

Systems and Control

Stanislaw H. Żak

*School of Electrical and Computer Engineering
Purdue University*

New York Oxford
OXFORD UNIVERSITY PRESS
2003

Oxford University Press

Oxford New York

Auckland Bangkok Buenos Aires Cape Town Chennai
Dar es Salaam Delhi Hong Kong Istanbul Karachi Kolkata
Kuala Lumpur Madrid Melbourne Mexico City Mumbai
Nairobi São Paulo Shanghai Taipei Tokyo Toronto

Copyright © 2003 by Oxford University Press, Inc.

Published by Oxford University Press, Inc.
198 Madison Avenue, New York, New York, 10016
<http://www.oup-usa.org>

Oxford is a registered trademark of Oxford University Press

All rights reserved. No part of this publication may be reproduced,
stored in a retrieval system, or transmitted, in any form or by any means,
electronic, mechanical, photocopying, recording, or otherwise,
without the prior permission of Oxford University Press.

Library of Congress Cataloging-in-Publication Data

Žak, Stanislaw H.

Systems and Control / Stanislaw H. Žak.

p. cm.

Includes bibliographical references and index.

ISBN 0-19-515011-2

1. Linear control systems. I. Title.

TJ220 .Z35 2002

629.8'32—dc21

2002066792

Printing number: 9 8 7 6 5 4 3 2 1

Printed in the United States of America
on acid-free paper

To
J. M. J.

PREFACE

This book is about modeling, analysis, and control of dynamical systems. Its objective is to familiarize the reader with the basics of dynamical system theory while, at the same time, equipping him or her with the tools necessary for control system design. The emphasis is on the design in order to show how dynamical system theory fits into practical applications. Various methods of a controller design are discussed.

This text can be used in a variety of ways. It is suitable for a senior- or graduate-level course on systems and control. The text can also be used in an introductory course on nonlinear systems. Finally, it can be used in an introductory course to modern automatic control. No special background is necessary to read this text beyond basic differential equations and elements of linear algebra. The book is also suitable for self-study. To accomplish this goal, the required mathematical techniques and terminology are reviewed in the Appendix. Depending on the background, the reader may refer to the material in the Appendix as needed.

Whenever one wishes to analyze a problem arising in the “real” world, the first order of business is the construction of a mathematical model of the problem. In this book we discuss mathematical modeling of problems from mechanical and electrical engineering, as well as from physics and biology. Constructing models of dynamical systems using a linguistic description that employs IF–THEN rules is also considered. Two main types of dynamical systems are common in applications: those for which the time variable is discrete and those for which the time variable is continuous. When the time variable is discrete, the dynamics of such systems are often modeled using difference equations. In the case when the time is continuous, ordinary differential equations are frequently chosen for modeling purposes. Both types of models are considered in the book. Nonlinear systems are emphasized to acknowledge the critical role that nonlinear phenomena play in science and technology. The models presented are the ones that can be used to design controllers. These models are constructed from the control engineering point of view. Most of the time, engineers develop models of man-made objects. This is in contrast with physicists, who are more interested in modeling nature. Models in the last chapter of this book, however, are of interest to engineers, physicists, or biologists. This chapter is an invitation to control engineering of the future where control designers will be challenged to tackle highly complex, nonlinear, dynamical systems using multidisciplinary tools.

As can be concluded from the above remarks, the scope of the book is quite broad. This is in order to show the multidisciplinary role of nonlinear dynamics and control. In particular, neural networks, fuzzy systems, and genetic algorithms are presented and a self-contained introduction to chaotic systems is provided. Lyapunov’s stability theory and its extensions are used in the analysis of dynamical system models. The objective of the stability analysis is to determine the

system behavior without solving the differential, or difference, equations modeling the system. In fact, the Lyapunov theory is used as a unifying medium for different types of dynamical systems analyzed. In particular, optimal, fuzzy, sliding mode, and chaotic controllers are all constructed with the aid of the Lyapunov method and its variants. Furthermore, a class of neural networks is analyzed using Lyapunov's method. Classical methods and the techniques of post-modern control engineering, such as fuzzy logic, neural networks, and genetic algorithms, are presented in a unified fashion. It is shown how these new tools of a control engineer supplement the classical ones.

In addition to worked examples, this book also contains exercises at the end of each chapter. A number of solved problems, as well as the exercises, require the use of a computer. Although there are now many quality software packages available that could be used with this book, MATLAB* was chosen because it appears to be the most widely used math-tools program in engineering circles, and it is also popular among science students. Student editions of MATLAB and SIMULINK are sufficient for the worked examples and exercises in this book. The solutions manual, with complete solutions to all the exercises in the book, is available from the publisher to the instructors who adopt this book for their courses.

There is a companion web site for this book, where the current errata are available as well as the MATLAB m-files of examples solved in the text. The URL of the site is <http://dynamo.ecn.purdue.edu/~zak/systems/Index.html>.

Acknowledgments

I am indebted to my wife, Mary Ann, and to my children, Michaela, George, Elliot, and Nick, for their support during the writing of this book.

I express my gratitude to numerous students on whom this text was tested. I have benefited tremendously while working with Professor Stefen Hui. His insights and discussions helped clarify many concepts discussed in this book. Dr. Yonggon Lee contributed many simulations included in the book and in the solutions manual. It was a pleasant experience working with Peter Gordon, Justin Collins, and the staff of Oxford University Press during the production of the book. Finally, I wish to thank the National Science Foundation for the much appreciated support during the last stages of this project.

It is most unfortunate that I could not find a co-author, because now I must bear the final responsibility for all the errors and omissions. So dear Reader if you find an error, no matter how minor, please let me know about it. Errors found by the author or by readers are corrected every printing.

Stanislaw H. Żak

*MATLAB is a registered trademark of the MathWorks, Inc., 24 Prime Park Way, Natick, MA 01760-1500. Phone: 508-647-7000, fax: 508-647-7001, e-mail: info@mathworks.com, <http://www.mathworks.com>.

CONTENTS

Preface xv

1 Dynamical Systems and Modeling 1

- 1.1 What Is a System? 1
- 1.2 Open-Loop Versus Closed-Loop 3
- 1.3 Axiomatic Definition of a Dynamical System 5
- 1.4 Mathematical Modeling 8
- 1.5 Review of Work and Energy Concepts 9
- 1.6 The Lagrange Equations of Motion 12
- 1.7 Modeling Examples 21
 - 1.7.1 Centrifugal Governor 21
 - 1.7.2 Ground Vehicle 22
 - 1.7.3 Permanent Magnet Stepper Motor 28
 - 1.7.4 Stick Balancer 35
 - 1.7.5 Population Dynamics 38
- Notes 40
- Exercises 40

2 Analysis of Modeling Equations 48

- 2.1 State-Plane Analysis 48
 - 2.1.1 Examples of Phase Portraits 49
 - 2.1.2 The Method of Isoclines 53
- 2.2 Numerical Techniques 59
 - 2.2.1 The Method of Taylor Series 59
 - 2.2.2 Euler's Methods 62
 - 2.2.3 Predictor–Corrector Method 64

2.2.4	Runge's Method	67
2.2.5	Runge–Kutta Method	70
2.3	Principles of Linearization	71
2.4	Linearizing Differential Equations	74
2.5	Describing Function Method	77
2.5.1	Scalar Product of Functions	77
2.5.2	Fourier Series	81
2.5.3	Describing Function in the Analysis of Nonlinear Systems	82
	Notes	88
	Exercises	88

3 Linear Systems 94

3.1	Reachability and Controllability	94
3.2	Observability and Constructability	107
3.3	Companion Forms	111
3.3.1	Controller Form	112
3.3.2	Observer Form	117
3.4	Linear State-Feedback Control	120
3.5	State Estimators	126
3.5.1	Full-Order Estimator	127
3.5.2	Reduced-Order Estimator	130
3.6	Combined Controller–Estimator Compensator	131
	Notes	141
	Exercises	141

4 Stability 149

4.1	Informal Introduction to Stability	149
4.2	Basic Definitions of Stability	150
4.3	Stability of Linear Systems	154
4.4	Evaluating Quadratic Indices	159
4.5	Discrete-Time Lyapunov Equation	165
4.6	Constructing Robust Linear Controllers	168
4.7	Hurwitz and Routh Stability Criteria	171
4.8	Stability of Nonlinear Systems	176
4.9	Lyapunov's Indirect Method	188

4.10	Discontinuous Robust Controllers	194
4.11	Uniform Ultimate Boundedness	201
4.12	Lyapunov-Like Analysis	210
4.13	LaSalle's Invariance Principle	213
	Notes	217
	Exercises	217

5 Optimal Control 225

5.1	Performance Indices	225
5.2	A Glimpse at the Calculus of Variations	227
	5.2.1 Variation and Its Properties	229
	5.2.2 Euler–Lagrange Equation	234
5.3	Linear Quadratic Regulator	244
	5.3.1 Algebraic Riccati Equation (ARE)	244
	5.3.2 Solving the ARE Using the Eigenvector Method	248
	5.3.3 Optimal Systems with Prescribed Poles	251
	5.3.4 Optimal Saturating Controllers	266
	5.3.5 Linear Quadratic Regulator for Discrete Systems on an Infinite Time Interval	270
5.4	Dynamic Programming	273
	5.4.1 Discrete-Time Systems	273
	5.4.2 Discrete Linear Quadratic Regulator Problem	277
	5.4.3 Continuous Minimum Time Regulator Problem	280
	5.4.4 The Hamilton–Jacobi–Bellman Equation	283
5.5	Pontryagin's Minimum Principle	292
	5.5.1 Optimal Control with Constraints on Inputs	292
	5.5.2 A Two-Point Boundary-Value Problem	304
	Notes	307
	Exercises	307

6 Sliding Modes 315

6.1	Simple Variable Structure Systems	315
6.2	Sliding Mode Definition	320
6.3	A Simple Sliding Mode Controller	324
6.4	Sliding in Multi-Input Systems	327

6.5	Sliding Mode and System Zeros	328
6.6	Nonideal Sliding Mode	329
6.7	Sliding Surface Design	333
6.8	State Estimation of Uncertain Systems	336
6.8.1	Discontinuous Estimators	336
6.8.2	Boundary Layer Estimators	338
6.9	Sliding Modes in Solving Optimization Problems	340
6.9.1	Optimization Problem Statement	340
6.9.2	Penalty Function Method	342
6.9.3	Dynamical Gradient Circuit Analysis	343
	Notes	356
	Exercises	357

7 Vector Field Methods 367

7.1	A Nonlinear Plant Model	368
7.2	Controller Form	369
7.3	Linearizing State-Feedback Control	380
7.4	Observer Form	382
7.5	Asymptotic State Estimator	387
7.6	Combined Controller–Estimator Compensator	388
	Notes	389
	Exercises	389

8 Fuzzy Systems 393

8.1	Motivation and Basic Definitions	393
8.2	Fuzzy Arithmetic and Fuzzy Relations	396
8.2.1	Interval Arithmetic	396
8.2.2	Manipulating Fuzzy Numbers	398
8.2.3	Fuzzy Relations	403
8.2.4	Composition of Fuzzy Relations	405
8.3	Standard Additive Model	407
8.4	Fuzzy Logic Control	410
8.5	Stabilization Using Fuzzy Models	420
8.5.1	Fuzzy Modeling	421

8.5.2	Constructing a Fuzzy Design Model Using a Nonlinear Model	422
8.5.3	Stabilizability of Fuzzy Models	427
8.5.4	A Lyapunov-Based Stabilizer	431
8.6	Stability of Discrete Fuzzy Models	437
8.7	Fuzzy Estimator	440
8.7.1	The Comparison Method for Linear Systems	442
8.7.2	Stability Analysis of the Closed-Loop System	444
8.8	Adaptive Fuzzy Control	447
8.8.1	Plant Model and Control Objective	447
8.8.2	Background Results	447
8.8.3	Controllers	449
8.8.4	Examples	456
	Notes	463
	Exercises	464

9 Neural Networks 469

9.1	Threshold Logic Unit	470
9.2	Identification Using Adaptive Linear Element	478
9.3	Backpropagation	484
9.4	Neural Fuzzy Identifier	489
9.5	Radial-Basis Function (RBF) Networks	494
9.5.1	Interpolation Using RBF Networks	494
9.5.2	Identification of a Single-Input, Single-State System	497
9.5.3	Learning Algorithm for the RBF Identifier	498
9.5.4	Growing RBF Network	500
9.5.5	Identification of Multivariable Systems	503
9.6	A Self-Organizing Network	506
9.7	Hopfield Neural Network	508
9.7.1	Hopfield Neural Network Modeling and Analysis	508
9.7.2	Analog-to-Digital Converter	514
9.8	Hopfield Network Stability Analysis	517
9.8.1	Hopfield Network Model Analysis	518
9.8.2	Single-Neuron Stability Analysis	520
9.8.3	Stability Analysis of the Network	522

9.9	Brain-State-in-a-Box (BSB) Models	529
9.9.1	Associative Memories	529
9.9.2	Analysis of BSB Models	533
9.9.3	Synthesis of Neural Associative Memory	535
9.9.4	Learning	546
9.9.5	Forgetting	547
	Notes	549
	Exercises	551
10	Genetic and Evolutionary Algorithms	553
10.1	Genetics as an Inspiration for an Optimization Approach	553
10.2	Implementing a Canonical Genetic Algorithm	555
10.3	Analysis of the Canonical Genetic Algorithm	559
10.4	Simple Evolutionary Algorithm (EA)	561
10.5	Evolutionary Fuzzy Logic Controllers	563
10.5.1	Vehicle Model and Control Objective	563
10.5.2	Case 1: EA Tunes Fuzzy Rules	565
10.5.3	Case 2: EA Tunes Fuzzy Membership Functions	567
10.5.4	Case 3: EA Tunes Fuzzy Rules and Membership Functions	568
	Notes	574
	Exercises	575
11	Chaotic Systems and Fractals	578
11.1	Chaotic Systems Are Dynamical Systems with Wild Behavior	578
11.2	Chaotic Behavior of the Logistic Equation	579
11.2.1	The Logistic Equation—An Example from Ecology	579
11.2.2	Stability Analysis of the Logistic Map	581
11.2.3	Period Doubling to Chaos	588
11.2.4	The Feigenbaum Numbers	592
11.3	Fractals	593
11.3.1	The Mandelbrot Set	595
11.3.2	Julia Sets	595
11.3.3	Iterated Function Systems	596
11.4	Lyapunov Exponents	599
11.5	Discretization Chaos	605

11.6	Controlling Chaotic Systems	607
11.6.1	Ingredients of Chaotic Control Algorithm	607
11.6.2	Chaotic Control Algorithm	608
	Notes	619
	Exercises	620

Appendix: Math Review 624

A.1	Notation and Methods of Proof	624
A.2	Vectors	626
A.3	Matrices and Determinants	629
A.4	Quadratic Forms	633
A.5	The Kronecker Product	637
A.6	Upper and Lower Bounds	640
A.7	Sequences	641
A.8	Functions	648
A.9	Linear Operators	651
A.10	Vector Spaces	654
A.11	Least Squares	656
A.12	Contraction Maps	659
A.13	First-Order Differential Equation	662
A.14	Integral and Differential Inequalities	663
	A.14.1 The Bellman-Gronwall Lemma	663
	A.14.2 A Comparison Theorem	664
A.15	Solving the State Equations	665
	A.15.1 Solution of Uncontrolled System	665
	A.15.2 Solution of Controlled System	666
A.16	Curves and Surfaces	667
A.17	Vector Fields and Curve Integrals	670
A.18	Matrix Calculus Formulas	674
	Notes	675
	Exercises	675
	Bibliography	679
	Index	693



CHAPTER 1

Dynamical Systems and Modeling

To catch a mouse alive, try this trick: When you see the beast running around a room, open a black umbrella and hold it on the floor at an angle that allows the mouse to run in (it helps if someone else can chase the mouse). Your prey, drawn to dark places, will probably head straight for the umbrella. Then quickly close it, go outside, and set your prisoner free.

—*Reader's Digest Practical Problem Solver* [241]

1.1 What Is a System?

A *system* is a collection of interacting components. An electric motor, an airplane, and a biological unit such as the human arm are examples of systems. A system is characterized by two properties, which are as follows:

1. The interrelations between the components that are contained within the system
2. The system boundaries that separate the components within the system from the components outside

The system boundaries can be real or imagined. They are elastic in the sense that we may choose, at any stage of the system analysis, to consider only a part of the original system as a system on its own. We call it a *subsystem* of the original system. On the other hand, we may decide to expand the boundaries of the original system to include new components. In Figure 1.1, we represent a system's boundaries using a box. The interactions between the system components may be governed, for example, by physical, biological, or economical laws. In dealing with systems, we are interested in the effects of external quantities upon the behavior of the system quantities. We refer to the external quantities acting on the system as the *inputs* to the system. The condition or the state of the system is described by the *state variables*. The state variables provide the information that, together with the knowledge of the system inputs, enables us to determine the future state of the system. “A *dynamical system* consists of a set of possible states, together with a rule that determines the present state in terms of past states” [7, p. 1]. An axiomatic description of a dynamical system is presented in Section 1.3. In practice it is often not possible or too expensive to measure or determine the values of all of the state variables. Instead, only their subset or combination can be measured. The system quantities whose behavior can be measured or observed are referred to as the system *outputs*.

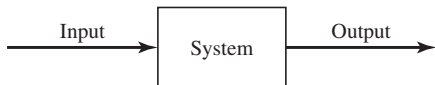


Figure 1.1 Representation of a system.

In engineering applications, when dealing with dynamical systems, we are interested in specifying the system inputs that force the system states or outputs to behave with time in some prespecified manner. That is, we are interested in *controlling* the system states or outputs. This is accomplished by means of a *controller* whose task is to produce the required system's inputs that in turn result in the desired system's outputs. An interconnection of the system and a controller is called a *control system*. In Figures 1.2 and 1.3, we show two different types of control systems that we discuss in the next section. Constructing a controller is a part of the *control problem*. The essential elements of the control problem, as described by Owens [218, p. 181] are as follows:

1. A specified objective for the system
2. A model of a dynamical system to be controlled
3. A set of admissible controllers
4. A means of measuring the performance of any given control strategy to evaluate its effectiveness

We now examine these elements one by one. The objective of a control system is to complete some specified task. This can be expressed as a combination of (a) constraints on the output or state variables and (b) limits on the time available to complete the control objective. For example, the objective of a controller might be to force the system output to settle within a certain percentage of the given value after prespecified time.

We now discuss the modeling part of the control problem. The first step in the controller design procedure is the construction of a *truth model* of the dynamics of the process to be controlled. The truth model is a *simulation model* that includes all the relevant characteristics of the process. The truth model is too complicated for use in the controller design. Thus, we need to develop a simplified model that can be used to design a controller. Such a simplified model is labeled by Friedland [91] as the *design model*. The design model should capture the essential features of the process. A common model of a dynamical system is the finite set of ordinary differential equations of the form

$$\begin{aligned}\dot{\mathbf{x}}(t) &= \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t)), & \mathbf{x}(t_0) &= \mathbf{x}_0, \\ \mathbf{y}(t) &= \mathbf{h}(t, \mathbf{x}(t), \mathbf{u}(t)),\end{aligned}$$

where the state $\mathbf{x} \in \mathbb{R}^n$, the input $\mathbf{u} \in \mathbb{R}^m$, the output $\mathbf{y} \in \mathbb{R}^p$, and \mathbf{f} and \mathbf{h} are vector-valued functions with $\mathbf{f} : \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ and $\mathbf{h} : \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^p$. Another common model of a dynamical system is the finite set of difference equations,

$$\begin{aligned}\mathbf{x}(k+1) &= \mathbf{f}(k, \mathbf{x}(k), \mathbf{u}(k)), & \mathbf{x}(k_0) &= \mathbf{x}_0, \\ \mathbf{y}(k) &= \mathbf{h}(k, \mathbf{x}(k), \mathbf{u}(k)),\end{aligned}$$

where $\mathbf{x}(k) = \mathbf{x}(kh)$, $\mathbf{u}(k) = \mathbf{u}(kh)$, h is the sampling interval, and $k \geq 0$ is an integer.

In our subsequent discussion, we regard the dynamical system and its model as equivalent. During the modeling process, properties of the system, physical constraints, and technical requirements should be taken into account. By itself the pair $(\mathbf{f}(t, \mathbf{x}, \mathbf{u}), \mathbf{h}(t, \mathbf{x}, \mathbf{u}))$ has little

engineering significance until constraints on the state \mathbf{x} , the control \mathbf{u} , and the output \mathbf{y} are also given. Control input signals are obtained from physical devices capable of providing only a limited amount of energy. We refer to the class of controllers that can be considered for the given control design problem as the set of admissible controllers.

We can evaluate the performance of any given control law by visual inspection of the transient characteristics of the controlled system after the design is complete. But this method of assessment of the effectiveness of the given control strategy is inaccurate. We therefore try to quantify the means of evaluating performance in the form of a performance index or cost functional at the beginning of the design process. We construct a numerical index of performance whose value reflects the quality of any admissible controller in accomplishing the system objective. This performance index is assumed to have the property that its numerical value decreases as the quality of the controller increases. This is the reason why we sometimes refer to the performance index as the penalty functional. The “best” controller then is the one that produces the smallest value of the performance index. We call an admissible controller that simultaneously ensures the completion of the system objective and the minimization of the performance index an *optimal controller*.

1.2 Open-Loop Versus Closed-Loop

We distinguish between two types of control systems. They are:

- Open-loop control systems
- Closed-loop control systems

An open-loop control system usually contains the following:

1. A process to be controlled, labeled *plant*
2. The controlling variable of the plant, called the *plant input*, or just *input* for short
3. The controlled variable of the plant, called the *plant output*, or just *output* for short
4. A *reference input*, which dictates the desired value of the output
5. A *controller* that acts upon the reference input in order to form the system input forcing the output behavior in accordance with the reference signal

Note that the plant and controller themselves can be considered as systems on their own. A schematic representation of an open-loop system is depicted in Figure 1.2. In an open-loop control system the output has no influence on the input or reference signal. The controller operates without taking into account the output. Thus, the plant input is formed with no influence of the output. A household appliance such as an iron is a simple example of an open-loop control system. In this example, we consider the iron itself as the plant or process to be controlled. Here we set the temperature control mechanism for the desired temperature. This setting—for example, very hot—is the reference input. We can view the potentiometer, inside the iron, as a controller, or a control actuator. It regulates the temperature of the iron by allowing the necessary amount of resistance to produce the desired temperature. The plant’s output is the temperature of the heated iron. Not every control system is an open-loop system. Another type of a control system is a closed-loop system. We can convert an open-loop system into a closed-loop system by adding the following components:

6. The *feedback loop* where the output signal is measured with a sensor and then the measured signal is fed back to the *summing junction*
7. The *summing junction*, where the measured output signal is subtracted from the reference (command) input signal in order to generate an *error signal*, also labeled as an *actuating signal*

A schematic of a closed-loop control system is shown in Figure 1.3. In a closed-loop system the error signal causes an appropriate action of the controller, which in turn “instructs” the plant to behave in a certain way in order to approach the desired output, as specified by the reference input signal. Thus, in the closed-loop system, the plant output information is fed back to the controller, and the controller then appropriately modifies the plant output behavior. A controller, also called a *compensator*, can be placed either in the *forward loop*, as in Figure 1.3, or in the *feedback loop*.

If the goal of the controller is to maintain the output at a constant value, then we have a *regulator* control system. A “centrifugal speed regulator,” or the “flyball governor,” commonly known as Watt’s governor, is an example of a regulator in a closed-loop system. A schematic of Watt’s governor is shown in Figure 1.4. The task of the centrifugal speed regulator is to control the setting of a throttle valve to maintain the desired engine speed. The nominal speed is set by adjusting the valve in the throttle. As the engine speeds up, the weights are thrown outwards, the throttle partially closes, and the engine slows down. As the engine slows down below its nominal speed, the throttle is opened up and the engine gains its speed. The operation of an engine equipped with a centrifugal speed regulator can be represented using the block diagram of Figure 1.3. The reference input to the closed-loop system is the desired engine speed. The desired engine speed is obtained by appropriately setting the valve. The system output is the actual engine speed. The flyball governor is our controller, while the engine represents the plant. The control signal generated by the controller is the setting of the throttle valve.

We mention at this juncture that Watt’s application of the centrifugal speed governor in 1788, to regulate the speed of the steam engine, marks the starting point for the development of automatic control—and, in particular, feedback control—as a science [192]. Before 1788,

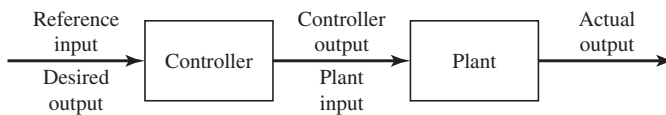


Figure 1.2 Open-loop control system.

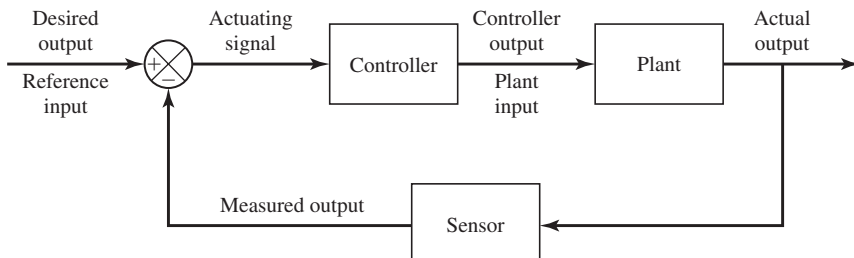


Figure 1.3 Closed-loop control system.

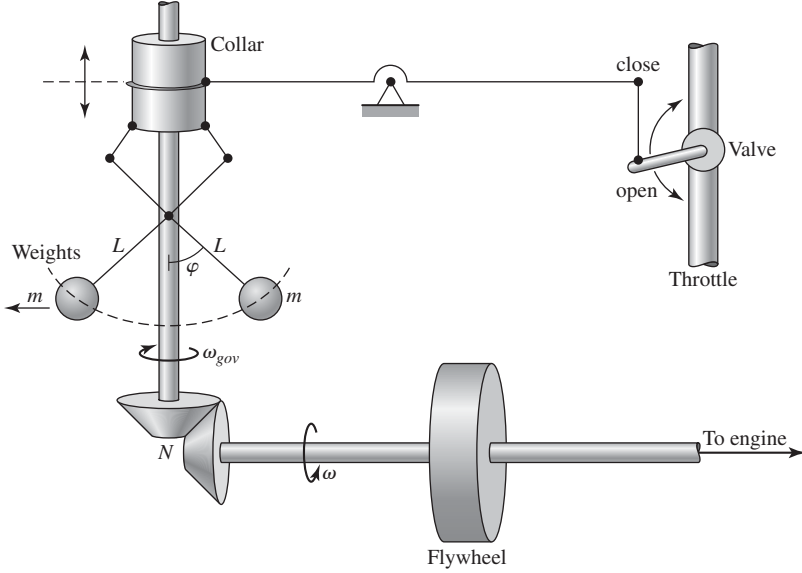


Figure 1.4 A simplified schematic of the flyball governor.

steam engines were controlled manually using the throttle valve. The flyball governor spread widely and quickly. It has been estimated that by 1868 there were about 75,000 flyball governors working in England alone [192]. The flyball governor is credited for the hatching of the industrial revolution. We devise a mathematical model of this revolutionary invention in Subsection 1.7.1, and in Section 4.7 we perform a quantitative analysis of its dynamical behavior.

If the goal of the controller is to force the output to follow a desired trajectory, then we have a *servomechanism* or *tracking* control system. As we mentioned before, to design a controller, we need first to analyze the plant quantitatively. The analysis requires a mathematical, or linguistic, description of the interrelations between the system quantities themselves as well as the interrelations between system quantities and system inputs. The mathematical description of these interrelations is called *mathematical modeling*. We devote Section 1.4 to mathematical modeling. Constructing models using linguistic description of the process to be controlled will be discussed in Subsection 8.5.1.

1.3 Axiomatic Definition of a Dynamical System

Following Kalman [147], we can define a dynamical system formally using the following axioms:

1. There is given a state space \mathcal{X} ; there is also T , an interval in the real line representing time.
2. There is given a space \mathcal{U} of functions on T that represent the inputs to the system.
3. For any initial time t_0 in T , any initial state \mathbf{x}_0 in \mathcal{X} , and any input \mathbf{u} in \mathcal{U} defined for $t \geq t_0$, the future states of the system are determined by the *transition mapping*

$$\Phi : T \times \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{X}$$

written as

$$\Phi(t_1; t_0, \mathbf{x}(t_0), \mathbf{u}(t)) = \mathbf{x}(t_1).$$

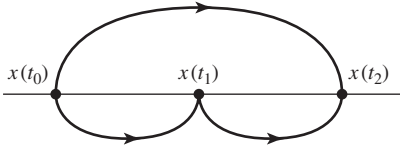


Figure 1.5 Illustration of the semigroup axiom.

Moreover, the following axioms are satisfied:

4. The identity property of the transition mapping, that is,

$$\Phi(t_0; t_0, \mathbf{x}(t_0), \mathbf{u}(t_0)) = \mathbf{x}(t_0).$$

5. The semigroup property of the transition mapping, that is,

$$\Phi(t_2; t_0, \mathbf{x}(t_0), \mathbf{u}(t)) = \Phi(t_2; t_1, \Phi(t_1; t_0, \mathbf{x}(t_0), \mathbf{u}(t)), \mathbf{u}(t)).$$

The semigroup property axiom states that it is irrelevant whether the system arrives at the state at time t_2 by a direct transition from the state at time t_0 , or by first going to an intermediate state at time t_1 , and then having been restarted from the state at time t_1 and moving to the state at time t_2 . In either case the system, satisfying the semigroup axiom, will arrive at the same state at time t_2 . This property is illustrated in Figure 1.5.

6. The causality property, that is,

$$\begin{aligned} \Phi(t; t_0, \mathbf{x}(t_0), \mathbf{u}_1(t)) &= \Phi(t; t_0, \mathbf{x}(t_0), \mathbf{u}_2(t)) && \text{for } t_0, t \in T \\ &\text{if and only if } \mathbf{u}_1(t) = \mathbf{u}_2(t) && \text{for all } t \in T. \end{aligned}$$

7. Every output of the system is a function of the form

$$\mathbf{h} : T \times \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{Y},$$

where \mathcal{Y} is the output space.

8. The transition mapping Φ and the output mapping \mathbf{h} are continuous functions.

9. If in addition the system satisfies the following axiom, then it is said to be *time-invariant*:

$$\Phi(t_1; t_0, \mathbf{x}(t_0), \mathbf{u}(t)) = \Phi(t_1 + \tau; t_0 + \tau, \mathbf{x}(t_0), \mathbf{u}(t)),$$

where $t_0, t_1 \in T$.

Thus, a dynamical system can be defined formally as a quintuple

$$\{T, \mathcal{X}, \mathcal{U}, \Phi, \mathcal{Y}\}$$

satisfying the above axioms.

Our attention will be focused on dynamical systems modeled by a set of ordinary differential equations

$$\dot{x}_i = f_i(t, x_1, x_2, \dots, x_n, u_1, u_2, \dots, u_m), \quad x_i(t_0) = x_{i0}, \quad i = 1, 2, \dots, n,$$

together with p functions

$$y_j = h_j(t, x_1, x_2, \dots, x_n, u_1, u_2, \dots, u_m), \quad j = 1, 2, \dots, p.$$

The system model state is $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^T \in \mathbb{R}^n$. The system input is $\mathbf{u} = [u_1 \ u_2 \ \dots \ u_m]^T \in \mathbb{R}^m$, and the system output is $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_p]^T \in \mathbb{R}^p$. In vector notation the above system

model has the form

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{f}(t, \mathbf{x}, \mathbf{u}), & \mathbf{x}(t_0) &= \mathbf{x}_0, \\ \mathbf{y} &= \mathbf{h}(t, \mathbf{x}, \mathbf{u}),\end{aligned}$$

where $\mathbf{f} : \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ and $\mathbf{h} : \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^p$ are vector-valued functions. In further considerations, we regard a dynamical system under considerations and its model represented by the above equations as equivalent. We now illustrate dynamical system axioms on a simple example.

◆ Example 1.1

We consider a dynamical system with no external inputs modeled by the first-order linear differential equation

$$\dot{x}(t) = ax(t), \quad x_0 = x(t_0), \quad (1.1)$$

where a is a constant, $x \in \mathcal{X} = \mathbb{R}$, and $T = \mathbb{R}$. The solution to the above differential equation has the form

$$x(t) = e^{a(t-t_0)} x(t_0).$$

The above system model can be characterized by the triple

$$\{T, \mathcal{X}, \Phi\},$$

where

$$\Phi : T \times \mathcal{X} \rightarrow \mathcal{X}$$

is the transition mapping. In this example,

$$\Phi(t; t_0, x(t_0)) = e^{a(t-t_0)} x(t_0).$$

It is clear that Φ is a continuous function. The transition function $\Phi = e^{a(t-t_0)} x(t_0)$ satisfies the semigroup property. Indeed,

$$\begin{aligned}\Phi(t_2; t_1, \Phi(t_1; t_0, x(t_0))) &= e^{a(t_2-t_1)} (e^{a(t_1-t_0)} x(t_0)) \\ &= e^{a(t_2-t_0)} x(t_0) \\ &= \Phi(t_2; t_0, x(t_0)).\end{aligned}$$

The identity property axiom is also satisfied:

$$\begin{aligned}\Phi(t_0; t_0, x(t_0)) &= e^{a(t_0-t_0)} x(t_0) \\ &= x(t_0).\end{aligned}$$

Finally, observe that the system model given by (1.1) is time invariant because

$$\begin{aligned}\Phi(t_1; t_0, x(t_0)) &= e^{a(t_1-t_0)} x(t_0) \\ &= e^{a(t_1+\tau-(t_0+\tau))} x(t_0) \\ &= \Phi(t_1 + \tau; t_0 + \tau, x(t_0)).\end{aligned}$$

1.4 Mathematical Modeling

The goal of mathematical modeling is to provide a mathematical description of the interrelations between the system quantities as well as the relations between these quantities and external inputs. A mathematical model can be viewed as “a mathematical representation of the significant, relevant aspects of an existing or proposed physical system. (Significance and relevance being in relation to an application where the model is to be used)” [181, p. xi]. Another definition of mathematical model that we found in reference 76, p. 1, is: “A mathematical model is a system of equations that purports to represent some phenomenon in a way that gives insight into the origins and consequences of its behavior.” (For the controller design purposes, we first construct a truth model of the given plant and then its design model.) In general, the more accurate the mathematical model is, the more complex the modeling equations are. This means that during the mathematical modeling process we may be forced to find a middle course between accuracy of the model and its complexity. Our goal is to generate the most accurate possible description of the given system while still being able to carry out the model analysis and computations in a reasonable amount of time. At first, we tend to make oversimplifying assumptions in order to obtain a very simple model. However, we should follow the principle of maximum simplicity, which states that “everything should be as simple as possible—but no simpler.” After extracting some information about the system behavior based on this first, possibly crude model, it may be desirable to construct a more accurate model. We try to create a mathematical description of a given physical process that incorporates realistic assumptions and physical constraints. We then extract information from this model and compare it with the experimental data. We perform this comparative study to determine the value of the obtained model. If the model appears to be inadequate for our needs, then it must be adjusted. The process of mathematical modeling is depicted in Figure 1.6. In the process of mathematical modeling, one can use different types of equations and their combinations like algebraic, ordinary and partial differential equations,

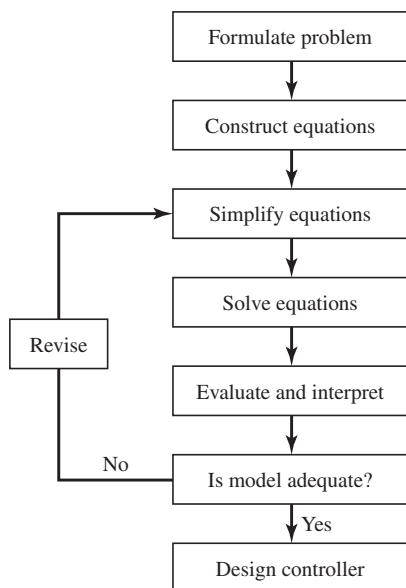


Figure 1.6 Diagram of mathematical modeling process and its use in the controller design process.

difference, integral, or functional equations. In this book, only difference and ordinary differential equations are used to model dynamical systems.

In the following, we discuss the Lagrange equations of motion that are useful in mathematical modeling of dynamical systems. We first review relevant concepts that are needed to derive the Lagrange equations. We then construct mathematical models of a few dynamical systems.

1.5 Review of Work and Energy Concepts

The purpose of this section is to review and acquaint ourselves with some basic notions from physics: work and energy. After a brief review of work and energy concepts, we derive the Lagrange equations of motion and use them to devise mathematical models of dynamical systems.

