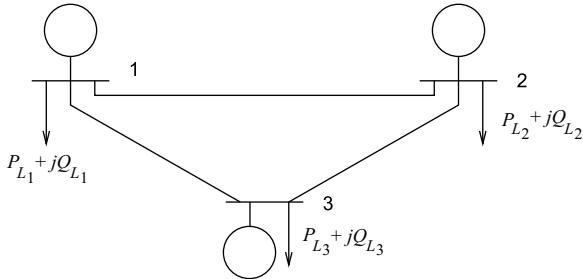
$$\begin{aligned} x_1 \cos(x_2) + x_2^2 - u_1 \cos(x_1) &= 1 \\ x_1 - x_2 + 3u_2 &= -3 \\ u_2 &\geq -0.8 \end{aligned}$$

(c) subject to

$$\begin{aligned} x_1 \cos(x_2) + x_2^2 - u_1 \cos(x_1) &= 1 \\ x_1 - x_2 + 3u_2 &= -3 \\ x_2 &\leq 0.30 \end{aligned}$$

using the penalty function $f(x_2) = ae^{b(x_2 - c)}$, where a and b are positive constants and c is the function offset.

Use an initial guess of $x^0 = [0 \ 0]'$ and $u^0 = [0 \ 0]'$ and $\gamma = 0.05$. You might want to experiment with other values of γ as well. Your stopping criterion should be $\|\nabla C\| \leq 0.01$.

**FIGURE 6.9**

Three-bus system

9. Consider the system shown in Figure 6.9. The bus and line data are given below:

Line	R	X	B
1-2	0.01	0.1	0.050
1-3	0.05	0.1	0.025
2-3	0.05	0.1	0.025
Bus	V	P _L	Q _L
1	1.00	0.35	0.10
2	1.02	0.40	0.25
3	1.02	0.25	0.10

The fuel costs for the generators are

$$F_1 : P_{G_1} + 1.5P_{G_1}^2$$

$$F_2 : 2P_{G_2} + P_{G_2}^2$$

$$F_3 : 2.5P_{G_3} + 0.5P_{G_3}^2$$

- (a) Using the equal incremental cost criterion, find the optimal scheduling for the units (remember that this method neglects system losses).
- (b) Using your answer for part (a) as the initial control vector, use the steepest descent method to find the optimal scheduling for this system, which considers system losses.
- (c) Now assume the following limits are imposed:

$$F_1 : P_{G_1} + 1.5P_{G_1}^2 \quad 0 \leq P_{G_1} \leq 0.6$$

$$F_2 : 2P_{G_2} + P_{G_2}^2 \quad 0 \leq P_{G_2} \leq 0.4$$

$$F_3 : 2.5P_{G_3} + 0.5P_{G_3}^2 \quad 0 \leq P_{G_3} \leq 0.1$$

Repeat part (b).

- (d) Interpret your results relating the generator settings to the cost functions.
10. For the system shown in Figure 6.9, the following measurements were obtained:

V_2	1.04
V_3	0.98
P_{G1}	0.58
P_{G2}	0.30
P_{G3}	0.14
P_{12}	0.12
P_{32}	-0.04
P_{13}	0.10

where $\sigma_V^2 = (0.01)^2$, $\sigma_{P_G}^2 = (0.015)^2$, and $\sigma_{P_{ij}}^2 = (0.02)^2$.

Estimate the system states, the error, and test for bad data using the chi-square test with $\alpha = 0.01$.

Eigenvalue Problems

Small signal stability is the ability of a system to maintain stability when subjected to small disturbances. Small signal analysis provides valuable information about the inherent dynamic characteristics of the system and assists in its design, operation, and control. Time domain simulation and eigenanalysis are the two main approaches to study system stability.

Eigenanalysis methods are widely used to perform small signal stability studies. The dynamic behavior of a system in response to small perturbations can be determined by computing the eigenvalues and eigenvectors of the system matrix. The locations of the eigenvalues can be used to investigate the system's performance. In addition, eigenvectors can be used to estimate the relative participation of the respective states in the corresponding disturbance modes.

A scalar λ is an eigenvalue of an $n \times n$ matrix A if there exists a nonzero $n \times 1$ vector v such that

$$Av = \lambda v \quad (7.1)$$

where v is the corresponding right eigenvector. If there exists a nonzero vector w such that

$$w^T A = \lambda w^T \quad (7.2)$$

then w is a left eigenvector. The set of all eigenvalues of A is called the spectrum of A . Normally the term "eigenvector" refers to the right eigenvector unless denoted otherwise. The eigenvalue problem in Equation (7.1) is called the standard eigenvalue problem. Equation (7.1) can be written as

$$(A - \lambda I) v = 0 \quad (7.3)$$

and thus is a homogeneous system of equations for x . This system has a nontrivial solution only if the determinant

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

The determinant equation is also called the characteristic equation for A and is an n th degree polynomial in λ . The eigenvalues of an $n \times n$ matrix A are the roots of the characteristic equation

$$\lambda^n + c_{n-1}\lambda^{n-1} + c_{n-2}\lambda^{n-2} + \dots + c_0 = 0 \quad (7.4)$$

Therefore, there are n roots (possibly real or complex) of the characteristic equation. Each one of these roots is also an eigenvalue of A .

7.1 The Power Method

The power method is one of the most common methods of finding the *dominant* eigenvalue of the $n \times n$ matrix A . The dominant eigenvalue is the largest eigenvalue in absolute value. Therefore, if $\lambda_1, \lambda_2, \dots, \lambda_n$ are eigenvalues of A , then λ_1 is the dominant eigenvalue of A if

$$|\lambda_1| > |\lambda_i| \quad (7.5)$$

for all $i = 2, \dots, n$.

The power method is actually an approach to finding the eigenvector v_1 corresponding to the dominant eigenvalue of the matrix A . Once the eigenvector is obtained, the eigenvalue can be extracted from the *Rayleigh quotient*:

$$\lambda = \frac{\langle Av, v \rangle}{\langle v, v \rangle} \quad (7.6)$$

The approach to finding the eigenvector v_1 is an iterative approach. Therefore, from an initial guess vector v^0 , a sequence of approximations v^k is constructed which hopefully converges as k goes to ∞ . The iterative algorithm for the power method is straightforward:

The Power Method

1. Let $k = 0$ and choose v^0 to be a nonzero $n \times 1$ vector.
2. $w^{k+1} = Av^k$
3. $\alpha_{k+1} = \|w^{k+1}\|$
4. $v^{k+1} = \frac{w^{k+1}}{\alpha^{k+1}}$
5. If $\|v^{k+1} - v^k\| < \varepsilon$, then done. Else, $k = k + 1$, go to Step 2.

The division by the norm of the vector in Step 4 is not a necessary step, but it keeps the size of the values of the eigenvector close to 1. Recall that a scalar times an eigenvector of A is still an eigenvector of A ; therefore, scaling has no adverse consequence. However, without Step 4 and $\alpha^k = 1$ for all k , the values of the updated vector may increase or decrease to the extent that the computer accuracy is affected.

Example 7.1

Use the power method to find the eigenvector corresponding to the dominant eigenvalue of the following matrix:

$$A = \begin{bmatrix} 6 & -2 \\ -8 & 3 \end{bmatrix}$$

Solution 7.1 Start with the initial guess $v^0 = [1 \ 1]^T$. Then

$$\begin{aligned} w^1 &= Av^0 = \begin{bmatrix} 4 \\ -5 \end{bmatrix} \\ \alpha^1 &= \|w^1\| = 6.4031 \\ v^1 &= \begin{bmatrix} 0.6247 \\ -0.7809 \end{bmatrix} \end{aligned}$$

The second iteration follows:

$$\begin{aligned} w^2 &= Av^1 = \begin{bmatrix} 5.3099 \\ -7.3402 \end{bmatrix} \\ \alpha^2 &= \|w^2\| = 9.0594 \\ v^2 &= \begin{bmatrix} 0.5861 \\ -0.8102 \end{bmatrix} \end{aligned}$$

Continuing to convergence yields the eigenvector:

$$v^* = \begin{bmatrix} 0.5851 \\ -0.8110 \end{bmatrix}$$

From the eigenvector, the corresponding eigenvalue is calculated

$$\lambda = \frac{\begin{bmatrix} 0.5851 & -0.8110 \end{bmatrix} \begin{bmatrix} 6 & -2 \\ -8 & 3 \end{bmatrix} \begin{bmatrix} 0.5851 \\ -0.8110 \end{bmatrix}}{\begin{bmatrix} 0.5851 & -0.8110 \end{bmatrix} \begin{bmatrix} 0.5851 \\ -0.8110 \end{bmatrix}} = \frac{8.7720}{1} = 8.7720$$

which is the largest eigenvalue of A (the smaller eigenvalue is 0.2280). ■

To see why the power method converges to the dominant eigenvector, let the initial guess vector v^0 be expressed as the linear combination

$$v^0 = \sum_{i=1}^n \beta_i v_i \tag{7.7}$$

where v_i are the actual eigenvectors of A and the coefficients β_i are chosen to make Equation (7.7) hold. Then, applying the power method (without loss of generality it can be assumed that $\alpha^k = 1$ for all k) yields

$$\begin{aligned} v^{k+1} &= Av^k = A^2 v^{k-1} = \dots = A^{k+1} v^0 \\ &= \sum_{i=1}^n \lambda_i^{k+1} \beta_i v_i = \lambda_1^{k+1} \left(\beta_1 v_1 + \sum_{i=2}^n \left(\frac{\lambda_i}{\lambda_1} \right)^{k+1} \beta_i v_i \right) \end{aligned}$$

Because

$$\left| \frac{\lambda_i}{\lambda_1} \right| < 1 \quad i = 2, \dots, n$$

then, as k successively increases, these terms go to zero, and only the component corresponding to v_1 is left.

The power method will fail if there are two largest eigenvalues of the same absolute magnitude. Recall that the eigenvalues of a real matrix A are in general complex and occur in conjugate pairs (which necessarily have the same absolute value). Therefore, if the largest eigenvalue of A is not real, then the power method will certainly fail to converge. For this reason, it is sensible to apply the power method only to matrices whose eigenvalues are known to be real. One class of real eigenvalue matrices is symmetric matrices.

There are also cases in which the power method may not converge to the eigenvector corresponding to the dominant eigenvalue. This would occur in the case in which $\beta_1 = 0$. This implies that the initial guess v^0 contains no component of the eigenvector v_1 . In this case, the method will converge to the eigenvector contained in the decomposition of v^0 of the next largest eigenvalue.

The rate of convergence of the power method is determined by the ratio $|\frac{\lambda_2}{\lambda_1}|$. Thus, if $|\lambda_2|$ is only slightly smaller than $|\lambda_1|$, then the power method will converge slowly and a large number of iterations will be required to meet the required accuracy.

There are several extensions to the power method. For example, if, instead of the dominant eigenvalue, the smallest eigenvalue was desired, then the power method can be applied to A^{-1} . Since the eigenvalues of A^{-1} are $\frac{1}{\lambda_n}, \dots, \frac{1}{\lambda_1}$, the inverse power method should converge to $\frac{1}{\lambda_n}$.

Another extension is the spectral shift. This approach uses the fact that the eigenvalues of $A - aI$ are $\lambda_1 - a, \dots, \lambda_n - a$. Thus, having computed the first eigenvalue λ_1 , the power method can be reapplied using the shifted matrix $A - \lambda_1 I$. This reduces the first eigenvalue to zero, and the power method now converges to the largest in absolute value of $\lambda_2 - \lambda_1, \dots, \lambda_n - \lambda_1$.

7.2 The QR Algorithm

Many methods for solving the eigenvalue problem are based on a sequence of similarity transformations with orthogonal matrices. Thus, if P is any non-singular matrix, then the matrices A and PAP^{-1} have the same eigenvalues. Furthermore, if v is an eigenvector of A , then Pv is an eigenvector of PAP^{-1} . If the matrix P is orthogonal, then the condition of the eigenproblem is not affected. This is the basis for the similarity transformation methods.

The QR method [21], [57], [58] is one of the most widely used decomposition methods for calculating eigenvalues of matrices. It uses a sequence of orthogonal similarity transformations [14] [31] such that $A = A_0, A_1, A_2, \dots$ is computed by

$$A_i = Q_i R_i, \quad R_i Q_i = A_{i+1}, \quad i = 0, 1, 2, \dots,$$

Similar to the LU factorization, the matrix A can also be factored into two matrices such that

$$A = QR \tag{7.8}$$

where Q is a unitary matrix and R is an upper triangular matrix. The matrix Q is *unitary* if

$$QQ^* = Q^*Q = I \tag{7.9}$$

where $(*)$ denotes a complex conjugate transpose.

Examples of unitary matrices are

$$Q_1 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad Q_2 = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

It also follows that the inverse of a unitary matrix is also its conjugate transpose, i.e.,

$$Q^{-1} = Q^*$$

This decomposition yields the column vectors $[a_1, a_2, \dots, a_n]$ of A and column vectors $[q_1, q_2, \dots, q_n]$ of Q such that

$$a_k = \sum_{i=1}^k r_{ik} q_i, \quad k = 1, \dots, n \tag{7.10}$$

The column vectors a_1, a_2, \dots, a_n must be orthonormalized from left to right into an orthonormal basis q_1, q_2, \dots, q_n .

In the implementation of the QR algorithm, it is common practice to transform A into a Hessenberg matrix H having the same eigenvalues and then apply the QR matrix to H . In the end, the matrix becomes upper triangular and the eigenvalues can be read off the diagonal. A Hessenberg matrix is essentially an upper triangular matrix with one extra set of nonzero elements directly below the diagonal. The reason for reducing A to a Hessenberg matrix is that this greatly reduces the total number of operations required for

the QR algorithm. A Hessenberg matrix has the form

$$H = \begin{bmatrix} * & * & * & * & \cdots & * & * & * & * \\ * & * & * & * & \cdots & * & * & * & * \\ 0 & * & * & * & \cdots & * & * & * & * \\ 0 & 0 & * & * & \cdots & * & * & * & * \\ \vdots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \ddots & * & * & * & * \\ 0 & 0 & 0 & 0 & \cdots & * & * & * & * \\ 0 & 0 & 0 & 0 & \cdots & 0 & * & * & * \\ 0 & 0 & 0 & 0 & \cdots & 0 & 0 & * & * \end{bmatrix}$$

where * indicates a nonzero entry.

The Householder method is one method used to reduce A to a Hessenberg matrix. For each $n \times n$ matrix A , there exist $n - 2$ Householder matrices H_1, H_2, \dots, H_{n-2} , such that for

$$Q = H_{n-2} \dots H_2 H_1$$

the matrix

$$P = Q^* A Q$$

is a Hessenberg matrix [29]. A matrix H is a Householder matrix if

$$H = I - 2 \frac{vv^*}{v^*v}$$

Note that Householder matrices are also unitary matrices. The vector v is chosen to satisfy

$$v_i = a_i \pm e_i \|a_i\|_2 \quad (7.11)$$

where the choice of sign is based upon the requirement that $\|v\|_2$ should not be too small, e_i is the i th column of I , and a_i is the i th column of A .

Example 7.2

Find the QR decomposition of the matrix A

$$A = \begin{bmatrix} 1 & 3 & 4 & 8 \\ 2 & 1 & 2 & 3 \\ 4 & 3 & 5 & 8 \\ 9 & 2 & 7 & 4 \end{bmatrix}$$

Solution 7.2 The first transformation will be applied to zero out the first column of A below the subdiagonal; thus

$$v_1 = a_1 + e_1 \|a_1\|_2$$

$$\begin{aligned}
 &= \begin{bmatrix} 1 \\ 2 \\ 4 \\ 9 \end{bmatrix} + 10.0995 \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \\
 &= \begin{bmatrix} 11.0995 \\ 2.0000 \\ 4.0000 \\ 9.0000 \end{bmatrix}
 \end{aligned}$$

leading to

$$\begin{aligned}
 H_1 &= I - 2 \frac{v_1 v_1^*}{(v_1^* v_1)} \\
 &= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} - \frac{2}{224.1990} \begin{bmatrix} 11.0995 \\ 2.0000 \\ 4.0000 \\ 9.0000 \end{bmatrix} \begin{bmatrix} 11.0995 & 2.0000 & 4.0000 & 9.0000 \end{bmatrix} \\
 &= \begin{bmatrix} -0.0990 & -0.1980 & -0.3961 & -0.8911 \\ -0.1980 & 0.9643 & -0.0714 & -0.1606 \\ -0.3961 & -0.0714 & 0.8573 & -0.3211 \\ -0.8911 & -0.1606 & -0.3211 & 0.2774 \end{bmatrix}
 \end{aligned}$$

and

$$H_1 A = \begin{bmatrix} -10.0995 & -3.4655 & -9.0103 & -8.1192 \\ 0 & -0.1650 & -0.3443 & 0.0955 \\ 0 & 0.6700 & 0.3114 & 2.1910 \\ 0 & -3.2425 & -3.5494 & -9.0702 \end{bmatrix}$$

The second iteration will operate on the part of the transformed matrix that excludes the first column and row. Therefore,

$$\begin{aligned}
 v_2 &= a_2 + e_2 \|a_2\|_2 \\
 &= \begin{bmatrix} -0.1650 \\ 0.6700 \\ -3.2425 \end{bmatrix} + 3.3151 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\
 &= \begin{bmatrix} 3.1501 \\ 0.6700 \\ -3.2425 \end{bmatrix}
 \end{aligned}$$

which results in

$$\begin{aligned}
 H_2 &= I - 2 \frac{v_2 v_2^*}{(v_2^* v_2)} \\
 &= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0.0498 & -0.2021 & 0.9781 \\ 0 & -0.2021 & 0.9570 & 0.2080 \\ 0 & 0.9781 & 0.2080 & -0.0068 \end{bmatrix}
 \end{aligned}$$

and

$$H_2 H_1 A = \begin{bmatrix} -10.0995 & -3.4655 & -9.0103 & -8.1192 \\ 0 & -3.3151 & -3.5517 & -9.3096 \\ 0 & 0 & -0.3708 & 0.1907 \\ 0 & 0 & -0.2479 & 0.6108 \end{bmatrix}$$

Continuing the process yields

$$v_3 = \begin{bmatrix} 0.0752 \\ -0.2479 \end{bmatrix}$$

$$H_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0.8313 \\ 0 & 0 & 0.5558 \end{bmatrix}$$

which results in

$$R = H_3 H_2 H_1 A = \begin{bmatrix} -10.0995 & -3.4655 & -9.0103 & -8.1192 \\ 0 & -3.3151 & -3.5517 & -9.3096 \\ 0 & 0 & -0.4460 & 0.4980 \\ 0 & 0 & 0 & -0.4018 \end{bmatrix}$$

and

$$Q = H_1 H_2 H_3 = \begin{bmatrix} -0.0990 & -0.8014 & -0.5860 & -0.0670 \\ -0.1980 & -0.0946 & 0.2700 & -0.9375 \\ -0.3961 & -0.4909 & 0.7000 & 0.3348 \\ -0.8911 & 0.3283 & -0.3060 & 0.0670 \end{bmatrix}$$

It can be verified that $A = QR$ and further that $Q^* = Q^{-1}$. ■

Elimination by QR decomposition can be considered as an alternative to Gaussian elimination. However, the number of multiplications and divisions required is more than twice the number required for Gaussian elimination. Therefore, QR decomposition is seldom used for the solution of linear systems, but it does play an important role in the calculation of eigenvalues.

Although the eigenvalue problem gives rise to a simple set of algebraic equations to determine the solution to

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

the practical problem of solving this equation is difficult. Computing the roots of the characteristic equation or the nullspace of a matrix is a process that is not well suited for computers. In fact, no generalized direct process exists for solving the eigenvalue problem in a finite number of steps. Therefore, iterative methods for calculation must be relied upon to produce a series of successively improved approximations to the eigenvalues of a matrix.

The QR method is commonly used to calculate the eigenvalues and eigenvectors of full matrices. As developed by Francis [14], the QR method produces

a series of similarity transformations

$$A_k = Q_k^* A_{k-1} Q_k \quad Q_k^* Q_k = I \quad (7.12)$$

where the matrix A_k is similar to A . The QR decomposition is repeatedly performed and applied to A as the subdiagonal elements are iteratively driven to zero. At convergence, the eigenvalues of A in descending order by magnitude appear on the diagonal of A_k .

The next step is to find the eigenvectors associated with each eigenvalue. Recall that

$$Av_i = \lambda_i v_i \quad (7.13)$$

for each eigenvalue and corresponding eigenvector $i = 1, \dots, n$. Equation (7.13) may also be written as

$$Av_i - \lambda_i v_i = 0$$

In other words, the matrix defined by $A - \lambda_i I$ is singular; thus only three of its rows (or columns) are independent. This fact can be used to determine the eigenvectors once the eigenvalues are known. Since $A - \lambda_i I$ is not of full rank, one of the elements of the eigenvector v_i can be chosen arbitrarily. To start, partition $A - \lambda_i I$ as

$$A - \lambda_i I = \begin{bmatrix} a_{11} & a_{1,2n} \\ a_{2n,1} & a_{2n,2n} \end{bmatrix}$$

where a_{11} is a scalar, $a_{1,2n}$ is a $1 \times (n-1)$ vector, $a_{2n,1}$ is a $(n-1) \times 1$ vector, and $a_{2n,2n}$ is an $(n-1) \times (n-1)$ matrix of rank $(n-1)$. Then let $v_i(1) = 1$ and solve for the remaining portion of the eigenvector as

$$\begin{bmatrix} v_i(2) \\ v_i(3) \\ \vdots \\ v_i(n) \end{bmatrix} = -a_{2n,2n}^{-1} a_{2n,1} v_i(1) \quad (7.14)$$

Now update $v_i(1)$ from

$$v_i(1) = -\frac{1}{a_{11}} a_{1,2n} * \begin{bmatrix} v_i(2) \\ v_i(3) \\ \vdots \\ v_i(n) \end{bmatrix}$$

Then the eigenvector corresponding to λ_i is

$$v_i = \begin{bmatrix} v_i(1) \\ v_i(2) \\ \vdots \\ v_i(n) \end{bmatrix}$$

The last step is to normalize the eigenvector; therefore,

$$v_i = \frac{v_i}{\|v_i\|}$$

Example 7.3

Find the eigenvalues and eigenvectors of the matrix of Example 7.2.

Solution 7.3 The first objective is to find the eigenvalues of the matrix A using the QR method. From Example 7.2, the first QR factorization yields the Q_0 matrix

$$Q_0 = \begin{bmatrix} -0.0990 & -0.8014 & -0.5860 & -0.0670 \\ -0.1980 & -0.0946 & 0.2700 & -0.9375 \\ -0.3961 & -0.4909 & 0.7000 & 0.3348 \\ -0.8911 & 0.3283 & -0.3060 & 0.0670 \end{bmatrix}$$

Using the given A matrix as A_0 , the first update A_1 is found by

$$\begin{aligned} A_1 &= Q_0^* A_0 Q_0 \\ &= \begin{bmatrix} 12.4902 & 10.1801 & 1.1599 & 0.3647 \\ 10.3593 & -0.9987 & -0.5326 & 1.2954 \\ -0.2672 & 0.3824 & -0.4646 & -0.1160 \\ 0.3580 & -0.1319 & 0.1230 & -0.0269 \end{bmatrix} \end{aligned}$$

The QR factorization of A_1 yields

$$Q_1 = \begin{bmatrix} -0.7694 & -0.6379 & -0.0324 & -0.0006 \\ -0.6382 & 0.7660 & 0.0733 & -0.0252 \\ 0.0165 & -0.0687 & 0.9570 & 0.2812 \\ -0.0221 & 0.0398 & -0.2786 & 0.9593 \end{bmatrix}$$

and the A_2 matrix becomes

$$\begin{aligned} A_2 &= Q_1^* A_1 Q_1 \\ &= \begin{bmatrix} 17.0913 & 4.8455 & -0.2315 & -1.0310 \\ 4.6173 & -5.4778 & -1.8116 & 0.6064 \\ -0.0087 & 0.0373 & -0.5260 & -0.1757 \\ 0.0020 & -0.0036 & 0.0254 & -0.0875 \end{bmatrix} \end{aligned}$$

Note that the elements below the diagonals are slowly decreasing to zero. This process is carried out until the final A matrix is obtained:

$$A_* = \begin{bmatrix} 18.0425 & 0.2133 & -0.5180 & -0.9293 \\ 0 & -6.4172 & -1.8164 & 0.6903 \\ 0 & 0 & -0.5269 & -0.1972 \\ 0 & 0 & 0 & -0.0983 \end{bmatrix} \quad (7.15)$$

The eigenvalues are on the diagonals of A_* and are in decreasing order by magnitude. Thus the eigenvalues are

$$\lambda_{1,\dots,4} = \begin{bmatrix} 18.0425 \\ -6.4172 \\ -0.5269 \\ -0.0983 \end{bmatrix}$$

The corresponding eigenvectors are:

$$\begin{bmatrix} 0.4698 \\ 0.2329 \\ 0.5800 \\ 0.6234 \end{bmatrix}, \begin{bmatrix} 0.6158 \\ 0.0539 \\ 0.2837 \\ -0.7330 \end{bmatrix}, \begin{bmatrix} 0.3673 \\ -0.5644 \\ -0.5949 \\ 0.4390 \end{bmatrix}, \begin{bmatrix} 0.0932 \\ 0.9344 \\ -0.2463 \\ -0.2400 \end{bmatrix}$$

■

7.2.1 Deflation

Once the elements of the last row are reduced to zero, the last row and column of the matrix may be neglected. This implies that the smallest eigenvalue is “deflated” by removing the last row and column. The procedure can then be repeated on the remaining $(n-1) \times (n-1)$ matrix.

7.2.2 Shifted QR

The speed of convergence of the QR method for calculating eigenvalues depends greatly on the location of the eigenvalues with respect to one another. The matrix $A - \sigma I$ has the eigenvalues $\lambda_i - \sigma$ for $i = 1, \dots, n$. If σ is chosen as an approximate value of the smallest eigenvalue λ_n , then $\lambda_n - \sigma$ becomes small. This will speed up the convergence in the last row of the matrix, since

$$\frac{|\lambda_n - \sigma|}{|\lambda_{n-1} - \sigma|} \ll 1$$

The QR iterations can converge very slowly in many instances. However, if some information about one or more of the eigenvalues is known a priori, then a variety of techniques can be applied to speed up convergence of the iterations. One such technique is the *shifted* QR method, in which a shift σ is introduced at each iteration such that the QR factorization at the k th iteration is performed on

$$A_k - \sigma I = Q_k R_k \quad (7.16)$$

and

$$A_{k+1} = Q_k^* (A_k - \sigma I) Q_k + \sigma I \quad (7.17)$$

If σ is a good estimate of an eigenvalue, then the $(n, n - 1)$ entry of A_k will converge rapidly to zero, and the (n, n) entry of A_{k+1} will converge to the eigenvalue closest to σ_k . The n th row and column can be removed (deflated), and an alternate shift can be applied.