Suppose we are given a particle of a constant mass m subjected to a force \mathbf{F} . Then, by Newton's second law we have

$$\mathbf{F} = m \frac{d\mathbf{v}}{dt} = m\mathbf{a}. \quad (1.2)$$

Here the force \mathbf{F} and the velocity \mathbf{v} are vectors. Therefore, they can be represented as

$$\mathbf{F} = [F_{x_1} \quad F_{x_2} \quad F_{x_3}]^T \quad \text{and} \quad \mathbf{v} = [v_{x_1} \quad v_{x_2} \quad v_{x_3}]^T = [\dot{x}_1 \quad \dot{x}_2 \quad \dot{x}_3]^T.$$

Using the above notation, we represent equation (1.2) as

$$\begin{bmatrix} F_{x_1} \\ F_{x_2} \\ F_{x_3} \end{bmatrix} = m \begin{bmatrix} \frac{dv_{x_1}}{dt} \\ \frac{dv_{x_2}}{dt} \\ \frac{dv_{x_3}}{dt} \end{bmatrix}. \quad (1.3)$$

Suppose now that a force \mathbf{F} is acting on a particle located at a point A and the particle moves to a point B . The work δW done by \mathbf{F} along infinitesimally small distance $\delta \mathbf{s}$ is

$$\begin{aligned} \delta W &= \mathbf{F}^T \delta \mathbf{s} \\ &= [F_{x_1} \quad F_{x_2} \quad F_{x_3}] \begin{bmatrix} \delta x_1 \\ \delta x_2 \\ \delta x_3 \end{bmatrix} \\ &= F_{x_1} \delta x_1 + F_{x_2} \delta x_2 + F_{x_3} \delta x_3, \end{aligned}$$

where $\delta \mathbf{s} = [\delta x_1 \quad \delta x_2 \quad \delta x_3]^T$. The work W_{AB} done on the path from A to B is obtained by integrating the above equation:

$$W_{AB} = \int_A^B (F_{x_1} dx_1 + F_{x_2} dx_2 + F_{x_3} dx_3).$$

We now would like to establish a relation between work and kinetic energy. For this observe that

$$\ddot{x} dx = \frac{d\dot{x}}{dt} dx = \dot{x} d\dot{x}.$$

Using the above relation and Newton's second law, we express W_{AB} in a different way as

$$\begin{aligned}
 W_{AB} &= \int_A^B (m\ddot{x}_1 dx_1 + m\ddot{x}_2 dx_2 + m\ddot{x}_3 dx_3) \\
 &= \int_A^B m(\dot{x}_1 d\dot{x}_1 + \dot{x}_2 d\dot{x}_2 + \dot{x}_3 d\dot{x}_3) \\
 &= m \left(\frac{\dot{x}_1^2}{2} + \frac{\dot{x}_2^2}{2} + \frac{\dot{x}_3^2}{2} \right) \Big|_A^B \\
 &= m \left(\frac{\|\mathbf{v}_B\|^2}{2} - \frac{\|\mathbf{v}_A\|^2}{2} \right) \\
 &= \frac{m\|\mathbf{v}_B\|^2}{2} - \frac{m\|\mathbf{v}_A\|^2}{2}, \tag{1.4}
 \end{aligned}$$

where $\mathbf{v}_A = [\dot{x}_{1A} \ \dot{x}_{2A} \ \dot{x}_{3A}]^T$, \mathbf{v}_B is defined in an analogous manner, and $\|\mathbf{v}\| = \mathbf{v}^T \mathbf{v}$. Observing that $m\|\mathbf{v}_B\|^2/2$ is the kinetic energy of the particle at the point B and $m\|\mathbf{v}_A\|^2/2$ is its kinetic energy at the point A , we obtain

$$W_{AB} = K_B - K_A.$$

The above relation can be used to define the kinetic energy of a particle as the work required to change its velocity from some value \mathbf{v}_A to a final value \mathbf{v}_B . The relation $W_{AB} = K_B - K_A = \Delta K$ is also known as the *work–energy theorem for a particle*.

A force is conservative if the work done by the force on a particle that moves through any round trip is zero. In other words, a force is conservative if the work done by it on a particle that moves between two points depends only on these points and not on the path followed.

We now review the notion of *potential energy*, or the *energy of configuration*. Recall that if the kinetic energy K of a particle changes by ΔK , then the potential energy U must change by an equal but opposite amount so that the sum of the two changes is zero; that is,

$$\Delta K + \Delta U = 0. \tag{1.5}$$

This is equivalent to saying that any change in the kinetic energy of the particle is compensated for by an equal but opposite change in the potential energy U of the particle so that their sum remains constant; that is,

$$K + U = \text{constant}.$$

See Exercise 1.1 for a rigorous proof of the above relation. The potential energy of a particle represents a form of stored energy that can be recovered and converted into the kinetic energy. If we now use the work–energy theorem and equation (1.5), then we obtain

$$W = \Delta K = -\Delta U.$$

The work done by a conservative force depends only on the starting and the end points of motion and not on the path followed between them. Therefore, for motion in one dimension, we obtain

$$\Delta U = -W = -\int_{x_0}^x F(s) ds. \tag{1.6}$$

We thus can write

$$F(x) = -\frac{dU(x)}{dx}. \quad (1.7)$$

A generalization of equation (1.6) to motion in three dimensions yields

$$\Delta U = -\int_{x_{10}}^{x_1} F_{x_1} ds - \int_{x_{20}}^{x_2} F_{x_2} ds - \int_{x_{30}}^{x_3} F_{x_3} ds,$$

and a generalization of equation (1.7) to motion in three dimensions is

$$\mathbf{F}(\mathbf{x}) = -\left[\frac{\partial U}{\partial x_1} \quad \frac{\partial U}{\partial x_2} \quad \frac{\partial U}{\partial x_3}\right]^T = -\nabla U(\mathbf{x}), \quad (1.8)$$

where $\mathbf{x} = [x_1 \ x_2 \ x_3]^T$.

For a formal proof of the property of a conservative vector field that the work done by it on a particle that moves between two points depends only on these points and not on the path followed, see Exercise 1.2.

We now represent Newton's equation in an equivalent format that establishes a connection with the Lagrange equations of motion, which are discussed in the following section. To proceed, note that

$$K = \frac{m \|\dot{\mathbf{x}}\|^2}{2} = \frac{m \dot{\mathbf{x}}^T \dot{\mathbf{x}}}{2}. \quad (1.9)$$

Hence,

$$\frac{\partial K}{\partial \dot{x}_i} = m \dot{x}_i, \quad i = 1, 2, 3. \quad (1.10)$$

Combining (1.8) and (1.10), we represent Newton's equations as

$$\frac{d}{dt} \left(\frac{\partial K}{\partial \dot{x}_i} \right) + (\nabla U)_i = 0, \quad i = 1, 2, 3. \quad (1.11)$$

Let

$$L = K - U. \quad (1.12)$$

The function L defined above is called the *Lagrangian function* or just the *Lagrangian*. Note that

$$\frac{\partial L}{\partial \dot{x}_i} = \frac{\partial K}{\partial \dot{x}_i}, \quad (1.13)$$

and

$$\frac{\partial L}{\partial x_i} = -\frac{\partial U}{\partial x_i}. \quad (1.14)$$

Substituting (1.13) and (1.14) into (1.11), we obtain

$$\frac{d}{dt} \left(\frac{\partial L}{\partial \dot{x}_i} \right) - \frac{\partial L}{\partial x_i} = 0, \quad i = 1, 2, 3. \quad (1.15)$$

Equations (1.15) are called the Lagrange equations of motion, in Cartesian coordinates, for a single particle. They are just an equivalent representation of Newton's equations given by (1.3).

1.6 The Lagrange Equations of Motion

A single particle with no constraints is free to take up any position in the three-dimensional space. This is the reason we need three coordinates to describe its position. In the case of a system of p particles, there are $3p$ degrees of freedom. In general, if there are constraints imposed on the particles, then the number of the degrees of freedom is

$$n = 3p - \text{number of constraints.}$$

For example, if the system consists of only one particle in space and if the particle is restricted to move along a line, then $n = 3 - 2 = 1$. We can use many different coordinate systems to describe a location of a particle in space. We will use the letter q as a symbol for a coordinate regardless of its nature. A generalized coordinate may represent, for example, angular displacement or electric charge. This indicates that the results we are going to obtain will be of general nature and will not be restricted to mechanical systems. We refer to q as a generalized coordinate. Thus, if

$$x = x(q_1, q_2, q_3),$$

$$y = y(q_1, q_2, q_3),$$

$$z = z(q_1, q_2, q_3),$$

then

$$\dot{x} = \frac{\partial x}{\partial q_1} \dot{q}_1 + \frac{\partial x}{\partial q_2} \dot{q}_2 + \frac{\partial x}{\partial q_3} \dot{q}_3,$$

and so on. We assume that the functions $x, y, z \in \mathcal{C}^2$; that is, they are twice continuously differentiable, which implies that their second partial derivatives are equal. Suppose now that a particle is confined to move on a surface. Therefore, to describe the location of the particle, we need only two generalized coordinates. The position of the particle on a surface in a three-dimensional space is then described as

$$x = x(q_1, q_2),$$

$$y = y(q_1, q_2),$$

$$z = z(q_1, q_2),$$

and hence

$$\begin{aligned} \delta x &= \frac{\partial x}{\partial q_1} \delta q_1 + \frac{\partial x}{\partial q_2} \delta q_2, \\ \delta y &= \frac{\partial y}{\partial q_1} \delta q_1 + \frac{\partial y}{\partial q_2} \delta q_2, \\ \delta z &= \frac{\partial z}{\partial q_1} \delta q_1 + \frac{\partial z}{\partial q_2} \delta q_2. \end{aligned} \tag{1.16}$$

We know that the work δW done by the force \mathbf{F} is $\delta W = F_x \delta x + F_y \delta y + F_z \delta z$. On the other hand, this work is equal to the change in the kinetic energy. Therefore,

$$\boxed{m(\ddot{x} \delta x + \ddot{y} \delta y + \ddot{z} \delta z) = F_x \delta x + F_y \delta y + F_z \delta z} \tag{1.17}$$

The above equation is known as *D'Alembert's equation*.

If we have a system of particles acted upon by forces F_1, \dots, F_p and if each of the particles undergoes displacements $\delta s_1, \dots, \delta s_p$, then the total work done by the forces is

$$\delta W = \sum_{i=1}^p (F_{x_i} \delta x_i + F_{y_i} \delta y_i + F_{z_i} \delta z_i). \quad (1.18)$$

Substituting (1.16) into D'Alembert's equation and rearranging, we obtain

$$\begin{aligned} \delta W &= m \left(\ddot{x} \frac{\partial x}{\partial q_1} + \ddot{y} \frac{\partial y}{\partial q_1} + \ddot{z} \frac{\partial z}{\partial q_1} \right) \delta q_1 + m \left(\ddot{x} \frac{\partial x}{\partial q_2} + \ddot{y} \frac{\partial y}{\partial q_2} + \ddot{z} \frac{\partial z}{\partial q_2} \right) \delta q_2 \\ &= \left(F_x \frac{\partial x}{\partial q_1} + F_y \frac{\partial y}{\partial q_1} + F_z \frac{\partial z}{\partial q_1} \right) \delta q_1 + \left(F_x \frac{\partial x}{\partial q_2} + F_y \frac{\partial y}{\partial q_2} + F_z \frac{\partial z}{\partial q_2} \right) \delta q_2. \end{aligned} \quad (1.19)$$

Let

$$F_{q_1} = F_x \frac{\partial x}{\partial q_1} + F_y \frac{\partial y}{\partial q_1} + F_z \frac{\partial z}{\partial q_1} \quad \text{and} \quad F_{q_2} = F_x \frac{\partial x}{\partial q_2} + F_y \frac{\partial y}{\partial q_2} + F_z \frac{\partial z}{\partial q_2}.$$

Then, (1.19) can be represented as

$$\delta W = F_{q_1} \delta q_1 + F_{q_2} \delta q_2.$$

To proceed further, note that

$$\frac{d}{dt} \left(\dot{x} \frac{\partial x}{\partial q_1} \right) = \ddot{x} \frac{\partial x}{\partial q_1} + \dot{x} \frac{d}{dt} \left(\frac{\partial x}{\partial q_1} \right).$$

Hence

$$\boxed{\ddot{x} \frac{\partial x}{\partial q_1} = \frac{d}{dt} \left(\dot{x} \frac{\partial x}{\partial q_1} \right) - \dot{x} \frac{d}{dt} \left(\frac{\partial x}{\partial q_1} \right)} \quad (1.20)$$

The time derivative of $x = x(q_1, q_2)$ is

$$\dot{x} = \frac{\partial x}{\partial q_1} \dot{q}_1 + \frac{\partial x}{\partial q_2} \dot{q}_2. \quad (1.21)$$

Differentiating the above partially with respect to \dot{q}_1 yields

$$\boxed{\frac{\partial \dot{x}}{\partial \dot{q}_1} = \frac{\partial x}{\partial q_1}} \quad (1.22)$$

Finally, we will show that

$$\boxed{\frac{d}{dt} \left(\frac{\partial x}{\partial q_1} \right) = \frac{\partial \dot{x}}{\partial q_1}} \quad (1.23)$$

Indeed, because $x = x(q_1, q_2)$ and the partial derivative $\partial x / \partial q_1$ is in general a function of q_1 and q_2 , we have

$$\frac{d}{dt} \left(\frac{\partial x}{\partial q_1} \right) = \frac{\partial}{\partial q_1} \left(\frac{\partial x}{\partial q_1} \right) \dot{q}_1 + \frac{\partial}{\partial q_2} \left(\frac{\partial x}{\partial q_1} \right) \dot{q}_2. \quad (1.24)$$

Taking the partial derivative of (1.21) with respect to q_1 yields

$$\begin{aligned} \frac{\partial \dot{x}}{\partial q_1} &= \frac{\partial}{\partial q_1} \left(\frac{\partial x}{\partial q_1} \right) \dot{q}_1 + \frac{\partial}{\partial q_1} \left(\frac{\partial x}{\partial q_2} \right) \dot{q}_2 \\ &= \frac{\partial}{\partial q_1} \left(\frac{\partial x}{\partial q_1} \right) \dot{q}_1 + \frac{\partial}{\partial q_2} \left(\frac{\partial x}{\partial q_1} \right) \dot{q}_2, \end{aligned} \quad (1.25)$$

because by assumption $x, y, z \in C^2$ and hence

$$\frac{\partial}{\partial q_1} \left(\frac{\partial x}{\partial q_2} \right) = \frac{\partial}{\partial q_2} \left(\frac{\partial x}{\partial q_1} \right).$$

With the above in mind, comparing the right-hand sides of (1.24) and (1.25), we obtain (1.23). We now use (1.20), (1.22), and (1.23) to arrive at the Lagrange equation of motion for a single particle. We first substitute (1.22) and (1.23) into (1.20) to get

$$\ddot{x} \frac{\partial x}{\partial q_1} = \frac{d}{dt} \left(\dot{x} \frac{\partial x}{\partial \dot{q}_1} \right) - \dot{x} \frac{\partial \dot{x}}{\partial q_1}. \quad (1.26)$$

Note that

$$\frac{\partial}{\partial \dot{q}_1} \left(\frac{\dot{x}^2}{2} \right) = \dot{x} \frac{\partial \dot{x}}{\partial \dot{q}_1} \quad (1.27)$$

and

$$\frac{\partial}{\partial q_1} \left(\frac{\dot{x}^2}{2} \right) = \dot{x} \frac{\partial \dot{x}}{\partial q_1}. \quad (1.28)$$

Substituting (1.27) and (1.28) into (1.26) gives

$$\ddot{x} \frac{\partial x}{\partial q_1} = \frac{d}{dt} \left(\frac{\partial(\dot{x}^2/2)}{\partial \dot{q}_1} \right) - \frac{\partial(\dot{x}^2/2)}{\partial q_1}.$$

We can obtain similar expressions for y and z . Taking this into account, when $\delta q_2 = 0$, equation (1.19) becomes

$$\begin{aligned} \delta W &= \left(\frac{d}{dt} \left(\frac{\partial}{\partial \dot{q}_1} \frac{m(\dot{x}^2 + \dot{y}^2 + \dot{z}^2)}{2} \right) - \frac{\partial}{\partial q_1} \frac{m(\dot{x}^2 + \dot{y}^2 + \dot{z}^2)}{2} \right) \delta q_1 \\ &= \left(F_x \frac{\partial x}{\partial q_1} + F_y \frac{\partial y}{\partial q_1} + F_z \frac{\partial z}{\partial q_1} \right) \delta q_1 \\ &= F_{q_1} \delta q_1. \end{aligned} \quad (1.29)$$

Let $K = \frac{1}{2}m(\dot{x}^2 + \dot{y}^2 + \dot{z}^2)$ denote the kinetic energy of the particle. Then, we can represent (1.29) as

$$\boxed{\frac{d}{dt} \left(\frac{\partial K}{\partial \dot{q}_1} \right) - \frac{\partial K}{\partial q_1} = F_{q_1}} \quad (1.30)$$

The equation above is called the *Lagrange equation of motion* for the q_1 coordinate. Using the same arguments as above, we can derive the Lagrange equation of motion for the q_2 coordinate. In general there are as many the Lagrange equations of motion as there are degrees of freedom of the particle.

Recall that the expression $\delta W = F_x \delta x + F_y \delta y + F_z \delta z$ is a work done by a force $\mathbf{F} = [F_x \ F_y \ F_z]^T$ for a general displacement $\delta \mathbf{s} = [\delta x \ \delta y \ \delta z]^T$. In terms of generalized coordinates and when, say, $\delta q_2 = 0$, we have

$$\begin{aligned} \delta W &= \left(F_x \frac{\partial x}{\partial q_1} + F_y \frac{\partial y}{\partial q_1} + F_z \frac{\partial z}{\partial q_1} \right) \delta q_1 \\ &= F_{q_1} \delta q_1 \\ &= \delta W_{q_1}. \end{aligned}$$

A generalized force F_{q_r} is of such nature that the product $F_{q_r} \delta q_r$ is the work done by driving forces when q_r alone is changed by δq_r . We mention that a generalized force does not have to be a force in the usual sense. For example, if q_r is an angle, then F_{q_r} must be a torque in order that $F_{q_r} \delta q_r$ be work.

We now derive a more general version of the Lagrange equations of motion for a single particle. In our derivation we will need the notion of *nonconservative* forces, which we discuss next. First recall that we interpret the kinetic energy of a particle as its ability to do work by virtue of its motion. This motion is the result of conservative and nonconservative forces acting upon the particle. A force is nonconservative if the work done by the force on a particle that moves through any round trip is not zero. Thus a force is nonconservative if the work done by that force on a particle that moves between two points depends on the path taken between these points.

Let us now suppose that in addition to the conservative forces, a single nonconservative force due to friction acts on the particle. Denote by $\sum W_c$ the sum of the work done by the conservative forces, and denote by W_f the work done by friction. We then have

$$W_f + \sum W_c = \Delta K. \quad (1.31)$$

Thus, if a nonconservative force, like friction, acts on a particle, then the total mechanical energy is not constant, but changes by the amount of work done by a nonconservative force according to the relation (1.31). Let us now denote by E the final mechanical energy of the particle and denote by E_0 its initial mechanical energy. Then, we can write

$$E - E_0 = \Delta E = W_f.$$

The work done by friction, the nonconservative force, is always negative. Hence, the final mechanical energy E is less than the initial mechanical energy. The “lost” mechanical energy is transformed into *internal energy* U_{int} . In other words, the loss of mechanical energy equals the

gain in internal energy and we write

$$\Delta E + U_{int} = 0.$$

In general, we have

$$\Delta K + \sum \Delta U + U_{int} + (\text{change in other forms of energy}) = 0.$$

This means that the total energy, kinetic plus potential plus internal plus all other forms of energy, does not change. This is the *principle of the conservation of energy*. To proceed further, recall that any generalized force F_{q_i} acting on our particle can be expressed as

$$F_{q_i} = F_x \frac{\partial x}{\partial q_i} + F_y \frac{\partial y}{\partial q_i} + F_z \frac{\partial z}{\partial q_i}.$$

Assuming that the forces are conservative, we get

$$F_{q_i} = - \left(\frac{\partial U}{\partial x} \frac{\partial x}{\partial q_i} + \frac{\partial U}{\partial y} \frac{\partial y}{\partial q_i} + \frac{\partial U}{\partial z} \frac{\partial z}{\partial q_i} \right).$$

We thus have

$$F_{q_i} = - \frac{\partial U}{\partial q_i}.$$

Therefore, the Lagrange equations of motion for a single particle for conservative forces can be represented as

$$\frac{d}{dt} \left(\frac{\partial K}{\partial \dot{q}_i} \right) - \frac{\partial K}{\partial q_i} = - \frac{\partial U}{\partial q_i}.$$

We can rearrange the above equations as

$$\frac{d}{dt} \left(\frac{\partial K}{\partial \dot{q}_i} \right) - \frac{\partial}{\partial q_i} (K - U) = 0.$$

We now introduce the *Lagrangian function* L , defined as

$$L = K - U.$$

Note that because U is not a function of \dot{q}_i , we have

$$\frac{\partial K}{\partial \dot{q}_i} = \frac{\partial L}{\partial \dot{q}_i}.$$

Taking the above into account, we can represent the Lagrange equations of motion for conservative forces as

$$\frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}_i} \right) - \frac{\partial L}{\partial q_i} = 0.$$

If some of the forces acting on the particle are nonconservative, the Lagrange equations of motion for our particle will take the form

$$\boxed{\frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}_i} \right) - \frac{\partial L}{\partial q_i} = \tilde{F}_{q_i}} \quad (1.32)$$

where \tilde{F}_{q_i} are nonconservative forces acting on the particle.

◆ Example 1.2

Consider the simple pendulum shown in Figure 1.7. The simple pendulum is an idealized body consisting of a point mass M , suspended by a weightless inextensible cord of length l . The simple pendulum is an example of a one-degree-of-freedom system with a generalized coordinate being the angular displacement θ . The pendulum kinetic energy is

$$K = \frac{1}{2} M l^2 \dot{\theta}^2.$$

We assume that the potential energy of the pendulum is zero when the pendulum is at rest. Hence, its potential energy is

$$U = Mgl(1 - \cos \theta),$$

where $g = 10 \text{ m/sec}^2$ is the acceleration due to gravity. The Lagrangian function is

$$\begin{aligned} L &= K - U \\ &= \frac{1}{2} M l^2 \dot{\theta}^2 - Mgl(1 - \cos \theta) \\ &= 6\dot{\theta}^2 - 60(1 - \cos \theta). \end{aligned}$$

The Lagrange equation describing the pendulum motion is

$$\frac{d}{dt} \left(\frac{\partial L}{\partial \dot{\theta}} \right) - \frac{\partial L}{\partial \theta} = 0.$$

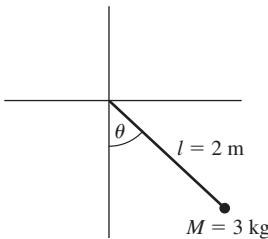


Figure 1.7 A simple pendulum.

In our example the above equation evaluates to

$$12\ddot{\theta} + 60 \sin \theta = 0,$$

or, equivalently,

$$\ddot{\theta} = -5 \sin \theta.$$

◆ Example 1.3

In Figure 1.8, the θ - r robot manipulator is shown. This example was adapted from Snyder [266, Chapter 10]. A lumped mass representation of the θ - r robot manipulator is shown in Figure 1.9. The mass $m_1 = 10$ kg represents the mass of the outer cylinder positioned at its center of mass. The constant distance $r_1 = 1$ m designates the fixed distance between the center of the mass of the outer cylinder and the center of rotation. The mass of the load is represented by $m_2 = 3$ kg and is assumed to be located at the end of a piston of a telescoping arm that is a variable radial distance r from the hub or center of rotation. The angle of rotation of the manipulator arm is θ . The inputs to the system are assumed to be (a) a torque T_θ

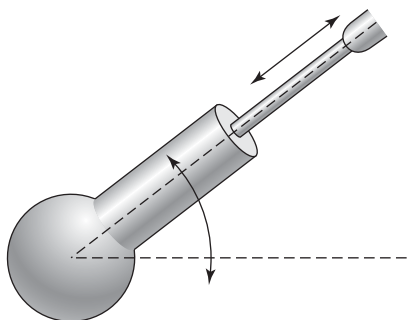


Figure 1.8 The θ - r robot manipulator.

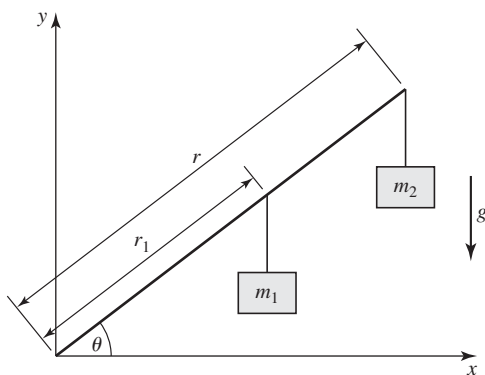


Figure 1.9 A lumped mass representation of the θ - r robot manipulator.

applied at the hub in the direction of θ and (b) a translational force F_r applied in the direction r . Finally, $g = 10 \text{ m/sec}^2$ is the gravitational constant. We construct the Lagrange equations of motion for the system, and then we represent the obtained model in a state-space format. We first find the total kinetic energy of the system. The position of the mass m_1 is

$$\begin{aligned}x_1 &= r_1 \cos \theta, \\y_1 &= r_1 \sin \theta.\end{aligned}$$

Note that r_1 is a constant. Differentiating x_1 and y_1 with respect to time yields

$$\begin{aligned}\dot{x}_1 &= -r_1 \dot{\theta} \sin \theta, \\ \dot{y}_1 &= r_1 \dot{\theta} \cos \theta.\end{aligned}$$

The magnitude squared of the velocity vector of the mass m_1 is

$$\begin{aligned}v_1^2 &= \dot{x}_1^2 + \dot{y}_1^2 \\ &= r_1^2 \dot{\theta}^2 \sin^2 \theta + r_1^2 \dot{\theta}^2 \cos^2 \theta \\ &= r_1^2 \dot{\theta}^2.\end{aligned}$$

Therefore, the kinetic energy of the mass m_1 is

$$K_1 = \frac{1}{2} m_1 v_1^2 = \frac{1}{2} m_1 r_1^2 \dot{\theta}^2.$$

We now derive an expression for the kinetic energy of the second mass. The position of the mass m_2 is

$$\begin{aligned}x_2 &= r \cos \theta, \\ y_2 &= r \sin \theta.\end{aligned}$$

Note that the distance r is not a constant. Differentiating x_2 and y_2 with respect to time yields

$$\begin{aligned}\dot{x}_2 &= \dot{r} \cos \theta - r \dot{\theta} \sin \theta, \\ \dot{y}_2 &= \dot{r} \sin \theta + r \dot{\theta} \cos \theta.\end{aligned}$$

The magnitude squared of the velocity vector of the mass m_2 is

$$\begin{aligned}v_2^2 &= (\dot{x}_2)^2 + (\dot{y}_2)^2 \\ &= (\dot{r} \cos \theta - r \dot{\theta} \sin \theta)^2 + (\dot{r} \sin \theta + r \dot{\theta} \cos \theta)^2 \\ &= \dot{r}^2 + r^2 \dot{\theta}^2.\end{aligned}$$

Hence, the kinetic energy of the mass m_2 is

$$K_2 = \frac{1}{2} m_2 v_2^2 = \frac{1}{2} m_2 (\dot{r}^2 + r^2 \dot{\theta}^2).$$

The total kinetic energy of the system is

$$K = K_1 + K_2 = \frac{1}{2} m_1 r_1^2 \dot{\theta}^2 + \frac{1}{2} m_2 (\dot{r}^2 + r^2 \dot{\theta}^2).$$

We now find the potential energy of the system. The potential energy of the mass m_1 is

$$U_1 = m_1 g r_1 \sin \theta.$$

The potential energy of the mass m_2 is

$$U_2 = m_2 g r \sin \theta.$$

The total potential energy of the system is

$$U = U_1 + U_2 = m_1 g r_1 \sin \theta + m_2 g r \sin \theta.$$

The Lagrangian function for the system is

$$L = K - U = \frac{1}{2} m_1 r_1^2 \dot{\theta}^2 + \frac{1}{2} m_2 \dot{r}^2 + \frac{1}{2} m_2 r^2 \dot{\theta}^2 - m_1 g r_1 \sin \theta - m_2 g r \sin \theta.$$

The manipulator has two degrees of freedom. Hence, we have two Lagrange equations of motion for this system:

$$\begin{aligned} \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{\theta}} \right) - \frac{\partial L}{\partial \theta} &= T_\theta, \\ \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{r}} \right) - \frac{\partial L}{\partial r} &= F_r. \end{aligned}$$

The first Lagrange equation yields

$$m_1 r_1^2 \ddot{\theta} + m_2 r^2 \ddot{\theta} + 2 m_2 r \dot{r} \dot{\theta} + g \cos \theta (m_1 r_1 + m_2 r) = T_\theta.$$

The second Lagrange equation gives

$$m_2 \ddot{r} - m_2 r \dot{\theta}^2 + m_2 g \sin \theta = F_r.$$

We next represent the above modeling equations in state-space format. (One can use, for example, MATLAB's Symbolic Math Toolbox to perform needed calculations.) Let the state variables be defined as

$$z_1 = \theta, \quad z_2 = \dot{\theta}, \quad z_3 = r, \quad z_4 = \dot{r}.$$

Then, we can write

$$\begin{bmatrix} \dot{z}_1 \\ \dot{z}_2 \\ \dot{z}_3 \\ \dot{z}_4 \end{bmatrix} = \begin{bmatrix} z_2 \\ \frac{-2m_2 z_2 z_3 z_4 - g \cos z_1 (m_1 r_1 + m_2 z_3) + u_1}{m_1 r_1^2 + m_2 z_3^2} \\ z_4 \\ \frac{z_2^2 z_3 - g \sin z_1 + u_2 / m_2}{z_3} \end{bmatrix}, \quad (1.33)$$

where $u_1 = T_\theta$ and $u_2 = F_r$.

In the following section, we develop models of a few more dynamical systems.

1.7 Modeling Examples

1.7.1 Centrifugal Governor

In this subsection, we model the centrifugal governor. In our modeling, we follow Pontryagin [235, pp. 213–220]. For an alternative approach to modeling the governor see Maxwell [197] and Fuller [92]. Some—for example, MacFarlane [192, p. 250]—claim that Maxwell founded the theory of automatic control systems in his 1868 paper “On governors.” However, this landmark paper is written in a terse and cryptic style. Fuller [92] does a wonderful job elucidating, correcting, and expanding Maxwell’s paper. Mayr [199, Chapter X] gives a fascinating historical account of the application of the centrifugal governor to the speed regulation of the steam engine and tracks the origins of the flyball governor.

A schematic of the governor is depicted in Figure 1.4, while forces acting on the governor are shown in Figure 1.10. For simplicity of notation, we assume that the length of the arms, on which the weights are fastened, are of unity length. If the weights rotate with angular velocity ω_{gov} , then the arms form an angle φ with the vertical position. Let v_{weight} denote the linear velocity of the weights, and let r be the radius of the circle along which the weights rotate. Then, $v_{weight} = \omega_{gov}r$. Note that $r = \sin \varphi$. The resulting centrifugal force acting on each weight is

$$\begin{aligned} m \frac{v_{weight}^2}{r} &= m \omega_{gov}^2 r \\ &= m \omega_{gov}^2 \sin \varphi. \end{aligned}$$

Simultaneously, a gravitational force mg acts on each weight. In addition, a frictional force in the hinge joints acts on each weight. We assume that the frictional force is proportional to $\dot{\varphi}$

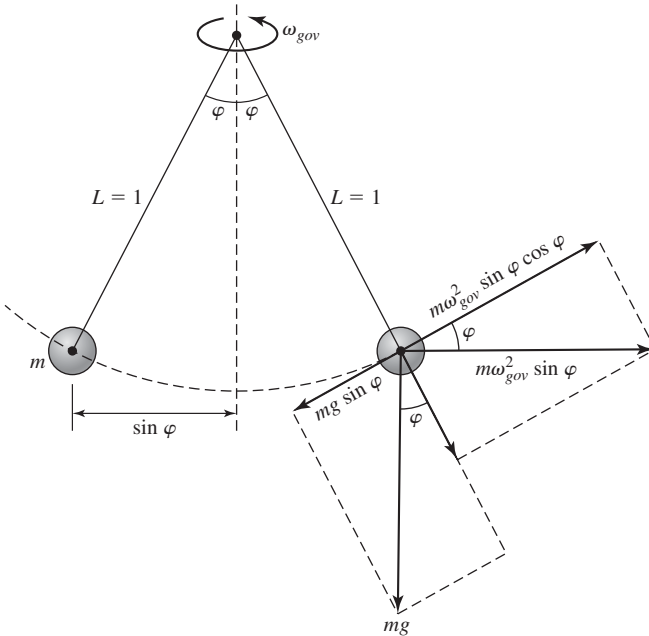


Figure 1.10 Forces acting on the weights in the governor.

with the coefficient of friction b . Summing forces along the axis perpendicular to an arm yields an equation of motion of the weight:

$$\begin{aligned} m\ddot{\varphi}L &= m\omega_{gov}^2 \sin \varphi \cos \varphi - mg \sin \varphi - b\dot{\varphi} \\ &= \frac{1}{2}m\omega_{gov}^2 \sin 2\varphi - mg \sin \varphi - b\dot{\varphi}. \end{aligned} \quad (1.34)$$

Note that $m\ddot{\varphi}L = m\ddot{\varphi}$ because by assumption, $L = 1$. The governor is connected to the engine via a gear with the gear ratio N such that

$$\omega_{gov} = N\omega, \quad (1.35)$$

where ω is the angular velocity of the flywheel. Substituting (1.35) into (1.34) gives

$$m\ddot{\varphi} = \frac{1}{2}mN^2\omega^2 \sin 2\varphi - mg \sin \varphi - b\dot{\varphi}.$$

Let I be the moment of inertia of the flywheel, τ_1 the torque delivered by the engine, and τ_l the load torque. Then, the equation modeling the motion of the flywheel is

$$I\dot{\omega} = \tau_1 - \tau_l. \quad (1.36)$$

The collar of the governor is connected with the valve so that at the nominal value of $\varphi = \varphi_0$ the engine delivers the nominal torque τ_d that results in the desired angular velocity $\omega = \omega_d$ of the flywheel. We model the interaction of the collar and the valve using the equation

$$\tau_1 = \tau_d + \kappa(\cos \varphi - \cos \varphi_0), \quad (1.37)$$

where $\kappa > 0$ is a proportionality constant that depends on the type of the governor. Substituting (1.37) into (1.36) yields

$$\begin{aligned} I\dot{\omega} &= \tau_1 - \tau_l \\ &= \tau_d + \kappa(\cos \varphi - \cos \varphi_0) - \tau_l \\ &= \kappa \cos \varphi - \tau, \end{aligned}$$

where $\tau = \tau_l - \tau_d + \kappa \cos \varphi_0$. We thus have two equations that model the governor:

$$\begin{aligned} m\ddot{\varphi} &= \frac{1}{2}mN^2\omega^2 \sin 2\varphi - mg \sin \varphi - b\dot{\varphi}, \\ I\dot{\omega} &= \kappa \cos \varphi - \tau. \end{aligned}$$

We use the state variables

$$x_1 = \varphi, \quad x_2 = \dot{\varphi}, \quad \text{and} \quad x_3 = \omega$$

to represent the above model in state-space format, i.e., as a system of first-order differential equations:

$$\begin{aligned} \dot{x}_1 &= x_2, \\ \dot{x}_2 &= \frac{1}{2}N^2x_3^2 \sin 2x_1 - g \sin x_1 - \frac{b}{m}x_2, \\ \dot{x}_3 &= \frac{\kappa}{I} \cos x_1 - \frac{\tau}{I}. \end{aligned} \quad (1.38)$$

1.7.2 Ground Vehicle

In this subsection, we model a ground vehicle in a turning maneuver, such as a step lane-change maneuver. Our modeling process follows that of Wong [303] and of Will and Žak [299]. During

such a step lane-change maneuver, the vehicle undergoes translational as well as rotational motion. When modeling a moving vehicle, it is convenient to use a reference coordinate frame attached to and moving with the vehicle. This is because with respect to the vehicle coordinate frame, moments of inertia of the vehicle are constant. We use the Society of Automotive Engineers (SAE) standard coordinate system as shown in Figure 1.11. We use lowercase letters to label axes of the vehicle coordinate system. Following the SAE convention, we define the coordinates as follows:

- The x coordinate is directed forward and on the longitudinal plane of symmetry of the vehicle.
- The y coordinate has a lateral direction out the right-hand side of the vehicle.
- The z coordinate is directed downward with respect to the vehicle.
- The yaw velocity about the z axis is denoted $\dot{\theta}$.

In developing our simple model, we neglect the effects of roll and pitch. The origin of the vehicle coordinate system is located at the vehicle center of gravity (CG). We also use a global, or fixed, coordinate system. Its axes are labeled using uppercase letters as shown in Figure 1.12. To form

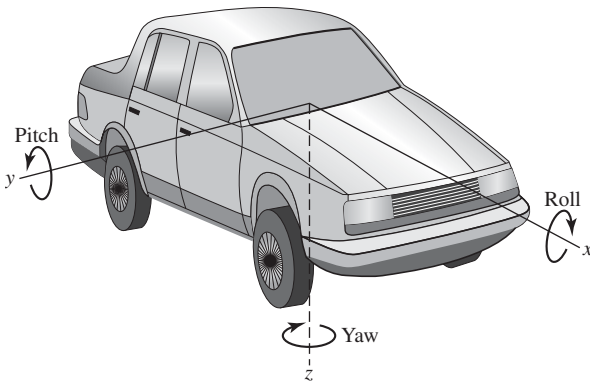


Figure 1.11 The Society of Automotive Engineers (SAE) coordinate system attached to the vehicle.

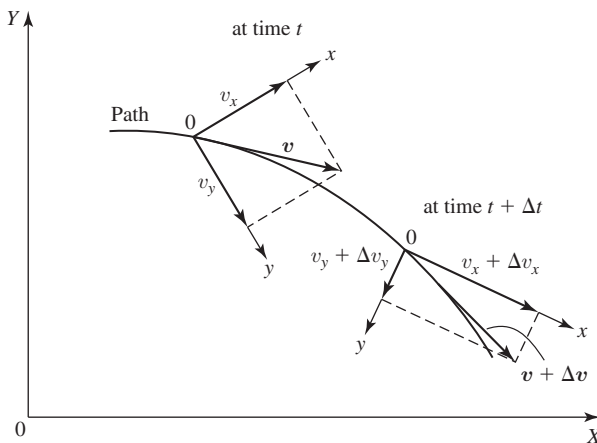


Figure 1.12 Analyzing vehicle motion using coordinate frame attached to the vehicle; snapshots of the vehicle velocity at times t and $t + \Delta t$.

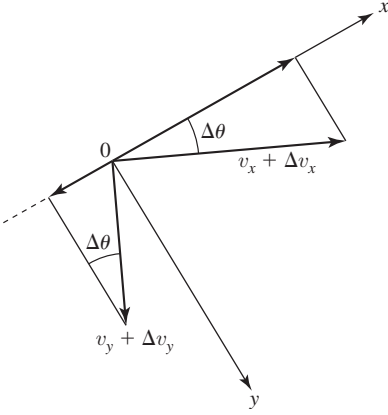


Figure 1.13 Computing the change in the longitudinal and lateral components of the vehicle velocity.

the equations of motion of the vehicle, we need an expression for the vehicle CG acceleration. We first compute the change in the longitudinal component of the velocity. Using Figures 1.12 and 1.13 and assuming small angles, we obtain

$$\begin{aligned} (v_x + \Delta v_x) \cos \Delta\theta - (v_y + \Delta v_y) \sin \Delta\theta - v_x &\approx \Delta v_x - v_y \Delta\theta \\ &= \Delta \dot{x} - \dot{y} \Delta\theta. \end{aligned}$$

We form the Newton quotient using the above. Then, taking the limit, we obtain the total longitudinal acceleration component a_x :

$$\begin{aligned} a_x &= \lim_{\Delta t \rightarrow 0} \frac{\Delta v_x - v_y \Delta\theta}{\Delta t} \\ &= \dot{v}_x - v_y \dot{\theta} \\ &= \ddot{x} - \dot{y} \omega_z, \end{aligned} \tag{1.39}$$

where $\omega_z = \dot{\theta}$. Using an argument similar to the one above, we obtain the lateral acceleration component a_y to be

$$\begin{aligned} a_y &= \dot{v}_y + v_x \dot{\theta} \\ &= \ddot{y} + \dot{x} \omega_z. \end{aligned} \tag{1.40}$$

We assume that the vehicle is symmetric with respect to the x - z plane. We use the so-called bicycle model that is shown in Figure 1.14, where δ_f denotes the front wheels' steering angle and δ_r is the rear wheels' steering angle. We denote by $2F_{xf}$ lumped longitudinal force at the front axle and denote by $2F_{xr}$ lumped longitudinal force at the rear axle. The lumped lateral force at the front axle is denoted $2F_{yf}$, while lumped lateral force at the rear axle is denoted $2F_{yr}$. Summing the forces along the x axis and assuming small angles, we obtain

$$\begin{aligned} m(\ddot{x} - \dot{y} \omega_z) &= 2F_{xf} \cos \delta_f - 2F_{yf} \sin \delta_f + 2F_{xr} \cos \delta_r + 2F_{yr} \sin \delta_r \\ &\approx 2F_{xf} + 2F_{xr}. \end{aligned} \tag{1.41}$$

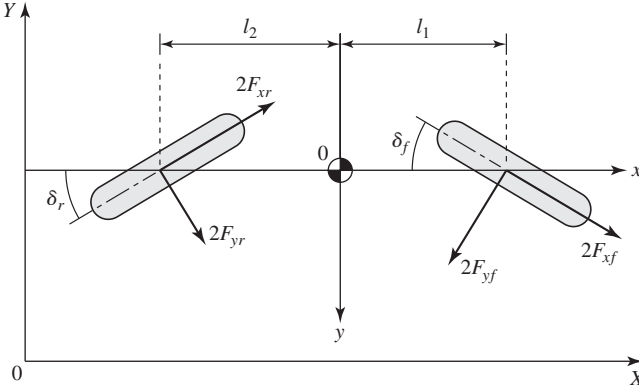


Figure 1.14 Bicycle model of a vehicle.

Next, summing the forces along the y axis and assuming small angles yields

$$\begin{aligned} m(\ddot{y} + \dot{x}\omega_z) &= 2F_{xf} \sin \delta_f + 2F_{yf} \cos \delta_f + 2F_{yr} \cos \delta_r - 2F_{xr} \sin \delta_r \\ &\approx 2F_{yf} + 2F_{yr}. \end{aligned} \quad (1.42)$$

The third equation of motion is obtained by summing the torques with respect to the z axis. Let I_z denote the moment of inertia of the vehicle about the z axis. Then, referring again to Figure 1.14 and assuming small angles, we obtain

$$\begin{aligned} I_z \dot{\omega}_z &= 2l_1 F_{yf} \cos \delta_f + 2l_1 F_{xf} \sin \delta_f - 2l_2 F_{yr} \cos \delta_r + 2l_2 F_{xr} \sin \delta_r \\ &\approx 2l_1 F_{yf} - 2l_2 F_{yr}. \end{aligned} \quad (1.43)$$

In our further analysis, we assume $v_x = \dot{x} = \text{constant}$; that is, the driver is neither accelerating nor braking the vehicle. This means that we need to use only the last two equations of motion in our model. To proceed, we need to define tire slip angles and cornering stiffness coefficients. The tire slip angle is the difference between the tire's heading direction, which is the steer angle, and the tire's actual travel direction. In other words, the tire slip angle is the angle between the desired direction of motion of the tire specified by the steer angle and the actual direction of travel of the center of the tire contact. We denote by $C_{\alpha f}$ the cornering stiffness coefficient of the front tires and denote by $C_{\alpha r}$ the cornering stiffness coefficient of the rear tires. The coefficients $C_{\alpha f}$ and $C_{\alpha r}$ are measured experimentally. A typical relation between a tire slip angle and a cornering force is shown in Figure 1.15. It is a nonlinear relation, where the cornering coefficient is

$$C_{\alpha w} = \frac{dF_{yw}}{d\alpha_w},$$

where $w = f$ or $w = r$. The lateral forces F_{yf} and F_{yr} that appear in (1.41)–(1.43) are functions of the corresponding tire slip angles α_f and α_r and cornering stiffness coefficients $C_{\alpha f}$ and $C_{\alpha r}$. In our simple linearized model, we assume that $C_{\alpha f} = C_{\alpha r} = \text{constant}$, and model dependencies between the lateral forces and tire slip angles are expressed as

$$F_{yf} = C_{\alpha f} \alpha_f \quad (1.44)$$

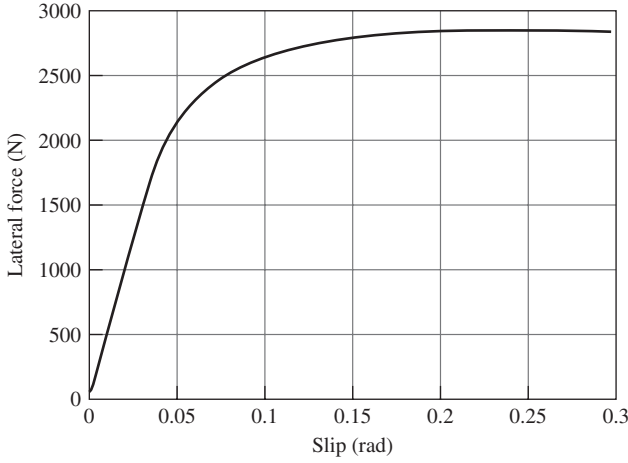


Figure 1.15 Typical cornering characteristic.

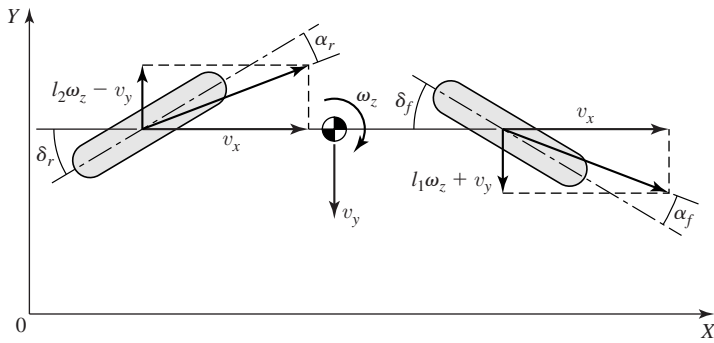


Figure 1.16 Computing tire slip angles.

and

$$F_{yr} = C_{ar}\alpha_r. \quad (1.45)$$

We formulate expressions for the tire slip angles using Figure 1.16. Assuming small angles, we form an expression for the front tire slip angle:

$$\begin{aligned} \alpha_f &= \delta_f - \tan^{-1} \left(\frac{l_1\omega_z + v_y}{v_x} \right) \\ &\approx \delta_f - \frac{l_1\omega_z + v_y}{v_x} \\ &= \delta_f - \frac{l_1\omega_z + \dot{y}}{\dot{x}}. \end{aligned} \quad (1.46)$$

The rear tire slip angle is

$$\begin{aligned}
 \alpha_r &= \delta_r - \tan^{-1} \left(-\frac{l_2 \omega_z - v_y}{v_x} \right) \\
 &= \delta_r + \tan^{-1} \left(\frac{l_2 \omega_z - v_y}{v_x} \right) \\
 &\approx \delta_r + \frac{l_2 \omega_z - v_y}{v_x} \\
 &= \delta_r + \frac{l_2 \omega_z - \dot{y}}{\dot{x}}.
 \end{aligned} \tag{1.47}$$

Hence

$$F_{yf} = C_{\alpha f} \left(\delta_f - \frac{l_1 \omega_z + \dot{y}}{\dot{x}} \right), \tag{1.48}$$

and

$$F_{yr} = C_{\alpha r} \left(\delta_r + \frac{l_2 \omega_z - \dot{y}}{\dot{x}} \right). \tag{1.49}$$

Substituting the above two equations into (1.42) and assuming small angles, we obtain

$$m\ddot{y} + m\dot{x}\omega_z = 2C_{\alpha f}\delta_f - 2C_{\alpha f}\frac{l_1\omega_z + \dot{y}}{v_x} + 2C_{\alpha r}\delta_r + 2C_{\alpha r}\frac{l_2\omega_z - \dot{y}}{v_x}.$$

Rearranging terms and performing some manipulations, we get

$$\ddot{y} = - \left(\frac{2C_{\alpha f} + 2C_{\alpha r}}{mv_x} \right) \dot{y} + \left(-v_x - \frac{2C_{\alpha f}l_1 - 2C_{\alpha r}l_2}{mv_x} \right) \omega_z + \frac{2C_{\alpha f}}{m}\delta_f + \frac{2C_{\alpha r}}{m}\delta_r. \tag{1.50}$$

Substituting the expressions for F_{yf} and F_{yr} into (1.43) and assuming small angles, we obtain

$$I_z \dot{\omega}_z = 2l_1 C_{\alpha f} \left(\delta_f - \frac{l_1 \omega_z + \dot{y}}{v_x} \right) \delta_f - 2l_2 C_{\alpha r} \left(\delta_r + \frac{l_2 \omega_z - \dot{y}}{v_x} \right).$$

Rearranging and performing manipulations yields

$$\dot{\omega}_z = - \left(\frac{2l_1 C_{\alpha f} - 2l_2 C_{\alpha r}}{I_z v_x} \right) \dot{y} - \left(\frac{2l_1^2 C_{\alpha f} + 2l_2^2 C_{\alpha r}}{I_z v_x} \right) \omega_z + \frac{2l_1 C_{\alpha f}}{I_z} \delta_f - \frac{2l_2 C_{\alpha r}}{I_z} \delta_r. \tag{1.51}$$

Later, we will analyze a vehicle performance in the fixed coordinates. In particular, we will be interested in an emergency lane-change maneuver. This requires that we transform the vehicle lateral variable into the fixed coordinates. Assuming small angles and referring to Figure 1.17, we obtain

$$\begin{aligned}
 \dot{Y} &= -\dot{x} \sin \psi - \dot{y} \cos \psi \\
 &\approx -v_x \psi - \dot{y}.
 \end{aligned} \tag{1.52}$$

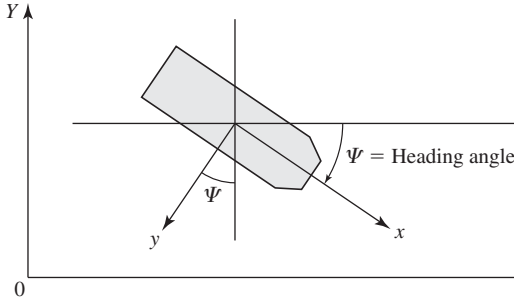


Figure 1.17 Computing the velocity of a vehicle with respect to the fixed coordinates.

Note that

$$\dot{X} = \dot{x} \cos \psi - \dot{y} \sin \psi.$$

For the step lane change, the yaw rate of the vehicle is equal to the yaw rate error; that is,

$$\omega_z = \dot{\psi}.$$

We further assume that

$$\theta = \psi. \quad (1.53)$$

We can now represent the vehicle modeling equations in a state-space format. Let

$$x_1 = \dot{y}, \quad x_2 = \theta, \quad x_3 = \omega_z, \quad x_4 = Y.$$

Then, using (1.50), (1.51), (1.52), and (1.53), we obtain

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \\ \dot{x}_4 \end{bmatrix} = \begin{bmatrix} -\frac{2C_{\alpha f} + 2C_{\alpha r}}{mv_x} & 0 & -v_x - \frac{2C_{\alpha f}l_1 - 2C_{\alpha r}l_2}{mv_x} & 0 \\ 0 & 0 & 1 & 0 \\ -\frac{2l_1C_{\alpha f} - 2l_2C_{\alpha r}}{I_z v_x} & 0 & -\frac{2l_1^2C_{\alpha f} + 2l_2^2C_{\alpha r}}{I_z v_x} & 0 \\ -1 & -v_x & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} \frac{2C_{\alpha f}}{m} & \frac{2C_{\alpha r}}{m} \\ 0 & 0 \\ \frac{2l_1C_{\alpha f}}{I_z} & -\frac{2l_2C_{\alpha r}}{I_z} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \delta_f \\ \delta_r \end{bmatrix}. \quad (1.54)$$

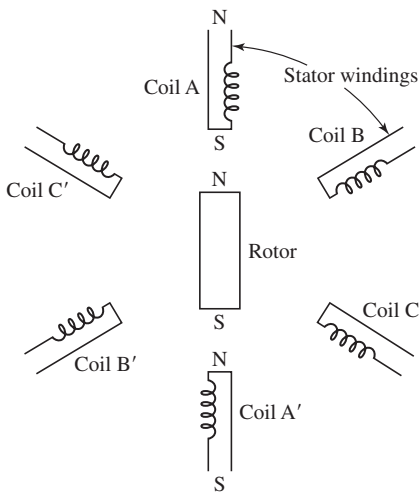
A possible vehicle parameters are given in Table 1.1.

1.7.3 Permanent Magnet Stepper Motor

In this subsection, we develop a model of a permanent magnet (PM) stepper motor. Stepper motors can be viewed as digital positioning devices. Specifically, stepper motors are electromechanical motion devices that convert input information in digital form into output that has

Table 1.1 Typical Vehicle Data

$C_{\alpha f} = C_{\alpha r}$	30,000.00 N/rad
l_1	1.20 m
l_2	1.22 m
m	1280 kg
I_z	2500 kg·m ²
v_x	18.3 m/sec

**Figure 1.18** A simplified schematic of an operation of a stepper motor.

mechanical form. The stepper motor can be driven by digital pulses of the following form: move clockwise or move counterclockwise. We now describe an idea of an operation of a PM stepper motor and then derive its modeling equations.

A simplified organization of a stepper motor is shown in Figure 1.18. The armature, or rotor, is a magnet. At any instant of the operation, exactly one pair of windings of the stator is energized. A pair of windings represents a phase. Each pair of the stator windings is wound in such a way that their magnetic fields are collinear—for example, windings A and A'. The rotor will then rotate until its magnetic field is aligned with the fields of the energized pair of the stator windings. Figure 1.18 shows the situation where the pair of windings A and A' is energized. If the rotor is at rest, as shown in Figure 1.18, and the pair A, A' is de-energized while the pair B, B' is energized, the rotor will rotate to the new position so that its magnetic field is aligned with the magnetic field of the pair B, B'. In this example of the three-phase motor, a single step results in the rotor motion of $360^\circ/6 = 60^\circ$. We next give a more detailed description of the electrical organization of the stepper motor and construct its model.

The basic PM stepper motor consists of a slotted stator and a rotor. There can be two sets of teeth on the rotor. These two sets of teeth are out of alignment with each other by a tooth width, as shown in Figure 1.19. Specifically, the rotor teeth of the left end-cap are displaced one-half pitch from the teeth on the right end-cap. A magnet, or an electromagnet, that is mounted on

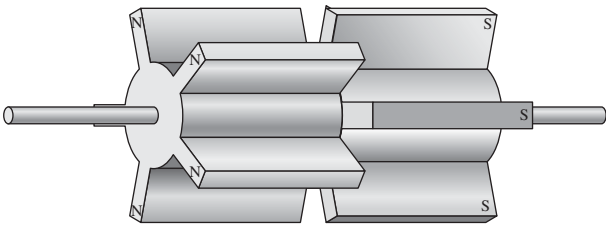
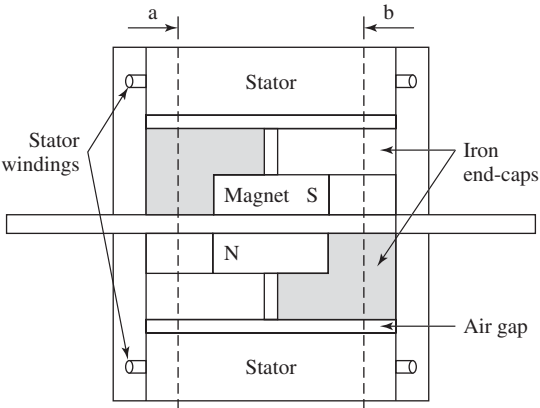
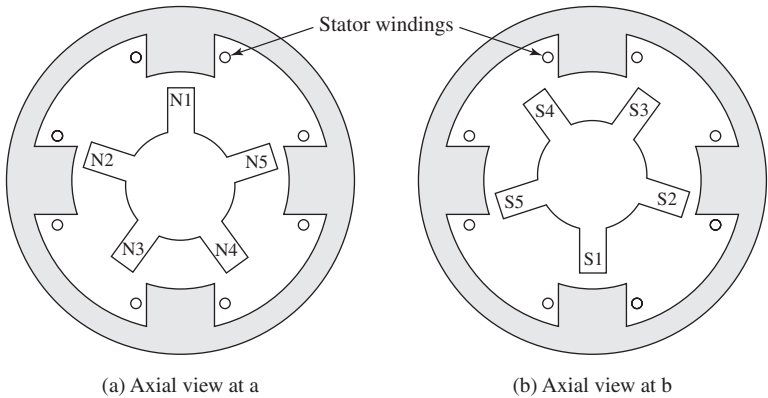


Figure 1.19 A possible rotor structure of a PM stepper motor.



(c) Side, cross-sectional view

Figure 1.20 Construction features of a four-pole stator and a five-pole rotor PM stepper motor.

the rotor magnetizes the iron end-caps. The left end-cap is magnetized as a north pole, while the right end-cap is magnetized as a south pole. The slotted stator is equipped with two or more individual coils. Figure 1.20 shows a four-pole stator and a five-pole rotor structure.

We now discuss the way in which the PM stepper motor can be used to perform precise positioning. Assume that the starting point is as represented in Figure 1.21. Here, coil A–A' is energized so that the upper stator pole is given an S-pole polarity and the corresponding

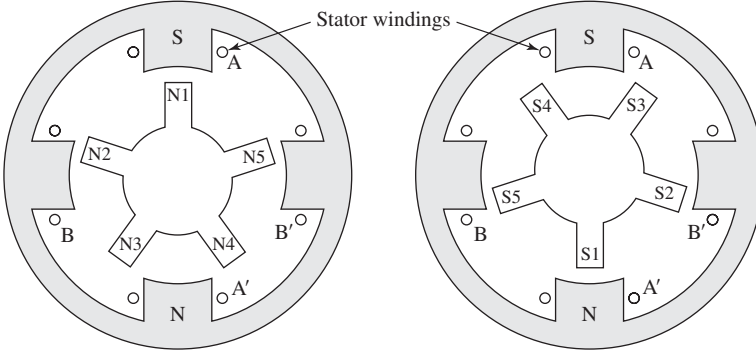


Figure 1.21 Coil A–A′ energized for an S–N orientation, while coil B–B′ de-energized.

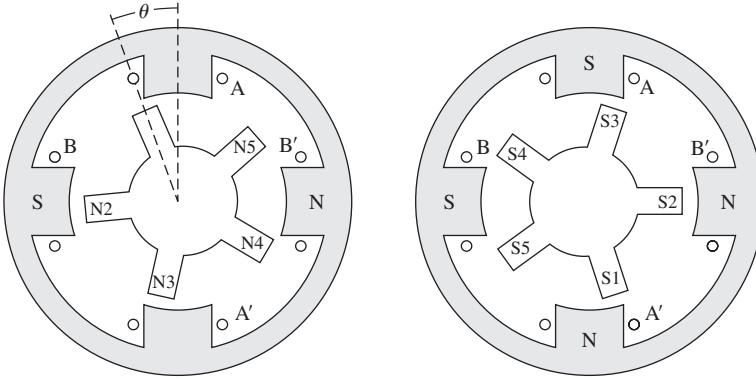


Figure 1.22 Coil A–A′ de-energized, while coil B–B′ energized for an S–N orientation.

lower stator is given an N-pole polarity. The rotor remains in this position as long as stator coil voltage—and hence stator polarities—remains unchanged. This is because the net torque acting on the rotor is zero. If we now de-energize coil A–A′ and energize coil B–B′ so that the left pole becomes an S-pole and the right pole assumes an N-polarity, then this results in the new steady-state position shown in Figure 1.22. We can proceed by alternatively energizing the two stator coils. The PM stepper motor can thus be controlled by digital pulses applied to its two inputs: move clockwise or move counterclockwise.

Our next objective is to devise a mathematical model of the PM stepper motor. In our discussion we use the following notation. We denote by R_t the number of rotor teeth. The angular displacement between rotor teeth is called the tooth pitch and denoted T_p . The angular displacement of the rotor as we change the input voltage from one pair of windings to the other is called the step length and denoted S_l . The number of phases—that is, the number of pairs of stator windings—is denoted N_p . In our example, we have $N_p = 2$, $R_t = 5$, and $T_p = 2\pi/R_t = 72^\circ$. In this example, it takes four input voltage switchings, and hence steps, to advance the rotor one tooth pitch. Indeed, because

$$T_p = 2N_p S_l = \frac{2\pi}{R_t}, \quad (1.55)$$

we obtain

$$S_l = \frac{\pi}{R_r N_p} = 18^\circ. \quad (1.56)$$

We denote the rotor angular displacement by θ .

Following Krause and Wasynczuk [169, pp. 350–351], we now trace the main path of the flux linking the A–A' windings for the rotor position shown in Figure 1.23. The flux leaves the left end-cap through the rotor tooth at the top that is aligned with the stator tooth having the A part of the A–A' winding. The flux then travels up through the stator tooth in the stator iron. Next, the flux splits and travels around the circumference of the stator, as shown in Figure 1.23(a). The flux returns to the south pole of the rotor through the stator tooth positioned at the bottom of the A' part of the A–A' winding, as shown in Figure 1.23(b). We denote the peak flux linkage produced by the rotor permanent magnet in the A–A' stator winding by $n\phi_m$, where n is the number of turns of the winding. For $\theta = 0^\circ$, corresponding to the rotor position in Figure 1.23, the flux linkage in the A–A' winding of the stator is maximal and is approximately equal to $n\phi_m$.

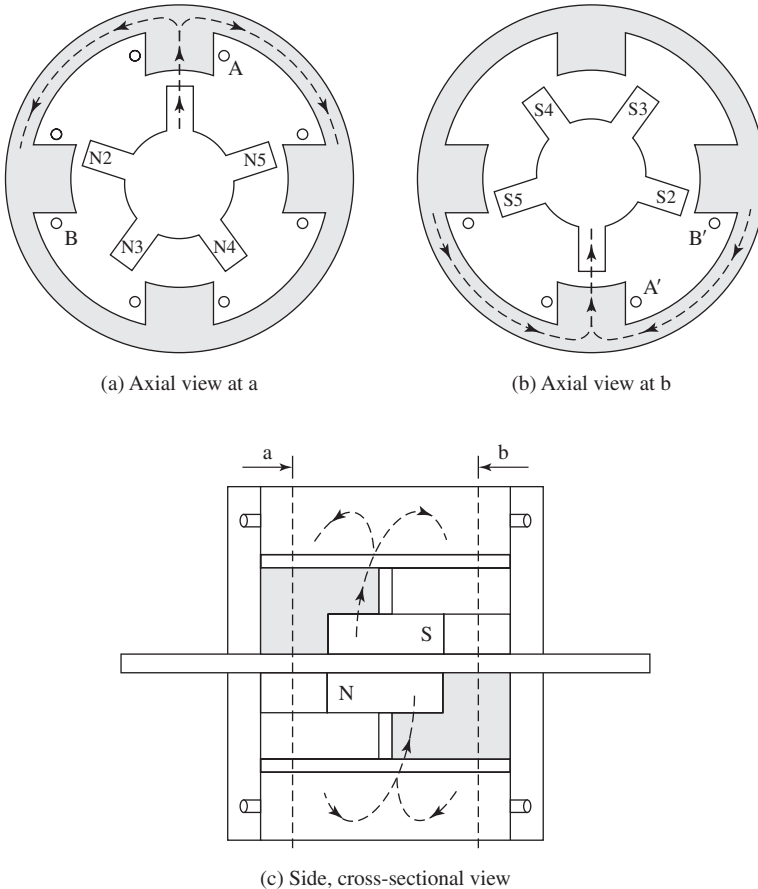


Figure 1.23 The main path of the flux linking the A–A' winding.

Table 1.2 Symbols Used in Modeling Equations of the PM Stepper Motor

T_p	Tooth pitch
S_l	Step length
R_t	Number of teeth on rotor
N_p	Number of phases, that is, stator winding pairs
θ	Angular displacement of the rotor
ω	Angular velocity of the rotor
$v_a(t)$	Voltage applied to phase a , that is, A–A' stator winding
$v_b(t)$	Voltage applied to phase b , that is, B–B' stator winding
$i_a(t)$	Current in phase a
$i_b(t)$	Current in phase b
R	Stator winding resistance per phase
L	Stator winding inductance per phase
B	Viscous friction coefficient
I	Rotor inertia
K_b	Back emf constant

For $\theta = 90^\circ / R_t$, corresponding to the rotor position in Figure 1.22, the flux linkage in the A–A' winding of the stator is minimal, approximately equal to zero. We model the flux linkage, ψ_a , produced by the permanent magnet in the A–A' stator winding as

$$\psi_a = n\phi_m \cos(R_t\theta). \quad (1.57)$$

The (main) flux, linking the B–B' winding for the rotor position shown in Figure 1.22, would enter the stator tooth on which the B part of the B–B' winding is located from the rotor tooth on the left. The flux then would travel around the circumference of the stator and return to the rotor through the stator pole on which the B' part of the B–B' winding is wound. For $\theta = 90^\circ / R_t$, corresponding to the rotor position in Figure 1.22, the flux linkage in the B–B' winding of the stator is maximal and is approximately equal to $n\phi_m$. For $\theta = 0^\circ$, corresponding to the rotor position in Figure 1.23, the flux linkage in the B–B' winding is minimal, approximately equal to zero. We model the flux linkage in the B–B' phase of the stator by

$$\psi_b = n\phi_m \sin(R_t\theta). \quad (1.58)$$

In the following development, we use symbols that are described in Table 1.2. Applying Kirchhoff's voltage law to the circuit of phase a yields

$$v_a(t) = Ri_a(t) + L \frac{di_a(t)}{dt} + e_a(t),$$

where e_a is the back emf induced in phase a ; that is,

$$\begin{aligned} e_a &= \frac{d\psi_a}{dt} \\ &= \frac{d\psi_a}{d\theta} \frac{d\theta}{dt} \\ &= -nR_t\phi_m \sin(R_t\theta) \frac{d\theta}{dt} \\ &= -K_b \omega \sin(R_t\theta), \end{aligned}$$

where $K_b = nR_t\phi_m$. Hence, the voltage equation for phase a can be represented as

$$v_a = Ri_a + L \frac{di_a}{dt} - K_b \omega \sin(R_t\theta).$$

Applying now Kirchhoff's voltage law to the circuit of phase b yields

$$v_b(t) = Ri_b(t) + L \frac{di_b(t)}{dt} + e_b(t),$$

where e_b is the back emf induced in phase b ; that is,

$$\begin{aligned} e_b &= \frac{d\psi_b}{dt} \\ &= \frac{d\psi_b}{d\theta} \frac{d\theta}{dt} \\ &= nR_t\phi_m \cos(R_t\theta) \frac{d\theta}{dt} \\ &= K_b \omega \cos(R_t\theta), \end{aligned}$$

where, as before, $K_b = nR_t\phi_m$. Hence, the voltage equation for phase b can be represented as

$$v_b = Ri_b + L \frac{di_b}{dt} + K_b \omega \cos(R_t\theta).$$

Let T_a and T_b denote torques produced by the currents i_a and i_b , respectively. Then, applying Newton's second law to the rotor yields

$$I \frac{d^2\theta}{dt^2} + B \frac{d\theta}{dt} = T_a + T_b.$$

The PM stepper motor is an electromechanical device in which electric power is transformed into mechanical power. The mechanical power, P_a , due to the torque T_a is

$$P_a = T_a \frac{d\theta}{dt}. \quad (1.59)$$

Assuming that there are no losses, this mechanical power is equal to the power developed in phase a :

$$P_a = i_a \frac{d\psi_a}{dt} = i_a \frac{d\psi_a}{d\theta} \frac{d\theta}{dt}. \quad (1.60)$$

Equating the right-hand sides of (1.59) and (1.60) and performing simple manipulations, we obtain

$$T_a = i_a \frac{d\psi_a}{d\theta} = -i_a K_b \sin(R_t\theta).$$

Applying similar arguments to phase b ; we get

$$T_b = i_b \frac{d\psi_b}{d\theta} = i_b K_b \cos(R_t\theta).$$

Taking the above equations into account, we represent the equation of motion of the rotor as

$$I \frac{d^2\theta}{dt^2} + B \frac{d\theta}{dt} = T_a + T_b = -i_a K_b \sin(R_t\theta) + i_b K_b \cos(R_t\theta).$$

Combining the above equations, we obtain a model of a two-phase (i.e., a two-input) PM stepper motor. We have

$$\begin{aligned}\frac{d\theta}{dt} &= \omega, \\ \frac{d\omega}{dt} &= -\frac{K_b}{I}i_a \sin(R_t\theta) + \frac{K_b}{I}i_b \cos(R_t\theta) - \frac{B}{I}\omega, \\ \frac{di_a}{dt} &= -\frac{R}{L}i_a + \frac{K_b}{L}\omega \sin(R_t\theta) + \frac{1}{L}v_a, \\ \frac{di_b}{dt} &= -\frac{R}{L}i_b - \frac{K_b}{L}\omega \cos(R_t\theta) + \frac{1}{L}v_b.\end{aligned}$$

We now define the state variables: $x_1 = \theta$, $x_2 = \omega$, $x_3 = i_a$, and $x_4 = i_b$. We let $K_1 = B/I$, $K_2 = K_b/I$, $K_3 = R_t$, $K_4 = K_b/L$, $K_5 = R/L$, $u_1 = \frac{1}{L}v_a$, and $u_2 = \frac{1}{L}v_b$. Using the above notation, we represent the PM stepper motor model in state space format:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \\ \dot{x}_4 \end{bmatrix} = \begin{bmatrix} x_2 \\ -K_1x_2 - K_2x_3 \sin(K_3x_1) + K_2x_4 \cos(K_3x_1) \\ K_4x_2 \sin(K_3x_1) - K_5x_3 \\ -K_4x_2 \cos(K_3x_1) - K_5x_4 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}. \quad (1.61)$$

1.7.4 Stick Balancer

In this subsection we develop a nonlinear model of the stick balancer, also known as the cart with an inverted pendulum. In our derivation of the model, we use Newton's laws, in a fashion similar to that of Kwakernaak and Sivan [173] and Ogata [216], where linearized models of the stick balancer were developed. An alternative derivation, using the method of D'Alembert, can be found in Cannon [41, Section 22.4].

A free-body diagram of a cart on which a stick (i.e., an inverted pendulum) is mounted is shown in Figure 1.24. Let $H = H(t)$ and $V = V(t)$ be horizontal and vertical reaction forces, respectively, in the pivot. That is, the force H is the horizontal reaction force that the cart exerts upon the pendulum, whereas $-H$ is the force exerted by the pendulum on the cart. Similar convention applies to the forces V and $-V$. The x and y are the coordinates of the fixed, nonrotating coordinate frame x - y . The angular displacement of the stick from the vertical position is $\theta = \theta(t)$. The mass of the cart is M , while the mass of the stick is m . The length of the stick is $2l$, and its center of gravity is at its geometric center. The control force applied to the cart is labeled u . We assume that the wheels of the cart do not slip. We model the frictional force of the cart wheels on the track by

$$f_c = \mu_c \text{sign}(\dot{x}), \quad (1.62)$$

where μ_c is the cart friction coefficient. Friedland [91, p. 201] refers to this model of friction as the classical Coulomb friction model. Let (x_G, y_G) be the coordinates of the center of gravity of the stick. Then,

$$\begin{aligned}x_G &= x + l \sin(\theta), \\ y_G &= l \cos(\theta).\end{aligned} \quad (1.63)$$

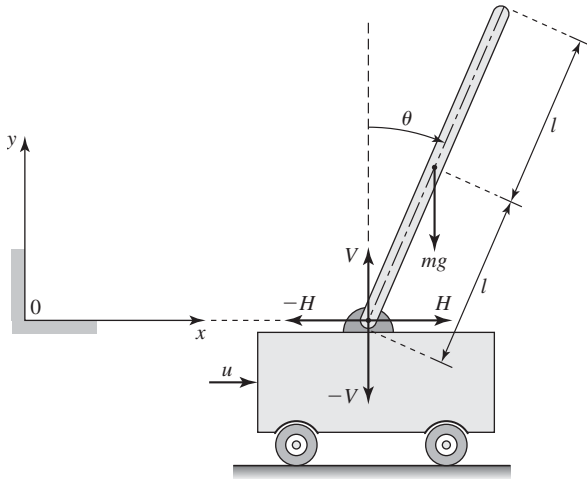


Figure 1.24 Free-body diagram of a cart with an inverted pendulum.

We are now ready to write down equations modeling the system. The equation that describes the rotational motion of the stick about its center of gravity is obtained by applying the rotational version of Newton's second law. Summing the moments about the center of gravity of the stick, we obtain

$$I \frac{d^2\theta}{dt^2} = Vl \sin(\theta) - Hl \cos(\theta), \quad (1.64)$$

where

$$I = \int_{-l}^l r^2 dm = \frac{ml^2}{3} \quad (1.65)$$

is the moment of inertia of the stick with respect to its center of gravity.