Using a shift and deflation in combination can significantly improve convergence. Additionally, if only one eigenvalue is desired of a particular magnitude, this eigenvalue can be isolated via the shift method. After the last row has been driven to zero, the eigenvalue can be obtained and the remainder of the QR iterations abandoned.

Example 7.4

Find the largest eigenvalue of Example 7.3 using shifts and deflation.

Solution 7.4 Start with using a shift of $\sigma = 15$. This is near the 18.0425 eigenvalue, so convergence to that particular eigenvalue should be rapid. Starting with the original A matrix as A_0 , the QR factorization of $A_0 - \sigma I$ yields

$$Q_0 = \begin{bmatrix} 0.8124 & -0.0764 & -0.2230 & -0.5334 \\ -0.1161 & 0.9417 & 0.0098 & -0.3158 \\ -0.2321 & -0.2427 & 0.7122 & -0.6164 \\ -0.5222 & -0.2203 & -0.6655 & -0.4856 \end{bmatrix}$$

and the update $A_1 = Q_0^*(A_0 - \sigma I)Q_0 + \sigma I$

$$A_1 = \begin{bmatrix} -4.9024 & 0.8831 & -1.6174 & 2.5476 \\ -0.2869 & 0.0780 & -0.1823 & 1.7775 \\ -2.9457 & 0.5894 & -1.5086 & 2.3300 \\ 2.5090 & 1.0584 & 3.1975 & 17.3330 \end{bmatrix}$$

The eigenvalue of interest ($\lambda = 18.0425$) will now appear in the lower right corner since, as the iterations progress, $A_{k+1}(n, n) - \sigma$ will be the smallest diagonal in magnitude. Recall that the eigenvalues are ordered on the diagonal from largest to smallest, and, since the largest eigenvalue is “shifted” by σ , it will now have the smallest magnitude. The convergence can be further increased by updating σ at each iteration, such that $\sigma_{k+1} = A_{k+1}(n, n)$. The iterations proceed as in Example 7.3. ■

7.2.3 Double Shifted QR

One problem that arises is if the eigenvalues sought are complex. As eigenvalues come in conjugate pairs, the QR may appear to fail since no dominant eigenvalue exists. However, instead of generating a single eigenvalue estimation, QR produces a 2×2 matrix “containing” the conjugate pair. By using deflation, the 2×2 matrix can be extracted and the eigenvalues easily computed. The remaining eigenvalues can then be computed as before.

Another approach is that, since complex eigenvalues always appear in conjugate pairs, it is natural to search for the pair simultaneously. This can be accomplished by collapsing two shifted QR steps in one double step with the two shifts being complex conjugates of each other.

Let σ_1 and σ_2 be two eigenvalues of a real matrix, where $\sigma_1^* = \sigma_2$. Performing two QR steps using Equations (7.16) and (7.17) yields

$$A_0 = Q_0 R_0 + \sigma_1 I \quad (7.18)$$

$$A_1 = Q_0^* (A_0 - \sigma_1 I) Q_0 + \sigma_1 I \quad (7.19)$$

$$= R_0 Q_0 + \sigma_1 I \quad (7.20)$$

and

$$A_1 = Q_1 R_1 + \sigma_2 I \quad (7.21)$$

$$A_2 = Q_1^* (A_1 - \sigma_2 I) Q_1 + \sigma_2 I \quad (7.22)$$

$$(7.23)$$

Combining Equations (7.20) and (7.21) yields

$$R_0 Q_0 + (\sigma_1 - \sigma_2) I = Q_1 R_1 \quad (7.24)$$

Multiplying by Q_0 on the left and R_0 on the right results in

$$Q_0 R_0 Q_0 R_0 + Q_0 (\sigma_1 - \sigma_2) R_0 = Q_0 Q_1 R_1 R_0 \quad (7.25)$$

Thus

$$Q_0 Q_1 R_1 R_0 = Q_0 R_0 (Q_0 R_0 + (\sigma_1 - \sigma_2) I) \quad (7.26)$$

$$= (A_0 - \sigma_1 I) (A_0 - \sigma_2 I) \quad (7.27)$$

$$= A_0^2 - 2\operatorname{Re}(\sigma) A_0 + |\sigma|^2 I \quad (7.28)$$

The right-hand side of Equation (7.28) is a real matrix; thus $(Q_0 Q_1) (R_1 R_0)$ is the QR factorization of this real matrix:

$$A_2 = (Q_0 Q_1)^* A_0 (Q_0 Q_1) \quad (7.29)$$

Example 7.5

Find the eigenvalues of

$$A = \begin{bmatrix} 9 & 3 & 13 & 7 & 17 \\ 4 & 6 & 15 & 9 & 1 \\ 19 & 9 & 5 & 11 & 19 \\ 20 & 12 & 3 & 2 & 15 \\ 9 & 6 & 6 & 6 & 10 \end{bmatrix}$$

Solution 7.5 Applying the QR as in Example 7.3 yields the following updated A_{k+1} after 12 iterations:

$$A_{13} = \left[\begin{array}{c|ccc|cc} 46.6670 & 5.4606 & 8.0469 & -13.9039 & -7.0580 \\ \hline 0.0000 & -5.8266 & 1.6519 & 2.1452 & 0.8102 \\ 0.0000 & -10.4370 & -12.3090 & 5.6634 & 4.1548 \\ \hline 0.0000 & 0.0000 & -0.0000 & -0.8085 & -1.3722 \\ 0.0000 & 0.0000 & -0.0000 & 9.3775 & 4.2771 \end{array} \right] \quad (7.30)$$

This matrix is now in upper triangular block form. The eigenvalues of the two 2×2 matrices can be computed. The eigenvalues of

$$\begin{bmatrix} -5.8266 & 1.6519 \\ -10.4370 & -12.3090 \end{bmatrix}$$

are $-9.0678 \pm j2.5953$ and the eigenvalues of

$$\begin{bmatrix} -0.8085 & -1.3722 \\ 9.3775 & 4.27710 \end{bmatrix}$$

are $1.7343 \pm j2.5303$. The real eigenvalue 46.6670 completes the set of eigenvalues of A . ■

7.3 Arnoldi Methods

In large interconnected systems, it is either impractical or intractable to find all of the eigenvalues of the system state matrix due to restrictions on computer memory and computational speed. The Arnoldi method has been developed as an algorithm that iteratively computes k eigenvalues of an $n \times n$ matrix A , where k is typically much smaller than n . This method therefore bypasses many of the constraints imposed by large matrix manipulation required by methods such as the QR decomposition. If the k eigenvalues are chosen selectively, they can yield rich information about the system under consideration, even without the full set of eigenvalues. The Arnoldi method was first developed in [3], but suffered from poor numerical properties such as loss of orthogonality and slow convergence. Several modifications to the Arnoldi method have overcome these shortcomings. The Modified Arnoldi Method (MAM) has been used frequently in solving eigenvalue problems in power system applications [32], [56]. This approach introduced preconditioning and explicit restart techniques to retain orthogonality. Unfortunately, however, an explicit restart will often discard useful information. The restart problem was solved by using implicitly shifted QR steps [49] in the Implicitly Restarted Arnoldi (IRA) method. Several commercial software packages have

been developed around the IRA method, including the well-known ARPACK and the MATLAB `speig` routines.

The basic approach of the Arnoldi method is to iteratively update a low-order matrix H whose eigenvalues successively approximate the selected eigenvalues of the larger A matrix, such that

$$AV = VH; \quad V^*V = I \quad (7.31)$$

where V is an $n \times k$ matrix and H is a $k \times k$ Hessenberg matrix. As the method progresses, the eigenvalues of A are approximated by the diagonal entries of H , yielding

$$HV_i = V_i D \quad (7.32)$$

where V_i is a $k \times k$ matrix whose columns are the eigenvectors of H (approximating the eigenvectors of A) and D is a $k \times k$ matrix whose diagonal entries are the eigenvalues of H (approximating the eigenvalues of A). The Arnoldi method is an orthogonal projection method onto a *Krylov* subspace.

The Arnoldi procedure is an algorithm for building an orthogonal basis of the Krylov subspace. One approach is given as:

The k -step Arnoldi Factorization

Starting with a vector v_1 of unity norm, for $j = 1, \dots, k$ compute:

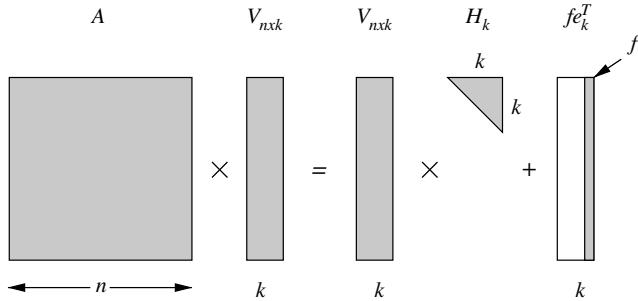
1. $H(i, j) = v_i^T A v_j$ for $i = 1, \dots, j$
2. $w_j = Av_j - \sum_{i=1}^j H(i, j)v_i$
3. $H(j+1, j) = \|w_j\|_2$
4. If $H(j+1, j) = 0$, then stop
5. $v_{j+1} = \frac{w_j}{H(j+1, j)}$

At each step, the algorithm multiplies the previous Arnoldi vector v_j by A and then orthonormalizes the resulting vector w_j against all previous v_i 's. The k -step Arnoldi factorization is shown in Figure 7.1 and is given by

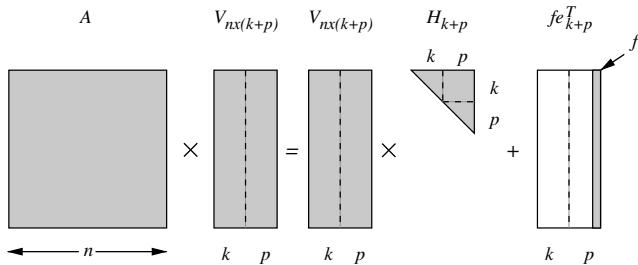
$$AV_k = V_k H_k + w_k e_k^T \quad (7.33)$$

The columns $V = [v_1, v_2, \dots, v_k]$ form an orthonormal basis for the Krylov subspace and H is the orthogonal projection of A onto this space. It is desirable for $\|w_k\|$ to become small because this indicates that the eigenvalues of H are accurate approximations to the eigenvalues of A . However, this “convergence” often comes at the price of numerical orthogonality in V . Therefore, the k -step Arnoldi factorization is “restarted” to preserve orthogonality.

Implicit restarting provides a means to extract rich information from very large Krylov subspaces while avoiding the storage and poor numerical properties associated with the standard approach. This is accomplished by continually compressing the information into a fixed size k -dimensional subspace, by

**FIGURE 7.1**

A \$k\$-step Arnoldi factorization

**FIGURE 7.2**

A \$(k+p)\$-step Arnoldi factorization

using a shifted QR mechanism. A \$(k+p)\$-step Arnoldi factorization

$$AV_{k+p} = V_{k+p}H_{k+p} + w_{k+p}e_{k+p}^T \quad (7.34)$$

is compressed to a factorization of length \$k\$ that retains the eigen-information of interest. This is accomplished using QR steps to apply \$p\$ shifts to yield

$$A\hat{V}_{k+p} = \hat{V}_{k+p}\hat{H}_{k+p} + \hat{w}_{k+p} \quad (7.35)$$

where \$\hat{V}_{k+p} = V_{k+p}Q\$, \$\hat{H}_{k+p} = Q^*H_{k+p}Q\$, and \$\hat{w}_{k+p} = w_{k+p}e_{k+p}^TQ\$. It may be shown that the first \$k-1\$ entries of the vector \$e_{k+p}^TQ\$ are zero [50]. Equating the first \$k\$ columns on both sides yields an updated \$k\$-step Arnoldi factorization. This now provides the “restart” vectors for extending the \$k\$-step Arnoldi factorization to the \$(k+p)\$-step Arnoldi factorization, shown in Figure 7.2.

The implicitly restarted Arnoldi algorithm consists of three main steps: initialization, iteration/refinement, and final calculation of the eigenvalues

and eigenvectors.

Implicitly Restarted Arnoldi Algorithm

1. Initialization

Using the vector v_1 as a starting vector, generate a k -step Arnoldi factorization. At each step k of the factorization, the vector V_k is augmented by a vector v_k satisfying Equation (7.33). Note that H_k is a Hessenberg matrix. The shaded regions in Figure 7.1 represent nonzero entries. The unshaded region of $w_k e_k^T$ is a zero matrix of $(k - 1)$ columns. The Arnoldi factorization is entirely dependent on the choice of initial vector v_1 .

2. Iteration/Refinement

(a) Extend the k -step factorization by p steps.