We next write the equation that describes the horizontal motion of the center of gravity of the stick. Applying Newton's second law along the x axis yields

$$m \frac{d^2}{dt^2} (x + l \sin(\theta)) = H. \quad (1.66)$$

Performing the required differentiation, we get

$$m(\ddot{x} + l[-\dot{\theta}^2 \sin(\theta) + \ddot{\theta} \cos(\theta)]) = H. \quad (1.67)$$

The equation that describes the vertical motion of the center of gravity of the stick is obtained by applying Newton's second law along the y axis. We have

$$m \frac{d^2}{dt^2} (l \cos(\theta)) = V - mg. \quad (1.68)$$

Performing the needed differentiation yields

$$ml(-\dot{\theta}^2 \cos(\theta) - \ddot{\theta} \sin(\theta)) = V - mg. \quad (1.69)$$

Finally, we apply Newton's second law to the cart to get

$$M \frac{d^2 x}{dt^2} = u - H - f_c. \quad (1.70)$$

Substituting (1.70) into (1.67) gives

$$m\ddot{x} + ml\ddot{\theta} \cos(\theta) - ml\dot{\theta}^2 \sin(\theta) + f_c = u - M\ddot{x}. \quad (1.71)$$

Now substituting equations (1.69) and (1.70) into (1.64) yields

$$I\ddot{\theta} = (mg - ml\dot{\theta}^2 \cos(\theta) - ml\ddot{\theta} \sin(\theta))l \sin(\theta) + f_c l \cos(\theta) - (u - M\ddot{x})l \cos(\theta). \quad (1.72)$$

We next substitute $u - M\ddot{x}$ from (1.71) into (1.72) and perform manipulations to get

$$I\ddot{\theta} = mgl \sin(\theta) - ml^2\ddot{\theta} - m\ddot{x}l \cos(\theta). \quad (1.73)$$

Let

$$a = \frac{1}{m + M}.$$

Then, we can represent (1.71) as

$$\ddot{x} = -mal\ddot{\theta} \cos(\theta) + mal\dot{\theta}^2 \sin(\theta) - af_c + au. \quad (1.74)$$

We substitute (1.74) into (1.73) to obtain

$$\ddot{\theta} = \frac{mgl \sin(\theta) - m^2 l^2 a \dot{\theta}^2 \sin(2\theta)/2 + mal \cos(\theta) f_c - mal \cos(\theta) u}{I - m^2 l^2 a \cos^2(\theta) + ml^2}. \quad (1.75)$$

Let $x_1 = \theta$ and $x_2 = \dot{\theta}$. Using the expression for I , given by (1.65), we represent (1.75) in state-space format as follows:

$$\begin{aligned} \dot{x}_1 &= x_2, \\ \dot{x}_2 &= \frac{g \sin(x_1) - mlax_2^2 \sin(2x_1)/2 + a \cos(x_1) f_c}{4l/3 - mla \cos^2(x_1)} - \frac{mal \cos(x_1) u}{4ml^2/3 - m^2 l^2 a \cos^2(x_1)}. \end{aligned} \quad (1.76)$$

We next substitute $\ddot{\theta}$, obtained from (1.73), into (1.74) to get

$$\ddot{x} = \frac{-mag \sin(2x_1)/2 + ax_2^2 \sin(x_1)4ml/3 + (u - f_c)4a/3}{4/3 - ma \cos^2(x_1)}. \quad (1.77)$$

Let $x_3 = x$ and $x_4 = \dot{x}$. Then, by combining (1.76) and (1.77) we obtain a state-space model of the inverted pendulum on a cart of the form

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \\ \dot{x}_4 \end{bmatrix} = \begin{bmatrix} x_2 \\ \frac{g \sin(x_1) - mlax_2^2 \sin(2x_1)/2}{4l/3 - mla \cos^2(x_1)} \\ x_4 \\ \frac{-mag \sin(2x_1)/2 + alx_2^2 \sin(x_1)4m/3}{4/3 - ma \cos^2(x_1)} \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{-a \cos(x_1)}{4l/3 - mla \cos^2(x_1)} \\ 0 \\ \frac{4a/3}{4/3 - ma \cos^2(x_1)} \end{bmatrix} (u - f_c). \quad (1.78)$$

While designing a controller for the stick balancer system, we will use an equivalent representation of (1.78):

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \\ \dot{x}_4 \end{bmatrix} = \begin{bmatrix} x_2 \\ \frac{g \sin(x_1)}{4l/3 - mla \cos^2(x_1)} \\ x_4 \\ \frac{-mag \sin(2x_1)/2}{4/3 - ma \cos^2(x_1)} \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{-a \cos(x_1)}{4l/3 - mla \cos^2(x_1)} \\ 0 \\ \frac{4a/3}{4/3 - ma \cos^2(x_1)} \end{bmatrix} (u - f_c + mlx_2^2 \sin(x_1)). \quad (1.79)$$

1.7.5 Population Dynamics

We first model dynamics of a single population. Then, we use the obtained model to investigate two interacting populations.

Let $x = \phi(t)$ be the population of the given species at a time t . The simplest model of the population dynamics of the given species is the one resulting from the hypothesis that the net rate of change of x is proportional to the current value of x . Let r be the rate of growth of the population. Then, the simple model of the population dynamics is

$$\dot{x} = rx. \quad (1.80)$$

Solving (1.80) subject to the initial condition, $x(0) = x_0$, yields

$$x = e^{rt} x_0. \quad (1.81)$$

In the above linear model, the growth rate, r , is constant. It does not account for the limitation of space, food supply, and so on. Therefore, we modify the unrestricted growth rate r to take into account the fact that the environment of the population can only support a certain number, say K , of the species. The constant K is called the *carrying capacity* of the environment or the *saturation level*. Specifically, if $x > K$, then because there is not enough food or space to support x , more species will die than are born, which results in a negative growth rate. On the other hand, if $x < K$, then there is an extra food and space and the population growth is positive. One function that satisfies the above requirements has the form

$$r \left(1 - \frac{x}{K} \right).$$

Using the above model of the growth rate, we obtain the following model of the population growth:

$$\dot{x} = r \left(1 - \frac{x}{K} \right) x. \quad (1.82)$$

Let $c = r/K$. Then, we represent (1.82) as

$$\boxed{\dot{x} = (r - cx)x} \quad (1.83)$$

Equation (1.83) is known as the *Verhulst equation* or the *logistic equation*. According to Boyce and DiPrima [32, p. 60], Verhulst, a Belgian mathematician, introduced equation (1.83) in 1838 to model human population growth. For reasons somewhat unclear, Verhulst referred to it as the logistic growth; hence (1.83) is often called the logistic equation.

We will use the logistic equation to analyze interacting species. We first model two interacting species that do not prey upon each other, but compete for the same resources—for example, for a common food supply. Let x_1 and x_2 be the populations, at time t , of the two competing species. We assume that the population of each of the species in the absence of the other is governed by a logistic equation. Thus, a model of two noninteracting species has the form

$$\begin{aligned}\dot{x}_1 &= (r_1 - c_1 x_1)x_1, \\ \dot{x}_2 &= (r_2 - c_2 x_2)x_2,\end{aligned}\tag{1.84}$$

where r_1 and r_2 are the unrestricted growth rates of the two populations, while r_1/c_1 and r_2/c_2 are their saturation levels. When both species are present, each will affect the growth rate of the other. Because x_1 and x_2 use the same resources, the growth rate of x_1 will be lowered by an amount proportional to the size of x_2 . We modify the growth equation for species x_1 as

$$\dot{x}_1 = (r_1 - c_1 x_1 - a_1 x_2)x_1,\tag{1.85}$$

where a_1 is a measure by which species x_2 lowers the growth rate of x_1 . Similarly, the growth rate of x_2 is reduced by $a_2 x_1$, where a_2 is a measure of the degree by which species x_1 lowers the growth rate of x_2 . Therefore, the growth equation for species x_2 becomes

$$\dot{x}_2 = (r_2 - c_2 x_2 - a_2 x_1)x_2.\tag{1.86}$$

In summary, equations (1.85) and (1.86) model two species that interact by competing for a common resource like food or other supply.

We will now study the situation in which one species, the predator, preys upon the other species, the prey, and there is no competition between the predator and prey for a common resource. Let x_1 and x_2 denote the populations of the prey and predator, respectively. We will make the following modeling assumptions:

1. When there is no predator—that is, $x_2 = 0$ —the prey enjoys an exponential growth governed by

$$\dot{x}_1 = r_1 x_1, \quad r_1 > 0.\tag{1.87}$$

2. When there is no prey—that is, $x_1 = 0$ —the predator perishes. We model this scenario using the equation

$$\dot{x}_2 = -r_2 x_2, \quad r_2 > 0.\tag{1.88}$$

3. The larger the population of the predator, the more prey is eaten. The larger the population of the prey, the easier hunting is for the predator. We model the growth rate decrease of the prey by including in (1.87) the term $-dx_1 x_2$, $d > 0$. The growth rate increase of the predator is modeled by adding the term $ex_1 x_2$ in (1.88), where $e > 0$.

The resulting model of the predator–prey relationship is

$$\begin{aligned}\dot{x}_1 &= (r_1 - dx_2)x_1 \\ \dot{x}_2 &= (-r_2 + ex_1)x_2\end{aligned}$$

(1.89)

Equations (1.89) are known as the *Lotka–Volterra equations*. They were published by Lotka in 1925 and by Volterra in 1926.

Notes

A landmark paper on mathematical description of linear dynamical systems from a controls point of view is by Kalman [147]. The first chapter of Sontag's book [269] is a nice, easy to read, and comprehensive introduction to the subject of mathematical control theory. Sontag's chapter is easily available on the web. For more examples of dynamical system models from mechanical and electrical engineering, the reader is referred to Truxal [283], Banks [19], Raven [240], Mohler [207], and D'Azzo and Houpis [57]. Wertz [298] has many excellent real-life models from space science. A well-written paper on vehicle modeling and control for broad audience is by Ackermann [2]. In addition to references 303 and 299, vehicle and tire dynamics modeling are investigated in detail by Nalecz [211], Nalecz and Bindemann [212, 213], and Smith and Starkey [263, 264]. For more on stepper motors, the reader is referred to Acarnley [1], Kenjo [153], Kenjo and Nagamori [154], Del Toro [61, Chapter 9], Krause and Wasynczuk [169, Chapter 8], and Zribi and Chiasson [319]. For further analysis of the equations modeling competing species and the predator–prey equations, we recommend Boyce and DiPrima [32, Sections 9.4 and 9.5]. Our derivation of the logistic equation follows that of Sandefur [253]. The logistic equation as well as equations modeling interacting species can be used for mathematical modeling of tumor dynamics and interactions between tumor and immune system. A survey of models for tumor–immune system dynamics can be found in the well-edited book by Adam and Bellomo [4].

Mayr writes on page 109 in reference 199, “It is still widely believed that the steam-engine governor is the oldest feedback device, and that James Watt had not only invented but also patented it. While both errors are easily refuted, we are still not able to reconstruct the history of this invention in all desired completeness.” Watt did not patent the governor. He did not invent it either. On page 112 of his book [199], Mayr adds the following: “But the application of the centrifugal pendulum in a system of speed regulation of steam engines was a new breakthrough for which the firm of Boulton & Watt, if not James Watt himself, clearly deserves the credit.”

EXERCISES

1.1 Prove the law of conservation of energy as stated in the following theorem:

Theorem 1.1 Given a conservative vector field \mathbf{F} —that is, $\mathbf{F} = -\nabla U$, where U is the potential energy—assume that a particle of mass m , satisfying Newton's law, moves on a differentiable curve \mathbf{C} . Then, the sum of the potential energy and kinetic energy is constant.

Hint Denote by $\mathbf{C}(t)$ a differentiable curve along which a particle of mass m moves. Hence,

$$\mathbf{F}(\mathbf{C}(t)) = m\ddot{\mathbf{C}}(t).$$

Prove that

$$U(\mathbf{C}(t)) + \frac{1}{2}m\|\dot{\mathbf{C}}(t)\|^2 = \text{constant}.$$

To show that the above relation holds, differentiate $U(\mathbf{C}(t)) + \frac{1}{2}m\|\dot{\mathbf{C}}(t)\|^2$ with respect to time and show that the time derivative is zero; hence the sum of the energies is constant.

- 1.2** Show that in a conservative vector field the work done by the field on a particle that moves between two points depends only on these points and not on the path followed.

Hint Note that the work done against the conservative force field \mathbf{F} along a curve \mathbf{C} between points $\mathbf{C}(t_1)$ and $\mathbf{C}(t_2)$ is

$$\int_{t_1}^{t_2} \mathbf{F}(\mathbf{C}(t))^T \dot{\mathbf{C}}(t) dt.$$

By assumption, \mathbf{F} is a conservative vector field. Hence,

$$\mathbf{F} = -\nabla U. \quad (1.90)$$

Combining the above yields

$$\int_{t_1}^{t_2} \mathbf{F}(\mathbf{C}(t))^T \dot{\mathbf{C}}(t) dt = - \int_{t_1}^{t_2} \nabla U(\mathbf{C}(t))^T \dot{\mathbf{C}}(t) dt.$$

Use the chain rule to represent the expression inside the integral on the right-hand side as

$$\nabla U(\mathbf{C}(t))^T \dot{\mathbf{C}}(t) = \frac{dU(\mathbf{C}(t))}{dt}.$$

- 1.3** An example of a closed-loop feedback system is a toilet flushing device connected to a water tank. A toilet system is shown in Figure 1.25. The control objective is to maintain water in the tank at a constant level. Draw a block diagram of the system.

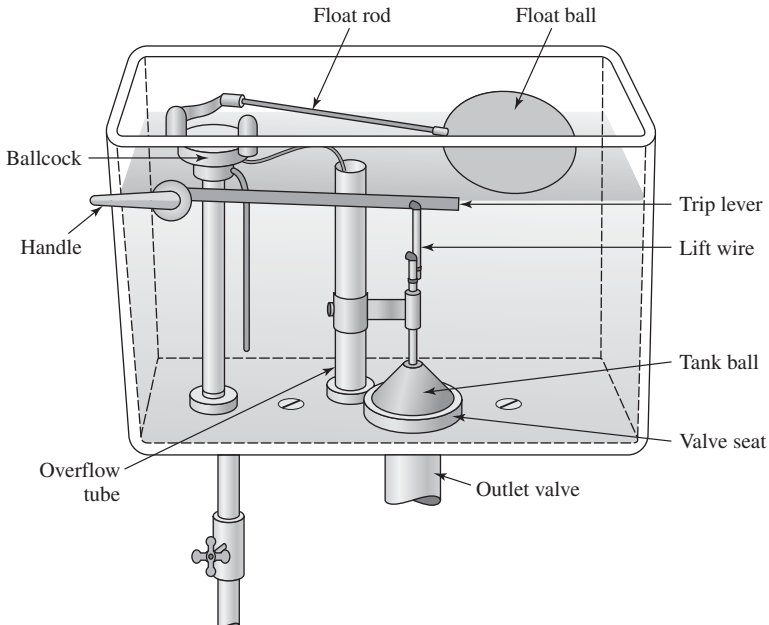
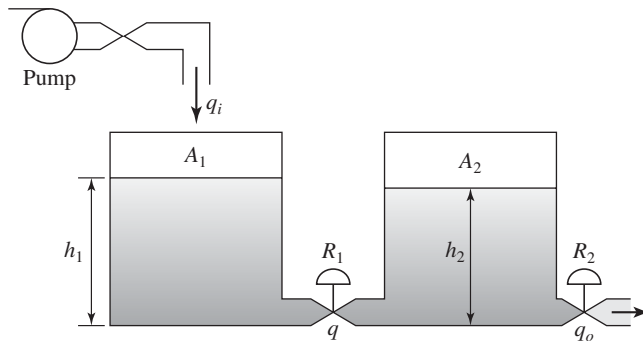


Figure 1.25 Toilet-flushing system of Exercise 1.3. (Adapted from reference 241, p. 211.)

Table 1.3 Numerical Data for the Two-Tank Pumped Storage System Model for Exercise 1.4

Tank 1 diameter	2 m
Tank 2 diameter	3 m
R_1	10 sec/m ²
R_2	20 sec/m ²

**Figure 1.26** A two-tank pumped storage system for Exercise 1.4.

- 1.4** The goal of this exercise is to develop a model of a two-tank pumped storage system. Figure 1.26 illustrates such a system. A pumped storage system can be used to supply water to a turbine producing electric power. The two storage tanks are connected together through a valve that may be modeled as a linear resistance R_1 . The resistance of the turbine supply pump is modeled by a linear resistance R_2 . Suppose that we wish to design a controller that would control the inflow rate q_i in order to maintain the output flow rate q_o at some desired value q_d . We assume that the linearized flow rates $q(t)$ and $q_o(t)$ are proportional to the corresponding water heights; that is,

$$q(t) = \frac{h_1(t) - h_2(t)}{R_1},$$

$$q_o(t) = \frac{h_2(t)}{R_2}.$$

Let A_1 and A_2 represent the surface areas in the two tanks. Then, the principle of conservation of mass for an incompressible fluid yields

$$q_i(t) - q(t) = A_1 \dot{h}_1(t),$$

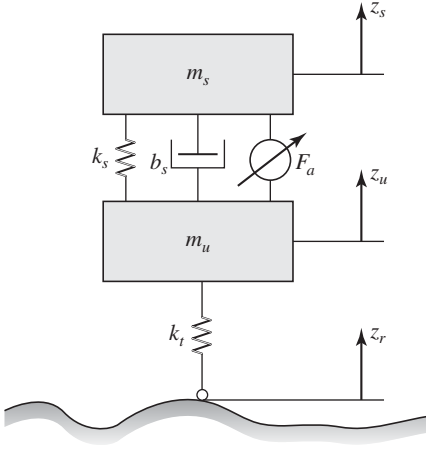
$$q(t) - q_o(t) = A_2 \dot{h}_2(t).$$

Numerical data of a two-tank pumped storage system is given in Table 1.3. Find the transfer function $\mathcal{L}(q_o(t))/\mathcal{L}(q_i(t))$.

- 1.5** The objective of this exercise is to derive a model of an automotive suspension. The role of an automotive suspension is to support the vehicle body on the axles. A simple linear model of the suspension assembly associated with a single wheel, called the quarter-car model, is shown in Figure 1.27. The vehicle body is represented by the “sprung” mass, m_s , while the tire and axle are represented by the “unsprung” mass, m_u . Model the tire

Table 1.4 Parameters of the Vehicle Suspension of Exercise 1.5

m_s	355 kg
m_u	46 kg
k_s	14384 N/m
b_s	1860 N/m/sec
k_t	300,000 N/m

**Figure 1.27** Linear quarter-car model for Exercise 1.5.

as a spring with spring coefficient k_t , and assume its damping to be negligible. The suspension consists of the spring, the shock absorber, and the variable force element. The role of the variable force element is to generate a force that compensates for the uneven road, which is the source of the system disturbances. Assume that the variable force element F_a can instantaneously provide any desired force. The above described model of an automotive suspension is referred to as a linear active suspension. Write down equations modeling the dynamics of the two-degree-of-freedom system depicted in Figure 1.27 by applying Newton's law to each mass of the system and taking into account the positive direction of the z coordinates. Define the state variables to be

$$x_1 = z_s - z_u = \text{suspension stroke,}$$

$$x_2 = \dot{z}_s = \text{sprung mass velocity,}$$

$$x_3 = z_u - z_r = \text{tire deflection,}$$

$$x_4 = \dot{z}_u = \text{unsprung mass velocity.}$$

Let $u = F_a$ be the system input, and let $d(t) = \dot{z}_r(t)$ be the disturbance input due to the road roughness. Represent the obtained linear model of active suspension in state-space format. Use data of a commercial vehicle given in Table 1.4. The data come from Sunwoo and Cheok [272]. Substitute the vehicle data into the modeling equations. A design goal of an automotive engineer may be to design a control law to isolate the sprung mass from road disturbances. In other words, he or she may wish to design a control law such that the vertical velocity of the sprung mass x_2 is as close to 0 as

possible in spite of road disturbances. Assume that we can measure the vertical velocity of the sprung mass. Write the output equation.

To investigate the effect of road disturbances on the vertical velocity of the sprung mass, set $u = 0$ and $\delta(t) = d \sin(6t)$. Generate a plot of y versus t for t in the range $[0 \text{ sec}, 5 \text{ sec}]$ and $d = 0.5$.

- 1.6** The objective of this exercise is to develop a model of a single-axis machine tool drive in which a motor drives a lead screw through a gear. Figure 1.28 illustrates such a system. This model comes from Driels [68, pp. 613–619]. The system operates as follows. The lead screw moves the table with a workpiece on it under the cutter to perform the desired operation. The table is positioned by a control system. We wish to place the table in a desirable position as rapidly and accurately as possible. This means that while modeling the machine tool drive system we should, in particular, take into account the flexibility of the lead screw. The length of the lead screw implies that its rotation at the gear box will not be the same as its rotation at the table. Thus, the angular twist along the length of the lead screw cannot be neglected. Model the lead screw flexibility using a torsional spring as illustrated in Figure 1.29. The system input is the motor torque τ_m , and the system output is the rotation of the end of the lead screw labeled θ_l . Assume that the damping torque, due to viscous friction, acting on the motor is proportional to its angular velocity $\dot{\theta}_m = \omega_m$ with the coefficient of viscous friction denoted c_m . Denote the motor moment of inertia by I_m , and denote by τ_1 the load torque on the motor gear due to the rest of the gear train. Write the equation modeling the dynamics of the motor shaft.

Let τ_2 be the torque transmitted to the load shaft, and let $N = N_1/N_2$ be the gear ratio, where N_1 denotes the number of teeth of the first gear, and N_2 is the number of teeth of the second gear. The work done by the first gear is equal to that of the second.

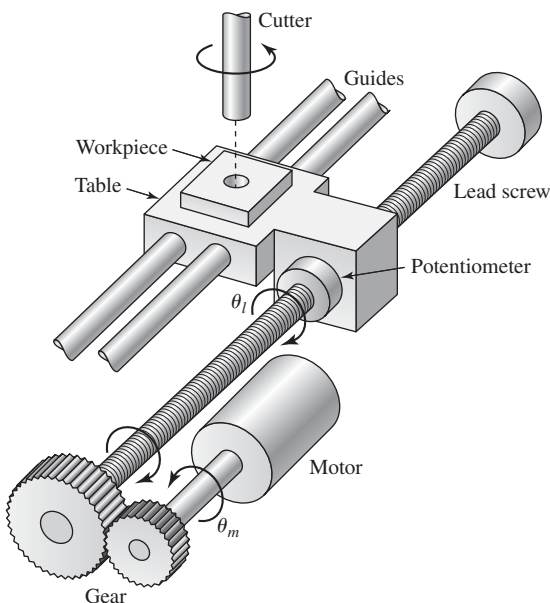


Figure 1.28 A schematic of a machine tool drive system of Exercise 1.6. (Adapted from Driels [68].)

Table 1.5 Data for the Machine Tool Drive System Model for Exercise 1.6

I_m	0.0001 Nm·sec ² /rad
I_l	0.001 Nm·sec ² /rad
$c_m = c_l$	0.01 Nm·sec/rad
k	1 Nm/rad
N	1/5

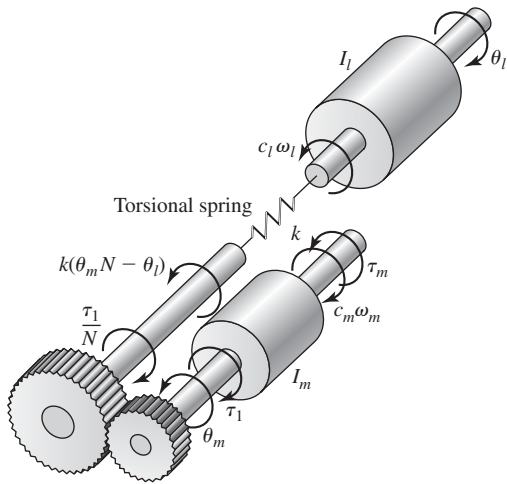


Figure 1.29 A mechanical model of the machine tool drive system for Exercise 1.6.

Model the load shaft flexibility with a torsional spring. The torsional spring coefficient is k . Assume that the moment of inertia of the lead screw is negligible compared to that of the load. Denote the moment of inertia of the load by I_l . The torque due to viscous friction of the load is $c_l\omega_l$, where c_l is the coefficient of the viscous friction of the load. The angular displacement of the end of the lead screw is θ_l . Write the equation modeling the dynamics of the load. Define the state variables to be

$$x_1 = \theta_m, \quad x_2 = \dot{\theta}, \quad x_3 = \theta_l, \quad x_4 = \dot{\theta}_l.$$

Represent the modeling equations in state-space format. Typical numerical data for a machine tool drive system are shown in Table 1.5.

- 1.7** Consider a bicycle model of the front-wheel-steering vehicle shown in Figure 1.30. Construct a two-dimensional state-space model of the vehicle, where $x_1 = \dot{y}$ and $x_2 = \dot{\theta}$. The input $u = \delta_f$ and the output $y = x_2$. In your derivation assume small angles. Vehicle data are given in Table 1.6.

Table 1.6 Vehicle Data for Exercise 1.7

$l_1 = l_2$	1 m
v_x	10 m/sec
m	1000 kg
I_z	1000 kg·m ²
$C_f = C_r$	10,000 N/rad

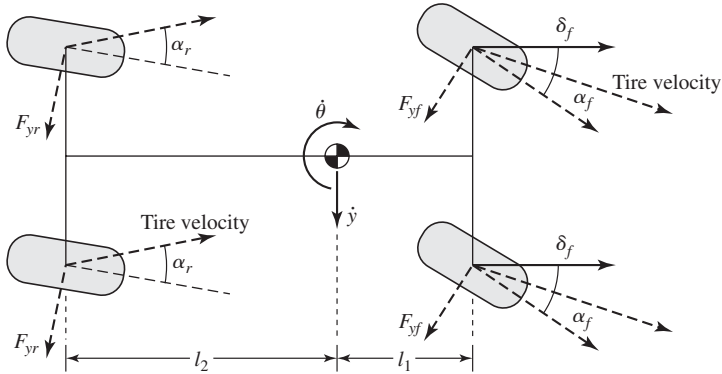


Figure 1.30 Bicycle model of a vehicle for Exercise 1.7.

- 1.8** Consider the two-mass–spring system depicted in Figure 1.31.
- Form the Lagrangian of the system.
 - Write the Lagrange equations of motion for the system.
 - Represent the obtained Lagrange equations of motion in state-space format.

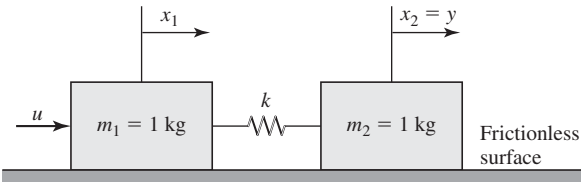


Figure 1.31 The two-mass–spring system for Exercise 1.8.

- 1.9** Consider the mechanical system shown in Figure 1.32.
- Form the Lagrangian for the system.
 - Write the Lagrange equation of motion for the system.
 - Represent the system model in state space format.

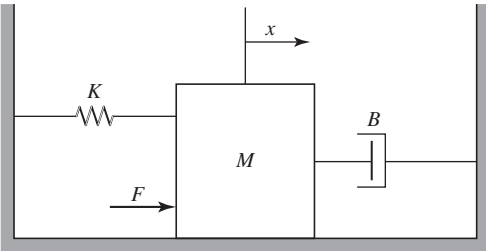


Figure 1.32 Mechanical system for Exercise 1.9.

1.10 Consider a dynamical system that consists of a cart with an inverted pendulum (point mass on a massless shaft) attached to it as depicted in Figure 1.33.

- Write the Lagrange equations of motion for the system.
- Represent the obtained model in state-space format using the state variables

$$x_1 = x, \quad x_2 = \dot{x}, \quad x_3 = \theta, \quad x_4 = \dot{\theta}.$$

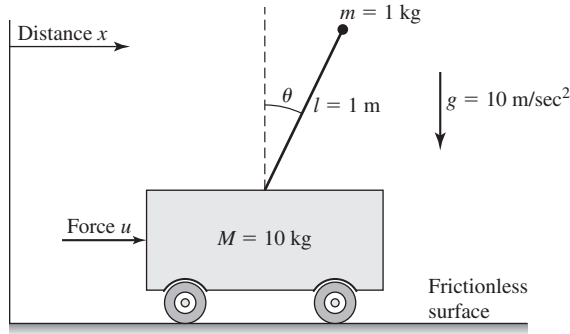


Figure 1.33 An inverted pendulum moving on a cart for Exercise 1.10.

1.11 Consider a schematic of a field-controlled DC motor shown in Figure 1.34. Let $x_1 = \theta$, $x_2 = \dot{\theta}$, $x_3 = i_f$, and $y = x_1$. Treat the field inductance L_f as a parasitic element and denote it as μ .

- Form the actual plant model P_μ in state-space format.
- Find the reduced-order design model P_0 .

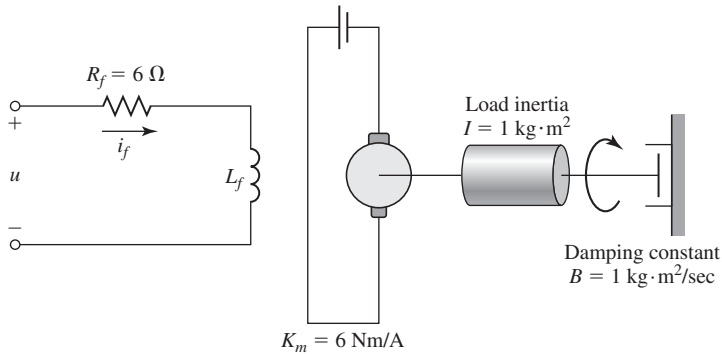


Figure 1.34 A diagram of a field-controlled DC motor for Exercise 1.11.



CHAPTER 2

Analysis of Modeling Equations

To be sure, mathematics can be extended to any branch of knowledge, including economics, provided the concepts are so clearly defined as to permit accurate symbolic representation. That is only another way of saying that in some branches of discourse it is desirable to know what you are talking about.

—James R. Newman

Most of dynamical systems analyzed in this book are modeled by ordinary differential equations. A main use of a mathematical model is to predict the system transient behavior. Unlike linear systems, where closed-form solutions can be written in terms of the system's eigenvalues and eigenvectors, finding analytical, exact, solutions to nonlinear differential equations can be very difficult or impossible. However, we can approximate the solutions of differential equations with difference equations whose solutions can be obtained easily using a computer. In this chapter we discuss a number of methods for solving differential equations. The first class of methods allows us to graphically determine solutions to second-order differential equations. Then, we discuss numerical techniques for solving differential equations. After that, we present two methods of linear approximation of nonlinear systems.

2.1 State-Plane Analysis

The *state plane*, which is a state space of two-dimensional systems, is also called the *phase plane*. Analysis in the state plane is applicable to linear and nonlinear systems modeled by second-order ordinary differential equations. The state-plane methods are graphical procedures for solving such equations. Using state-plane methods, one can graphically determine the transient response of a second-order dynamical system. We now introduce relevant terms. Consider a class of second-order systems whose dynamics are modeled by the differential equation

$$\ddot{x} = f(x, \dot{x}). \quad (2.1)$$

Define the state variables

$$x_1 = x \quad \text{and} \quad x_2 = \dot{x}.$$

Then, equation (2.1) can be represented as

$$\begin{aligned}\dot{x}_1 &= x_2, \\ \dot{x}_2 &= f(x_1, x_2).\end{aligned}\tag{2.2}$$

The plane with coordinates x_1 and x_2 is labeled as the state plane or phase plane. A solution of (2.1) can be illustrated by plotting x versus t , or by plotting x_2 versus x_1 using t as a parameter. To each value of the state $\mathbf{x}(t) = [x_1(t) \ x_2(t)]^T$ there corresponds a point in the (x_1, x_2) plane called the *representative point* (RP). As t varies, the RP describes a curve in the state plane called a *trajectory*. A family of trajectories is called a *phase portrait*.

2.1.1 Examples of Phase Portraits

◆ Example 2.1

Consider the spring–mass system shown in Figure 2.1, where x is the displacement of the mass M from equilibrium, and k is the linear spring constant. The spring is assumed to be linear; that is, it obeys Hooke’s law. Applying Newton’s law, we obtain the equation of motion of the mass:

$$M\ddot{x} + kx = 0.$$

Let $\omega^2 = k/M$. Then, the above equation can be represented in standard form:

$$\ddot{x} + \omega^2 x = 0.\tag{2.3}$$

Observe that

$$\ddot{x} = \dot{x} \frac{\ddot{x}}{\dot{x}} = \dot{x} \frac{d\dot{x}/dt}{dx/dt} = \dot{x} \frac{d\dot{x}}{dx}.$$

Using the above, we represent equation (2.3) as

$$\dot{x} \frac{d\dot{x}}{dx} + \omega^2 x = 0.$$

Separating the variables yields

$$\frac{\dot{x}}{\omega^2} d\dot{x} + x dx = 0.$$

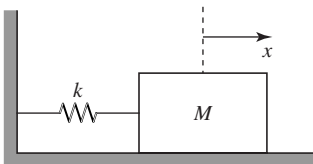


Figure 2.1 A spring–mass system of Example 2.1.

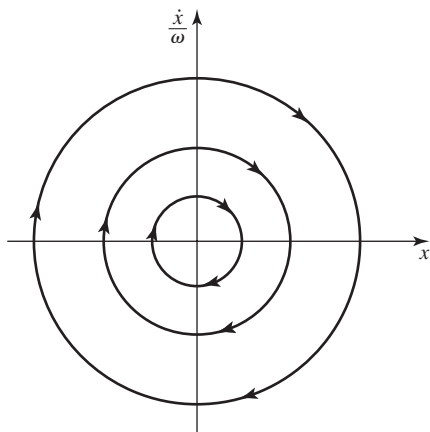


Figure 2.2 A phase portrait for the spring–mass system of Example 2.1.

Integrating the above gives

$$\left(\frac{\dot{x}}{\omega}\right)^2 + x^2 = A^2,$$

where A is the integration constant determined by the initial conditions. We see that a phase portrait of equation (2.3) consists of a family of circles centered at the origin of the $(x, \dot{x}/\omega)$ plane as illustrated in Figure 2.2. Different circular trajectories, in this plane, correspond to different values of the constant A representing initial conditions. Note that

$$x = \omega \int \frac{\dot{x}}{\omega} dt,$$

where $\omega > 0$. Hence, for positive values of \dot{x}/ω the displacement x increases. This means that for the values $\dot{x}/\omega > 0$ the arrows on the trajectories indicate x moving to the right. For negative values of \dot{x}/ω —that is, for the values of \dot{x}/ω below the horizontal axis—the arrows on the trajectories indicate x decreasing by going to the left. Therefore, as time increases, the RP moves in a clockwise direction along the trajectory.

◆ Example 2.2

In this example, we consider a model of a satellite shown in Figure 2.3. We assume that the satellite is rigid and is in a frictionless environment. It can rotate about an axis perpendicular to the page as a result of torque τ applied to the satellite by firing the thrusters. The system input is the applied torque, and the system output is the attitude angle θ . The satellite's moment of inertia is I . Applying the rotational

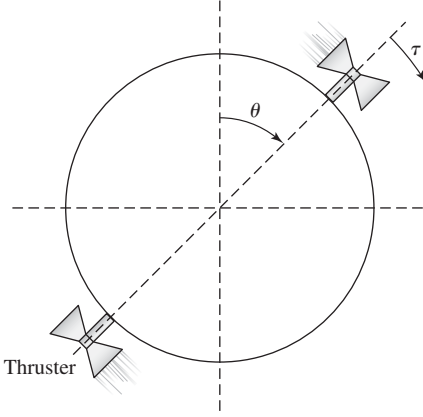


Figure 2.3 A model of a rigid satellite of Example 2.2.

version of Newton's law, we obtain the satellite's equation of motion:

$$I \frac{d^2\theta}{dt^2} = \tau,$$

which can be represented in the following equivalent form:

$$\frac{d^2\theta}{dt^2} = \frac{1}{I}\tau.$$

Let

$$\frac{1}{I}\tau = u,$$

then the equation describing the satellite's motion becomes

$$\frac{d^2\theta}{dt^2} = u. \quad (2.4)$$

We assume that when the thrusters fire, the thrust is constant; that is, $u = \pm U$. For simplicity, let $U = 1$. Using the state variables $x_1 = \theta$ and $x_2 = \dot{\theta}$, we obtain the satellite's state-space model:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u, \quad (2.5)$$

where $u = \pm 1$. We first consider the case when $u = 1$. With this control, the system equations are

$$\begin{aligned} \dot{x}_1 &= x_2, \\ \dot{x}_2 &= 1. \end{aligned}$$

Solving the second equation yields

$$x_2 = t + c_1,$$

where c_1 is an integration constant. We substitute $x_2 = t + c_1$ into $\dot{x}_1 = x_2$ to get

$$\dot{x}_1 = t + c_1.$$

Solving the above equation, we obtain

$$\begin{aligned} x_1 &= \int (t + c_1) dt \\ &= \frac{1}{2}(t + c_1)^2 + c_2 \\ &= \frac{1}{2}x_2^2 + c_2, \end{aligned}$$

where c_2 is an integration constant. In the state plane the above equation represents a family of parabolas open toward the positive x_1 axis. The representative point moves upward along these parabolas because $\dot{x}_2 = 1$; that is, $\dot{x}_2 > 0$. Typical system trajectories for this case are shown in Figure 2.4.

Similarly, when $u = -1$, we have

$$\dot{x}_2 = -1.$$

Solving the above differential equation yields

$$x_2 = -t + \tilde{c}_1,$$

where \tilde{c}_1 is an integration constant. Substituting $x_2 = -t + \tilde{c}_1$ into $\dot{x}_1 = x_2$ and then solving the resulting equation yields

$$x_1 = -\frac{1}{2}x_2^2 + \tilde{c}_2,$$

where \tilde{c}_2 is an integration constant. In the state plane, this equation represents

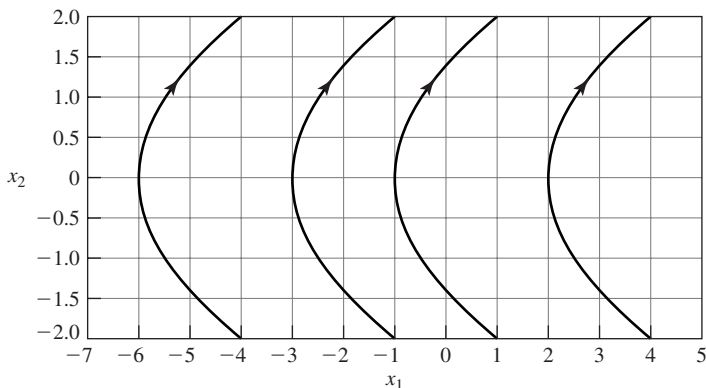


Figure 2.4 Trajectories of the system (2.5) for $u = 1$.