Each of the p additions represents an eigenvalue/eigenvector that can be discarded at the end of the iteration if it does not meet the chosen criteria. In general, the choice of p is a trade-off between the length of factorization that may be tolerated and the rate of convergence. For most problems, the size of p is determined experimentally. The only requirement is that $1 \leq p \leq n - k$.

(b) Calculate eigenvalues of H_{k+p}

After the p -step extension has been completed, the eigenvalues of H_{k+p} are calculated by the QR method and sorted according to a predetermined sort criterion S and ordered from best to worst. The p worst eigenvalues ($\sigma_1, \sigma_2, \dots, \sigma_p$) are used as shifts to perform p shifted QR factorizations. Since the matrix H_{k+p} in the Arnoldi factorization

$$AV_{k+p} = V_{k+p}H_{k+p} + w_{k+p}e_{k+p}^T \quad (7.36)$$

is relatively small, the shifted QR factorization can be used efficiently to calculate the eigenvalues of H .

(c) Update the Arnoldi matrices

$$\begin{aligned}\hat{V}_{k+p} &= V_{k+p}Q \\ \hat{H}_{k+p} &= Q^*H_{k+p}Q \\ \hat{w}_{k+p} &= w_{k+p}e_{k+p}^T Q\end{aligned}$$

Note that the updated matrix \hat{V}_{k+p} has orthonormal columns since it is the product of V and an orthogonal matrix Q .

(d) Obtain a new k -step Arnoldi factorization by equating the first k columns on each side of Equation (7.35) and discarding the last p equations:

$$A\hat{V}_k = \hat{V}_k \hat{H}_k + \hat{w}_k e_k^T$$

The vector \hat{w} is the new residual vector that is being driven to zero.

(e) If

$$\|AV_k - V_k H_k\| \leq \varepsilon$$

where ε is the preselected convergence tolerance, then the iteration/refinement terminates. Otherwise the process is repeated until tolerance is achieved.

3. Eigenvalue/Eigenvector Calculation

The last step in the Arnoldi method is to compute the eigenvalues and eigenvectors of the reduced matrix H_k from

$$H_k V_k + V_h D_k \quad (7.37)$$

The eigenvectors of A are then calculated as

$$V_k = V_h V_k \quad (7.38)$$

and the desired eigenvalues of A may be obtained from the diagonal entries of D_k :

$$AV_k = V_k D_k \quad (7.39)$$

Example 7.6

Using a three-step Arnoldi factorization, find the two smallest (in magnitude) eigenvalues and corresponding eigenvectors of the matrix of Example 7.2.

Solution 7.6 Since the two smallest eigenvalues are desired, the value of k is two. After the initialization step, the two-step Arnoldi method will be extended up to three steps; therefore, p is one. Thus, at each step, three eigenvalues will be calculated and the worst eigenvalue will be discarded.

The factorization can be initialized with an arbitrary nonzero vector. In many software implementations, the starting vector is chosen randomly such that all of the entries have absolute value less than 0.5. The starting vector for this example will be

$$v_0 = \begin{bmatrix} 0.2500 \\ 0.2500 \\ 0.2500 \\ 0.2500 \end{bmatrix}$$

To satisfy the requirement that the initial vector have unity norm, the starting vector is normalized to yield

$$\begin{aligned} v_1 &= \frac{Av_0}{\|Av_0\|} \\ &= \begin{bmatrix} 0.4611 \\ 0.2306 \\ 0.5764 \\ 0.6340 \end{bmatrix} \end{aligned}$$

After the initial vector has been chosen, the Arnoldi factorization is applied for k steps; thus

$$h_{2,1}v_2 = Av_1 - h_{1,1}v_1 \quad (7.40)$$

where v_2 produces the second column of the matrix V_k and $h_{1,1}$ is chosen such that

$$h_{1,1} = \langle v_1, Av_1 \rangle = v_1^T Av_1 \quad (7.41)$$

where $\langle \cdot \rangle$ denotes inner product. Thus, solving Equation (7.41) yields $h_{1,1} = 18.0399$. Applying the Arnoldi factorization for w_1 yields

$$\begin{aligned} w_1 &= h_{2,1}v_2 = Av_1 - h_{1,1}v_1 \\ &= \begin{bmatrix} 1 & 3 & 4 & 8 \\ 2 & 1 & 2 & 3 \\ 4 & 3 & 5 & 8 \\ 9 & 2 & 7 & 4 \end{bmatrix} \begin{bmatrix} 0.4611 \\ 0.2306 \\ 0.5764 \\ 0.6340 \end{bmatrix} - (18.0399) \begin{bmatrix} 0.4611 \\ 0.2306 \\ 0.5764 \\ 0.6340 \end{bmatrix} \\ &= \begin{bmatrix} 0.2122 \\ 0.0484 \\ 0.0923 \\ -0.2558 \end{bmatrix} \end{aligned}$$

The factor $h_{2,1}$ is chosen to normalize v_2 to unity; thus $h_{2,1} = 0.3483$, and

$$v_2 = \begin{bmatrix} 0.6091 \\ 0.1391 \\ 0.2650 \\ -0.7345 \end{bmatrix}$$

Calculating the remaining values of the Hessenberg matrix yields

$$\begin{aligned} h_{1,2} &= v_1^* Av_2 = 0.1671 \\ h_{2,2} &= v_2^* Av_2 = -6.2370 \end{aligned}$$

and

$$w_2 = h_{3,2}v_2 = Av_2 - h_{1,2}v_1 - h_{2,2}v_2 = \begin{bmatrix} -0.0674 \\ 0.5128 \\ -0.1407 \\ -0.0095 \end{bmatrix}$$

These values can be checked to verify that they satisfy Equation (7.33) for $i = 2$:

$$AV_2 = V_2H_2 + w_2 [0 \ 1]$$

where

$$V_2 = [v_1 \ v_2] = \begin{bmatrix} 0.4611 & 0.6091 \\ 0.2306 & 0.1391 \\ 0.5764 & 0.2650 \\ 0.6340 & -0.7345 \end{bmatrix}$$

and

$$H_2 = \begin{bmatrix} h_{1,1} & h_{1,2} \\ h_{2,1} & h_{2,2} \end{bmatrix} = \begin{bmatrix} 18.0399 & 0.1671 \\ 0.3483 & -6.2370 \end{bmatrix}$$

This completes the initialization stage.

After the initial k -step Arnoldi sequence has been generated, it can be extended to $k + p$ steps. In this example $p = 1$, so only one more extension is required. From the initialization, $w_2 = h_{3,2}v_2$ from which $h_{3,2}$ and v_2 can be extracted (recalling that $\|v_2\| = 1.0$) to yield $h_{3,2} = 0.5361$ and

$$v_3 = \begin{bmatrix} -0.1257 \\ 0.9565 \\ -0.2625 \\ -0.0178 \end{bmatrix}$$

The Hessenberg matrix H_3 becomes

$$H_3 = \begin{bmatrix} h_{1,1} & h_{1,2} & h_{1,3} \\ h_{2,1} & h_{2,2} & h_{2,3} \\ 0 & h_{3,2} & h_{3,3} \end{bmatrix} = \begin{bmatrix} 18.0399 & 0.1671 & 0.5560 \\ 0.3483 & -6.2370 & 2.0320 \\ 0 & 0.5361 & -0.2931 \end{bmatrix}$$

where

$$\begin{aligned} h_{1,3} &= v_1^T A v_3 \\ h_{2,3} &= v_2^T A v_3 \\ h_{3,3} &= v_3^T A v_3 \end{aligned}$$

and

$$w_3 = Av_3 - h_{1,3}v_1 - h_{2,3}v_2 - h_{3,3}v_3 = \begin{bmatrix} 0.0207 \\ -0.0037 \\ -0.0238 \\ 0.0079 \end{bmatrix}$$

The next step is to compute (using QR factorization) and sort the eigenvalues (and eigenvectors) of the small matrix H_3 . The eigenvalues of H_3 are

$$\sigma = \begin{bmatrix} 18.0425 \\ -6.4166 \\ -0.1161 \end{bmatrix}$$

Since the smallest two eigenvalues are desired, the eigenvalues are sorted such that the desired eigenvalues are at the bottom (which they already are). The undesired eigenvalue estimate is $\sigma_1 = 18.0425$. Applying the shifted QR factorization to $H_3 - \sigma_1 I$ yields

$$H_3 - \sigma_1 I = \begin{bmatrix} -0.0026 & 0.1671 & 0.5560 \\ 0.3483 & -24.2795 & 2.0320 \\ 0 & 0.5361 & -18.3356 \end{bmatrix}$$

and

$$QR = \begin{bmatrix} -0.0076 & -0.0311 & 0.9995 \\ 1.0000 & -0.0002 & 0.0076 \\ 0 & 0.9995 & 0.0311 \end{bmatrix} \begin{bmatrix} 0.3483 & -24.2801 & 2.0277 \\ 0 & 0.5364 & -18.3445 \\ 0 & 0 & 0 \end{bmatrix}$$

From Q , the update \hat{H} can be found:

$$\begin{aligned} \hat{H}_3 &= Q^* H_3 Q \\ &= \begin{bmatrix} -6.2395 & 2.0216 & 0.2276 \\ 0.5264 & -0.2932 & -0.5673 \\ 0 & 0 & 18.0425 \end{bmatrix} \end{aligned}$$

Note that the $\hat{H}(3, 2)$ element is now zero. Continuing the algorithm yields the update for \hat{V} :

$$\hat{V} = V_3 Q = \begin{bmatrix} 0.6056 & -0.1401 & 0.4616 \\ 0.1373 & 0.9489 & 0.2613 \\ 0.2606 & -0.2804 & 0.5699 \\ -0.7392 & -0.0374 & 0.6276 \end{bmatrix}$$

where $V_3 = [v_1 \ v_2 \ v_3]$, and

$$\hat{w}_3 e^T = w_3 e_3^T Q = \begin{bmatrix} 0 & 0.0207 & 0.0006 \\ 0 & -0.0037 & -0.0001 \\ 0 & -0.0238 & -0.0007 \\ 0 & 0.0079 & 0.0002 \end{bmatrix}$$

Note that the first column of $\hat{w}e^T$ is zeros, so that a new k -step Arnoldi factorization can be obtained by equating the first k columns on each side such that

$$A\hat{V}_2 = \hat{V}_2 \hat{H}_2 + \hat{w}_2 e_2^T \quad (7.42)$$

The third columns of \hat{V} and \hat{H} are discarded.

This iteration/refinement procedure is continued until

$$\|AV - VH\| = \|we^T\| < \varepsilon$$

at which time the calculated eigenvalues will be obtained within order of ε accuracy. ■

7.4 Singular Value Decomposition

Singular value decomposition (SVD) produces three matrices whose product is the (possibly rectangular) matrix A . In matrix form, the SVD is

$$A = U\Sigma V^T \quad (7.43)$$

where U satisfies $U^T U = I$ and the columns of U are the orthonormal eigenvectors of AA^T , V satisfies $V^T V = I$ and the columns of V are the orthonormal eigenvectors of $A^T A$, and Σ is a diagonal matrix containing the square roots of the eigenvalues corresponding to U (or V) in descending order.

The SVD decomposition can be found by applying either the QR or Arnoldi method to the matrices $A^T A$ and AA^T to compute the eigenvalues and eigenvectors. Once the eigenvalues are found, the singular values are the square roots. The condition number of a matrix is a measure of the “invertibility” of a matrix and is defined as the ratio of the largest singular value to the smallest singular value. A large condition number indicates a nearly singular matrix.

Example 7.7

Find the singular value decomposition of

$$A = \begin{bmatrix} 1 & 2 & 4 & 9 & 3 \\ 3 & 1 & 3 & 2 & 6 \\ 4 & 2 & 5 & 7 & 7 \\ 8 & 3 & 8 & 4 & 10 \end{bmatrix}$$

Solution 7.7 The matrix A is 4×5 ; therefore U will be a 4×4 matrix, Σ will be a 4×5 matrix with the four singular values on the diagonal followed by a column of zeros, and V will be a 5×5 matrix.

Starting with

$$\hat{A} = A^T A$$

the QR method for finding the eigenvalues and eigenvectors of \hat{A} yields

$$\hat{A} = \begin{bmatrix} 90 & 37 & 97 & 75 & 129 \\ 37 & 18 & 45 & 46 & 56 \\ 97 & 45 & 114 & 109 & 145 \\ 75 & 46 & 109 & 150 & 128 \\ 129 & 56 & 145 & 128 & 194 \end{bmatrix}$$

$$D = \begin{bmatrix} 507.6670 & 0 & 0 & 0 & 0 \\ 0 & 55.1644 & 0 & 0 & 0 \\ 0 & 0 & 3.0171 & 0 & 0 \\ 0 & 0 & 0 & 0.1516 & 0 \\ 0 & 0 & 0 & 0 & 0.0000 \end{bmatrix}$$

$$V^T = \begin{bmatrix} -0.3970 & 0.4140 & -0.3861 & -0.7189 & -0.0707 \\ -0.1865 & -0.0387 & -0.2838 & 0.3200 & -0.8836 \\ -0.4716 & 0.0610 & -0.5273 & 0.5336 & 0.4595 \\ -0.4676 & -0.8404 & 0.0688 & -0.2645 & 0.0177 \\ -0.6054 & 0.3421 & 0.6983 & 0.1615 & -0.0530 \end{bmatrix}$$

The matrix Σ has the squares of the singular values on the diagonals and must be the same dimension as A ; thus

$$\Sigma = \begin{bmatrix} 22.5315 & 0 & 0 & 0 \\ 0 & 7.4273 & 0 & 0 \\ 0 & 0 & 1.7370 & 0 \\ 0 & 0 & 0 & 0.3893 \end{bmatrix}$$

To find U , repeat with $\hat{A} = AA^T$ to obtain

$$U = \begin{bmatrix} -0.3853 & 0.8021 & -0.2009 & 0.4097 \\ -0.3267 & -0.2367 & 0.7502 & 0.5239 \\ -0.5251 & 0.2161 & 0.3574 & -0.7415 \\ -0.6850 & -0.5039 & -0.5188 & 0.0881 \end{bmatrix}$$

■

In addition to condition number, another common use of the SVD is to calculate the pseudoinverse A^+ of a non-square matrix A . The most commonly encountered pseudoinverse is the Moore–Penrose matrix inverse, which is a special case of a general type of pseudoinverse known as a matrix 1-inverse. The pseudoinverse is commonly used to solve the least-squares problem $Ax = b$ when A is nonsingular or nonsquare. From Section 6.1, the least-squares problem solution is given by

$$\begin{aligned} x &= (A^T A)^{-1} Ab \\ &= A^+ b \end{aligned}$$

The matrix A^+ can be found through LU factorization, but a much more common approach is to use SVD. In this case, the pseudoinverse is given by

$$A^+ = V\Sigma^+U^T \quad (7.44)$$

where Σ^+ is a matrix of the same dimension as A^T with the reciprocal of the singular values on the diagonal.

Example 7.8

Repeat Example 6.1 using a pseudoinverse.

Solution 7.8 The system of equations is repeated here for convenience:

$$\begin{bmatrix} 4.27 \\ -1.71 \\ 3.47 \\ 2.50 \end{bmatrix} = \begin{bmatrix} 0.4593 & -0.0593 \\ 0.0593 & -0.4593 \\ 0.3111 & 0.0889 \\ 0.0889 & 0.3111 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} \quad (7.45)$$

Use the pseudoinverse to solve for

$$\begin{bmatrix} V_1 \\ V_2 \end{bmatrix}$$

Applying the SVD

$$U = \begin{bmatrix} -0.5000 & -0.6500 & -0.5505 & -0.1567 \\ 0.5000 & -0.6500 & 0.1566 & 0.5505 \\ -0.5000 & -0.2785 & 0.8132 & -0.1059 \\ -0.5000 & 0.2785 & -0.1061 & 0.8131 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 0.5657 & 0 \\ 0 & 0.5642 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$V = \begin{bmatrix} -0.7071 & -0.7071 \\ -0.7071 & 0.7071 \end{bmatrix}$$

The matrix Σ^+ is

$$\Sigma^+ = \begin{bmatrix} 1.7677 & 0 & 0 & 0 \\ 0 & 1.7724 & 0 & 0 \end{bmatrix}$$

leading to A^+

$$A^+ = \begin{bmatrix} 1.4396 & 0.1896 & 0.9740 & 0.2760 \\ -0.1896 & -1.4396 & 0.2760 & 0.9740 \end{bmatrix}$$

and

$$\begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = \begin{bmatrix} 1.4396 & 0.1896 & 0.9740 & 0.2760 \\ -0.1896 & -1.4396 & 0.2760 & 0.9740 \end{bmatrix} \begin{bmatrix} 4.27 \\ -1.71 \\ 3.47 \\ 2.50 \end{bmatrix}$$

$$= \begin{bmatrix} 9.8927 \\ 5.0448 \end{bmatrix}$$

which is the same solution as Example 6.1. ■

7.5 Modal Identification

Although many systems are inherently nonlinear, in some instances they may respond to well-tuned linear controls. In order to implement linear feedback control, the system designer must have an accurate model of sufficiently low

order from which to design the control. Several approaches to developing such lower-order models have included dynamic equivalencing, eigenanalysis, and pole/zero cancellation. Frequently, however, the original system is too complex or the parameters are not known with enough accuracy to produce an adequate reduced-order model. In practice, the system may have parameters that drift with time or operating condition, which compromises the accuracy of the mathematical model. In these cases, it is desirable to extract the modal information directly from the system response to a perturbation. Using this approach, it may be possible to replace the actual dynamic model with an estimated linear model that is derived from the system output waveform. The time-varying dynamic response of a power system to a disturbance may be composed of numerous modes that must be identified. Several methods have been proposed to extract the pertinent modal information from time-varying responses. An appropriate method must consider the inclusion of nonlinearities, the size of the model that can be effectively utilized, and the reliability of the results.

Methods that are applied directly to the nonlinear system simulation or field measurements include the effects of nonlinearities. In full-state eigenvalue analysis, the size of the system model is typically limited to several hundred states with present computing capabilities. This means that a typical system containing several thousand nodes must be reduced using dynamic equivalencing. Modal analysis techniques that operate directly on system output are not limited by system size. This means that standard time-domain-analysis results are directly usable. This eliminates the possibility of losing some of system modal content due to reduction. The estimated linear model may then be used for control design applications or other linear analysis techniques. The estimated model may be chosen to be of lower order than the original model, but still retain the dominant modal characteristics.

This problem may be posed such that, given a set of measurements that vary with time, it is desired to fit a time-varying waveform of prespecified form to the actual waveform (i.e., minimize the error between the actual measured waveform and the proposed waveform). The coefficients of the prespecified waveform yield the dominant modal characteristics of the underlying linear system. Consider the following linear system:

$$\dot{x}(t) = Ax(t) \quad x(t_0) = x_0 \quad (7.46)$$

where

$$x_i(t) = \sum_{k=1}^n a_k e^{(b_k t)} \cos(\omega_k t + \theta_k) \quad (7.47)$$

is one of the n states. The parameters a_k and θ_k are derived from the influence of the initial conditions, whereas the parameters b_k and ω_k are derived from the eigenvalues of A . The estimation of these responses yields modal information about the system that can be used to predict possible unstable

behavior, controller design, parametric summaries for damping studies, and modal interaction information.

Any time-varying function can be fit to a series of complex exponential functions over a finite time interval. However, it is not practical to include a large number of terms in the fitting function. The problem then becomes one of minimizing the error between the actual time-varying function and the proposed function by estimating the magnitude, phase, and damping parameters of the fitting function. In the problem of estimating a nonlinear waveform by a series of functions, the minimization function is given by

$$\min f = \sum_{i=1}^N \left[\sum_{k=1}^n \left[a_k e^{(b_k t_i)} \cos(\omega_k t_i + \theta_k) \right] - y_i \right]^2 \quad (7.48)$$

where n is the number of desired modes of the approximating waveform, N is the number of data samples, y_i is the sampled waveform, and

$$[a_1 \ b_1 \ \omega_1 \ \theta_1, \dots, a_n \ b_n \ \omega_n \ \theta_n]^T$$

are the parameters to be estimated.

There are several approaches to estimating the modal content of a time-varying waveform. The Prony method is well known and widely used in power systems applications. The matrix pencil approach was introduced for extracting poles from antennas' electromagnetic transient responses. The Levenberg–Marquardt iteratively updates the modal parameters by an analytic optimization to minimize the error between the resulting waveform and the input data.

7.5.1 Prony Method

One approach to estimating the various parameters is the *Prony method* [23]. This method is designed to directly estimate the parameters for the exponential terms by fitting the function

$$\hat{y}(t) = \sum_{i=1}^n A_i e^{\sigma_i t} \cos(\omega_i t + \phi_i) \quad (7.49)$$

to an observed measurement for $y(t)$, where $y(t)$ consists of N samples

$$y(t_k) = y(k), \quad k = 0, 1, \dots, N - 1$$

that are evenly spaced by a time interval Δt . Since the measurement signal $y(t)$ may contain noise or dc offset, it may have to be conditioned before the fitting process is applied.

Note that Equation (7.49) can be recast in complex exponential form as

$$\hat{y}(t) = \sum_{i=1}^n B_i e^{\lambda_i t} \quad (7.50)$$

which can be translated to

$$\hat{y}(k) = \sum_{i=1}^n B_i z_i^k \quad (7.51)$$

where

$$z_i = e^{(\lambda_i \Delta t)} \quad (7.52)$$

The system eigenvalues λ can be found from the discrete modes by

$$\lambda_i = \frac{\ln(z_i)}{\Delta t} \quad (7.53)$$

The z_i are the roots of the n th order polynomial

$$z^n - (a_1 z^{n-1} + a_2 z^{n-2} + \dots + a_n z^0) = 0 \quad (7.54)$$

where the a_i coefficients are unknown and must be calculated from the measurement vector.

The basic Prony method is summarized as

Prony Method

1. Assemble selected elements of the record into a Toeplitz data matrix.
2. Fit the data with a discrete linear prediction model, such as a least squares solution.
3. Find the roots of the characteristic polynomial (7.54) associated with the model of Step 1.
4. Using the roots of Step 3 as the complex modal frequencies for the signal, determine the amplitude and initial phase for each mode.

These steps are performed in the z -domain, translating the eigenvalues to the s -domain as a final step.

The approach to the Toeplitz (or the closely related Hankel) matrix assembly of Step 1 has received the most attention in the literature. The problem can be formulated in many different ways. If the initial (i.e., $i < 0$) and post (i.e., $i > N$) conditions are assumed to be zero, then the subscript ranges X_1 through X_4 in Figure 7.3 represent four such formulations. The range X_1 is termed the covariance problem. Because the initial and post conditions are not used, no assumption is required on their value. The range X_4 is called the correlation problem; it incorporates both initial and post conditions. The remaining problems, X_2 and X_3 , are termed the prewindowed and postwindowed methods.

In the majority of practical cases, the Toeplitz matrix is nonsquare with more rows than columns. A Toeplitz matrix is a matrix with a constant diagonal in which each descending diagonal from left to right is constant.

$$\begin{array}{c}
 \begin{array}{ccccc}
 & \uparrow & \uparrow & & \\
 & X_2 & & & \\
 & \downarrow & \downarrow & & \\
 X_4 & & X_1 & & \\
 & \downarrow & \downarrow & & \\
 & X_3 & & & \\
 & \downarrow & & & \\
 \end{array}
 \left[\begin{array}{ccccc}
 y_0 & 0 & \dots & 0 & -1 \\
 y_1 & y_0 & \dots & 0 & 0 \\
 \vdots & \vdots & \ddots & \vdots & \vdots \\
 y_n & y_{n-1} & \dots & y_0 & a_1 \\
 y_{n+1} & y_n & \dots & y_1 & 0 \\
 \vdots & \vdots & \ddots & \vdots & \vdots \\
 y_N & y_{N-1} & \dots & y_{N-n-1} & a_n \\
 0 & \vdots & \ddots & \vdots & 0 \\
 0 & \dots & 0 & y_N & 0
 \end{array} \right] = \left[\begin{array}{c} 0 \\ 0 \\ \vdots \\ 0 \end{array} \right]
 \end{array}$$

FIGURE 7.3
Toeplitz matrix

The system of equations requires a least-squares solution to find the factors a_1 through a_n . After the a_i coefficients are obtained, the n roots z_i of the polynomial in Equation (7.52) can be found by factoring.

Once z_i has been computed from the roots of Equation (7.54), then the eigenvalues λ_i can be calculated from Equation (7.53). The next step is to find the B_i that produces $\hat{y}(k) = y(k)$ for all k . This leads to the following relationship:

$$\left[\begin{array}{cccc}
 z_1^0 & z_2^0 & \dots & z_n^0 \\
 z_1^1 & z_2^1 & \dots & z_n^1 \\
 \vdots & \vdots & \vdots & \vdots \\
 z_1^{N-1} & z_2^{N-1} & \dots & z_n^{N-1}
 \end{array} \right] \left[\begin{array}{c} B_1 \\ B_2 \\ \vdots \\ B_n \end{array} \right] = \left[\begin{array}{c} y(0) \\ y(1) \\ \vdots \\ y(N-1) \end{array} \right] \quad (7.55)$$

which can be succinctly expressed as

$$ZB = Y \quad (7.56)$$

Note that the matrix Z is $n \times N$; therefore, Equation (7.56) must also be solved by the least-squares method. The estimating waveform $\hat{y}(t)$ is then calculated from Equation (7.50). The reconstructed signal $\hat{y}(t)$ will usually not fit $y(t)$ exactly. An appropriate measure for the quality of this fit is a “signal to noise ratio (SNR)” given by

$$\text{SNR} = 20 \log \frac{\|\hat{y} - y\|}{\|y\|} \quad (7.57)$$

where the SNR is given in decibels (dB).

Since the fit for this method may be inexact, it is desirable to have control over the level of error between the fitting function and the original waveform. In this case, a nonlinear least squares can provide improved results.