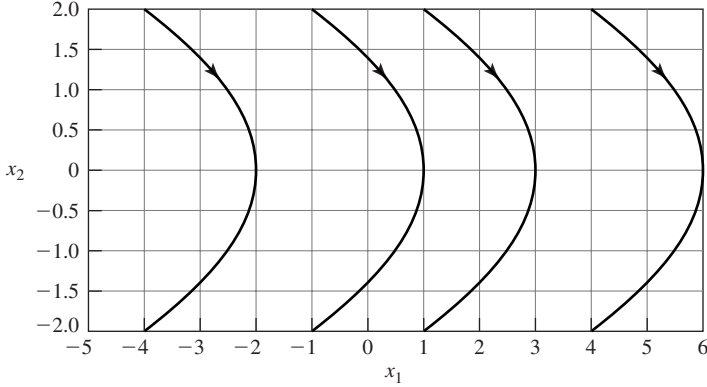


Figure 2.5 Trajectories of the system (2.5) for $u = -1$.

a family of parabolas open toward the negative x_1 axis. The representative point moves downward along these parabolas because $\dot{x}_2 = -1$; that is, $\dot{x} < 0$. Typical system trajectories for this case are shown in Figure 2.5.

In the above examples, systems' trajectories were constructed analytically. Analytical methods are useful for systems modeled by differential equations that can be easily solved. If the system of differential equations cannot be easily solved analytically, we can use graphical or numerical methods. We now describe a graphical method for solving second-order differential equations or a system of two first-order differential equations.

2.1.2 The Method of Isoclines

An *isocline* is a curve along which the trajectories' slope is constant. In other words, an isocline is a locus (a set of points), of the trajectories' constant slope. To obtain equations that describe isoclines, we consider a system of two first-order differential equations of the form

$$\begin{aligned}\dot{x}_1 &= \frac{dx_1}{dt} = f_1(x_1, x_2), \\ \dot{x}_2 &= \frac{dx_2}{dt} = f_2(x_1, x_2).\end{aligned}\tag{2.6}$$

Dividing the second of the above equations by the first one, we obtain

$$\frac{dx_2}{dx_1} = \frac{f_2(x_1, x_2)}{f_1(x_1, x_2)}.\tag{2.7}$$

Thus, we eliminated the independent variable t from the set of the first-order differential equations given by (2.6). In equation (2.7), we consider x_1 and x_2 as the independent and

dependent variables, respectively. Let

$$m = m(x_1, x_2) = \frac{dx_2}{dx_1}.$$

Note that m is just the slope of the tangent to the trajectory passing through the point $[x_1 \ x_2]^T$. The locus of the trajectory constant slope,

$$\frac{dx_2}{dx_1} = m(x_1, x_2) = \text{constant},$$

can be obtained from the equation

$$f_2(x_1, x_2) = m f_1(x_1, x_2).$$

The curve that satisfies the above equation is an isocline corresponding to the trajectories' slope m because a trajectory crossing the isocline will have its slope equal to m . The idea of the method of isoclines is to construct several isoclines in the state plane. This then allows one to construct a field of local tangents m . Then, the trajectory passing through any given point in the state plane is obtained by drawing a continuous curve following the directions of the field.

◆ Example 2.3

Consider the spring–mass system from Example 2.1. Defining the state variables $x_1 = x$ and $x_2 = \dot{x}$, we obtain the state equations

$$\begin{aligned}\dot{x}_1 &= x_2, \\ \dot{x}_2 &= -\omega^2 x_1.\end{aligned}$$

Let, for simplicity, $\omega = 1$. Then,

$$\frac{dx_2}{dx_1} = m = -\frac{x_1}{x_2}.$$

From the above we obtain the equation of isoclines:

$$x_2 = -\frac{1}{m}x_1.$$

This equation describes a family of straight lines through the origin. In Figure 2.6, we show some of the isoclines and a typical trajectory for the above spring–mass system. Note that in this example the trajectory slope is m , while the isocline slope is $-1/m$. This means that trajectories of the spring–mass system are perpendicular to the isoclines; that is, the trajectories are circles centered at the origin of the state plane.

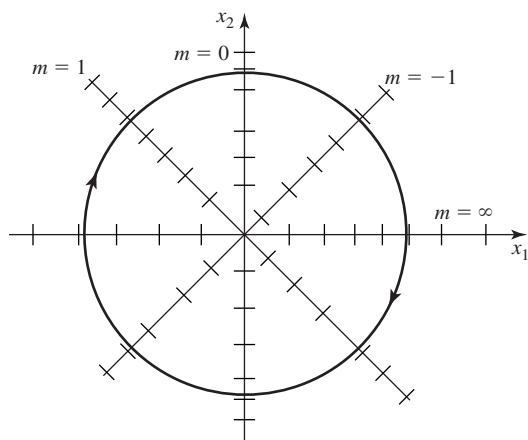


Figure 2.6 Isoclines and a typical trajectory for the spring-mass system of Example 2.3.

◆ Example 2.4

Given a second-order differential equation

$$\ddot{x} = -4\eta\dot{x} - 4x, \quad (2.8)$$

where η is a parameter. We will find values of η for which there are isoclines along which the trajectory slope and the isocline slope are equal. Such isoclines, if they exist, are called the *asymptotes*.

Let $x_1 = x$ and $x_2 = \dot{x}$. Then, the second-order differential equation given by (2.8) can be represented as a system of first-order differential equations:

$$\begin{aligned} \dot{x}_1 &= x_2, \\ \dot{x}_2 &= -4\eta x_2 - 4x_1. \end{aligned}$$

The isocline equation is

$$\frac{dx_2}{dx_1} = \frac{-4\eta x_2 - 4x_1}{x_2} = m,$$

or, equivalently,

$$x_2 = -\left(\frac{4}{m+4\eta}\right)x_1. \quad (2.9)$$

A trajectory slope is equal to an isocline slope if and only if

$$-\frac{4}{m+4\eta} = m;$$

that is,

$$m^2 + 4\eta m + 4 = 0. \quad (2.10)$$

The discriminant of (2.10) is

$$\Delta = 4\sqrt{\eta^2 - 1}.$$

Hence, the condition for the existence of asymptotes is

$$|\eta| \geq 1. \quad (2.11)$$

◆ Example 2.5

Consider the system of differential equations

$$\begin{aligned} \dot{x}_1 &= x_1 + x_2 + 1, \\ \dot{x}_2 &= -x_2 + 2. \end{aligned}$$

- (a) Find the equation of the isocline corresponding to the trajectory slope $m = 1$.
- (b) Does the given system have asymptotes? If yes, then write the equations of the asymptotes.

The trajectory slope is

$$\frac{dx_2}{dx_1} = m = \frac{-x_2 + 2}{x_1 + x_2 + 1}.$$

Hence, the equation of the isocline corresponding to the trajectory slope $m = 1$ is obtained by solving

$$\frac{-x_2 + 2}{x_1 + x_2 + 1} = 1.$$

We obtain

$$x_2 = -\frac{1}{2}x_1 + \frac{1}{2}.$$

The trajectory slope is

$$m = \frac{-x_2 + 2}{x_1 + x_2 + 1}. \quad (2.12)$$

From (2.12), we obtain an equation of isoclines of the form

$$x_2 = -\frac{m}{m+1}x_1 - \frac{m-2}{m+1}.$$

An asymptote is the isocline for which its slope and the trajectory slope are equal.

Therefore, we can obtain the equations of asymptotes by solving

$$m = -\frac{m}{m+1}.$$

We obtain

$$m^2 + 2m = m(m+2) = 0.$$

Hence, the equations of asymptotes are

$$x_2 = 2 \quad \text{and} \quad x_2 = -2x_1 - 4.$$

We can use a computer to draw phase portraits interactively. Below is an example of MATLAB's script that was used to sketch a phase portrait of the van der Pol equation shown in Figure 2.7.

```
t0=0; % initial time
tf=20; % final time
tspan=tf-t0;
x0=[-4 -4]'; % initial condition
button=1;
% we will now draw internal axes
p=4*[-1 0;1 0];
clf;plot(p(:,1),p(:,2))
hold on
plot(p(:,2),p(:,1))
axis(4*[-1 1 -1 1])
while(button==1)
% select a new initial condition and press the left button. Press
% the right button if you want to stop.
[t,x]=ode45(@my_xdot,tspan,x0); % invoke your own m.file
plot(x(:,1),x(:,2))
[x1,x2,button]=ginput(1);
x0=[x1 x2]';
end
```

The van der Pol equation has the form

$$\ddot{x} - \mu(1 - x^2)\dot{x} + x = 0, \quad (2.13)$$

where μ is a parameter. We can equivalently represent the van der Pol equation as

$$\begin{aligned} \dot{x}_1 &= x_2, \\ \dot{x}_2 &= \mu(1 - x_1^2)x_2 - x_1. \end{aligned} \quad (2.14)$$

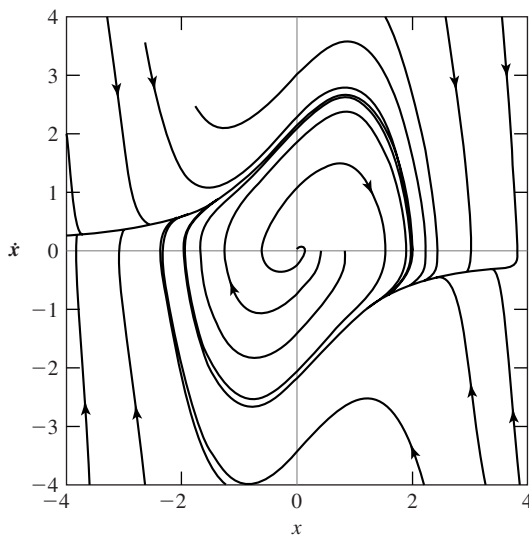


Figure 2.7 A phase-plane portrait of the van der Pol equation for $\mu = 1$.

An inspection of the phase portrait of the van der Pol equation reveals the existence of a closed trajectory in the state plane. This solution is an example of a limit cycle.

Definition 2.1 A limit cycle is a closed trajectory in the state plane, such that no other closed trajectory can be found arbitrarily close to it.

Note that the differential equation modeling the spring–mass system does not have a limit cycle because the solutions are not isolated. Limit cycles can only appear in nonlinear systems; they do not exist in linear systems.

Theorem 2.1 (Bendixson) If f_1 and f_2 in (2.7) have continuous partial derivatives in a domain \mathcal{A} bounded by a closed path \mathbf{C} and if $\partial f_1/\partial x_1 + \partial f_2/\partial x_2$ does not vanish and is of constant sign in that domain, then (2.7) can have no limit cycle in that domain.

Proof We prove the theorem by contradiction. Equation (2.7) is satisfied on any of its trajectory. Therefore, if (2.7) has a limit cycle, then clearly

$$f_1 dx_2 - f_2 dx_1 = 0. \quad (2.15)$$

We now apply Green's theorem, stated on page 674, to the vector field

$$[-f_2 \quad f_1]^T$$

to obtain

$$\int_{\mathbf{C}} f_1 dx_2 - f_2 dx_1 = \iint_{\mathcal{A}} \left(\frac{\partial f_1}{\partial x_1} + \frac{\partial f_2}{\partial x_2} \right) dx_1 dx_2. \quad (2.16)$$

By (2.15), the right-hand side of (2.16) must vanish, which contradicts the hypothesis. Thus, if $\partial f_1/\partial x_1 + \partial f_2/\partial x_2$ does not vanish and is of constant sign in a domain \mathcal{A} bounded by a closed path \mathbf{C} , then no limit cycle can exist in \mathcal{A} . The proof of the theorem is complete.

◆ Example 2.6

We apply Bendixson's theorem to the van der Pol equation represented by (2.14). We have

$$f_1 = x_2 \quad \text{and} \quad f_2 = \mu(1 - x_1^2)x_2 - x_1.$$

We calculate $\partial f_1 / \partial x_1 + \partial f_2 / \partial x_2$ to get

$$\partial f_1 / \partial x_1 + \partial f_2 / \partial x_2 = \mu(1 - x_1^2).$$

It follows from the Bendixson theorem that there are no limit cycles in the region $|x_1| < 1$ because $\partial f_1 / \partial x_1 + \partial f_2 / \partial x_2 > 0$ in that region.

Bendixson's theorem gives us a necessary condition for a closed curve to be a limit cycle. It is useful for the purpose of establishing the nonexistence of a limit cycle. In general, it is very difficult to establish the existence of a limit cycle. Sufficient conditions for the existence of a limit cycle are discussed, for example, by Hochstadt [123, Chapter 7]. A method for predicting limit cycles, called the describing function method, is presented in Section 2.5.

2.2 Numerical Techniques

In this section we discuss numerical techniques for solving differential equations. We begin our presentation of numerical techniques that allow us to determine approximate solutions to the given system of differential equations with a method based on a Taylor series expansion. This section is based on the Class Notes of Prof. V. B. Haas [109].

2.2.1 The Method of Taylor Series

The basis for many numerical methods for solving differential equations is the Taylor formula. Suppose we are given a state-space model of the form

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t)), \quad \mathbf{x}(t_0) = \mathbf{x}_0. \quad (2.17)$$

The method we will discuss depends on the fact that if $\mathbf{f}(t, \mathbf{x}, \mathbf{u})$ is differentiable with respect to its $n + m + 1$ variables $x_1, \dots, x_n, u_1, \dots, u_m, t$ in some $(n + m + 1)$ -dimensional region \mathcal{D} and if $\mathbf{x}(t)$ is a solution of (2.17), where $(t_0, \mathbf{x}_0, \mathbf{u}(t_0))$ lies in \mathcal{D} , then we can expand $\mathbf{x}(t)$ into a Taylor series to obtain

$$\begin{aligned} \mathbf{x}(t) &= \mathbf{x}_0 + \dot{\mathbf{x}}(t_0)(t - t_0) + \ddot{\mathbf{x}}(t_0)\frac{(t - t_0)^2}{2} + \dots \\ &= \mathbf{x}_0 + \mathbf{f}(t_0, \mathbf{x}_0, \mathbf{u}(t_0))(t - t_0) + \left[\frac{\partial \mathbf{f}(t_0, \mathbf{x}_0, \mathbf{u}(t_0))}{\partial t} + \sum_{i=1}^n \frac{\partial \mathbf{f}(t_0, \mathbf{x}_0, \mathbf{u}(t_0))}{\partial x_i} \dot{x}_i \right. \\ &\quad \left. + \sum_{j=1}^m \frac{\partial \mathbf{f}(t_0, \mathbf{x}_0, \mathbf{u}(t_0))}{\partial u_j} \dot{u}_j \right] \frac{(t - t_0)^2}{2} + \dots \end{aligned}$$

If we stop after the q th term of the Taylor series solution, then the remainder after q terms of a Taylor series gives us an estimate of the error. We can expect the Taylor series solution to give a fairly accurate approximation of the solution of the state-space equations only in a small interval about t_0 . Thus, if $\Delta t = t - t_0$ is some small time interval such that the Taylor series converges for $t \in [t_0, t_0 + \Delta t]$, we can use the expansion

$$\mathbf{x}(t) = \mathbf{x}_0 + \frac{d\mathbf{x}(t_0)}{dt}(t - t_0) + \frac{d^2\mathbf{x}(t_0)}{dt^2} \frac{(t - t_0)^2}{2} + \dots$$

to evaluate $\mathbf{x}(t_1) = \mathbf{x}(t_0 + \Delta t)$. Then, we can start over again solving $\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t))$ with the new initial condition $\mathbf{x}(t_1)$. We illustrate the method with simple examples.

◆ Example 2.7

Consider the first-order differential equation

$$\frac{dx}{dt} = t^2 - x^2, \quad x(1) = 1.$$

Expanding $x = x(t)$ into a Taylor series about $x(1) = 1$ yields

$$\begin{aligned} x(t) = x(1) &+ \frac{dx(1)}{dt}(t - 1) + \frac{d^2x(1)}{dt^2} \frac{(t - 1)^2}{2} + \frac{d^3x(1)}{dt^3} \frac{(t - 1)^3}{3!} \\ &+ \frac{d^4x(1)}{dt^4} \frac{(t - 1)^4}{4!} \dots \end{aligned}$$

Because $dx/dt = t^2 - x^2$, we obtain by successive differentiations

$$\begin{aligned} \frac{d^2x}{dt^2} &= 2t - 2x \frac{dx}{dt}, \\ \frac{d^3x}{dt^3} &= 2 - 2 \left(\frac{dx}{dt} \right)^2 - 2x \frac{d^2x}{dt^2}, \\ \frac{d^4x}{dt^4} &= -6 \frac{dx}{dt} \frac{d^2x}{dt^2} - 2x \frac{d^3x}{dt^3}, \\ &\vdots \end{aligned}$$

Evaluating the above at $t = 1$ yields

$$\begin{aligned} \frac{dx(1)}{dt} &= 0, \\ \frac{d^2x(1)}{dt^2} &= 2, \\ \frac{d^3x(1)}{dt^3} &= -2, \\ \frac{d^4x(1)}{dt^4} &= 4, \\ &\vdots \end{aligned}$$

Hence,

$$x(t) = 1 + (t-1)^2 - \frac{(t-1)^3}{3} + \frac{(t-1)^4}{6} + \dots$$

◆ Example 2.8

In this example, we consider the second-order differential equation

$$\frac{d^2 x}{dt^2} + 3 \frac{dx}{dt} + 6tx = 0,$$

subject to the initial conditions $x(0) = 1$ and $dx(0)/dt = 1$. Using $d^2 x/dt^2 = -3dx/dt - 6tx$, we obtain by successive differentiation

$$\begin{aligned} \frac{d^3 x}{dt^3} &= -3 \frac{d^2 x}{dt^2} - 6t \frac{dx}{dt} - 6x, \\ \frac{d^4 x}{dt^4} &= -3 \frac{d^3 x}{dt^3} - 6t \frac{d^2 x}{dt^2} - 12 \frac{dx}{dt}, \\ \frac{d^5 x}{dt^5} &= -3 \frac{d^4 x}{dt^4} - 6t \frac{d^3 x}{dt^3} - 18 \frac{d^2 x}{dt^2}, \\ &\vdots \end{aligned}$$

Taking into account the initial conditions, we evaluate all coefficients at $t = 0$ to get

$$\begin{aligned} \frac{d^2 x(0)}{dt^2} &= -3, \\ \frac{d^3 x(0)}{dt^3} &= 3, \\ \frac{d^4 x(0)}{dt^4} &= -21, \\ \frac{d^5 x(0)}{dt^5} &= 117, \\ &\vdots \end{aligned}$$

Hence,

$$\begin{aligned} x(t) &= x(0) + \frac{dx(0)}{dt}t + \frac{d^2 x(0)}{dt^2} \frac{t^2}{2} + \frac{d^3 x(0)}{dt^3} \frac{t^3}{3!} + \frac{d^4 x(0)}{dt^4} \frac{t^4}{4!} + \frac{d^5 x(0)}{dt^5} \frac{t^5}{5!} + \dots \\ &= 1 + t - \frac{3}{2}t^2 + \frac{1}{2}t^3 - \frac{7}{8}t^4 + \frac{117}{120}t^5 + \dots \end{aligned}$$

We now illustrate the method of Taylor series when applied to a system of first-order differential equations.

◆ Example 2.9

We consider the following system of two first-order equations:

$$\begin{aligned} dx_1/dt &= f_1(t, x_1, x_2), \\ dx_2/dt &= f_2(t, x_1, x_2). \end{aligned}$$

The initial conditions at $t = t_0$ are $x_1(t_0) = x_{10}$, $x_2(t_0) = x_{20}$. Having the initial conditions, we can determine the values $dx_1(t_0)/dt$ and $dx_2(t_0)/dt$ from the original equations. We then differentiate dx_1/dt and dx_2/dt to obtain

$$\begin{aligned} \frac{d^2 x_1}{dt^2} &= \frac{d}{dt} f_1 \left(t, x_1, x_2, \frac{dx_1}{dt}, \frac{dx_2}{dt} \right), \\ \frac{d^2 x_2}{dt^2} &= \frac{d}{dt} f_2 \left(t, x_1, x_2, \frac{dx_1}{dt}, \frac{dx_2}{dt} \right). \end{aligned}$$

Now, $d^2 x_1/dt^2$ and $d^2 x_2/dt^2$ can be calculated by substituting known quantities on the right-hand sides of the above equations. By successive differentiation and substitution, we can determine $d^3 x_1/dt^3$ and $d^3 x_2/dt^3$ as well as higher derivatives at $t = t_0$. The solution to the given system is then computed as

$$\begin{aligned} x_1(t) &= x_{10} + \frac{dx_1(t_0)}{dt}(t - t_0) + \frac{d^2 x_1(t_0)}{dt^2} \frac{(t - t_0)^2}{2} + \frac{d^3 x_1(t_0)}{dt^3} \frac{(t - t_0)^3}{3!} + \dots, \\ x_2(t) &= x_{20} + \frac{dx_2(t_0)}{dt}(t - t_0) + \frac{d^2 x_2(t_0)}{dt^2} \frac{(t - t_0)^2}{2} + \frac{d^3 x_2(t_0)}{dt^3} \frac{(t - t_0)^3}{3!} + \dots \end{aligned}$$

The Taylor series method can be programmed on a computer. However, this is not a very efficient method from a numerical point of view. In the following, we discuss other methods for numerical solution of the state equations that are related to the method of Taylor series. We first present two simple numerical techniques known as the forward and backward Euler methods.

2.2.2 Euler's Methods

Suppose that we wish to obtain the solution of the modeling equation

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t))$$

in the time interval $[t_0, t_f]$, subject to the initial condition $\mathbf{x}(t_0) = \mathbf{x}_0$ and the input vector $\mathbf{u}(t)$. We divide the interval $[t_0, t_f]$ into N equal subintervals of width

$$h = \frac{t_f - t_0}{N}.$$

We call h the *step length*. We set

$$t_k = t_0 + kh.$$

To proceed further, we approximate the derivative $d\mathbf{x}(t)/dt$ at time t_k by

$$\frac{d\mathbf{x}(t_k)}{dt} \approx \frac{\mathbf{x}(t_{k+1}) - \mathbf{x}(t_k)}{h}$$

and write

$$\frac{\mathbf{x}(t_{k+1}) - \mathbf{x}(t_k)}{h} = \mathbf{f}(t_k, \mathbf{x}(t_k), \mathbf{u}(t_k)).$$

Hence, for $k = 1, 2, \dots, N$, we have

$$\mathbf{x}(t_{k+1}) = \mathbf{x}(t_k) + h\mathbf{f}(t_k, \mathbf{x}(t_k), \mathbf{u}(t_k))$$

The above is known as the *forward Euler* algorithm. We see that if h is sufficiently small, we can approximately determine the state at time $t_1 = t_0 + h$ from the initial condition \mathbf{x}_0 to get

$$\mathbf{x}(t_1) = \mathbf{x}_0 + h\mathbf{f}(t_0, \mathbf{x}_0, \mathbf{u}(t_0)).$$

Once we have determined the approximate solution at time t_1 , we can determine the approximate solution at time $t_2 = t_1 + h$, and so on. The forward Euler method can also be arrived at when instead of considering the differential equation $\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t))$, we start with its equivalent integral representation:

$$\mathbf{x}(t) = \mathbf{x}_0 + \int_{t_0}^t \mathbf{f}(\tau, \mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau. \quad (2.18)$$

We then use the *rectangular rule* for numerical integration that can be stated as follows. The area under the curve $z = g(t)$ between $t = t_k$ and $t = t_k + h$ is approximately equal to the area of the rectangle $ABCD$ as shown in Figure 2.8; that is,

$$\int_{t_k}^{t_k+h} g(t) dt \approx hg(t_k). \quad (2.19)$$

Applying the rectangular rule of integration to (2.18), we obtain the forward Euler method.

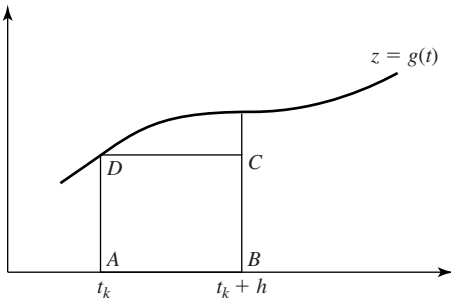


Figure 2.8 Illustration of the rectangular rule of integration.

The *backward Euler* method differs from the forward Euler method in the way we approximate the derivative:

$$\frac{d\mathbf{x}(t_k)}{dt} \approx \frac{\mathbf{x}(t_k) - \mathbf{x}(t_{k-1})}{h}.$$

This approximation gives the backward Euler formula:

$$\mathbf{x}(t_{k+1}) = \mathbf{x}(t_k) + h\mathbf{f}(t_{k+1}, \mathbf{x}(t_{k+1}), \mathbf{u}(t_{k+1}))$$

In the above formula, $\mathbf{x}(t_{k+1})$ is a function of itself. For this reason, we say that the backward Euler method is an implicit integration algorithm. Implicit integration algorithms require additional computation to solve for $\mathbf{x}(t_{k+1})$. For further discussion of this issue, we refer to Parker and Chua [222].

Euler's methods are rather elementary, and thus they are not as accurate as some of the more sophisticated techniques that we are going to discuss next.

2.2.3 Predictor–Corrector Method

The predictor–corrector method, also known as Heun's method, is based on the trapezoidal rule for numerical integration. The trapezoidal rule states that the area under the curve $z = g(t)$ between $t = t_k$ and $t = t_k + h$ is approximately equal to the area of the trapezoid $ABCD$ as depicted in Figure 2.9; that is,

$$\int_{t_k}^{t_k+h} g(t) dt \approx h \frac{g(t_k) + g(t_k + h)}{2}. \quad (2.20)$$

Suppose we wish to find $\mathbf{x}(t)$ at $t = t_1 > t_0$. Applying the trapezoid rule of integration to (2.17), we obtain

$$\mathbf{x}(t_1) = \mathbf{x}_0 + \frac{h}{2}(\mathbf{f}(t_0, \mathbf{x}_0, \mathbf{u}(t_0)) + \mathbf{f}(t_1, \mathbf{x}(t_1), \mathbf{u}(t_1))).$$

The above is an implicit algorithm, and it may be viewed as a merging of the forward and backward Euler algorithms. If the function \mathbf{f} is nonlinear, then we will, in general, not be able

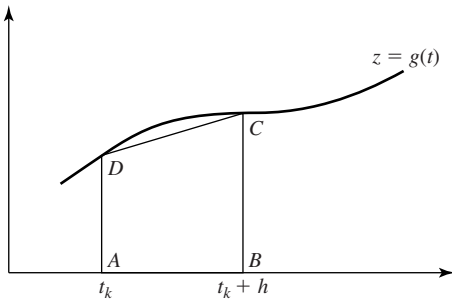


Figure 2.9 The trapezoidal rule of integration.

to solve the above equation for $\mathbf{x}(t_1)$ exactly. We attempt to obtain $\mathbf{x}(t_1)$ iteratively. Let

$$\begin{aligned} h &= t_{k+1} - t_k, \\ \mathbf{x}_1^* &= \mathbf{x}_0 + h \mathbf{f}(t_0, \mathbf{x}_0, \mathbf{u}(t_0)), \\ \mathbf{m}_0 &= \mathbf{f}(t_0, \mathbf{x}_0, \mathbf{u}(t_0)), \\ \mathbf{m}_1 &= \mathbf{f}(t_1, \mathbf{x}_1^*, \mathbf{u}(t_1)). \end{aligned}$$

A first approximation to $\mathbf{x}(t_1)$ is

$$\mathbf{x}_1^* = \mathbf{x}_0 + h \mathbf{f}(t_0, \mathbf{x}_0, \mathbf{u}(t_0))$$

and is called a *predictor*. The next approximation is

$$\mathbf{x}(t_1) = \mathbf{x}_0 + \frac{h}{2}(\mathbf{m}_0 + \mathbf{m}_1)$$

and is called a *corrector*.

Suppose now that we wish to find $\mathbf{x}(t)$ at $t = t_f > t_0$. We first divide the interval $[t_0, t_f]$ into N equal subintervals using evenly spaced subdivision points to obtain

$$t_0 < t_1 < t_2 < \cdots < t_N = t_f.$$

Once we have found $\mathbf{x}(t_1)$, we can apply the same procedure to find $\mathbf{x}(t_2)$, and so on.

◆ Example 2.10

We use the predictor–corrector method to find $x(1)$, where the step length h equals 0.5, for the following differential equation:

$$\frac{dx}{dt} = tx^2 - x, \quad x(0) = 1. \quad (2.21)$$

The calculations are summarized in Table 2.1.

Table 2.1 Solving the Differential Equation (2.21) Using the Predictor–Corrector Method

t_0	x_0	t_1	m_0	x_1^*	m_1	$m = (m_0 + m_1)/2$	Δx
0	1	0.5	−1.0	0.5	−0.375	−0.688	−0.344
0.5	0.656	1.0	−0.441	0.435	−0.246	−0.344	−0.172
1.0	0.484						

We could use $\mathbf{x}(t_1)$ obtained from the corrector formula to construct a next approximation to $\mathbf{x}(t_1)$. In general, we can use the formula

$$\mathbf{x}^{(k)}(t_1) = \mathbf{x}_0 + \frac{h}{2} (\mathbf{f}(t_0, \mathbf{x}_0, \mathbf{u}(t_0)) + \mathbf{f}(t_1, \mathbf{x}^{(k-1)}(t_1), \mathbf{u}(t_1)))$$

to evaluate $\mathbf{x}(t_1)$ until two successive iterations agree to the desired accuracy. This general algorithm is referred to as a second-order predictor–corrector method.

An n -dimensional system of first-order differential equations can be worked analogously. We illustrate the procedure for $n = 2$, where

$$\begin{aligned} \frac{dx}{dt} &= f(t, x, y), & x(t_0) &= x_0, \\ \frac{dy}{dt} &= g(t, x, y), & y(t_0) &= y_0. \end{aligned}$$

Let

$$\begin{aligned} m_0 &= f(t_0, x_0, y_0), \\ n_0 &= g(t_0, x_0, y_0). \end{aligned}$$

Then, the predictor equations are

$$\begin{aligned} x_1^* &= x_0 + m_0 h, \\ y_1^* &= y_0 + n_0 h. \end{aligned}$$

Let

$$\begin{aligned} m_1 &= f(t_1, x_1^*, y_1^*), \\ n_1 &= g(t_1, x_1^*, y_1^*). \end{aligned}$$

Then, the corrector equations are

$$\begin{aligned} x(t_1) &= x_0 + \frac{m_0 + m_1}{2} h, \\ y(t_1) &= y_0 + \frac{n_0 + n_1}{2} h. \end{aligned}$$

Using the above, we can generate the corresponding formulas for t_k .

We will now show that the predictor–corrector method yields a solution that agrees with the Taylor series solution through terms of degree two. We consider the scalar case, where f is assumed to be differentiable with respect to its arguments and $\partial^2 f / \partial t \partial x = \partial^2 f / \partial x \partial t$. First observe that $m_1 = f(t_1, x_1^*)$ can be written as

$$\begin{aligned} f(t_1, x_1^*) &= f(t_0 + h, x_0 + m_0 h) \\ &\approx f(t_0, x_0) + [f_t \quad f_x] \begin{bmatrix} h \\ m_0 h \end{bmatrix} + \frac{1}{2} [h \quad m_0 h] \begin{bmatrix} f_{tt} & f_{tx} \\ f_{xt} & f_{xx} \end{bmatrix} \begin{bmatrix} h \\ m_0 h \end{bmatrix}, \end{aligned}$$

where $f_t = \partial f / \partial t$, $f_x = \partial f / \partial x$, and so on. Substituting the above into the corrector equation

yields

$$\begin{aligned}
 x(t_1) &= x_0 + \frac{h}{2}(m_0 + m_1) \\
 &= x_0 + \frac{h}{2}(m_0 + m_0) + \frac{h}{2} \left([f_t \quad f_x] \begin{bmatrix} h \\ m_0 h \end{bmatrix} \right. \\
 &\quad \left. + \frac{1}{2} [h \quad m_0 h] \begin{bmatrix} f_{tt} & f_{tx} \\ f_{xt} & f_{xx} \end{bmatrix} \begin{bmatrix} h \\ m_0 h \end{bmatrix} \right) \\
 &= x_0 + m_0 h + (f_t + f_x m_0) \frac{h^2}{2} + (f_{tt} + 2f_{tx} m_0 + f_{xx} m_0^2) \frac{h^3}{4}.
 \end{aligned}$$

The Taylor series expansion, on the other hand, gives

$$\begin{aligned}
 x(t_1) &= x_0 + \dot{x} h + \ddot{x} \frac{h^2}{2} + x^{(3)} \frac{h^3}{6} + \cdots \\
 &= x_0 + f(t_0, x_0) h + (f_t + f_x \dot{x}) \frac{h^2}{2} \\
 &\quad + (f_{tt} + 2f_{tx} \dot{x} + f_{xx} (\dot{x})^2 + f_x \ddot{x}) \frac{h^3}{6} + \cdots,
 \end{aligned}$$

where f_t , f_x , \dot{x} , and so on, are all evaluated at $t = t_0$. Comparing the right-hand sides of the last two equations, we see that the solution obtained using the predictor–corrector method agrees with the Taylor series expansion of the true solution through terms of degree two.

2.2.4 Runge's Method

Runge's method is an improvement over the predictor–corrector method in the sense that it agrees with the Taylor series solution through terms of degree three. The method depends upon integration by means of Simpson's rule, which says that the area under the curve $x = g(t)$ between t_0 and $t_0 + h$ is given approximately by

$$\int_{t_0}^{t_0+h} g(t) dt \approx \frac{h}{6}(x_0 + 4x_1 + x_2), \quad (2.22)$$

where $x_0 = g(t_0)$, $x_1 = g(t_0 + h/2)$, and $x_2 = g(t_0 + h)$. The above formula is derived by passing a parabola through the three points (t_0, x_0) , $(t_0 + h/2, x_1)$, and $(t_0 + h, x_2)$ and assuming that the area under the given curve is approximately equal to the area under the parabola, as illustrated in Figure 2.10. To verify Simpson's formula, we translate the axes so that in the new coordinates (\tilde{t}, \tilde{x}) we have $t = t_0 + h/2 + \tilde{t}$, $x = \tilde{x}$. In the new coordinates the three points become

$$(-h/2, x_0), \quad (0, x_1), \quad (h/2, x_2).$$

The parabola through the above points is described by the equation

$$\tilde{x} = a\tilde{t}^2 + b\tilde{t} + c.$$

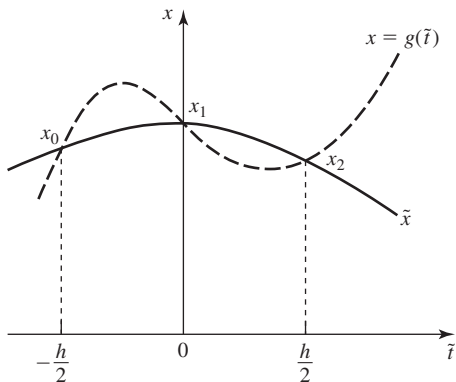


Figure 2.10 Integration using Simpson's rule.

Hence

$$\begin{aligned}x_0 &= a \frac{h^2}{4} - b \frac{h}{2} + c, \\x_1 &= c, \\x_2 &= a \frac{h^2}{4} + b \frac{h}{2} + c.\end{aligned}$$

Using the above three equations, we calculate

$$\begin{aligned}c &= x_1, \\b &= \frac{x_2 - x_0}{h}, \\a &= \frac{2}{h^2}(x_0 - 2x_1 + x_2).\end{aligned}$$

Then, the area under the parabola is

$$\begin{aligned}\int_{-h/2}^{h/2} (a\tilde{t}^2 + b\tilde{t} + c) d\tilde{t} &= \int_{-h/2}^{h/2} \left(\frac{2}{h^2}(x_0 - 2x_1 + x_2)\tilde{t}^2 + \frac{x_2 - x_0}{h}\tilde{t} + x_1 \right) d\tilde{t} \\&= \frac{h}{6}(x_0 + 4x_1 + x_2).\end{aligned}$$

If we are given the differential equation

$$\frac{dx}{dt} = f(t, x), \quad x(t_0) = x_0,$$

we integrate it to obtain

$$x(t_1) = x(t_0) + \int_{t_0}^{t_1} f(s, x(s)) ds.$$

Then, we approximate the integral

$$\int_{t_0}^{t_1} f(s, x(s)) ds$$

by

$$\frac{h}{6}(m_0 + 4m_1 + m_3),$$

where

$$\begin{aligned} h &= t_1 - t_0, \\ m_0 &= f(t_0, x_0), \\ m_1 &= f\left(t_0 + \frac{h}{2}, x_0 + m_0 \frac{h}{2}\right), \\ m_2 &= f(t_0 + h, x_0 + m_0 h), \\ m_3 &= f(t_0 + h, x_0 + m_2 h). \end{aligned}$$

Thus, we may write

$$x(t_1) = x_0 + \Delta x,$$

where $\Delta x = mh$, and $m = \frac{1}{6}(m_0 + 4m_1 + m_3)$.

◆ Example 2.11

We use Runge's method to find $x(1)$, where the step length h equals 0.5, for the differential equation (2.21). We found $x(1)$ using the predictor–corrector method in Example 2.10. The calculations for Runge's method are summarized in Table 2.2.

Table 2.2 Solving the Differential Equation (2.21) Using Runge's Method

t_0	x_0	m_0	m_1	m_2	m_3	m	Δx
0	1	−1	−0.610	−0.375	−0.482	−0.654	−0.327
0.5	0.673	−0.446	−0.325	−0.248	−0.248	−0.332	−0.166
1.0	0.507						

Runge's method can, of course, be used for systems of first-order differential equations. If, for example,

$$\begin{aligned} \frac{dx}{dt} &= f(t, x, y), & x(t_0) &= x_0, \\ \frac{dy}{dt} &= g(t, x, y), & y(t_0) &= y_0, \end{aligned}$$

then

$$\begin{aligned}x(t_1) &= x_0 + \Delta x, \\y(t_1) &= y_0 + \Delta y,\end{aligned}$$

where

$$\begin{aligned}\Delta x &= \frac{h}{6}(m_0 + 4m_1 + m_3), \\ \Delta y &= \frac{h}{6}(n_0 + 4n_1 + n_3).\end{aligned}$$

In the above expressions,

$$\begin{aligned}m_0 &= f(t_0, x_0, y_0), \\ m_1 &= f(t_0 + h/2, x_0 + m_0 h/2, y_0 + n_0 h/2), \\ m_2 &= f(t_0 + h, x_0 + m_0 h, y_0 + n_0 h), \\ m_3 &= f(t_0 + h, x_0 + m_2 h, y_0 + n_2 h), \\ n_0 &= g(t_0, x_0, y_0), \\ n_1 &= g(t_0 + h/2, x_0 + m_0 h/2, y_0 + n_0 h/2), \\ n_2 &= g(t_0 + h, x_0 + m_0 h, y_0 + n_0 h), \\ n_3 &= g(t_0 + h, x_0 + m_2 h, y_0 + n_2 h).\end{aligned}$$

2.2.5 Runge–Kutta Method

This method is an improvement over the Runge method, and it agrees with the Taylor series solution through terms of degree four; for this reason it is sometimes called the fourth-order Runge–Kutta algorithm. It is based on the use of the formula

$$x(t_1) = x(t_0) + \frac{h}{6}(m_0 + 2m_1 + 2m_2 + m_3), \quad (2.23)$$

where

$$\begin{aligned}m_0 &= f(t_0, x_0), \\ m_1 &= f(t_0 + h/2, x_0 + m_0 h/2), \\ m_2 &= f(t_0 + h/2, x_0 + m_1 h/2), \\ m_3 &= f(t_0 + h, x_0 + h m_2).\end{aligned}$$

Thus, we may write

$$x(t_1) = x_0 + \Delta x,$$

where $\Delta x = mh$, and $m = \frac{1}{6}(m_0 + 2m_1 + 2m_2 + m_3)$.

◆ Example 2.12

We apply the Runge–Kutta method to find $x(1)$, where the step length h equals 0.5, for the differential equation (2.21) that we used Examples 2.10 and 2.11. Our calculations for the Runge–Kutta method are summarized in Table 2.3.

Table 2.3 Solving the Differential Equation (2.21) Using the Runge–Kutta Method

t_0	x_0	m_0	m_1	m_2	m_3	m	Δx
0	1	−1.0	−0.609	−0.668	−0.444	−0.666	−0.333
0.5	0.667	−0.444	−0.324	−0.328	−0.250	−0.333	−0.167
1.0	0.500						

The above problem has been worked by the predictor–corrector, Runge’s, and the Runge–Kutta methods. This problem can be solved exactly by means of the Taylor series. The Taylor series solution is

$$x(t) = 1 - t + t^2 - t^3 + \cdots$$

The series converges to $1/(1+t)$ for $-1 < t < 1$. A check on the differential equation shows that

$$x(t) = \frac{1}{1+t}$$

is indeed a solution for all $t > -1$. Thus, $x(1) = 0.500$.

The algorithms presented above require only one input point, $x(t_k)$, at each step to compute $x(t_{k+1})$. Such algorithms are referred to as single-step algorithms. In the multistep algorithms, a number of previous points are used, along with the corresponding values of f , to evaluate $x(t_{k+1})$; that is, the multistep algorithms reuse past information about the trajectory when evaluating its new point. In general, it may be difficult to decide which type of algorithm is more efficient because the performance of a particular algorithm is problem dependent. Well-known multistep algorithms are: Adams–Bashforth, Adams–Moulton, and Gear. More information about the above-mentioned algorithms can be found in Parker and Chua [222, Section 4.1].

2.3 Principles of Linearization

The replacement of a nonlinear system model by its linear approximation is called *linearization*. The motivation for linearization is that the dynamical behavior of many nonlinear system models can be well approximated within some range of variables by linear system models. Then, we can use well-developed techniques for analysis and synthesis of linear systems to analyze a nonlinear system at hand. However, the results of analysis of nonlinear systems using their linearized models should be carefully interpreted in order to avoid unacceptably large errors due to the approximation in the process of linearization.

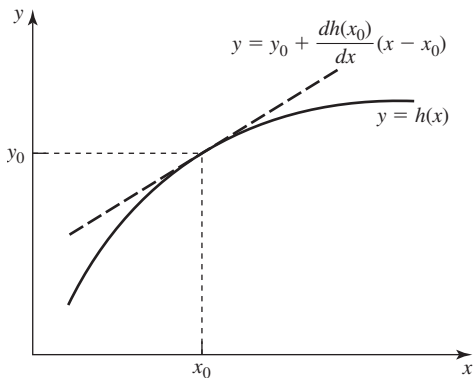


Figure 2.11 Linear approximation to a function $y = h(x)$.

We begin our discussion of linear approximations of nonlinear systems by considering a nonlinear element with a state variable x and an output variable y that are related by the equation

$$y = h(x),$$

where the function $h : \mathbb{R} \rightarrow \mathbb{R}$ is continuously differentiable; that is, $h \in \mathcal{C}^1$. Let x_0 be an operating point. We expand h into a Taylor series about the operating point x_0 to obtain

$$\begin{aligned} y &= h(x) \\ &= h(x_0) + \frac{dh(x_0)}{dx}(x - x_0) + \text{higher-order terms.} \end{aligned}$$

Linearization of $h(x)$ about the point x_0 consists of replacing h by the linear approximation

$$\begin{aligned} y &= h(x_0) + \frac{dh(x_0)}{dx}(x - x_0) \\ &= y_0 + \frac{dh(x_0)}{dx}(x - x_0), \end{aligned} \tag{2.24}$$

where $y_0 = h(x_0)$. Let $\Delta y = y - y_0$ and $\Delta x = x - x_0$. Then, we can represent (2.24) as

$$\Delta y = \frac{dh(x_0)}{dx} \Delta x.$$

Over a small range of Δx , the line (2.24) is a good approximation to the curve $y = h(x)$ in a neighborhood of the operating point x_0 ; see Figure 2.11 for an illustration of the above statement. We now illustrate the process of linearization on the simple pendulum.

◆ Example 2.13

Consider the simple pendulum shown in Figure 2.12. Let g be the gravity constant. The tangential force component, $mg \sin \theta$, acting on the mass m returns the pendulum to its equilibrium position. By Newton's second law we obtain

$$F = -mg \sin \theta.$$

The equilibrium point for the pendulum is $\theta = 0^\circ$. Note that $F(0^\circ) = 0$. Therefore,

$$\Delta F = F \quad \text{and} \quad \Delta \theta = \theta.$$

We also have

$$\frac{dF}{d\theta}(0^\circ) = -mg \cos 0^\circ = -mg.$$

Hence, the linearized model about the equilibrium position $\theta = 0^\circ$ has the form

$$F = -mg\theta.$$

Thus, for small angular displacements θ , the force F is proportional to the displacements. This approximation is quite accurate for $-\pi/4 \leq \theta \leq \pi/4$, as can be seen in Figure 2.13.

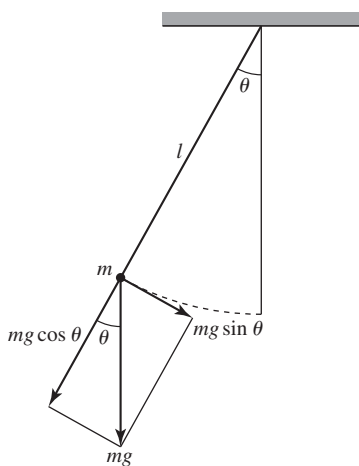


Figure 2.12 Forces acting on the simple pendulum.

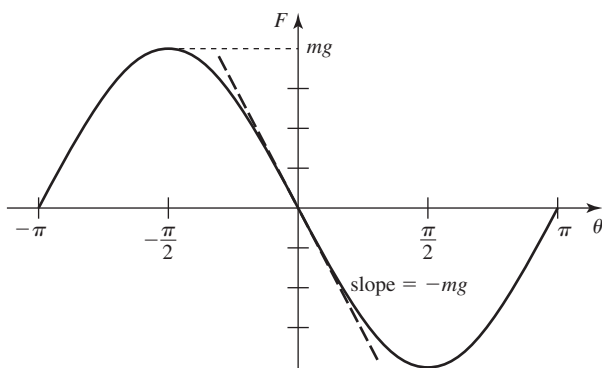


Figure 2.13 Linearizing the equation modeling the simple pendulum.

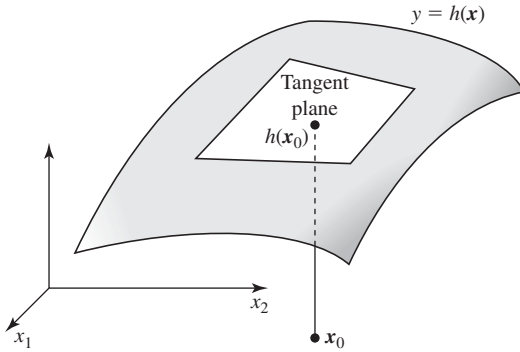


Figure 2.14 Tangent plane as a linear approximation for a function of two variables.

If $h : \mathbb{R}^n \rightarrow \mathbb{R}$ —that is, $y = h(x_1, x_2, \dots, x_n)$, which means that the dependent variable depends upon several variables—the same principle applies. Let

$$\mathbf{x}_0 = [x_{10} \quad x_{20} \quad \cdots \quad x_{n0}]^T$$

be the operating point. The Taylor series expansion of h about the operating point \mathbf{x}_0 yields

$$y - h(\mathbf{x}_0) = \nabla h(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) + \text{higher-order terms},$$

where

$$\nabla h(\mathbf{x}_0)^T = \left[\left. \frac{\partial h}{\partial x_1} \right|_{\mathbf{x}=\mathbf{x}_0} \quad \left. \frac{\partial h}{\partial x_2} \right|_{\mathbf{x}=\mathbf{x}_0} \quad \cdots \quad \left. \frac{\partial h}{\partial x_n} \right|_{\mathbf{x}=\mathbf{x}_0} \right].$$

Geometrically, the linearization of h about \mathbf{x}_0 can be thought of as placing a tangent plane onto the nonlinear surface at the operating point \mathbf{x}_0 as illustrated in Figure 2.14.

2.4 Linearizing Differential Equations

We now consider a dynamical system modeled by a set of nonlinear differential equations:

$$\begin{aligned} \frac{dx_1}{dt} &= f_1(x_1, x_2, \dots, x_n, u_1, u_2, \dots, u_m), \\ \frac{dx_2}{dt} &= f_2(x_1, x_2, \dots, x_n, u_1, u_2, \dots, u_m), \\ &\vdots \\ \frac{dx_n}{dt} &= f_n(x_1, x_2, \dots, x_n, u_1, u_2, \dots, u_m). \end{aligned}$$

We assume that the functions $f_i, i = 1, 2, \dots, n$, are continuously differentiable. The above set of differential equations can be represented in vector form as

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, \mathbf{u}). \quad (2.25)$$

Let $\mathbf{u}_e = [u_{1e} \quad u_{2e} \quad \cdots \quad u_{me}]^T$ be a constant input that forces the system (2.25) to settle into a constant equilibrium state $\mathbf{x}_e = [x_{1e} \quad x_{2e} \quad \cdots \quad x_{ne}]^T$; that is, \mathbf{u}_e and \mathbf{x}_e satisfy

$$\mathbf{f}(\mathbf{x}_e, \mathbf{u}_e) = \mathbf{0}.$$

We now perturb the equilibrium state by allowing

$$\mathbf{x} = \mathbf{x}_e + \Delta \mathbf{x}, \quad \mathbf{u} = \mathbf{u}_e + \Delta \mathbf{u}.$$

Taylor's expansion yields

$$\begin{aligned} \frac{d}{dt} \mathbf{x} &= \mathbf{f}(\mathbf{x}_e + \Delta \mathbf{x}, \mathbf{u}_e + \Delta \mathbf{u}) \\ &= \mathbf{f}(\mathbf{x}_e, \mathbf{u}_e) + \frac{\partial \mathbf{f}}{\partial \mathbf{x}}(\mathbf{x}_e, \mathbf{u}_e) \Delta \mathbf{x} + \frac{\partial \mathbf{f}}{\partial \mathbf{u}}(\mathbf{x}_e, \mathbf{u}_e) \Delta \mathbf{u} + \text{higher-order terms,} \end{aligned}$$

where

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}}(\mathbf{x}_e, \mathbf{u}_e) = \left[\begin{array}{ccc} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_n}{\partial x_1} & \cdots & \frac{\partial f_n}{\partial x_n} \end{array} \right]_{\substack{\mathbf{x}=\mathbf{x}_e \\ \mathbf{u}=\mathbf{u}_e}} \quad \text{and} \quad \frac{\partial \mathbf{f}}{\partial \mathbf{u}}(\mathbf{x}_e, \mathbf{u}_e) = \left[\begin{array}{ccc} \frac{\partial f_1}{\partial u_1} & \cdots & \frac{\partial f_1}{\partial u_m} \\ \vdots & & \vdots \\ \frac{\partial f_n}{\partial u_1} & \cdots & \frac{\partial f_n}{\partial u_m} \end{array} \right]_{\substack{\mathbf{x}=\mathbf{x}_e \\ \mathbf{u}=\mathbf{u}_e}}$$

are the Jacobian matrices of \mathbf{f} with respect to \mathbf{x} and \mathbf{u} , evaluated at the equilibrium point, $[\mathbf{x}_e^T \ \mathbf{u}_e^T]^T$. Note that

$$\frac{d}{dt} \mathbf{x} = \frac{d}{dt} \mathbf{x}_e + \frac{d}{dt} \Delta \mathbf{x} = \frac{d}{dt} \Delta \mathbf{x}$$

because \mathbf{x}_e is constant. Furthermore, $\mathbf{f}(\mathbf{x}_e, \mathbf{u}_e) = \mathbf{0}$. Let

$$\mathbf{A} = \frac{\partial \mathbf{f}}{\partial \mathbf{x}}(\mathbf{x}_e, \mathbf{u}_e) \quad \text{and} \quad \mathbf{B} = \frac{\partial \mathbf{f}}{\partial \mathbf{u}}(\mathbf{x}_e, \mathbf{u}_e).$$

Neglecting higher-order terms, we arrive at the linear approximation

$$\frac{d}{dt} \Delta \mathbf{x} = \mathbf{A} \Delta \mathbf{x} + \mathbf{B} \Delta \mathbf{u}.$$

Similarly, if the outputs of the nonlinear system model are of the form

$$\begin{aligned} y_1 &= h_1(x_1, x_2, \dots, x_n, u_1, u_2, \dots, u_m), \\ y_2 &= h_2(x_1, x_2, \dots, x_n, u_1, u_2, \dots, u_m), \\ &\vdots \\ y_p &= h_p(x_1, x_2, \dots, x_n, u_1, u_2, \dots, u_m) \end{aligned}$$

or in vector notation

$$\mathbf{y} = \mathbf{h}(\mathbf{x}, \mathbf{u}),$$

then Taylor's series expansion can again be used to yield the linear approximation of the above output equations. Indeed, if we let

$$\mathbf{y} = \mathbf{y}_e + \Delta \mathbf{y},$$

then we obtain

$$\Delta \mathbf{y} = \mathbf{C} \Delta \mathbf{x} + \mathbf{D} \Delta \mathbf{u},$$

where

$$C = \frac{\partial \mathbf{h}}{\partial \mathbf{x}}(\mathbf{x}_e, \mathbf{u}_e) = \left[\begin{array}{ccc} \frac{\partial h_1}{\partial x_1} & \cdots & \frac{\partial h_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial h_p}{\partial x_1} & \cdots & \frac{\partial h_p}{\partial x_n} \end{array} \right] \bigg|_{\substack{\mathbf{x}=\mathbf{x}_e \\ \mathbf{u}=\mathbf{u}_e}} \in \mathbb{R}^{p \times n}$$

and

$$D = \frac{\partial \mathbf{h}}{\partial \mathbf{u}}(\mathbf{x}_e, \mathbf{u}_e) = \left[\begin{array}{ccc} \frac{\partial h_1}{\partial u_1} & \cdots & \frac{\partial h_1}{\partial u_m} \\ \vdots & & \vdots \\ \frac{\partial h_p}{\partial u_1} & \cdots & \frac{\partial h_p}{\partial u_m} \end{array} \right] \bigg|_{\substack{\mathbf{x}=\mathbf{x}_e \\ \mathbf{u}=\mathbf{u}_e}} \in \mathbb{R}^{p \times m}$$

are the Jacobian matrices of \mathbf{h} with respect to \mathbf{x} and \mathbf{u} , evaluated at the equilibrium point $[\mathbf{x}_e^T \ \mathbf{u}_e^T]^T$.

◆ Example 2.14

In this example, we linearize the model of Watt's governor that we derived in Section 1.7.1. Recall that the equations modeling the governor have the form

$$\begin{aligned} \dot{x}_1 &= x_2, \\ \dot{x}_2 &= \frac{1}{2} N^2 x_3^2 \sin 2x_1 - g \sin x_1 - \frac{b}{m} x_2, \\ \dot{x}_3 &= \frac{\kappa}{I} \cos x_1 - \frac{\tau}{I}. \end{aligned}$$

We linearize the above model of an equilibrium state of the form

$$\mathbf{x}_e = [x_{1e} \ 0 \ x_{3e}]^T.$$

Equating the right-hand sides of the above engine–governor equations to zero yields

$$\begin{aligned} x_{2e} &= 0 \\ N^2 x_{3e}^2 &= \frac{g}{\cos x_{1e}} \\ \cos x_{1e} &= \frac{\tau}{\kappa} \end{aligned}$$

The Jacobian matrix of the engine–governor model evaluated at the equilibrium state is

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} \bigg|_{\mathbf{x}=\mathbf{x}_e} = \left[\begin{array}{ccc} 0 & 1 & 0 \\ N^2 x_{3e}^2 \cos 2x_{1e} - g \cos x_{1e} & -\frac{b}{m} & N^2 x_{3e} \sin 2x_{1e} \\ -\frac{\kappa}{I} \sin x_{1e} & 0 & 0 \end{array} \right].$$

We use the following trigonometric identities

$$\cos 2\alpha = \cos^2 \alpha - \sin^2 \alpha, \quad \sin 2\alpha = 2 \sin \alpha \cos \alpha$$

and the previously obtained relation, $N^2 x_{3e}^2 = g/\cos x_{1e}$, to eliminate $N^2 x_{3e}^2$ from the Jacobian matrix. Performing needed manipulations, we obtain

$$\left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{x}_e} = \begin{bmatrix} 0 & 1 & 0 \\ -\frac{g \sin^2 x_{1e}}{\cos x_{1e}} & -\frac{b}{m} & \frac{2g \sin x_{1e}}{x_{3e}} \\ -\frac{\kappa}{I} \sin x_{1e} & 0 & 0 \end{bmatrix}.$$

Hence, the linearized model of the engine–governor system can be represented as

$$\begin{bmatrix} \Delta \dot{x}_1 \\ \Delta \dot{x}_2 \\ \Delta \dot{x}_3 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ -\frac{g \sin^2 x_{1e}}{\cos x_{1e}} & -\frac{b}{m} & \frac{2g \sin x_{1e}}{x_{3e}} \\ -\frac{\kappa}{I} \sin x_{1e} & 0 & 0 \end{bmatrix} \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \Delta x_3 \end{bmatrix}. \quad (2.26)$$

2.5 Describing Function Method

The describing function method allows us to apply familiar frequency domain techniques, used in the linear system analysis, to the analysis of a class of nonlinear dynamical systems. The method can be used to predict limit cycles in nonlinear systems. The describing function method can be viewed as a “harmonic linearization” of a nonlinear element. The method provides a “linear approximation” to the nonlinear element based on the assumption that the input to the nonlinear element is a sinusoid of known, constant amplitude. The fundamental harmonic of the element’s output is compared with the input sinusoid to determine the steady-state amplitude and phase relation. This relation is the describing function for the nonlinear element. The describing function method is based on the Fourier series. In our review of Fourier series, we use the notion of a scalar product over a function space, which we discuss next. Our presentation of the scalar product in the context of spaces of functions follows that of Lang [175]. The reader familiar with the Fourier series may go directly to Subsection 2.5.3, where the describing function method analysis is presented.

2.5.1 Scalar Product of Functions

Let $\mathcal{C}([a, b])$ be the set of continuous functions on the interval $[a, b] \subset \mathbb{R}$.

Definition 2.2 The scalar product of two real-valued functions $u, v \in \mathcal{C}([a, b])$ is

$$\langle u, v \rangle = \int_a^b u(t)v(t) dt.$$

Using simple properties of the integral, we can verify that the above scalar product satisfies the following conditions:

1. For all $u, v \in \mathcal{C}([a, b])$, we have $\langle u, v \rangle = \langle v, u \rangle$.
2. If u, v, w are elements of $\mathcal{C}([a, b])$, then $\langle u, v + w \rangle = \langle u, v \rangle + \langle u, w \rangle$.
3. If x is a number, then $\langle xu, v \rangle = x \langle u, v \rangle = \langle u, xv \rangle$.
4. For all $v \in \mathcal{C}([a, b])$, we have $\langle v, v \rangle \geq 0$, and $\langle v, v \rangle > 0$ if $v \neq 0$.

The functions $u, v \in \mathcal{C}([a, b])$ are said to be *orthogonal* if $\langle u, v \rangle = 0$. A set of functions from $\mathcal{C}([a, b])$ is said to be *mutually orthogonal* if each distinct pair of functions in the set is orthogonal.

◆ Example 2.15

The functions $\psi_n(t) = \sin n\omega t$, $\phi_n(t) = \cos n\omega t$, $n = 1, 2, \dots$, and $\phi_0(t) = 1$ form a mutually orthogonal set of functions on the interval $[t_0, t_0 + 2\pi/\omega]$. This can be verified by direct integration. For example, for positive m and n such that $m \neq n$, we have

$$\begin{aligned}
 \langle \psi_m, \psi_n \rangle &= \int_{t_0}^{t_0+2\pi/\omega} \sin m\omega t \sin n\omega t \, dt \\
 &= \frac{1}{2} \int_{t_0}^{t_0+2\pi/\omega} (\cos(m-n)\omega t - \cos(m+n)\omega t) \, dt \\
 &= \frac{1}{2\omega} \left(\frac{\sin(m-n)\omega t}{m-n} - \frac{\sin(m+n)\omega t}{m+n} \right) \Bigg|_{t_0}^{t_0+2\pi\omega} \\
 &= 0.
 \end{aligned} \tag{2.27}$$

Note that $m+n \neq 0$ because m and n are positive. Using similar computations, we can show that for $m \neq n$ we obtain

$$\langle \phi_m, \phi_n \rangle = 0 \tag{2.28}$$

and that for all non-negative m and positive n we have

$$\langle \phi_m, \psi_n \rangle = 0. \tag{2.29}$$

We define the *norm* of $v \in \mathcal{C}([a, b])$ to be

$$\|v\| = \sqrt{\langle v, v \rangle}.$$

◆ Example 2.16

Consider $\psi_n(t) = \sin n\omega t$ on the interval $[t_0, t_0 + 2\pi/\omega]$. Then,

$$\begin{aligned}
 \|\psi_n\|^2 &= \int_{t_0}^{t_0+2\pi/\omega} \sin n\omega t \sin n\omega t \, dt \\
 &= \int_{t_0}^{t_0+2\pi/\omega} \sin^2 n\omega t \, dt \\
 &= \frac{1}{2} \int_{t_0}^{t_0+2\pi/\omega} (1 - \cos 2n\omega t) \, dt \\
 &= \frac{1}{2} \left(t - \frac{\sin 2n\omega t}{2n\omega} \right) \Big|_{t_0}^{t_0+2\pi/\omega} \\
 &= \frac{\pi}{\omega}.
 \end{aligned} \tag{2.30}$$

Similarly,

$$\|\phi_n\|^2 = \frac{\pi}{\omega}, \quad n = 1, 2, \dots, \tag{2.31}$$

and

$$\|\phi_0\|^2 = \frac{2\pi}{\omega}. \tag{2.32}$$

So far we discussed only continuous functions. In many applications we have to work with more general functions, specifically with piecewise continuous functions.

Definition 2.3 A function f is said to be *piecewise continuous* on an interval $[a, b]$ if the interval can be partitioned by a finite number of points, say t_i , so that:

1. $a = t_0 < t_1 < \dots < t_n = b$.
2. The function f is continuous on each open subinterval (t_{i-1}, t_i) .
3. The function f approaches a finite limit as the end points of each subinterval are approached from within the subinterval; that is, the limits

$$\lim_{\substack{h \rightarrow 0 \\ h > 0}} f(t_i - h) \quad \text{and} \quad \lim_{\substack{h \rightarrow 0 \\ h > 0}} f(t_i + h)$$

both exist.

The graph of a piecewise continuous function is shown in Figure 2.15. In further discussion, we assume that two piecewise functions f and g are equal if they have the same values at the points where they are continuous. Thus, the functions shown in Figure 2.16 are considered to be equal, though they differ at the points of discontinuity.

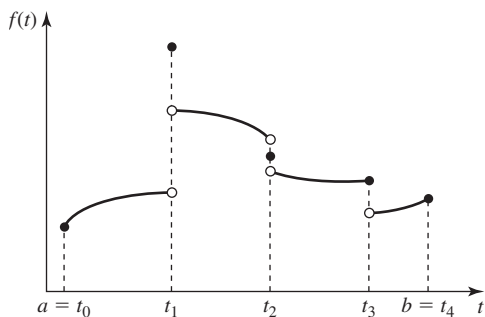


Figure 2.15 A graph of a piecewise continuous function.

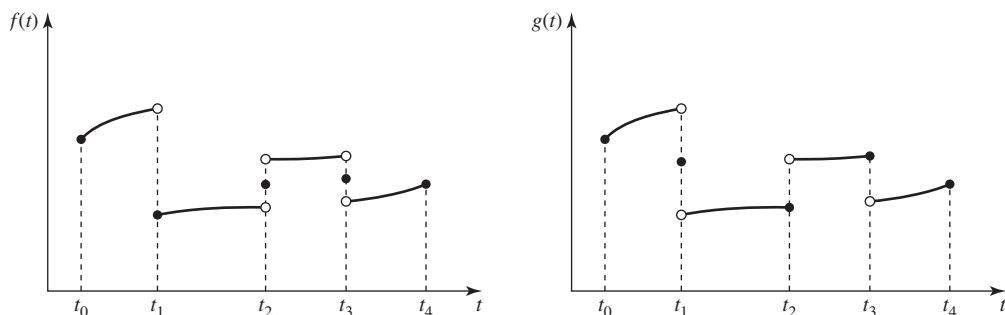


Figure 2.16 Two piecewise continuous functions that we consider to be equal.

Theorem 2.2 Let V be the space of functions that are piecewise continuous on the interval $[a, b]$. Let $f \in V$. Then, $\|f\| = 0$ if and only if $f(x) = 0$ for all but finite number of points x in the interval $[a, b]$.

Proof (\Leftarrow) It is clear that if $f(x) = 0$ except for a finite number of $x \in [a, b]$, then

$$\|f\|^2 = \int_a^b f^2(x) dx = 0.$$

(\Rightarrow) Suppose f is piecewise continuous on $[a, b]$, and let $a = t_0 < t_1 < \cdots < t_n = b$ be a partition of the interval $[a, b]$ such that f is continuous on each subinterval $[t_{i-1}, t_i]$ except possibly at the end points. Suppose that $\|f\| = 0$. Then, also $\|f\|^2 = \langle f, f \rangle = 0$. The above means that

$$\int_a^b f^2(x) dx = 0,$$

where the integral is the sum of the integrals over the intervals $[t_{i-1}, t_i]$; that is,

$$\int_a^b f^2(x) dx = \sum_{i=1}^n \int_{t_{i-1}}^{t_i} f^2(x) dx = 0.$$

Each integral satisfies

$$\int_{t_{i-1}}^{t_i} f^2(x) dx \geq 0.$$

Hence, each such integral must be equal to 0. The function f is continuous on (t_{i-1}, t_i) . Therefore,

$$f^2(x) = 0 \quad \text{for } t_{i-1} < x < t_i,$$

which means that

$$f(x) = 0 \quad \text{for } t_{i-1} < x < t_i.$$

Hence, $f(x) = 0$ except at a finite number of points. The proof is complete.

The above theorem implies that the scalar product on the space V of piecewise continuous functions is positive definite. In other words, for a piecewise continuous function $f \in V$, we have $\langle f, f \rangle \geq 0$, and $\langle f, f \rangle = 0$ if and only if $f = 0$ for all but a finite number of points in the interval $[a, b]$.

2.5.2 Fourier Series

We begin with the definition of a periodic function.

Definition 2.4 A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is said to be periodic with the period T if the domain of f contains $t + T$ whenever it contains t and if

$$f(t + T) = f(t)$$

for every value of t . The smallest positive value of T for which $f(t + T) = f(t)$ is called the fundamental period of f .

It follows from the above definition that if T is a period of f , then $2T$ is also a period, and so is any integral multiple of T .

A periodic function, satisfying certain assumptions spelled out later, may be represented by the series

$$f(t) = A_0 + \sum_{k=1}^{\infty} (A_k \cos k\omega t + B_k \sin k\omega t), \quad (2.33)$$

where $\omega = 2\pi/T$, with T being the fundamental period of f . On the set of points where the series (2.33) converges, it defines f , whose value at each point is the sum of the series for that value of t , and we say that the series (2.33) is the Fourier series for f .

Suppose that the series (2.33) converges. We use the results of the previous subsection to find expressions for the coefficients A_k and B_k . We first compute A_k coefficients. Multiplying (2.33) by $\phi_m = \cos m\omega t$, where $m > 0$ is a fixed positive integer, and then integrating both sides with respect to t from t_0 to $t_0 + 2\pi/\omega$ yields

$$\langle \phi_m, f \rangle = A_0 \langle \phi_m, \phi_0 \rangle + \sum_{k=1}^{\infty} A_k \langle \phi_m, \phi_k \rangle + \sum_{k=1}^{\infty} B_k \langle \phi_m, \psi_k \rangle, \quad (2.34)$$

where $\phi_0 = 1$. It follows from (2.27), (2.28), and (2.29) that the only nonvanishing term on the right-hand side of (2.34) is the one for which $k = m$ in the first summation. Using (2.31),

we obtain

$$\begin{aligned} A_m &= \frac{\langle \phi_m, f \rangle}{\langle \phi_m, \phi_m \rangle} \\ &= \frac{\omega}{\pi} \int_{t_0}^{t_0+2\pi/\omega} f(t) \cos m\omega t \, dt, \quad m = 1, 2, \dots \end{aligned} \quad (2.35)$$

To obtain A_0 , we integrate (2.33). Using (2.32), we get

$$A_0 = \frac{\langle \phi_0, f \rangle}{\langle \phi_0, \phi_0 \rangle} = \frac{\omega}{2\pi} \int_{t_0}^{t_0+2\pi/\omega} f(t) \, dt. \quad (2.36)$$

An expression for B_m may be obtained by multiplying (2.33) by $\sin m\omega t$ and then integrating term by term from t_0 to $t_0 + 2\pi/\omega$. Using the orthogonality relations (2.27)–(2.29) yields

$$B_m = \frac{\langle \psi_m, f \rangle}{\langle \psi_m, \psi_m \rangle} = \frac{\omega}{\pi} \int_{t_0}^{t_0+2\pi/\omega} f(t) \sin m\omega t \, dt, \quad m = 1, 2, \dots \quad (2.37)$$

Suppose now that a function f is given that is periodic with period $2\pi/\omega$ and integrable on the interval $[t_0, t_0 + 2\pi/\omega]$. Then, coefficients A_k and B_k can be computed using (2.35)–(2.37), and a series of the form (2.33) can be formally constructed. However, we do not know whether this series converges for each value of t and, if so, whether its sum is $f(t)$. The answer is given by the theorem that we state after explaining the notation used. We denote the limit of $f(t)$ as t approaches c from the right by $f(c+)$; that is,

$$f(c+) = \lim_{t \rightarrow c+} f(t).$$

Similarly, $f(c-) = \lim_{t \rightarrow c-} f(t)$ denotes the limit of $f(t)$ as $t \rightarrow c$ from the left. The mean value of the right- and left-hand limits at the point c is

$$(f(c+) + f(c-))/2.$$

Note that at any point c where f is continuous, we have

$$f(c+) = f(c-) = f(c).$$

Theorem 2.3 Assume that f and its derivative are piecewise continuous on the interval $[t_0, t_0 + 2\pi/\omega]$. Furthermore, suppose that f is defined outside the interval $[t_0, t_0 + 2\pi/\omega]$ so that it is periodic with period $2\pi/\omega$. Then, f has a Fourier series of the form (2.33) whose coefficients are given by (2.35)–(2.37). The Fourier series converges to $(f(t+) + f(t-))/2$ at all points of the interval $[t_0, t_0 + 2\pi/\omega]$.

With this knowledge of the Fourier series, it is now possible to discuss the describing function method for analyzing a class of nonlinear control systems.

2.5.3 Describing Function in the Analysis of Nonlinear Systems

The describing function method can be considered as an extension of the Nyquist stability criterion to nonlinear systems. In our further discussion, we sometimes use a double box to represent a nonlinear element as illustrated in Figure 2.17.

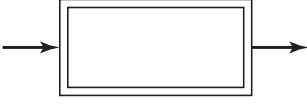


Figure 2.17 A symbol for a nonlinear element.

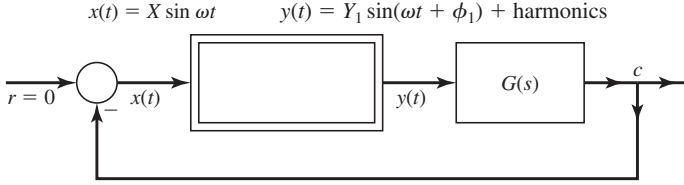


Figure 2.18 A nonlinear feedback system with a nonlinear element in the forward path.

Consider a feedback system that contains a nonlinear element as shown in Figure 2.18. Assume that the input to the nonlinear element is a sinusoidal signal $x(t) = X \sin \omega t$, where X is the amplitude of the input sinusoid. The output of the nonlinear element is, in general, not sinusoidal in response to a sinusoidal input. Suppose that the nonlinear element output is periodic with the same period as its input and that it may be expanded in the Fourier series,

$$\begin{aligned} y(t) &= A_0 + \sum_{k=1}^{\infty} (A_k \cos k\omega t + B_k \sin k\omega t) \\ &= A_0 + \sum_{k=1}^{\infty} Y_k \sin(k\omega t + \phi_k). \end{aligned} \quad (2.38)$$

Thus $y(t)$ consists of a fundamental component

$$Y_1 \sin(\omega t + \phi_1),$$

together with a mean level A_0 and harmonic components at frequencies $2\omega, 3\omega, \dots$. The describing function method relies on the assumption that only the fundamental harmonic is significant. This assumption is called the *filtering hypothesis*. The filtering hypothesis can be expressed as

$$|G(j\omega)| \gg |G(kj\omega)| \quad \text{for } k = 2, 3, \dots,$$

which requires that the linear element, modeled by the transfer function $G(s)$, be a low-pass filter. The filtering hypothesis is often valid because many, if not most, control systems are low-pass filters. This results in higher harmonics being more attenuated compared with the fundamental harmonic component. The nonlinear element is then approximately modeled by its describing function,

$$N(X) = \frac{Y_1}{X} \angle \phi_1 = \frac{\text{fundamental component of output}}{\text{input}}. \quad (2.39)$$

The describing function is, in general, a complex-valued function of the amplitude X of the input sinusoid.

The stability of the nonlinear closed-loop system depicted in Figure 2.18 can be analyzed using Nyquist's criterion applied to the associated linear system shown in Figure 2.19. The

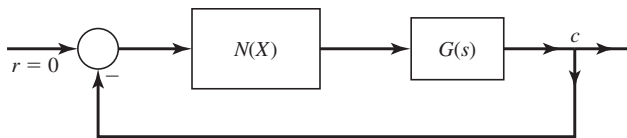


Figure 2.19 Associated linear system used in the describing function analysis of the nonlinear system shown in Figure 2.18.

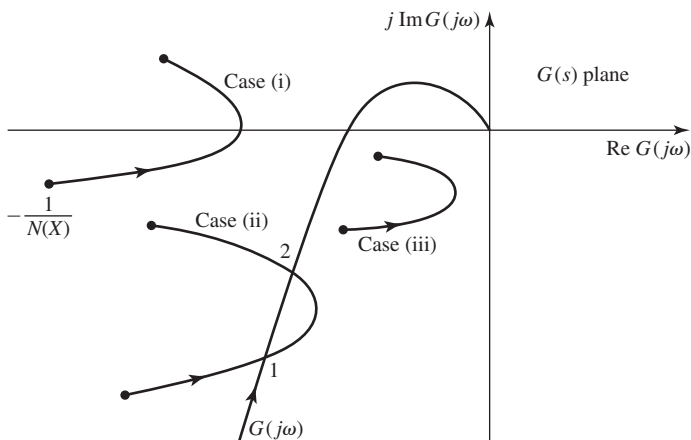


Figure 2.20 Stability analysis using describing function.

associated linear closed-loop system characteristic equation is

$$1 + N(X)G(s) = 0. \quad (2.40)$$

In the $G(j\omega)$ complex plane, the locus traced out by $-1/N$ as X varies is considered a generalization of the $(-1, j0)$ point in the linear case. The stability of the nonlinear system is determined by computing the encirclements of the $-1/N$ locus by the $G(j\omega)$ plot. We might have the situations shown in Figure 2.20. In case (i), the $-1/N$ locus is not encircled by the $G(j\omega)$ plot. Assuming that the poles of $G(s)$ are in the left-half complex plane, the closed-loop system is asymptotically stable. For a limit cycle to exist, we must have an intersection between the polar plot of $G(j\omega)$ and the $-1/N$ locus. We have two intersections between $-1/N$ and $G(j\omega)$ in case (ii). At intersection point 1, any increase in X will cause the operating point to be encircled and will consequently cause X to grow. If there is a decrease in X , then the oscillations will decay. Thus, point 1 represents an unstable limit cycle. At intersection point 2, an increase in X would shift the operating point outside the $G(j\omega)$ plot resulting in the decaying oscillations. Point 2 represents a stable limit cycle. In case (iii), any oscillation will grow without limit.