7.5.2 The Matrix Pencil Method

The Prony method described in the previous section is a “polynomial” method in that it includes the process of finding the poles z_i of a characteristic polynomial. The matrix pencil (MP) method produces a matrix whose roots provide z_i . The poles are found as the solution of a generalized eigenvalue problem [24], [46]. The matrix pencil is given by

$$[Y_2] - \lambda [Y_1] = [Z_1] [B] \{[Z_0] - \lambda [I]\} [Z_2] \quad (7.58)$$

where

$$[Y] = \begin{bmatrix} y(0) & y(1) & \dots & y(L) \\ y(1) & y(2) & \dots & y(L+1) \\ \vdots & \vdots & & \vdots \\ y(N-L) & y(N-L+1) & \dots & y(N) \end{bmatrix} \quad (7.59)$$

$$[Z_0] = \text{diag}[z_1, z_2, \dots, z_n] \quad (7.60)$$

$$[Z_1] = \begin{bmatrix} 1 & 1 & \dots & 1 \\ z_1 & z_2 & \dots & z_n \\ \vdots & \vdots & & \vdots \\ z_1^{(N-L-1)} & z_2^{(N-L-2)} & \dots & Z_n^{(N-L-1)} \end{bmatrix} \quad (7.61)$$

$$[Z_2] = \begin{bmatrix} 1 & z_1 & \dots & z_1^{L-1} \\ 1 & z_2 & \dots & z_2^{L-1} \\ \vdots & \vdots & & \vdots \\ 1 & Z_n & \dots & Z_n^{l-1} \end{bmatrix} \quad (7.62)$$

$[B]$ = matrix of residuals

$[I]$ = $n \times n$ identity matrix

n = desired number of eigenvalues

L = pencil parameter, such that $n \leq L \leq N - n$

Matrix Pencil Method

1. Choose L such that $n \leq L \leq N - n$.
2. Construct the matrix $[Y]$.

3. Perform a singular value decomposition of $[Y]$ to obtain

$$[Y] = [U][S][V]^T \quad (7.63)$$

where $[U]$ and $[V]$ are unitary matrices and contain the eigenvectors of $[Y][Y]^T$ and $[Y]^T[Y]$, respectively.

4. Construct the matrices $[V_1]$ and $[V_2]$ such that

$$V_1 = [v_1 \ v_2 \ v_3 \ \dots \ v_{n-1}] \quad (7.64)$$

$$V_2 = [v_2 \ v_3 \ v_4 \ \dots \ v_n] \quad (7.65)$$

where v_i is the i th right singular vector of V .

5. Construct $[Y_1]$ and $[Y_2]$:

$$[Y_1] = [V_1]^T[V_1]$$

$$[Y_2] = [V_2]^T[V_1]$$

6. The desired poles z_i may be found as the generalized eigenvalues of the matrix pair $\{[Y_2]; [Y_1]\}$.

From this point, the remainder of the algorithm follows that of the Prony method to calculate the eigenvalues λ and the residual matrix B .

If the pencil parameter L is chosen such that $L = N/2$, then the performance of the method is very close to the optimal bound [24].

It has been shown that, under noise, the statistical variance of the poles found from the matrix pencil method is always less than that of the Prony method [24].

7.5.3 The Levenberg–Marquardt Method

The nonlinear least squares for data fitting applications has the general form

$$\text{minimize } f(x) = \sum_{k=1}^N [\hat{y}(x, t_i) - y_i]^2 \quad (7.66)$$

where y_i is the output of the system at time t_i , and x is the vector of magnitudes, phases, and damping coefficients of Equation (7.49), which arise from the eigenvalues of the state matrix of the system.

To find the minimum of $f(x)$, the same procedure for developing the Newton–Raphson iteration is applied. The function $f(x)$ is expanded about some x_0 by the Taylor series

$$f(x) \approx f(x_0) + (x - x_0)^T f'(x_0) + \frac{1}{2}(x - x_0)^T f''(x_0)(x - x_0) + \dots \quad (7.67)$$

where

$$\begin{aligned} f'(x) &= \frac{\partial f}{\partial x_j} \quad \text{for } j = 1, \dots, n \\ f''(x) &= \frac{\partial^2 f}{\partial x_j \partial x_k} \quad \text{for } j, k = 1, \dots, n \end{aligned}$$

If the higher-order terms in the Taylor expansion are neglected, then minimizing the quadratic function on the right-hand side of Equation (7.67) yields

$$x_1 = x_0 - [f''(x_0)]^{-1} f'(x_0) \quad (7.68)$$

which yields an approximation for the minimum of the function $f(x)$. This is also one Newton–Raphson iteration update for solving the necessary minimization condition

$$f'(x) = 0$$

The Newton–Raphson Equation (7.68) may be rewritten as the iterative linear system

$$A(x_k)(x_{k+1} - x_k) = g(x_k) \quad (7.69)$$

where

$$\begin{aligned} g_j(x) &= -\frac{\partial f}{\partial x_j}(x) \\ a_{jk}(x) &= \frac{\partial^2 f}{\partial x_j \partial x_k}(x) \end{aligned}$$

and the matrix A is the system Jacobian (or similar iterative matrix).

The derivatives of Equation (7.66) are

$$\frac{\partial f}{\partial x_j}(x) = 2 \sum_{k=1}^N [\hat{y}_k - y_k] \frac{\partial \hat{y}_i}{\partial x_j}(x)$$

and

$$\frac{\partial^2 f}{\partial x_j \partial x_k}(x) = 2 \sum_{k=1}^N \left\{ \frac{\partial \hat{y}_k}{\partial x_j}(x) \frac{\partial \hat{y}_i}{\partial x_k}(x) + [\hat{y}_k - y_k] \frac{\partial^2 \hat{y}_k}{\partial x_j \partial x_k}(x) \right\}$$

In this case, the matrix element a_{jk} contains second derivatives of the functions \hat{y}_i . These derivatives are multiplied by the factor $[\hat{y}_i(x) - y_i]$ and will become small during the minimization of f . Therefore, the argument can be made that these terms can be neglected during the minimization process. Note that, if the method converges, it will converge regardless of whether the exact Jacobian is used in the iteration. Therefore, the iterative matrix A can be simplified as

$$a_{jk} = 2 \sum_{i=1}^N \frac{\partial \hat{y}_i}{\partial x_j}(x) \frac{\partial \hat{y}_i}{\partial x_k}(x) \quad (7.70)$$

and note that $a_{jj}(x) > 0$.

The Levenberg–Marquardt method modifies Equation (7.69) by introducing the matrix \hat{A} with entries

$$\begin{aligned}\hat{a}_{jj} &= (1 + \gamma) a_{jj} \\ \hat{a}_{jk} &= a_{jk} \quad j \neq k\end{aligned}$$

where γ is some positive parameter. Equation (7.69) becomes

$$\hat{A}(x_0)(x_1 - x_0) = g \quad (7.71)$$

For large γ , the matrix \hat{A} will become diagonally dominant. As γ approaches zero, Equation (7.71) will turn into the Newton–Raphson method. The Levenberg–Marquardt method has the basic feature of varying γ to select the optimal characteristics of the iteration. The basic Levenberg–Marquardt algorithm is summarized:

Levenberg–Marquardt Method

1. Set $k = 0$. Choose an initial guess x_0 , γ , and a factor α .
2. Solve the linear system of Equation (7.71) to obtain x_{k+1} .
3. If $f(x_{k+1}) > f(x_k)$, reject x_{k+1} as the new approximation, replace γ by $\alpha\gamma$, and repeat Step 2.
4. If $f(x_{k+1}) < f(x_k)$, accept x_{k+1} as the new approximation, replace γ by γ/α , set $k = k + 1$, and repeat Step 2.
5. Terminate the iteration when

$$\|x_{k+1} - x_k\| < \varepsilon$$

In the problem of estimating a nonlinear waveform by a series of functions, the minimization function is given by

$$\text{minimize } f = \sum_{i=1}^N \left[\sum_{k=1}^m \left[a_k e^{(b_k t_i)} \cos(\omega_k t_i + \theta_k) \right] - y_i \right]^2 \quad (7.72)$$

where m is the number of desired modes of the approximating waveform, and $x = [a_1 \ b_1 \ \omega_1 \ \theta_1 \ \dots \ a_m \ b_m \ \omega_m \ \theta_m]^T$.

As with all of the nonlinear iterative methods, the ability of the Levenberg–Marquardt method to converge to a solution depends on the choice of initial guess. In this case, it is wise to use the results of the matrix pencil or the Prony method to provide the initial values.

7.5.4 Eigensystem Realization Algorithm

The Eigensystem Realization Algorithm (ERA) is based on the singular value decomposition of the Hankel matrix H_0 associated with the linear ringdown of the system. A Hankel matrix is a square matrix with constant skew-diagonals. The Hankel matrices are typically assembled using all of the available data such that the top left-most element of H_0 is y_0 and the bottom right-most element of H_1 is y_N . The Hankel matrices are assembled such that

$$H_0 = \begin{bmatrix} y_0 & y_1 & \cdots & y_r \\ y_1 & y_2 & \cdots & y_{r+1} \\ \vdots & \vdots & \ddots & \vdots \\ y_r & y_{r+1} & \cdots & y_{N-1} \end{bmatrix} \quad (7.73)$$

$$H_1 = \begin{bmatrix} y_1 & y_2 & \cdots & y_{r+1} \\ y_2 & y_3 & \cdots & y_{r+2} \\ \vdots & \vdots & \ddots & \vdots \\ y_{r+1} & y_{r+2} & \cdots & y_N \end{bmatrix} \quad (7.74)$$

where $r = \frac{N}{2} - 1$. This choice of r assumes that the number of data points is sufficient such that $r > n$.

The ERA formulation begins by separating the singular value decomposition of H_0 into two components according to the relative size of the singular values:

$$H_0 = U\Sigma V^T = [U_n \ U_z] \begin{bmatrix} \Sigma_n & 0 \\ 0 & \Sigma_z \end{bmatrix} \begin{bmatrix} V_n^T \\ V_z^T \end{bmatrix} \quad (7.75)$$

where Σ_n and Σ_z are diagonal matrices with their elements ordered by magnitude:

$$\Sigma_n = \text{diag } (\sigma_1, \sigma_2, \dots, \sigma_n) \quad (7.76)$$

$$\Sigma_z = \text{diag } (\sigma_{n+1}, \sigma_{n+2}, \dots, \sigma_N) \quad (7.77)$$

and the singular values are ordered by magnitude such that

$$\sigma_1 > \sigma_2 > \dots > \sigma_n > \sigma_{n+1} > \sigma_{n+2} > \dots > \sigma_N$$

The SVD is a useful tool for determining an appropriate value for n . The ratio of the singular values contained in Σ can determine the best approximation of n . The ratio of each singular value σ_i to the largest singular value σ_{\max} is compared to a threshold value, where p is the number of significant decimal digits in the data:

$$\frac{\sigma_i}{\sigma_{\max}} \approx 10^{-p}$$

If p is set to 3, then any singular values with a ratio below 10^{-3} are assumed to be part of the noise and are not included in the reconstruction of the system. The value of n should be set to the number of singular values with a ratio

above the threshold 10^{-p} . It can be shown that, for a linear system of order n , the diagonal elements of Σ_z are zero (assuming that the impulse response is free of noise). The practical significance of this result is that the relative size of the singular values provides an indication of the identified system order. If the singular values exhibit a significant grouping such that $\sigma_n \geq \sigma_{n+1}$, then, from the partitioned representation, H_0 can be approximated by

$$H_0 \approx U_n \Sigma_n V_n^T \quad (7.78)$$

The method for obtaining the eigenvalue realization algorithm solution can be summarized as follows:

1. Assemble selected elements of the record into Hankel data matrices H_0 and H_1 .
2. Perform the singular value decomposition of H_0 and estimate the system order n based on the magnitude of the singular values.
3. Computer the discrete system matrices as follows:

$$A = \Sigma_n^{-\frac{1}{2}} U_n^T H_1 V_n \Sigma_n^{-\frac{1}{2}} \quad (7.79)$$

$$B = \Sigma_n^{-\frac{1}{2}} V_n^T (1 : n, 1 : n) \quad (7.80)$$

$$C = U_n (1 : N, 1 : n) \Sigma_n^{-\frac{1}{2}} \quad (7.81)$$

$$D = y_0 \quad (7.82)$$

4. Calculate continuous system matrices A_c, B_c assuming a zero order hold and sampling interval Δt :

$$A_c = \ln \left(\frac{A}{\Delta t} \right) \quad (7.83)$$

$$B_c = \left[\int_0^{\Delta t} e^{A\tau} d\tau \right]^{-1} B \quad (7.84)$$

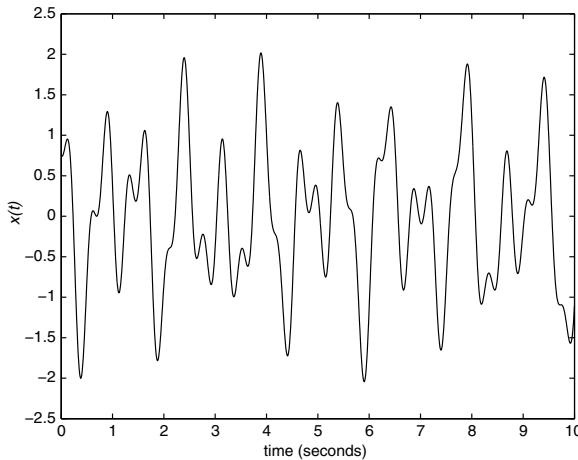
The reduced system response can then be computed from the continuous matrices.

7.5.5 Examples

The effectiveness of these methods will be illustrated with several examples ranging from a simple three-mode linear system to an actual power system oscillation.