◆ Example 2.17

Consider the nonlinear control system shown in Figure 2.21 that was analyzed in Power and Simpson [237, pp. 339–341]. We first derive the describing function of the ideal relay. When the input of the ideal relay is greater than zero, its output is M . When the relay's input is less than zero, its output is $-M$. Thus, the ideal relay

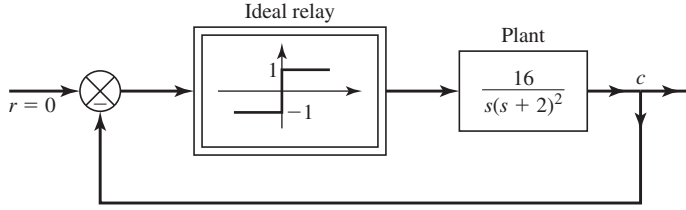


Figure 2.21 Nonlinear feedback control system analyzed in Example 2.17.

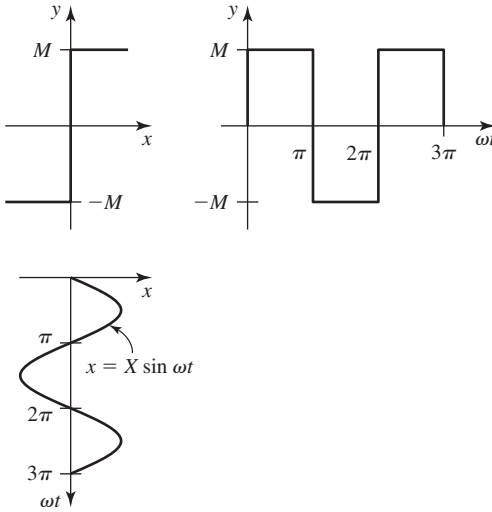


Figure 2.22 Input and output waveforms of an ideal relay.

can be described as

$$y = \begin{cases} M & \text{for } x > 0 \\ -M & \text{for } x < 0. \end{cases} \quad (2.41)$$

We assume that the input to the relay is a sinusoid, $x(t) = X \sin \omega t$, drawn vertically down from the relay characteristic in Figure 2.22. The output $y(t)$ of the relay is a rectangular wave. It is drawn horizontally to the right of the relay characteristic in Figure 2.22. This rectangular wave is an odd function. For an odd function, we have $A_k = 0$, $k = 0, 1, \dots$, in its Fourier series. Thus, the Fourier series of the relay's output has the form

$$y(t) = \sum_{k=1}^{\infty} B_k \sin k\omega t.$$

The fundamental harmonic of $y(t)$ is

$$B_1 \sin \omega t = Y_1 \sin \omega t,$$

where

$$\begin{aligned}
 B_1 &= \frac{\omega}{\pi} \int_0^{2\pi/\omega} y(t) \sin \omega t \, dt \\
 &= \frac{2\omega}{\pi} \int_0^{\pi/\omega} M \sin \omega t \, dt \\
 &= \frac{2M\omega}{\pi} \int_0^{\pi/\omega} \sin \omega t \, dt \\
 &= -\frac{2M}{\pi} \cos(\omega t) \Big|_0^{\pi/\omega} \\
 &= \frac{4M}{\pi}.
 \end{aligned} \tag{2.42}$$

Hence, the describing function of the ideal relay is

$$N(X) = \frac{Y_1}{X} \angle 0^\circ = \frac{4M}{\pi X}.$$

A normalized plot of the describing function of the ideal relay nonlinearity is depicted in Figure 2.23. In our example, $M = 1$ and thus the describing function of the relay is

$$N(X) = \frac{4}{\pi X}.$$

The sinusoidal transfer function of the linear plant model is

$$G(j\omega) = \frac{16}{j\omega(j\omega + 2)^2} = \frac{-64\omega + j16(\omega^2 - 4)}{\omega(4 + \omega^2)^2}.$$

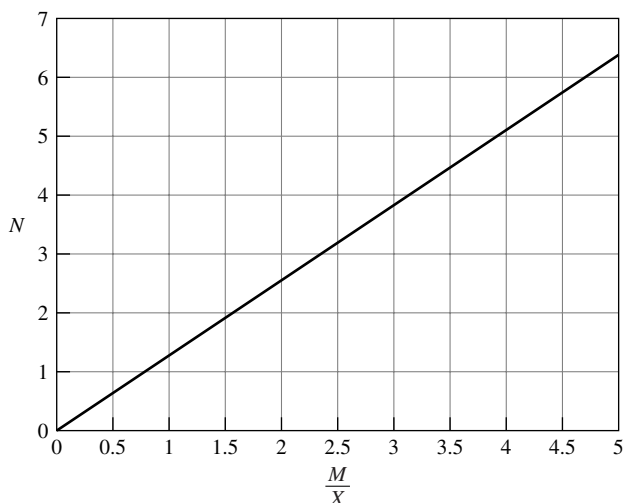


Figure 2.23 A plot of the describing function for the ideal relay.

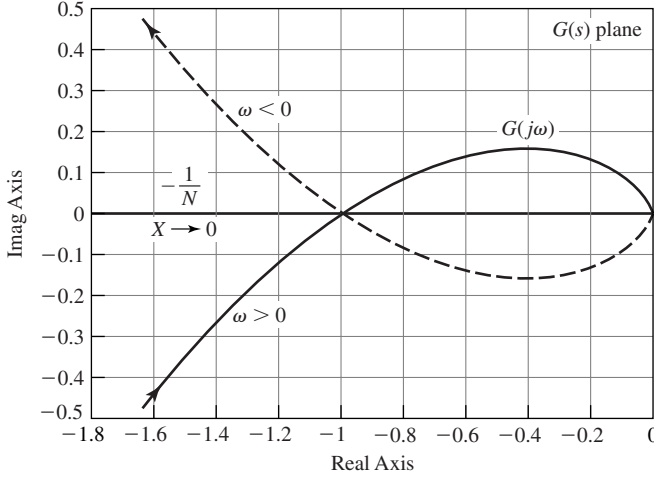


Figure 2.24 Plots of $G(j\omega)$ and $-1/N$ of Example 2.17.

At $\omega = 2$ rad/sec the polar plot of $G(j\omega)$ intersects the plot of $-1/N$ (see Figure 2.24), to form a stable limit cycle. The amplitude of the limit cycle is obtained from the relation $G(j2) = -1/N$; that is,

$$\frac{-64}{(4 + 4)^2} = -\frac{\pi X}{4}.$$

Hence, the amplitude of the limit cycle is

$$X = \frac{4}{\pi},$$

while the amplitude of the fundamental harmonic of the plant output, c , is

$$\frac{4}{\pi} |G(j2)| = \frac{4}{\pi}.$$

Note that the filtering hypothesis is satisfied in this example. Indeed, the Fourier series of the output of the relay is

$$\frac{4}{\pi} \left(\sin \omega t + \frac{1}{3} \sin 3\omega t + \frac{1}{5} \sin 5\omega t + \dots \right).$$

Thus, for example, the amplitude of the third harmonic of the plant output is

$$\frac{4}{3\pi} |G(j6)| = \frac{1}{45} \frac{4}{\pi}.$$

The higher harmonics are even more attenuated compared with the fundamental harmonic.

Notes

The system of differential equations (2.43) in Exercise 2.6 comes from Boyce and DiPrima [32, p. 523]. Al-Khafaji and Tooley [6] provide an introductory treatment of numerical solutions of differential equations. For further discussion of solving differential equations using numerical methods and analyses of their computational efficiency, the reader is referred to Salvadori and Baron [252], Conte and de Boor [53, Chapter 8], or Parker and Chua [222]. For an in-depth treatment of the describing function method, we refer to Graham and McRuer [105] and Chapters 6 and 7 of Hsu and Meyer [128]. Graham and McRuer [105] make an interesting remark about the origins of the describing function method. They note on page 92 of their 1961 book: “While the notion of representing a nonlinearity by an ‘equivalent’ linear element is quite old and has been used by many writers, the first instance of a major systematic exploitation of the technique was probably by N. M. Kryloff and N. N. Bogoliuboff, *Introduction to Non-linear Mechanics* (a free translation of the Russian edition of 1937, by S. Lefschetz), Princeton University Press, Princeton, N.J., 1943. See also N. Minorsky, *Introduction to Non-linear Mechanics*, J. W. Edwards, Ann Arbor, Mich., 1947.” The so-called sinusoidal describing function that we analyzed in the previous section, was developed almost simultaneously in several different countries during and just after World War II. The first to introduce the method in the open literature appears to be A. Tustin of England, who published his paper in the *Journal of the IEE* in 1947.

The behavior of numerous dynamical systems is often modeled by a set of linear, first-order, differential equations. Frequently, however, a linear model is a result of linearization of a nonlinear model. Linear models seem to dominate the controls’ literature. Yet nature is nonlinear. “All physical systems are nonlinear and have time-varying parameters in some degree” [105, p. 1]. How can one reconcile this paradox? Graham and McRuer [105, p. 1] have the following answer: “Where the effect of the nonlinearity is very small, or if the parameters vary only slowly with time, linear constant-parameter methods of analysis can be applied to give an approximate answer which is adequate for engineering purposes. The analysis and synthesis of physical systems, predicated on the study of linear constant-parameter mathematical models, has been, in fact, an outstandingly successful enterprise. A system represented as linear, however, is to some extent a mathematical abstraction that can never be encountered in a real world. Either by design or because of nature’s ways it is often true that experimental facts do not, or would not, correspond with any prediction of linear constant-parameter theory. In this case nonlinear or time-varying-parameter theory is essential to the description and understanding of physical phenomena.”

EXERCISES

- 2.1** Sketch a phase-plane portrait of the differential equation

$$\ddot{x} - \dot{x} + 2x + 2\text{sign}(2x + \dot{x}) = 0.$$

- 2.2** For the system of differential equations

$$\begin{aligned}\dot{x}_1 &= -\frac{7}{4}x_1 + \frac{1}{4}x_2, \\ \dot{x}_2 &= \frac{3}{4}x_1 - \frac{5}{4}x_2,\end{aligned}$$

find the phase-plane asymptotes.

2.3 For a dynamical system modeled by

$$\begin{aligned}\dot{x}_1 &= x_2, \\ \dot{x}_2 &= -x_1,\end{aligned}$$

find the time T it takes the system trajectory to move from the state $\mathbf{x}(0) = [0 \ 1]^T$ to the state $\mathbf{x}(T) = [1 \ 0]^T$.

2.4 Consider a trajectory in the state plane shown in Figure 2.25. Find the time it takes the representative point to move from A to B.

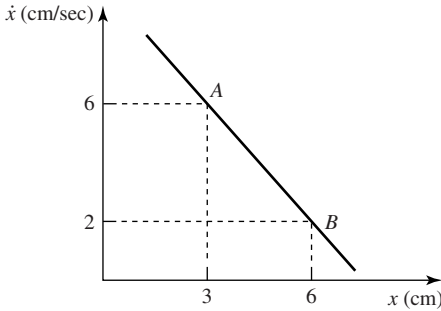


Figure 2.25 A trajectory for Exercise 2.4.

2.5 Given the following initial value problem:

$$\dot{x} = f(t, x), \quad x(t_0) = x_0.$$

- Derive the Euler formula for solving the above initial value problem using a Taylor series.
- Make use of a Taylor series with a remainder to compute the local error formula assuming that the data at the n th step are correct.

2.6 Use Bendixson's theorem to investigate the existence of limit cycles for the system

$$\begin{aligned}\dot{x}_1 &= x_1 + x_2 - x_1(x_1^2 + x_2^2), \\ \dot{x}_2 &= -x_1 + x_2 - x_2(x_1^2 + x_2^2).\end{aligned}\tag{2.43}$$

Sketch a phase-plane portrait for this system.

2.7 Use Bendixson's theorem to investigate the existence of limit cycles for the system

$$\begin{aligned}\dot{x}_1 &= x_2 - x_1(x_1^2 + x_2^2 - c), \\ \dot{x}_2 &= -x_1 - x_2(x_1^2 + x_2^2 - c),\end{aligned}\tag{2.44}$$

where c is a constant. Sketch phase-plane portraits for this system for $c = -0.2$ and $c = 0.2$. The above nonlinear system was analyzed in reference 222, p. 204.

2.8 Verify that the Runge method agrees with the Taylor series solution through terms of degree three.

2.9 Verify that the Runge–Kutta method agrees with the Taylor series solution through terms of degree four.

- 2.10** Consider the nonlinear mass–spring–damper system depicted in Figure 2.26. The spring coefficient $K = 0.5x^2$ N/m, the damping coefficient $B = 0.1\dot{x}^2$ N·sec/m, and $M = 1$ kg.

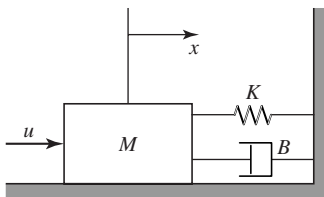


Figure 2.26 The mass–spring–damper system of Exercise 2.10.

- Write the second-order differential equation describing the motion of the system.
 - Choose $x_1 = x$ and $x_2 = \dot{x}$. Represent the equation of motion derived in (a) in state-space format.
 - For the model from (b), find the equilibrium state corresponding to $u = 0$. Then, linearize the model about the obtained equilibrium point.
- 2.11** (Based on Toptscheev and Tsypliyakov [282, pp. 29–31]) A cyclotron, whose operation is illustrated in Figure 2.27, accelerates charged particles into high energies so that they can be used in atom-smashing experiments. The operation of the cyclotron can be described as follows. Positive ions enter the cyclotron from the ion source, S, and are available to be accelerated. An electric oscillator establishes an accelerating potential difference across the gap of the two D-shaped objects, called dees. An emerging ion from the ion source, S, finds the dee that is negative. It will accelerate toward this dee and will enter it. Once inside the dee, the ion is screened from electric fields by the metal walls of the dees. The dees are immersed in a magnetic field so that the ion's trajectory bends in a circular path. The radius of the path depends on the ion's velocity. After some time, the ion emerges from the dee on the other side of the ion source. In

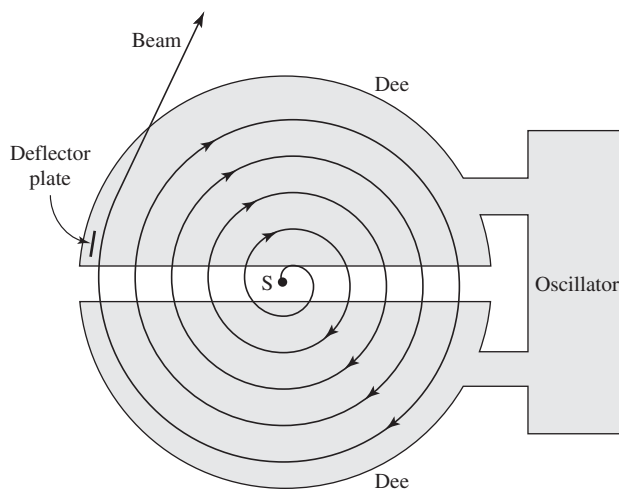


Figure 2.27 Schematic of a cyclotron of Exercise 2.11.

the meantime, the accelerating potential changes its sign. Thus the ion again faces a negative dee. The ion further accelerates and describes a semicircle of larger radius. This process goes on until the ion reaches the outer edge of one dee, where it is pulled out of the cyclotron by a negatively charged deflector plate. The moving ion is acted upon by a deflecting force of magnitude

$$q(v + \Delta v)(B + \Delta B),$$

where q is the charge of the ion, v is the ion's speed, Δv is the increase in the ion's speed, and B is the magnetic field acting on the ion. The nominal radius of the ion's orbit is r . The change in the magnetic field as the result of a change in the orbit of the ion is ΔB . As the ion speed increases, the relativistic mass m increases by Δm . The ion moving in a circular path undergoes centrifugal force of magnitude

$$\frac{(m + \Delta m)(v + \Delta v)^2}{r + \Delta r}.$$

By Newton's second law, the equation describing the motion of the ion is

$$(m + \Delta m) \frac{d^2}{dt^2}(r + \Delta r) = \frac{(m + \Delta m)(v + \Delta v)^2}{r + \Delta r} - q(v + \Delta v)(B + \Delta B).$$

Linearize the above equation. The motion of the ion along a nominal path is described by

$$m \frac{d^2 r}{dt^2} = \frac{mv^2}{r} - qvB.$$

Find the transfer function relating the change in radius of motion, Δr , of the ion and the change in the momentum of the ion:

$$\Delta p = m\Delta v + \Delta mv.$$

In other words, find the transfer function $\mathcal{L}(\Delta r)/\mathcal{L}(\Delta p)$, where \mathcal{L} is the Laplace transform operator.

2.12 Linearize the permanent-magnet stepper motor model, analyzed in Subsection 1.7.3, about $x_{1eq} = \theta_{eq} = \text{constant}$.

2.13 Let v_1, v_2, \dots, v_n be nonzero, mutually orthogonal elements of the space V and let v be an arbitrary element of V .

- Find the component of v along v_j .
- Let c_j denote the component of v along v_j . Show that the element

$$v - c_1 v_1 - \dots - c_n v_n$$

is orthogonal to v_1, v_2, \dots, v_n .

- Let c_j be the component of $v \in V$ along v_j . Show that

$$\sum_{i=1}^n c_i^2 \leq \|v\|^2.$$

The above relation is called the Bessel inequality.

2.14 Consider characteristics of two nonlinear elements that are shown in Figure 2.28. Let

$$N_S(x) = \frac{2}{\pi} \left[\arcsin \frac{1}{x} + \frac{1}{x} \cos \left(\arcsin \frac{1}{x} \right) \right]. \quad (2.45)$$

Then, the describing function of the limiter is

$$N_l = \begin{cases} K, & M \leq S \\ K N_S \left(\frac{M}{S} \right), & M > S. \end{cases}$$

The describing function of the dead zone is

$$N_{d-z} = \begin{cases} 0, & M \leq A, \\ K \left[1 - N_S \left(\frac{M}{A} \right) \right], & M > A. \end{cases}$$

Derive the describing function of the nonlinear element whose characteristic is shown in Figure 2.29.

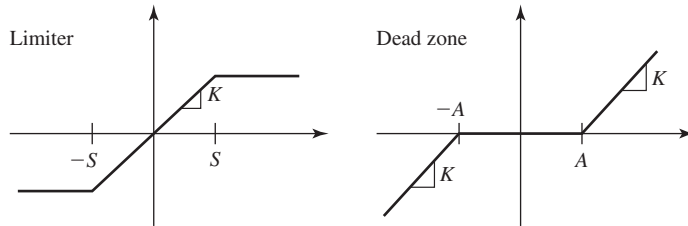


Figure 2.28 Characteristics of a limiter and dead zone nonlinear elements of Exercise 2.14.

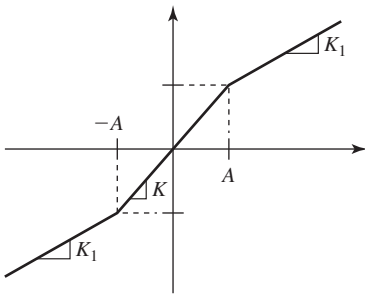


Figure 2.29 Characteristic of a nonlinear element of Exercise 2.14.

2.15 The describing function of the nonlinear element shown in Figure 2.30 is

$$N_1 = \begin{cases} 0, & M \leq A, \\ K \left[1 - N_S \left(\frac{M}{A} \right) \right] + \frac{4B}{\pi M} \sqrt{1 - \left(\frac{A}{M} \right)^2}, & M > A, \end{cases}$$

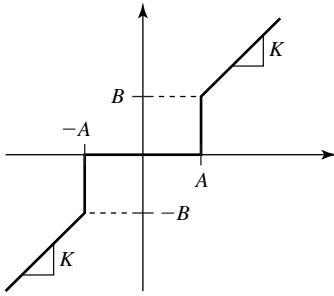


Figure 2.30 Nonlinear element with the describing function N_1 of Exercise 2.15.

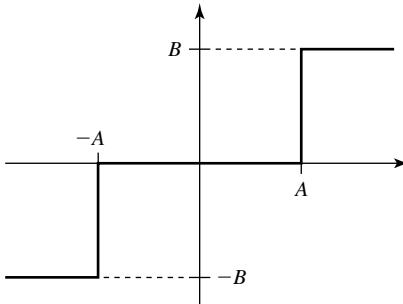


Figure 2.31 Characteristic of the ideal relay with dead zone.

where $N_S(x)$ is given by (2.45). The describing function of the dead-zone nonlinearity is given in Exercise 2.14 and its characteristic is shown in Figure 2.28. Derive the describing function of the ideal relay with dead zone whose characteristic is shown in Figure 2.31.



CHAPTER 3

Linear Systems

In bear country, wear bells to avoid surprising bruins. Don't approach a bear, but don't run away from one either. Movement may incite a charge, and a bear can outrun you.

If a bear approaches you, talk to it calmly in a low, firm voice. Slowly back away until the bear is out of sight. If you climb a tree, get at least 12 feet up.

If a bear senses that its food may be stolen or that its cub is in danger, it may charge. If this happens, drop to the ground and assume a fetal position, protecting your head with your arms, and play dead. Keep your pack on; it may serve as armor.

—*Reader's Digest Practical Problem Solver* [241]

We begin our discussion of linear system theory with a description of the concept of linearity by Mitchell J. Feigenbaum [227, p. 3], who is a researcher of nonlinear systems. “Linearity means that the rule that determines what a piece of a system is going to do next is not influenced by what it is doing now. More precisely, this is intended in a differential or incremental sense: For a linear spring, the increase of its tension is proportional to the increment whereby it is stretched, with the ratio of these increments exactly independent of how much it has already been stretched. Such a spring can be stretched arbitrarily far, and in particular will never snap or break. Accordingly, no real spring is linear.” The real world is nonlinear. However, there are a number of reasons to investigate linear systems. Linear models are often used to approximate nonlinear systems. Many times, this approximation is sufficient for the controller design for the underlying nonlinear system. Knowledge of linear system theory helps to understand the intricate theory of nonlinear systems.

In this chapter we acquaint ourselves with basic concepts from linear systems that will be used in subsequent chapters.

3.1 Reachability and Controllability

We first consider a linear, time-invariant, discrete-time system modeled by

$$\mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k), \quad \mathbf{x}(0) = \mathbf{x}_0, \quad (3.1)$$

where $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$. The solution to (3.1) is

$$\mathbf{x}(i) = A^i \mathbf{x}(0) + \sum_{k=0}^{i-1} A^{i-k-1} B \mathbf{u}(k),$$

which can be represented as

$$\mathbf{x}(i) = A^i \mathbf{x}(0) + [B \quad AB \quad \cdots \quad A^{i-1}B] \begin{bmatrix} \mathbf{u}(i-1) \\ \mathbf{u}(i-2) \\ \vdots \\ \mathbf{u}(0) \end{bmatrix}. \quad (3.2)$$

Let

$$U_i = [B \quad AB \quad \cdots \quad A^{i-1}B], \quad i = 1, 2, \dots$$

Note that $U_i \in \mathbb{R}^{n \times mi}$. If for some positive l we have

$$\text{rank } U_l = \text{rank } U_{l+1},$$

then all the columns of $A^l B$ are linearly dependent on those of U_l . This, in turn, implies that all the columns of $A^{l+1}B, A^{l+2}B, \dots$ must also be linearly dependent on those of U_l . Hence,

$$\text{rank } U_l = \text{rank } U_{l+1} = \text{rank } U_{l+2} = \dots \quad (3.3)$$

The above implies that the rank of U_i increases by at least one when i is increased by one, until the maximum value of $\text{rank } U_i$ is attained. The maximum value of the rank U_i is guaranteed to be achieved for $i = n$, which means that we always have

$$\text{rank } U_n = \text{rank } U_{n+1} = \text{rank } U_{n+2} = \dots \quad (3.4)$$

It follows from (3.3) that if $\text{rank } B = p$, then

$$\text{rank}[B \quad AB \quad \cdots \quad A^{n-1}B] = \text{rank}[B \quad AB \quad \cdots \quad A^{n-p}B]. \quad (3.5)$$

We say that the system (3.1) is *reachable* if for any \mathbf{x}_f there exists an integer $q > 0$ and a control sequence, $\{\mathbf{u}(i) : i = 0, 1, \dots, q-1\}$, that transfers $\mathbf{x}(0) = \mathbf{0}$ to $\mathbf{x}(q) = \mathbf{x}_f$.

Theorem 3.1 The system (3.1) is reachable if and only if

$$\text{rank}[B \quad AB \quad \cdots \quad A^{n-1}B] = n.$$

Proof Using (3.2) where $\mathbf{x}(0) = \mathbf{0}$ and $\mathbf{x}(q) = \mathbf{x}_f$, we obtain

$$\begin{aligned}\mathbf{x}_f &= \sum_{k=0}^{q-1} \mathbf{A}^{q-k-1} \mathbf{B} \mathbf{u}(k) \\ &= [\mathbf{B} \quad \mathbf{AB} \quad \dots \quad \mathbf{A}^{q-1} \mathbf{B}] \begin{bmatrix} \mathbf{u}(q-1) \\ \mathbf{u}(q-2) \\ \vdots \\ \mathbf{u}(0) \end{bmatrix}.\end{aligned}\quad (3.6)$$

It follows from the above that for any \mathbf{x}_f there exists a control sequence $\{\mathbf{u}(i) : i = 0, 1, \dots, q-1\}$ that transfers $\mathbf{x}(0) = \mathbf{0}$ to $\mathbf{x}(q) = \mathbf{x}_f$ if and only if

$$\text{rank}[\mathbf{B} \quad \mathbf{AB} \quad \dots \quad \mathbf{A}^{q-1} \mathbf{B}] = \text{rank } \mathbf{U}_q = n.$$

By (3.4) the maximum value of the rank of \mathbf{U}_q is guaranteed to be attained for $q = n$, which proves the theorem.

A notion closely related to reachability is the notion of *controllability*. The notion of controllability is related to the property of a controlled system being transferable from any given state to the origin $\mathbf{0}$ of the state space by means of an appropriately chosen control law.

Theorem 3.2 The system (3.1) is controllable if and only if

$$\text{range}(\mathbf{A}^n) \subset \text{range}[\mathbf{B} \quad \mathbf{AB} \quad \dots \quad \mathbf{A}^{n-1} \mathbf{B}].$$

Proof Using the solution formula for the system (3.1), for the case when $\mathbf{x}(0) = \mathbf{x}_0$ and $\mathbf{x}_f = \mathbf{x}(q) = \mathbf{0}$, we obtain

$$\begin{aligned}\mathbf{A}^q \mathbf{x}_0 &= - \sum_{k=0}^{q-1} \mathbf{A}^{q-k-1} \mathbf{B} \mathbf{u}(k) \\ &= -[\mathbf{B} \quad \mathbf{AB} \quad \dots \quad \mathbf{A}^{q-1} \mathbf{B}] \begin{bmatrix} \mathbf{u}(q-1) \\ \mathbf{u}(q-2) \\ \vdots \\ \mathbf{u}(0) \end{bmatrix}.\end{aligned}\quad (3.7)$$

The system (3.1) is controllable if and only if for some $q > 0$ and arbitrary initial condition \mathbf{x}_0 , the vector $\mathbf{A}^q \mathbf{x}_0$ is in the range of $[\mathbf{B} \quad \mathbf{AB} \quad \dots \quad \mathbf{A}^{q-1} \mathbf{B}] = \mathbf{U}_q$. Taking into account that the maximal range of \mathbf{U}_q is guaranteed to be attained for $q = n$, the system (3.1) is controllable if and only if the condition of the theorem is satisfied.

Note that the condition of the above theorem can be equivalently represented as

$$\text{rank}[\mathbf{B} \quad \mathbf{AB} \quad \dots \quad \mathbf{A}^{n-1} \mathbf{B}] = \text{rank}[\mathbf{B} \quad \mathbf{AB} \quad \dots \quad \mathbf{A}^{n-1} \mathbf{B} \quad \mathbf{A}^n]. \quad (3.8)$$

We illustrate the above discussion with a simple example that we found in Kaczorek [144, p. 482].

◆ Example 3.1

Consider a discrete-time system modeled by

$$\begin{aligned}\mathbf{x}(k+1) &= \mathbf{A}\mathbf{x}(k) + \mathbf{b}u(k) \\ &= \begin{bmatrix} 0 & a \\ 0 & 0 \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} 1 \\ 0 \end{bmatrix} u(k), \quad a \in \mathbb{R}.\end{aligned}\quad (3.9)$$

The above system model is not reachable because

$$\text{rank}[\mathbf{b} \quad \mathbf{A}\mathbf{b}] = \text{rank} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} = 1 < 2.$$

Using (3.6), it is easy to check that, for example, the state $\mathbf{x}_f = [0 \ 1]^T$ is not reachable from $\mathbf{x}(0) = \mathbf{0}$ because \mathbf{x}_f does not belong to the range of $[\mathbf{b} \quad \mathbf{A}\mathbf{b}]$. In fact, nonreachable states form the set

$$\left\{ c \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad c \in \mathbb{R} \right\}.$$

However, the system (3.9) is controllable because for any $a \in \mathbb{R}$,

$$\text{rank}[\mathbf{b} \quad \mathbf{A}\mathbf{b}] = \text{rank} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} = \text{rank}[\mathbf{b} \quad \mathbf{A}\mathbf{b} \quad \mathbf{A}^2\mathbf{b}] = \text{rank} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix};$$

that is,

$$\text{range}(\mathbf{A}^2) = \text{range}(\mathbf{O}) \subset \text{range}[\mathbf{b} \quad \mathbf{A}\mathbf{b}].$$

(See appendix, Section A.3, for an introduction to the \mathbf{O} notation.) Using (3.7), we can easily see that arbitrary initial state $\mathbf{x}(0) \in \mathbb{R}^2$ of the system (3.9) can be transferred to the origin of \mathbb{R}^2 using, for example, the zero control sequence, $u(0) = u(1) = 0$. This is because $\mathbf{A}^2 = \mathbf{O}$ and hence $\mathbf{A}^2\mathbf{x}(0) = \mathbf{0}$ for any $\mathbf{x}(0) \in \mathbb{R}^2$.

In summary, for discrete-time linear systems, reachability implies controllability, and the two notions are equivalent if the matrix \mathbf{A} of the given discrete system is nonsingular.

We next discuss the notions of reachability and controllability for continuous-time linear systems.

We say that the system $\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$, where $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times m}$, or equivalently the pair (\mathbf{A}, \mathbf{B}) , is reachable if for any $\mathbf{x}_f \in \mathbb{R}^n$ there is t_1 and a control law $\mathbf{u}(t)$ that transfers $\mathbf{x}(t_0) = \mathbf{0}$ to $\mathbf{x}(t_1) = \mathbf{x}_f$.

We say that the system $\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$, or equivalently the pair (\mathbf{A}, \mathbf{B}) , is *controllable* if there is a control law $\mathbf{u}(t)$ that transfers any initial state $\mathbf{x}_0 = \mathbf{x}(t_0)$ at any time t_0 to the origin $\mathbf{0}$ of the state space at some time $t_1 > t_0$.

It turns out that in the case of continuous, time-invariant systems the two notions are equivalent. Indeed, from the solution formula for the controlled system (A.151), for $\mathbf{x}(0) = \mathbf{0}$ and

$\mathbf{x}(t_1) = \mathbf{x}_f$, we obtain

$$\mathbf{x}_f = \int_0^{t_1} e^{A(t_1-t)} \mathbf{B} \mathbf{u}(t) dt.$$

Premultiplying the above equation by e^{-At_1} and rearranging, we get

$$\mathbf{v} = \int_0^{t_1} e^{-At} \mathbf{B} \mathbf{u}(t) dt, \quad (3.10)$$

where $\mathbf{v} = e^{-At_1} \mathbf{x}_f$. On the other hand, from the solution formula for the controlled system (A.151) for $\mathbf{x}(0) = \mathbf{x}_0$ and $\mathbf{x}(t_1) = \mathbf{0}$, we obtain

$$\mathbf{0} = e^{At_1} \mathbf{x}_0 + \int_0^{t_1} e^{A(t_1-t)} \mathbf{B} \mathbf{u}(t) dt.$$

Premultiplying the above equation by e^{-At_1} and rearranging yields

$$-\mathbf{x}_0 = \int_0^{t_1} e^{-At} \mathbf{B} \mathbf{u}(t) dt. \quad (3.11)$$

Comparing (3.10) and (3.11), we conclude that reachability of the system $\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$ is equivalent to its controllability. This is not true in the case of discrete-time systems, where reachability implies controllability, and the two notions are equivalent if the matrix \mathbf{A} of the given discrete system is nonsingular.

There are many criteria for checking if a given system is reachable/controllable or not. We present a few of them in this section.

Theorem 3.3 The following are equivalent:

1. The system $\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$ is reachable.
2. $\text{rank}[\mathbf{B} \ \mathbf{AB} \ \dots \ \mathbf{A}^{n-1}\mathbf{B}] = n$.
3. The n rows of the matrix $e^{-At} \mathbf{B}$ are linearly independent over the set of real numbers for all $t \in [0, \infty)$.
4. The matrix

$$\mathbf{W}(t_0, t_1) = \int_{t_0}^{t_1} e^{-At} \mathbf{B} \mathbf{B}^T e^{-A^T t} dt$$

is nonsingular for all $t_1 > t_0$.

Proof We first prove the implication $1 \Rightarrow 2$; that is, if the pair (\mathbf{A}, \mathbf{B}) is reachable, then the so-called controllability matrix

$$[\mathbf{B} \ \mathbf{AB} \ \dots \ \mathbf{A}^{n-1}\mathbf{B}] \in \mathbb{R}^{n \times mn}$$

is of full rank n . We prove this statement by contraposition. Thus, we assume that

$$\text{rank}[\mathbf{B} \ \mathbf{AB} \ \dots \ \mathbf{A}^{n-1}\mathbf{B}] < n;$$

then, there is a constant n -vector $\mathbf{q} \neq \mathbf{0}$ so that

$$\mathbf{q}^T \mathbf{B} = \mathbf{0}^T, \quad \mathbf{q}^T \mathbf{A} \mathbf{B} = \mathbf{0}^T, \dots, \mathbf{q}^T \mathbf{A}^{n-1} \mathbf{B} = \mathbf{0}^T.$$

By the Cayley–Hamilton theorem the matrix \mathbf{A} satisfies its own characteristic equation. Hence

$$\mathbf{A}^n = -a_{n-1} \mathbf{A}^{n-1} - \dots - a_1 \mathbf{A} - a_0 \mathbf{I}_n.$$

Therefore

$$\mathbf{q}^T \mathbf{A}^n \mathbf{B} = \mathbf{q}^T (-a_{n-1} \mathbf{A}^{n-1} \mathbf{B} - \dots - a_1 \mathbf{A} \mathbf{B} - a_0 \mathbf{B}) = \mathbf{0}^T.$$

By induction we obtain

$$\mathbf{q}^T \mathbf{A}^i \mathbf{B} = \mathbf{0}^T \quad \text{for } i = n+1, n+2, \dots$$

Now let $\mathbf{x}(0) = \mathbf{0}$ and $\mathbf{x}(t_1) = \mathbf{q}$. We will now show that there is no control law that can transfer the system from $\mathbf{x}(0) = \mathbf{0}$ to $\mathbf{x}(t_1) = \mathbf{q}$. From the solution formula for the controlled system, we obtain

$$\mathbf{q} = \int_0^{t_1} e^{\mathbf{A}(t_1-t)} \mathbf{B} \mathbf{u}(t) dt.$$

We now premultiply the resulting equation by \mathbf{q}^T to get

$$0 \neq \|\mathbf{q}\|^2 = \int_0^{t_1} \mathbf{q}^T e^{\mathbf{A}(t_1-t)} \mathbf{B} \mathbf{u}(t) dt = 0$$

because

$$\mathbf{q}^T e^{\mathbf{A}(t_1-t)} \mathbf{B} = \mathbf{q}^T \left(\mathbf{B} + (t_1 - t) \mathbf{A} \mathbf{B} + \frac{(t_1 - t)^2}{2!} \mathbf{A}^2 \mathbf{B} - \dots \right) = \mathbf{0}^T.$$

Thus the system cannot be transferred from $\mathbf{x}(0) = \mathbf{0}$ to $\mathbf{x}(t_1) = \mathbf{q}$, and hence the pair (\mathbf{A}, \mathbf{B}) is not reachable.