Simple Example

The application of the methods will initially consider the waveform shown in

**FIGURE 7.4**

Three-mode waveform

Figure 7.4, which is generated from

$$x(t) = \sum_{i=1}^3 a_i e^{b_i t} (\cos \omega_i t + \theta_i)$$

where

mode	a_i	b_i	ω_i	θ_i
1	1.0	-0.01	8.0	0.0
2	0.6	-0.03	17.0	π
3	0.5	0.04	4.7	$\pi/4$

Each of the methods described previously is used to estimate the signal parameters and to reconstruct the waveform.

Prony

mode	a_i	b_i	ω_i	θ_i
1	0.9927	-0.0098	7.9874	0.0617
2	0.6009	-0.0304	17.0000	3.1402
3	0.5511	0.0217	4.6600	0.9969

Matrix Pencil

mode	a_i	b_i	ω_i	θ_i
1	1.0121	-0.0130	8.0000	0.0008
2	0.6162	-0.0361	16.9988	3.1449
3	0.5092	0.0364	4.6953	0.7989

Levenberg–Marquardt

mode	a_i	b_i	ω_i	θ_i
1	1.0028	-0.0110	7.9998	0.0014
2	0.6010	-0.0305	16.9994	3.1426
3	0.5051	0.0378	4.6967	0.7989

The reconstruction error in each waveform is measured.

$$\text{error} = \sum_{i=1}^N \left[\sum_{k=1}^m \left[a_k e^{(b_k t_i)} \cos(\omega_k t_i + \theta_k) \right] - y_i \right]^2 \quad (7.85)$$

and the errors for each method are

Method	error
Matrix pencil	0.1411
Levenberg–Marquardt	0.0373
Prony	3.9749

Not surprisingly, the Levenberg–Marquardt yielded the best results since it is an iterative method, whereas the other estimation methods are linear noniterative methods.

Power System Example

In this example, the accuracy of the methods will be compared using the dynamic response of a time-domain simulation of a large Midwestern utility system, shown in Figure 7.5. This simulation contains several hundred states comprising a wide range of responses. The number of dominant modes is not known. The results of a fast Fourier transform (FFT) are shown in Figure 7.6. From this figure, it appears as if there are five dominant modes that contribute significantly to the original waveform, with several of the modes concentrated at low frequencies. Therefore, the estimation methods introduced earlier will be applied to extract five modes.

Extracting five modes, the results are shown in Figure 7.7 and summarized:

Prony

mode	a_i	b_i	ω_i	θ_i
1	1.7406	-0.5020	3.7835	-1.4870
2	1.5723	-0.1143	4.8723	-1.1219
3	1.0504	-0.0156	6.2899	-0.0331
4	2.1710	-0.2455	7.7078	2.2011
5	0.9488	-0.3515	8.3854	-1.6184

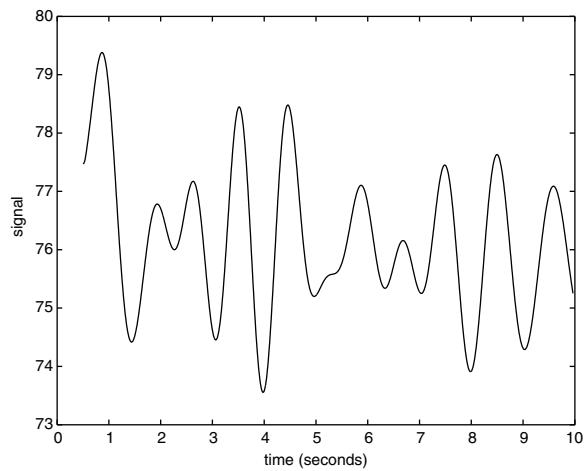


FIGURE 7.5
PSS/E waveform

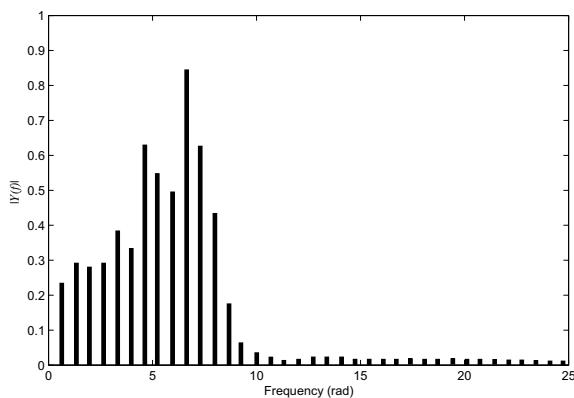
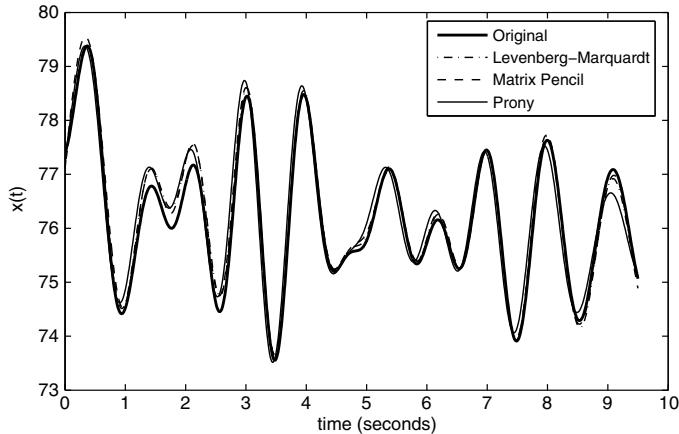


FIGURE 7.6
FFT of PSS/E waveform

**FIGURE 7.7**

Reconstruction of the PSS/E waveform using various methods

Matrix Pencil

mode	a_i	b_i	ω_i	θ_i
1	2.0317	-0.5610	3.7357	-1.5158
2	1.3204	-0.0774	4.8860	-1.3607
3	0.7035	0.0527	6.3030	-0.3093
4	1.2935	-0.2400	7.4175	2.9957
5	0.6718	-0.0826	8.0117	0.3790

Levenberg–Marquardt

mode	a_i	b_i	ω_i	θ_i
1	1.8604	-0.5297	3.6774	-1.3042
2	1.1953	-0.0578	4.8771	-1.3405
3	0.8164	0.0242	6.2904	-0.2537
4	1.9255	-0.2285	7.6294	2.1993
5	0.6527	-0.2114	8.3617	-1.5187

The error in each method as determined by Equation (7.85) and the relative computation times are

method	CPU	Error
Prony	0.068	228.16
Matrix pencil	39.98	74.81
Levenberg–Marquardt	42.82*	59.74

* Depends on initial condition

Note that the Prony method is the most computationally efficient since it only requires two least-squares solutions. The matrix pencil method is more computationally expensive since it requires a singular value decomposition of a relatively large matrix. It also requires an eigensolution, but, since the matrix itself is relatively small (the size of the number of required modes), it is not overly burdensome. Not surprisingly, the Levenberg–Marquardt method is the most computationally expensive since it is an iterative method. Its computational burden is directly related to the initial guess: the better the initial guess, the faster the method converges. It is wise to choose the initial guess as the parameters obtained from either the Prony or the matrix pencil methods.

Similarly, the level of error in each method varies with the complexity of the method. The Levenberg–Marquardt method yields the best results, but with the greatest computational effort. The Prony has the largest error, but this is offset by the relative speed of computation.

7.6 Power System Applications

7.6.1 Participation Factors

In the analysis of large-scale power systems, it is sometimes desirable to have a measure of the impact that a particular state has on a selected system mode (or eigenvalue). In some cases, it is desirable to know whether a set of physical states has influence over an oscillatory mode such that control of that component may mitigate the oscillations. Another use is to identify which system components contribute to an unstable mode. One tool for identifying which states significantly participate in a selected mode is the method of *participation factors* [62]. In large-scale power systems, participation factors can also be used to identify interarea oscillations versus those that persist only within localized regions (intraarea oscillations).

Participation factors provide a measure of the influence each dynamic state has on a given mode or eigenvalue. Consider a linear system

$$\dot{x} = Ax \quad (7.86)$$

The participation factor p_{ki} is a sensitivity measure of the i th eigenvalue to the (k, k) diagonal entry of the system A matrix. This is defined as

$$p_{ki} \frac{\partial \lambda_i}{\partial a_{kk}} \quad (7.87)$$

where λ_i is the i th eigenvalue and a_{kk} is the k th diagonal entry of A . The participation factor p_{ki} relates the k th state variable to the i th eigenvalue. An

equivalent, but more common expression for the participation factor is also defined as

$$p_{ki} = \frac{w_{ki}v_{ik}}{w_i^T v_i} \quad (7.88)$$

where w_{ki} and v_{ik} are the k th entries of the left and right eigenvectors associated with λ_i . As with eigenvectors, participation factors are frequently normalized to unity, such that

$$\sum_{k=1}^n p_{ki} = 1 \quad (7.89)$$

When the participation factors are normalized, they provide a straightforward measure of the percent of impact each state has on a particular mode. Participation factors for complex eigenvalues (and eigenvectors) are defined in terms of magnitudes, rather than complex quantities. In the case of complex eigenvalues, the participation factors are defined as

$$p_{ki} = \frac{|v_{ik}| |w_{ki}|}{\sum_{i=1}^n |v_{ik}| |w_{ki}|} \quad (7.90)$$

In some applications, it may be preferred to retain the complex nature of the participation factors to yield both phase and magnitude information [32].

7.7 Problems

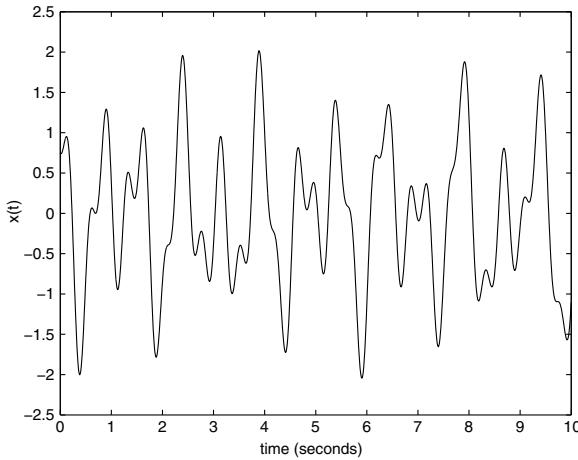
- Find the eigenvalues and eigenvectors of the following matrices.

$$A_1 = \begin{bmatrix} 5 & 4 & 1 & 1 \\ 4 & 5 & 1 & 1 \\ 1 & 1 & 4 & 2 \\ 1 & 1 & 2 & 4 \end{bmatrix}$$

$$A_2 = \begin{bmatrix} 2 & 3 & 4 \\ 7 & -1 & 3 \\ 1 & -1 & 5 \end{bmatrix}$$

- Find the complex eigenvalues of the follow matrix method.

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

**FIGURE 7.8**

Waveform for Problem 5

3. Using the implicitly restarted Arnoldi method, find the two largest eigenvalues and the corresponding eigenvectors of

$$A = \begin{bmatrix} 8 & 3 & 4 & 4 & 10 \\ 3 & 6 & 6 & 5 & 4 \\ 4 & 6 & 10 & 7 & 8 \\ 4 & 5 & 7 & 10 & 4 \\ 10 & 4 & 8 & 4 & 2 \end{bmatrix}$$

4. Find the singular values and singular vectors of the matrix:

$$A_0 = \begin{bmatrix} 13 & 4 & 7 & 6 & 20 \\ 6 & 3 & 20 & 9 & 7 \\ 9 & 8 & 19 & 11 & 15 \\ 1 & 4 & 2 & 19 & 14 \end{bmatrix}$$

5. Generate the waveform shown in Figure 7.8 on the interval $t \in [0, 10]$ with a time step of 0.01 seconds.

$$x(t) = \sum_{i=1}^3 a_i e^{b_i t} (\cos c_i t + d_i)$$

where

mode	a_i	b_i	c_i	d_i
1	1.0	-0.01	8.0	0.0
2	0.6	-0.03	17.0	π
3	0.5	0.04	4.7	$\pi/4$

- (a) Using 100 equidistant points on the interval [0, 10], estimate the six system eigenvalues using Prony analysis. How do these compare with the actual eigenvalues?
- (b) Using 100 equidistant points on the interval [0, 10], estimate the six system eigenvalues using Levenberg–Marquardt. How do these eigenvalues compare with the actual eigenvalues? With those obtained from the Prony analysis?
- (c) Using 100 equidistant points on the interval [0, 10], estimate the six system eigenvalues using the matrix pencil method. How do these eigenvalues compare with the actual eigenvalues? With those obtained from the Prony analysis?
- (d) Using all of the points, estimate the six system eigenvalues using Prony analysis. How do these compare with the actual eigenvalues?
- (e) Using all of the points, estimate the six system eigenvalues using Levenberg–Marquardt. How do these compare with the actual eigenvalues?
- (f) Using all of the points, estimate the six system eigenvalues using the matrix pencil method. How do these compare with the actual eigenvalues?
- (g) Using all of the points, estimate the two dominant modes (two complex eigenvalue pairs) of the system response using the matrix pencil method. Substitute the estimated parameters into

$$x(t) = \sum_{i=1}^2 a_i e^{b_i t} (\cos c_i t + d_i)$$

and plot this response versus the three-mode response. Discuss the differences and similarities.

References

- [1] R. A. M. van Amerongen, “A general-purpose version of the fast decoupled loadflow,” *IEEE Transactions on Power Systems*, vol. 4, no. 2, May 1989.
- [2] P. M. Anderson and A. A. Fouad, *Power System Control and Stability*. Ames, IA: Iowa State University Press, 1977.
- [3] W. E. Arnoldi, “The principle of minimized iterations in the solution of the matrix eigenvalue problem,” *Quart. Appl. Math.*, vol. 9, 1951.
- [4] M. S. Bazaraa, H. D. Sherali, and C. M. Shetty, *Nonlinear Programming: Theory and Algorithms*, 2nd ed., Hoboken, NJ: Wiley Press, 1993.
- [5] L. T. Biegler, T. F. Coleman, A. R. Conn, and F. N. Santosa, *Large-Scale Optimization with Applications – Part II: Optimal Design and Control*. New York: Springer-Verlag, Inc., 1997.
- [6] E. Boman and B. Hendrickson, “Support theory for preconditioning,” *SIAM J. Matrix Analysis*, vol. 25, no. 3, pp. 694–717, 2004.
- [7] K. E. Brenan, S. L. Campbell, and L. R. Petzold, *Numerical Solution of Initial-Value Problems in Differential-Algebraic Equations*. Philadelphia: Society for Industrial and Applied Mathematics, 1995.
- [8] L. O. Chua and P. Lin, *Computer Aided Analysis of Electronic Circuits: Algorithms and Computational Techniques*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1975.
- [9] G. Dahlquist and A. Bjorck, *Numerical Methods*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1974.
- [10] G. Dantzig, *Linear Programming: Introduction*. Secaucus, NJ: Springer-Verlag, Inc., 1997.
- [11] H. W. Dommel and W. F. Tinney, “Optimal power flow solutions,” *IEEE Transactions on Power Apparatus and Systems*, vol. 87, no. 10, pp. 1866 – 1874, October 1968.
- [12] S. Eisenstat, M. Gursky, M. Schultz, and A. Sherman, “The Yale Sparse Matrix Package I: The symmetric codes,” *International Journal of Numerical Methods Engineering*, vol. 18, 1982, pp. 1145 – 1151.

- [13] O. I. Elgerd, *Electric Energy System Theory, An Introduction*. New York: McGraw-Hill Book Company, 1982.
- [14] J. Francis, “The QR transformation: A unitary analogue to the LR Transformation,” *Comp. Journal*, vol. 4, 1961.
- [15] C. W. Gear, *Numerical Initial Value Problems in Ordinary Differential Equations*, Englewood Cliffs, NJ: Prentice-Hall, Inc., 1971.
- [16] C. W. Gear, “The simultaneous numerical solution of differential-algebraic equations,” *IEEE Transactions on Circuit Theory*, vol. 18, pp. 89 – 95, 1971.
- [17] C. W. Gear and L. R. Petzold, “ODE methods for the solution of differential/algebraic systems,” *SIAM Journal of Numerical Analysis*, vol. 21, no. 4, pp. 716 – 728, August 1984.
- [18] A. George and J. Liu, “A fast implementation of the minimum degree algorithm using quotient graphs,” *ACM Transactions on Mathematical Software*, vol. 6, no. 3, September 1980, pp. 337 – 358.
- [19] A. George and J. Liu, “The evaluation of the minimum degree ordering algorithm,” *SIAM Review*, vol. 31, March 1989, pp. 1 – 19.
- [20] H. Glavitsch and R. Bacher, “Optimal power flow algorithms,” *Control and Dynamic Systems*, vol. 41, part 1, *Analysis and Control System Techniques for Electric Power Systems*, New York: Academic Press, 1991.
- [21] G. H. Golub and C. F. Van Loan, *Matrix Computations*, Baltimore: Johns Hopkins University Press, 1983.
- [22] G. K. Gupta, C. W. Gear, and B. Leimkuhler, “Implementing linear multistep formulas for solving DAEs,” Report no. UIUCDCS-R-85-1205, University of Illinois, Urbana, IL, April 1985.
- [23] J. F. Hauer, C. J. Demeure, and L. L. Scharf, “Initial results in Prony analysis of power system response signals,” *IEEE Transactions on Power Systems*, vol. 5, no. 1, February 1990.
- [24] Y. Hua and T. Sarkar, “Generalized pencil-of-function method for extracting poles of an EM system from its transient response,” *IEEE Transactions on Antennas and Propagation*, vol 37, no. 2, February 1989.
- [25] M. Ilic and J. Zaborszky, *Dynamics and Control of Large Electric Power Systems*, New York: Wiley-Interscience, 2000.
- [26] D. Kahaner, C. Moler, and S. Nash, *Numerical Methods and Software*, Englewood Cliffs, NJ: Prentice-Hall, 1989.
- [27] C. T. Kelley, *Iterative Methods for Linear and Nonlinear Equations*, Frontiers in Applied Mathematics, Society for Industrial and Applied Mathematics (SIAM), vol 16, 1995.

- [28] J. A. Kelner, L. Orecchia, A. Sidford, and Z. A. Zhu, “A simple, combinatorial algorithm for solving SDD systems in nearly-linear time,” arXiv preprint arXiv:1301.6628, 2013.
- [29] R. Kress, *Numerical Analysis*, New York: Springer-Verlag, 1998.
- [30] I. Koutis, G. L. Miller, and R. Peng, “Approaching Optimality For Solving SDD Linear Systems,” *IEEE 51st Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 235 – 244, 2010.
- [31] V. N. Kublanovskaya, “On some algorithms for the solution of the complete eigenvalue problem,” *USSR Comp. Math. Phys.*, vol. 3, pp. 637 – 657, 1961.
- [32] P. Kundur, *Power System Stability and Control*. New York: McGraw-Hill, 1994.
- [33] J. Liu, “Modification of the minimum-degree algorithm by multiple elimination,” *ACM Transactions on Mathematical Software*, vol. 11, no. 2, June 1985, pp. 141 – 153.
- [34] B. M. Maggs, G. L. Miller, O. Parekh, R. Ravi, S. Leung, and M. Woo, “Solving Symmetric Diagonally-Dominant Systems by Preconditioning,” *Proc. IEEE 38th Annual Symposium on Foundations of Computer Science*, Miami Beach, FL, 1997.
- [35] H. Markowitz, “The elimination form of the inverse and its application to linear programming,” *Management Science*, vol. 3, 1957, pp. 255 – 269.
- [36] A. Monticelli, “Fast decoupled load flow: Hypothesis, derivations, and testing,” *IEEE Transactions on Power Systems*, vol. 5, no. 4, pp. 1425 – 1431, 1990.
- [37] J. Nanda, P. Bijwe, J. Henry, and V. Raju, “General purpose fast decoupled power flow,” *IEEE Proceedings-C*, vol. 139, no. 2, March 1992.
- [38] J. M. Ortega and W. C. Rheinboldt, *Iterative Solution of Nonlinear Equations in Several Variables*, San Diego: Academic Press, Inc., 1970.
- [39] A. F. Peterson, S. L. Ray, and R. Mittra, *Computational Methods for Electromagnetics*, New York: IEEE Press, 1997.
- [40] M. J. Quinn, *Designing Efficient Algorithms for Parallel Computers*, New York: McGraw-Hill Book Company, 1987.
- [41] Y. Saad and M. Schultz, “GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems,” *SIAM J. Sci. Stat. Comput.*, vol. 7, no. 3, July 1986.
- [42] Y. Saad, *Iterative Methods for Sparse Linear Systems*, Philadelphia: Society for Industrial and Applied Mathematics Press, 2003.

- [43] O. R. Saavedra, A. Garcia, and A. Monticelli, “The representation of shunt elements in fast decoupled power flows,” *IEEE Transactions on Power Systems*, vol. 9, no. 3, August 1994.
- [44] J. Sanchez-Gasca and J. Chow, “Performance comparison of three identification methods for the analysis of electromechanical oscillations,” *IEEE Transactions on Power Systems*, vol. 14, no. 3, August 1999.
- [45] J. Sanchez-Gasca, K. Clark, N. Miller, H. Okamoto, A. Kurita, and J. Chow, “Identifying linear models from time domain simulations,” *IEEE Computer Applications in Power*, April 1997.
- [46] T. Sarkar and O. Pereira, “Using the matrix pencil method to estimate the parameters of a sum of complex exponentials,” *IEEE Antennas and Propagation*, vol. 37, no. 1, February 1995.
- [47] P. W. Sauer and M. A. Pai, *Power System Dynamics and Stability*, Upper Saddle River, NJ: Prentice-Hall, 1998.
- [48] G. Soderlind, “DASP3—A program for the numerical integration of partitioned stiff ODEs and differential/algebraic systems,” Report TRITA-NA-8008, The Royal Institute of Technology, Stockholm, Sweden, 1980.
- [49] D. C. Sorensen, “Implicitly restarted Arnoldi/Lanzcos methods for large scale eigenvalue calculations,” in D. E. Keyes, A. Sameh, and V. Venkatakrishnan, editors, *Parallel Numerical Algorithms: Proceedings of an ICASE/LaRC Workshop, May 23–25, 1994, Hampton, VA*, Kluwer, 1995.
- [50] D. C. Sorensen, “Implicit application of polynomial filters in a k -step Arnoldi method,” *SIAM J. Mat. Anal. Appl.*, vol. 13, no. 1, 1992.
- [51] D. A. Spielman, N. Srivastava, “Graph sparsification by effective resistances,” *SIAM Journal on Computing*, vol. 40, no. 6, 2011, pp. 1913 – 1926.
- [52] P. A. Stark, *Introduction to Numerical Methods*, London, UK: The Macmillan Company, 1970.
- [53] B. Stott and O. Alsac, “Fast decoupled load flow,” *IEEE Transactions on Power Apparatus and Systems*, vol. 93, pp. 859 – 869, 1974.
- [54] G. Strang, *Linear Algebra and Its Applications*, San Diego: Harcourt Brace Javanonich, 1988.
- [55] W. Tinney and J. Walker, “Direct solutions of sparse network equations by optimally ordered triangular factorizations,” *Proceedings of the IEEE*, vol. 55, no. 11, November 1967, pp. 1801 – 1809.
- [56] L. Wang and A. Semlyen, “Application of sparse eigenvalue techniques to the small signal stability analysis of large power systems,” *IEEE Transactions on Power Systems*, vol. 5, no. 2, May 1990.

- [57] D. S. Watkins, *Fundamentals of Matrix Computations*. New York: John Wiley and Sons, 1991.
- [58] J. H. Wilkinson, *The Algebraic Eigenvalue Problem*. Oxford, England: Clarendon Press, 1965.
- [59] M. Yannakakis, “Computing the minimum fill-in is NP-complete,” *SIAM Journal of Algebraic Discrete Methods*, vol. 2, 1981, pp. 77 – 79.
- [60] T. Van Cutsem and C. Vournas, *Voltage Stability of Electric Power Systems*, Boston: Kluwer Academic Publishers, 1998.
- [61] R. S. Varga, *Matrix Iterative Analysis*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1962.
- [62] G. C. Verghese, I. J. Perez-Arriaga, and F. C. Schweppe, “Selective modal analysis with applications to electric power systems,” *IEEE Transactions on Power Systems*, vol. 101, pp. 3117 – 3134, Sept. 1982.
- [63] W. Zangwill and C. Garcia, *Pathways to Solutions, Fixed Points, and Equilibria*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1981.

Engineering – Electrical

“... presents a nonconventional approach to teach or understand power system analysis: mathematics first, then each topic is related to power system applications. ... This approach is ideal for researchers and graduate students, and can immediately lead them into the power system field. ... Algorithms, however sophisticated, are explained with clarity, along with numerical examples to help the reader get the point.”

—**Lingling Fan**, University of South Florida, Tampa, USA

“... an excellent combination of topics regarding computational aspects and numerical algorithms for power system analysis, operations, and control. ... very useful to teach on analysis techniques for large-scale energy systems.”

—**Hao Zhu**, University of Illinois, Urbana-Champaign, USA

“... an excellent textbook ... for a graduate-level course in electric power engineering. ... covers a broad range of topics related to computational methods for power systems. ... contains very good problems for students' homework. I highly recommend this book for graduate teaching in electric power.”

—**Fangxing Li**, University of Tennessee, Knoxville, USA

Computational Methods for Electric Power Systems introduces computational methods that form the basis of many analytical studies in power systems. The book provides the background for a number of widely used algorithms that underlie several commercial software packages, linking concepts to power system applications. By understanding the theory behind many of the algorithms, the reader can make better use of the software and make more informed decisions (e.g., choice of integration method and step size in simulation packages).

This **Third Edition** contains new material on preconditioners for linear iterative methods, Broyden's method, and Jacobian-free Newton–Krylov methods. It includes additional problems and examples, as well as updated examples on sparse lower-upper (LU) factorization. It also adds coverage of the eigensystem realization algorithm and the double-shift method for computing complex eigenvalues.



CRC Press

Taylor & Francis Group
an Informa business

www.crcpress.com

6000 Broken Sound Parkway, NW
Suite 300, Boca Raton, FL 33487
711 Third Avenue
New York, NY 10017
2 Park Square, Milton Park
Abingdon, Oxon OX14 4RN, UK

K25073

ISBN: 978-1-4987-1159-3



9 781498 711593