We will now prove the implication $2 \Rightarrow 3$; that is, if the controllability matrix is of full rank, then the n rows of $e^{-\mathbf{A}t} \mathbf{B}$ are linearly independent over \mathbb{R} for all $t \in [0, \infty)$. We prove this statement by contraposition. So assume that for some nonzero n -vector \mathbf{q} we have

$$\mathbf{q}^T e^{-\mathbf{A}t} \mathbf{B} = \mathbf{0}^T, \quad \text{for all } t \in [0, \infty).$$

If we set $t = 0$, then we obtain

$$\mathbf{q}^T \mathbf{B} = \mathbf{0}^T.$$

Differentiating $\mathbf{q}^T e^{-\mathbf{A}t} \mathbf{B} = \mathbf{0}$ and letting $t = 0$ yields

$$\mathbf{q}^T \mathbf{A} \mathbf{B} = \mathbf{0}^T.$$

Continuing this process gives

$$\mathbf{q}^T [\mathbf{B} \quad \mathbf{A} \mathbf{B} \quad \dots \quad \mathbf{A}^{n-1} \mathbf{B}] = \mathbf{0}^T,$$

which means that $\text{rank}[\mathbf{B} \quad \mathbf{A} \mathbf{B} \quad \dots \quad \mathbf{A}^{n-1} \mathbf{B}] < n$.

We will now prove, by contraposition, the implication $3 \Rightarrow 4$; that is, if the n rows of $e^{-\mathbf{A}t} \mathbf{B}$ are linearly independent over \mathbb{R} for all $t \in [0, \infty)$, then $\det \mathbf{W}(t_0, t_1) \neq 0$ for any $t_1 > t_0$. So assume that $\det \mathbf{W}(t_0, t_1) = 0$. Then, there is a nonzero n -vector

$\mathbf{q} \in \mathbb{R}^n$ such that

$$\mathbf{q}^T \int_{t_0}^{t_1} e^{-\mathbf{A}t} \mathbf{B} \mathbf{B}^T e^{-\mathbf{A}^T t} dt = \mathbf{0}^T$$

for all $t \in [t_0, t_1]$. We postmultiply the above equation by \mathbf{q} to get

$$\int_{t_0}^{t_1} (\mathbf{q}^T e^{-\mathbf{A}t} \mathbf{B}) (\mathbf{B}^T e^{-\mathbf{A}^T t} \mathbf{q}) dt = 0.$$

Let $\mathbf{v}(t) = \mathbf{B}^T e^{-\mathbf{A}^T t} \mathbf{q}$; then the above equation can be written as

$$\int_{t_0}^{t_1} \|\mathbf{v}(t)\|^2 dt = 0.$$

However, the above is true if and only if $\mathbf{v}(t) = \mathbf{0}$ —that is, if and only if the rows of $e^{-\mathbf{A}t} \mathbf{B}$ are linearly dependent over the set of real numbers.

We now prove, constructively, the implication $4 \Rightarrow 1$. Let $\mathbf{x}(t_0) = \mathbf{x}_0$ be an arbitrary initial state and let $\mathbf{x}(t_1) = \mathbf{x}_f$ be an arbitrary final state. Then, we claim that the control law

$$\mathbf{u}(t) = \mathbf{B}^T e^{-\mathbf{A}^T t} \mathbf{W}^{-1}(t_0, t_1) (e^{-\mathbf{A}t_1} \mathbf{x}(t_1) - e^{-\mathbf{A}t_0} \mathbf{x}_0) \quad (3.12)$$

will transfer the system from \mathbf{x}_0 to \mathbf{x}_f . Note that to construct the above control law the matrix $\mathbf{W}(t_0, t_1)$ must be invertible. Substituting (3.12) into

$$\mathbf{x}(t_1) = e^{\mathbf{A}(t_1-t_0)} \mathbf{x}_0 + \int_{t_0}^{t_1} e^{\mathbf{A}(t_1-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau$$

yields

$$\begin{aligned} \mathbf{x}(t_1) &= e^{\mathbf{A}(t_1-t_0)} \mathbf{x}_0 \\ &+ e^{\mathbf{A}t_1} \int_{t_0}^{t_1} e^{-\mathbf{A}\tau} \mathbf{B} \mathbf{B}^T e^{-\mathbf{A}^T \tau} d\tau \mathbf{W}^{-1}(t_0, t_1) (e^{-\mathbf{A}t_1} \mathbf{x}(t_1) - e^{-\mathbf{A}t_0} \mathbf{x}_0). \end{aligned}$$

Performing easy manipulations gives $\mathbf{x}(t_1) = \mathbf{x}(t_1)$, which means that indeed the proposed control law does the required job. Thus, we have established equivalence between the four statements of the theorem.

◆ Example 3.2

Consider an armature-controlled DC motor system whose schematic is shown in Figure 3.1, where the system parameters are $R_a = 1 \Omega$, $L_a \approx 0 \text{ H}$, $K_b = 5 \text{ V/rad/sec}$, $K_i = 5 \text{ Nm/A}$, and the effective moment of inertia is $I_{eq} = 0.1 \text{ kg}\cdot\text{m}^2$. The gear inertia and friction are negligible.

Applying Kirchhoff's voltage law to the armature circuit gives

$$R_a i_a + K_b \dot{\theta} = e_a. \quad (3.13)$$

The armature-controlled DC motor uses a constant field current, and therefore the motor torque is

$$\tau_m = K_i i_a.$$

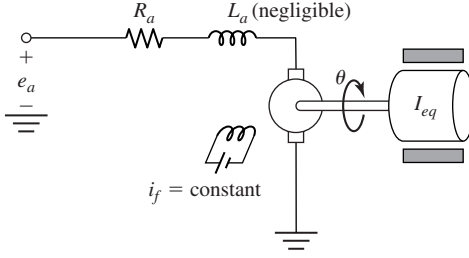


Figure 3.1 Schematic of an armature-controlled DC motor system of Example 3.2.

The motor torque is equal to the torque delivered to the load. Hence,

$$I_{eq} \ddot{\theta} = K_i i_a. \quad (3.14)$$

Substituting i_a from (3.13) into the above equation and dividing both sides by I_{eq} yields

$$\begin{aligned} \ddot{\theta} &= \frac{K_i}{I_{eq}} \left(\frac{e_a - K_b \dot{\theta}}{R_a} \right) \\ &= -\frac{K_i K_b}{I_{eq} R_a} \dot{\theta} + \frac{K_i}{I_{eq} R_a} e_a. \end{aligned}$$

Let $x_1 = \theta$, $x_2 = \dot{\theta}$, and $u = e_a$. Then taking into account the system parameters, we represent the above equation in a state-space format:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & -250 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 50 \end{bmatrix} u.$$

We will now find a control u that transfers the system from the initial state $\mathbf{x}(t_0) = \mathbf{x}(0) = [10 \ 1]^T$ to the final state $\mathbf{x}(2) = [0 \ 0]^T$. Recall that the solution of a controlled system has the form

$$\mathbf{x}(t) = e^{\mathbf{A}(t-t_0)} \mathbf{x}(t_0) + \int_{t_0}^t e^{\mathbf{A}(t-\tau)} \mathbf{B} u(\tau) d\tau.$$

Our goal is to find u so that

$$\mathbf{0} = e^{2\mathbf{A}} \mathbf{x}(0) + \int_0^2 e^{\mathbf{A}(2-\tau)} \mathbf{B} u(\tau) d\tau.$$

One can verify that one such u has the form

$$u(t) = -\mathbf{B}^T e^{\mathbf{A}^T(2-t)} \left(\int_0^2 e^{\mathbf{A}(2-\tau)} \mathbf{B} \mathbf{B}^T e^{\mathbf{A}^T(2-\tau)} d\tau \right)^{-1} e^{2\mathbf{A}} \mathbf{x}(0),$$

where

$$e^{\mathbf{A}t} = \begin{bmatrix} 1 & \frac{1}{250}(1 - e^{-250t}) \\ 0 & e^{-250t} \end{bmatrix}.$$

Hence,

$$\begin{aligned} \int_0^2 e^{A(2-\tau)} \mathbf{B} \mathbf{B}^T e^{A^T(2-\tau)} d\tau &= \int_0^2 \begin{bmatrix} 1 & \frac{1}{250} (1 - e^{-250(2-\tau)}) \\ 0 & e^{-250(2-\tau)} \end{bmatrix} \begin{bmatrix} 0 \\ 50 \end{bmatrix} [0 \ 50] \\ &\quad \times \begin{bmatrix} 1 & 0 \\ \frac{1}{250} (1 - e^{-250(2-\tau)}) & e^{-250(2-\tau)} \end{bmatrix} d\tau \\ &= \begin{bmatrix} 0.07976 & 0.02 \\ 0.02 & 5 \end{bmatrix}. \end{aligned}$$

The inverse of the above matrix is

$$\begin{bmatrix} 12.5502 & -0.0502 \\ -0.0502 & 0.2002 \end{bmatrix}.$$

Hence,

$$\begin{aligned} u(t) &= -[0 \ 50] \begin{bmatrix} 1 & 0 \\ \frac{1}{250} (1 - e^{-250(2-t)}) & e^{-250(2-t)} \end{bmatrix} \\ &\quad \times \begin{bmatrix} 12.5502 & -0.0502 \\ -0.0502 & 0.2002 \end{bmatrix} e^{2A} \begin{bmatrix} 10 \\ 1 \end{bmatrix} \\ &= 50.22e^{(250t-500)} - 25.11. \end{aligned}$$

The above control strategy transfers the system from $\mathbf{x}(0) = [10 \ 1]^T$ to $\mathbf{x}(2) = [0 \ 0]^T$. Plots of x_1 , x_2 versus time are shown in Figure 3.2. A plot of u versus time is given in Figure 3.3.

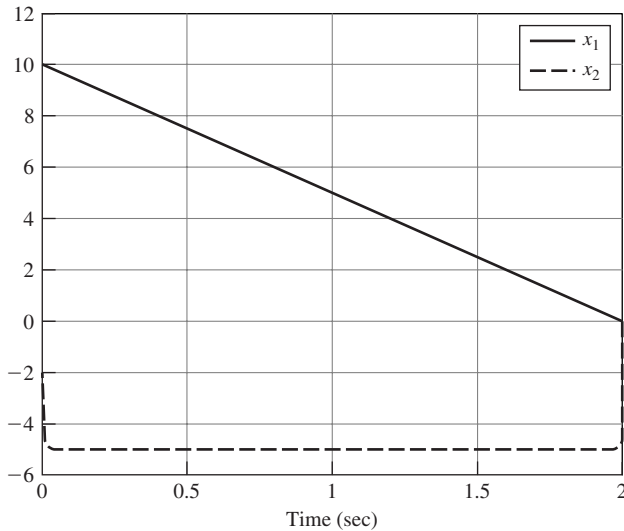


Figure 3.2 Plots of state variables versus time of Example 3.2.

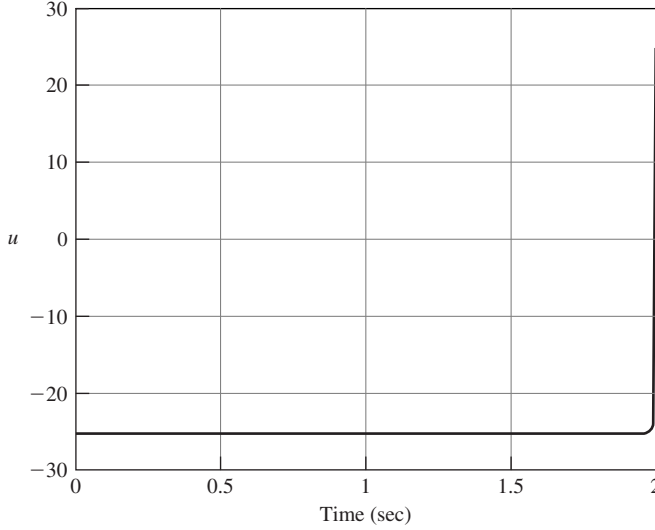


Figure 3.3 A plot of the control action versus time of Example 3.2.

We now give a necessary and sufficient condition for the pair (A, B) to be nonreachable.

Theorem 3.4 The pair (A, B) is nonreachable if and only if there is a similarity transformation $z = Tx$ such that

$$\tilde{A} = TAT^{-1} = \begin{bmatrix} A_1 & A_2 \\ O & A_4 \end{bmatrix}, \quad \tilde{B} = TB = \begin{bmatrix} B_1 \\ O \end{bmatrix},$$

where the pair (A_1, B_1) is reachable, $A_1 \in \mathbb{R}^{r \times r}$, $B_1 \in \mathbb{R}^{r \times m}$, and the rank of the controllability matrix of the pair (A, B) equals r .

Proof We first prove necessity (\Rightarrow). After forming the controllability matrix of the pair (A, B) , we use elementary row operations to get

$$T[B \quad AB \quad \cdots \quad A^{n-1}B] = [TB \quad (TAT^{-1})(TB) \quad \cdots \quad (TAT^{-1})^{n-1}(TB)]$$

$$= \begin{bmatrix} x & x & x & \cdots & x & x & x \\ 0 & x & x & \cdots & x & x & x \\ \vdots & & & \vdots & & & \vdots \\ 0 & 0 & 0 & \cdots & x & x & x \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \vdots & & & \vdots & & & \vdots \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} B_1 & A_1 B_1 & A_1^2 B_1 & \cdots & A_1^{n-1} B_1 \\ O & O & O & \cdots & O \end{bmatrix}, \quad (3.15)$$

where the symbol x denotes a “don’t care”—that is, an unspecified scalar. We now show that indeed

$$\tilde{\mathbf{A}} = \mathbf{T} \mathbf{A} \mathbf{T}^{-1} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{A}_3 & \mathbf{A}_4 \end{bmatrix}, \quad \tilde{\mathbf{B}} = \mathbf{T} \mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{O} \end{bmatrix},$$

where $\mathbf{A}_3 = \mathbf{O}$. It is clear that $\mathbf{T} \mathbf{B}$ must have the above form. If \mathbf{A}_3 were not a zero matrix, then we would have

$$\tilde{\mathbf{A}} \tilde{\mathbf{B}} = \begin{bmatrix} \mathbf{A}_1 \mathbf{B}_1 \\ \mathbf{A}_3 \mathbf{B}_1 \end{bmatrix}.$$

Comparing the above with (3.15), we conclude that $\mathbf{A}_3 \mathbf{B}_1 = \mathbf{O}$. Thus,

$$\tilde{\mathbf{A}} \tilde{\mathbf{B}} = \begin{bmatrix} \mathbf{A}_1 \mathbf{B}_1 \\ \mathbf{O} \end{bmatrix}.$$

Next we consider

$$\tilde{\mathbf{A}}^2 \tilde{\mathbf{B}} = \begin{bmatrix} \mathbf{A}_1^2 \mathbf{B}_1 \\ \mathbf{A}_3 \mathbf{A}_1 \mathbf{B}_1 \end{bmatrix}.$$

Again, comparing the above with (3.15), we conclude that

$$\mathbf{A}_3 \mathbf{A}_1 \mathbf{B}_1 = \mathbf{O},$$

and thus

$$\tilde{\mathbf{A}}^2 \tilde{\mathbf{B}} = \begin{bmatrix} \mathbf{A}_1^2 \mathbf{B}_1 \\ \mathbf{O} \end{bmatrix}.$$

Continuing in this manner, we conclude that $\mathbf{A}_3 \mathbf{A}_1^{n-1} \mathbf{B}_1 = \mathbf{O}$. Thus it follows that

$$\mathbf{A}_3 \begin{bmatrix} \mathbf{B}_1 & \mathbf{A}_1 \mathbf{B}_1 & \cdots & \mathbf{A}_1^{n-1} \mathbf{B}_1 \end{bmatrix} = \mathbf{O}.$$

Because $\text{rank}[\mathbf{B}_1 \ \mathbf{A}_1 \mathbf{B}_1 \ \cdots \ \mathbf{A}_1^{n-1} \mathbf{B}_1] = r$ —that is, the controllability matrix of the pair $(\mathbf{A}_1, \mathbf{B}_1)$ is of full rank—we have to have $\mathbf{A}_3 = \mathbf{O}$, and the proof of necessity is complete.

To prove sufficiency (\Leftarrow), note that if there is a similarity transformation such that

$$\tilde{\mathbf{A}} = \mathbf{T} \mathbf{A} \mathbf{T}^{-1} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{O} & \mathbf{A}_4 \end{bmatrix}, \quad \tilde{\mathbf{B}} = \mathbf{T} \mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{O} \end{bmatrix},$$

then the controllability matrix of the pair $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$ has the form

$$\begin{bmatrix} \mathbf{B}_1 & \mathbf{A}_1 \mathbf{B}_1 & \mathbf{A}_1^2 \mathbf{B}_1 & \cdots & \mathbf{A}_1^{n-1} \mathbf{B}_1 \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \cdots & \mathbf{O} \end{bmatrix}.$$

Hence, the pair $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$ is nonreachable. Because the similarity transformation preserves the reachability property, the pair (\mathbf{A}, \mathbf{B}) is nonreachable as well.

◆ Example 3.3

For the system model,

$$\dot{\mathbf{x}} = \begin{bmatrix} -1 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 4 & 0 \\ 1 & 1 & -1 & 2 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 1 & 1 \\ 0 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{u},$$

construct a state variable transformation so that in the new coordinates the reachable part is separated and distinct from the nonreachable part.

We first form the controllability matrix of the pair (\mathbf{A}, \mathbf{B}) to get

$$\mathbf{C} = [\mathbf{B} \quad \mathbf{AB} \quad \mathbf{A}^2\mathbf{B} \quad \mathbf{A}^3\mathbf{B}] = \begin{bmatrix} 1 & 1 & -1 & 0 & 3 & 4 & -1 & 4 \\ 0 & 0 & 1 & 1 & 0 & 1 & 3 & 5 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 3 & 2 & 7 & 7 & 19 \end{bmatrix}.$$

We then use the MATLAB command $[\mathbf{Q}, \mathbf{R}] = \text{qr}(\mathbf{C})$ to obtain

$$\mathbf{Q} = \begin{bmatrix} -1 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & -1 & 0 & 0 \end{bmatrix}, \quad \mathbf{R} = \begin{bmatrix} -1 & -1 & 1 & 0 & -3 & -4 & 1 & -4 \\ 0 & -1 & -1 & -3 & -2 & -7 & -7 & -19 \\ 0 & 0 & -1 & -1 & 0 & -1 & -3 & -5 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix},$$

where the matrices \mathbf{Q} and \mathbf{R} satisfy $\mathbf{C} = \mathbf{QR}$, and \mathbf{R} is upper row triangular such that $\text{rank } \mathbf{R} = \text{rank } \mathbf{C}$. Thus, premultiplying the controllability matrix, \mathbf{C} , by \mathbf{Q}^{-1} reduces this matrix to an upper row triangular matrix \mathbf{R} because $\mathbf{Q}^{-1}\mathbf{C} = \mathbf{R}$. Let $\mathbf{z} = \mathbf{Q}^{-1}\mathbf{x}$ be the state variable transformation. Then in the new coordinates the matrices \mathbf{A} and \mathbf{B} take the form

$$\mathbf{Q}^{-1}\mathbf{A}\mathbf{Q} = \left[\begin{array}{ccc|c} -1 & 1 & 1 & 0 \\ 1 & 2 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ \hline 0 & 0 & 0 & 4 \end{array} \right] \quad \text{and} \quad \mathbf{Q}^{-1}\mathbf{B} = \left[\begin{array}{cc} -1 & -1 \\ 0 & -1 \\ 0 & 0 \\ \hline 0 & 0 \end{array} \right].$$

The dimension of the nonreachable part is 1.

If the eigenvalues of the noncontrollable part of a given pair (\mathbf{A}, \mathbf{B}) are all in the open left-hand complex plane—that is, the nonreachable part is asymptotically stable—then the pair (\mathbf{A}, \mathbf{B}) is called *stabilizable*.

We now give yet another test for reachability. This test is also known as the Popov–Belevitch–Hautus (PBH) eigenvector test. We denote the set of eigenvalues of a matrix \mathbf{A} as $\text{eig}(\mathbf{A})$.

Theorem 3.5 The pair (\mathbf{A}, \mathbf{B}) is reachable if and only if

$$\text{rank}[s\mathbf{I}_n - \mathbf{A} \quad \mathbf{B}] = n \quad \text{for all } s \in \text{eig}(\mathbf{A}).$$

Proof We first prove the necessity part (\Rightarrow) by contraposition. If there is a vector $\mathbf{q} \neq \mathbf{0}$ such that for some $s \in \text{eig}(\mathbf{A})$,

$$\mathbf{q}^T [s\mathbf{I}_n - \mathbf{A} \quad \mathbf{B}] = \mathbf{0}^T,$$

that is,

$$\mathbf{q}^T \mathbf{A} = s\mathbf{q}^T, \quad \mathbf{q}^T \mathbf{B} = \mathbf{0}^T,$$

then

$$\begin{aligned} \mathbf{q}^T \mathbf{A} \mathbf{B} &= s\mathbf{q}^T \mathbf{B} = \mathbf{0}^T, \\ \mathbf{q}^T \mathbf{A}^2 \mathbf{B} &= s\mathbf{q}^T \mathbf{A} \mathbf{B} = \mathbf{0}^T, \end{aligned}$$

and so on, until

$$\mathbf{q}^T \mathbf{A}^{n-1} \mathbf{B} = s\mathbf{q}^T \mathbf{A}^{n-2} \mathbf{B} = \mathbf{0}^T.$$

We write the above equations in the form

$$\mathbf{q}^T [\mathbf{B} \quad \mathbf{A} \mathbf{B} \quad \cdots \quad \mathbf{A}^{n-1} \mathbf{B}] = \mathbf{0}^T,$$

which means that the pair (\mathbf{A}, \mathbf{B}) is nonreachable.

We now prove the sufficiency part (\Leftarrow) also by contraposition. So we assume that $\text{rank}[\mathbf{B} \quad \mathbf{A} \mathbf{B} \quad \cdots \quad \mathbf{A}^{n-1} \mathbf{B}] = r < n$. By Theorem 3.4, there exists a state variable transformation such that in the new basis the pair (\mathbf{A}, \mathbf{B}) takes the form

$$\tilde{\mathbf{A}} = \mathbf{T} \mathbf{A} \mathbf{T}^{-1} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{O} & \mathbf{A}_4 \end{bmatrix}, \quad \tilde{\mathbf{B}} = \mathbf{T} \mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{O} \end{bmatrix}.$$

We choose \mathbf{q} of the form

$$\mathbf{q}^T = [\mathbf{0} \quad \mathbf{z}],$$

where $\mathbf{z} \neq \mathbf{0}$ is a left eigenvector of \mathbf{A}_4 ; that is,

$$\mathbf{z} \mathbf{A}_4 = s\mathbf{z},$$

where $s \in \text{eig}(\mathbf{A}_4)$. Thus

$$\begin{aligned} \mathbf{q}^T \tilde{\mathbf{A}} &= [\mathbf{0} \quad \mathbf{z}] \tilde{\mathbf{A}} \\ &= [\mathbf{0} \quad s\mathbf{z}] \\ &= s\mathbf{q}^T. \end{aligned}$$

Clearly,

$$\mathbf{q}^T [s\mathbf{I}_n - \tilde{\mathbf{A}} \quad \tilde{\mathbf{B}}] = \mathbf{0}^T,$$

which means that the rank of the matrix $[s\mathbf{I}_n - \tilde{\mathbf{A}} \quad \tilde{\mathbf{B}}]$ is less than n , and therefore the proof of the sufficiency part is complete.

3.2 Observability and Constructability

As we have seen in the previous section, system reachability implies the ability to completely control the entire state of the system. Suppose now that the system state is not directly accessible. Instead, we have the output of the system,

$$\mathbf{y} = \mathbf{C}\mathbf{x} + \mathbf{D}\mathbf{u},$$

where $\mathbf{C} \in \mathbb{R}^{p \times n}$, $p \leq n$, and $\mathbf{D} \in \mathbb{R}^{p \times m}$. However, we still want to know the dynamical behavior of the entire state.

We say that the system

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}, \\ \mathbf{y} &= \mathbf{C}\mathbf{x} + \mathbf{D}\mathbf{u},\end{aligned}$$

or equivalently the pair (\mathbf{A}, \mathbf{C}) , is *observable* if there is a finite $t_1 > t_0$ such that for arbitrary $\mathbf{u}(t)$ and resulting $\mathbf{y}(t)$ over $[t_0, t_1]$, we can determine $\mathbf{x}(t_0)$ from complete knowledge of the system input \mathbf{u} and output \mathbf{y} .

Note that once $\mathbf{x}(t_0)$ is known, we can determine $\mathbf{x}(t)$ from knowledge of $\mathbf{u}(t)$ and $\mathbf{y}(t)$ over any finite time interval $[t_0, t_1]$. We now show how to constructively determine $\mathbf{x}(t_0)$, given $\mathbf{u}(t)$ and $\mathbf{y}(t)$. While doing so we will establish a sufficiency condition for the system to be observable. The condition is also necessary for the system observability, as we show next. To proceed with constructive determination of the state vector, we note that the solution $\mathbf{y}(t)$ is given by

$$\mathbf{y}(t) = \mathbf{C}e^{\mathbf{A}(t-t_0)}\mathbf{x}_0 + \int_{t_0}^t \mathbf{C}e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau) d\tau + \mathbf{D}\mathbf{u}(t).$$

We subtract $\int_{t_0}^t \mathbf{C}e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau) d\tau + \mathbf{D}\mathbf{u}(t)$ from both sides of the above equation. Let

$$\mathbf{g}(t) = \mathbf{y}(t) - \int_{t_0}^t \mathbf{C}e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau) d\tau - \mathbf{D}\mathbf{u}(t),$$

then we have

$$\mathbf{g}(t) = \mathbf{C}e^{\mathbf{A}(t-t_0)}\mathbf{x}(t_0),$$

where \mathbf{g} is known to us. Our goal now is to determine $\mathbf{x}(t_0)$. Having $\mathbf{x}(t_0)$, we can determine the entire $\mathbf{x}(t)$ for all $t \in [t_0, t_1]$ from the formula

$$\mathbf{x}(t) = e^{\mathbf{A}(t-t_0)}\mathbf{x}_0 + \int_{t_0}^t e^{\mathbf{A}(t-\tau)}\mathbf{B}\mathbf{u}(\tau) d\tau.$$

Premultiplying both side of $\mathbf{g}(t)$ by $e^{\mathbf{A}^T(t-t_0)}\mathbf{C}^T$ and integrating the resulting expression between the limits t_0 and t_1 , we get

$$\int_{t_0}^{t_1} e^{\mathbf{A}^T(t-t_0)}\mathbf{C}^T\mathbf{C}e^{\mathbf{A}(t-t_0)} dt \mathbf{x}_0 = \int_{t_0}^{t_1} e^{\mathbf{A}^T(t-t_0)}\mathbf{C}^T\mathbf{g}(t) dt.$$

Let

$$V(t_0, t_1) = \int_{t_0}^{t_1} e^{A^T t} \mathbf{C}^T \mathbf{C} e^{A t} dt.$$

Then, after performing some manipulations and assuming that the matrix $V(t_0, t_1)$ is invertible, we obtain

$$\mathbf{x}_0 = e^{A t_0} \mathbf{V}^{-1}(t_0, t_1) e^{A^T t_0} \int_{t_0}^{t_1} e^{A^T(t-t_0)} \mathbf{C}^T \mathbf{g}(t) dt.$$

Knowledge of \mathbf{x}_0 allows us to reconstruct the entire state $\mathbf{x}(t)$ over the interval $[t_0, t_1]$. We conclude that if the matrix $V(t_0, t_1)$ is invertible, then the system is observable.

◆ Example 3.4

For a dynamical system modeled by

$$\begin{aligned} \dot{\mathbf{x}} &= \begin{bmatrix} -1/2 & 1/2 \\ -1/2 & 1/2 \end{bmatrix} \mathbf{x} + \begin{bmatrix} -1 \\ 1 \end{bmatrix} u = \mathbf{A}\mathbf{x} + \mathbf{b}u, \\ y &= [1 \quad 1] \mathbf{x} = \mathbf{c}\mathbf{x}, \end{aligned}$$

we shall determine $\mathbf{x}(t)$ over the time interval $[0, 10]$ knowing that $u(t) = 1$ for $t \geq 0$ and $y(t) = t^2 + 2t$ for $t \geq 0$.

We first find $\mathbf{x}(0)$. Knowledge of $\mathbf{x}(0)$ will allow us to find the entire state vector for all $t \in [0, 10]$. We have $t_0 = 0$ and $t_1 = 10$. Hence,

$$\begin{aligned} \mathbf{x}(t_0) = \mathbf{x}(0) &= e^{A t_0} \mathbf{V}^{-1}(t_0, t_1) e^{A^T t_0} \int_{t_0}^{t_1} e^{A^T(t-t_0)} \mathbf{C}^T \mathbf{g}(t) dt \\ &= \mathbf{V}^{-1}(0, 10) \int_0^{10} e^{A^T t} \mathbf{C}^T \mathbf{g}(t) dt, \end{aligned}$$

where

$$\mathbf{g}(t) = y(t) - \mathbf{c} \int_0^t e^{A(t-\tau)} \mathbf{b}u(\tau) d\tau.$$

To evaluate the above expression, we first use formula (A.148) to compute

$$\begin{aligned} e^{A t} &= \mathcal{L}^{-1}([s\mathbf{I}_2 - \mathbf{A}]^{-1}) = \mathcal{L}^{-1} \left(\begin{bmatrix} \frac{1}{s} - \frac{1}{2} \frac{1}{s^2} & \frac{1}{2} \frac{1}{s^2} \\ -\frac{1}{2} \frac{1}{s^2} & \frac{1}{s} + \frac{1}{2} \frac{1}{s^2} \end{bmatrix} \right) \\ &= \begin{bmatrix} 1 - \frac{1}{2}t & \frac{1}{2}t \\ -\frac{1}{2}t & 1 + \frac{1}{2}t \end{bmatrix}. \end{aligned}$$

Next, we compute

$$\begin{aligned}
 \mathbf{V}(0, 10) &= \int_0^{10} e^{\mathbf{A}^T t} \mathbf{c}^T \mathbf{c} e^{\mathbf{A} t} dt \\
 &= \int_0^{10} \left(\begin{bmatrix} 1-t \\ 1+t \end{bmatrix} \begin{bmatrix} 1-t & 1+t \end{bmatrix} \right) dt = \int_0^{10} \begin{bmatrix} 1-2t+t^2 & 1-t^2 \\ 1-t^2 & 1+2t+t^2 \end{bmatrix} dt \\
 &= \left[\begin{array}{cc} t-t^2+\frac{t^3}{3} & t-\frac{t^3}{3} \\ t-\frac{t^3}{3} & t+t^2+\frac{t^3}{3} \end{array} \right] \bigg|_0^{10} = \begin{bmatrix} 243.333 & -323.333 \\ -323.333 & 443.333 \end{bmatrix}.
 \end{aligned}$$

The inverse of $\mathbf{V}(0, 10)$ is

$$\mathbf{V}^{-1}(0, 10) = \begin{bmatrix} 0.133 & 0.097 \\ 0.097 & 0.073 \end{bmatrix}.$$

We can now compute $g(t)$ to get

$$\begin{aligned}
 g(t) &= y(t) - \mathbf{c} \int_0^t e^{\mathbf{A}(t-\tau)} \mathbf{b} u(\tau) d\tau \\
 &= t^2 + 2t - [1 \quad 1] \int_0^t \begin{bmatrix} -1+t-\tau \\ 1+t+\tau \end{bmatrix} d\tau = t^2 + 2t - \int_0^t (2t-2\tau) d\tau \\
 &= t^2 + 2t - 2t(\tau) \bigg|_0^t + \tau^2 \bigg|_0^t = t^2 + 2t - t^2 \\
 &= 2t.
 \end{aligned}$$

Hence,

$$\int_0^{10} e^{\mathbf{A}^T t} \mathbf{c}^T g(t) dt = \int_0^{10} \begin{bmatrix} (1-t)2t \\ (1+t)2t \end{bmatrix} dt = \left[\begin{array}{c} t^2 - \frac{2}{3}t^3 \\ t^2 + \frac{2}{3}t^3 \end{array} \right] \bigg|_0^{10} = \begin{bmatrix} -566.6667 \\ 766.6667 \end{bmatrix}.$$

We can now compute

$$\mathbf{x}(0) = \mathbf{V}^{-1}(0, 10) \int_0^{10} e^{\mathbf{A}^T t} \mathbf{c}^T g(t) dt = \begin{bmatrix} -1 \\ 1 \end{bmatrix}.$$

Alternatively, we could use the Laplace transform method to find $\mathbf{x}(0)$. Knowing $\mathbf{x}(0)$, we can find $\mathbf{x}(t)$ for $t \geq 0$,

$$\begin{aligned}
 \mathbf{x}(t) &= e^{\mathbf{A} t} \mathbf{x}(0) + \int_0^t e^{\mathbf{A}(t-\tau)} \mathbf{b} u(\tau) d\tau \\
 &= \begin{bmatrix} -1+t \\ 1+t \end{bmatrix} + \begin{bmatrix} (-1+t)t - \frac{t^2}{2} \\ (1+t)t - \frac{t^2}{2} \end{bmatrix} = \begin{bmatrix} \frac{t^2}{2} - 1 \\ \frac{t^2}{2} + 2t + 1 \end{bmatrix}.
 \end{aligned}$$

We will now show that the invertibility of $V(t_0, t_1)$ is also necessary for observability. We first note that one can use similar arguments to those in the previous section to show that $V(t_0, t_1)$ is nonsingular if and only if the n columns of the matrix

$$Ce^{At}$$

are linearly independent for all $t \in [0, \infty)$ over the set of real numbers. Hence, if $V(t_0, t_1)$ is singular, there exists a nonzero constant vector, say \mathbf{x}_a , such that $V(t_0, t_1)\mathbf{x}_a = \mathbf{0}$. This implies that

$$\mathbf{g}(t) = Ce^{At}\mathbf{x}_0 = Ce^{At}(\mathbf{x}_0 + \mathbf{x}_a).$$

Thus, $\mathbf{x}(0) = \mathbf{x}_0 + \mathbf{x}_a$ yields the same response as $\mathbf{x}(0) = \mathbf{x}_0$, which means that we cannot determine the system state, or in other words, the state fails to be observable. In view of the above result, it follows that by repeating essentially the same arguments used in the previous section we can prove the following theorem:

Theorem 3.6 The following statements are equivalent:

1. The pair (\mathbf{A}, \mathbf{C}) is observable.
2. The matrix $V(t_0, t_1)$ is nonsingular for all $t_1 > t_0$.
3. The n columns of Ce^{At} are linearly independent for all $t \in [0, \infty)$ over the real numbers.
4. The observability matrix

$$\begin{bmatrix} \mathbf{C} \\ \mathbf{CA} \\ \vdots \\ \mathbf{CA}^{n-1} \end{bmatrix} \in \mathbb{R}^{pn \times n}$$

is of full rank n .

5.

$$\text{rank} \begin{bmatrix} s\mathbf{I}_n - \mathbf{A} \\ \mathbf{C} \end{bmatrix} = n \quad \text{for all } s \in \text{eig}(\mathbf{A}).$$

◆ Example 3.5

For the dynamical system model,

$$\begin{aligned} \dot{\mathbf{x}} &= \mathbf{Ax} = \begin{bmatrix} 2 & 0 \\ 3 & -1 \end{bmatrix} \mathbf{x}, \\ y &= \mathbf{cx} = [-1 \quad 1] \mathbf{x}, \end{aligned}$$

we will find a nonzero initial vector $\mathbf{x}(0) = \mathbf{x}_0$ such that $y(t) = 0$ for all $t \geq 0$.

We have $\mathbf{x}(t) = e^{\mathbf{A}t} \mathbf{x}_0$, where

$$e^{\mathbf{A}t} = \mathcal{L}^{-1}((s\mathbf{I} - \mathbf{A})^{-1}) = \begin{bmatrix} e^{2t} & 0 \\ e^{2t} - e^{-t} & e^{-t} \end{bmatrix}.$$

We want

$$0 = y(t) = [-1 \quad 1]e^{\mathbf{A}t} \mathbf{x}_0 \quad \text{for all } t \geq 0.$$

For example, if

$$\mathbf{x}_0 = [1 \quad 1]^T,$$

then $y(t) = 0$ for all $t \geq 0$. Note that we were able to find $\mathbf{x}_0 \neq \mathbf{0}$ such that $y(t) = 0$ for all $t \geq 0$, because for all $t \geq 0$ the columns of $\mathbf{c}e^{\mathbf{A}t}$ are linearly dependent over the real numbers because the pair (\mathbf{A}, \mathbf{c}) is nonobservable. The vector \mathbf{x}_0 is in the null space of the observability matrix

$$\begin{bmatrix} \mathbf{c} \\ \mathbf{c}\mathbf{A} \end{bmatrix}.$$

We can also prove a result analogous to that of Theorem 3.4 for the given pair (\mathbf{A}, \mathbf{C}) ; that is, using a similarity transformation we can separate the observable part from the nonobservable. If the eigenvalues of the nonobservable part of a given pair (\mathbf{A}, \mathbf{C}) are all in the open left-hand complex plane—that is, the nonobservable part is asymptotically stable—then we say that the pair (\mathbf{A}, \mathbf{C}) is *detectable*.

A notion related to observability is the notion of constructability. We say that the system

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u},$$

$$\mathbf{y} = \mathbf{C}\mathbf{x} + \mathbf{D}\mathbf{u},$$

or equivalently the pair (\mathbf{A}, \mathbf{C}) , is *constructable* if there is a finite $t_1 > t_0$ such that for arbitrary $\mathbf{u}(t)$ and resulting $\mathbf{y}(t)$ over $[t_0, t_1]$ we can determine $\mathbf{x}(t_1)$ from complete knowledge of the system input \mathbf{u} and output \mathbf{y} .

One can verify that in the case of continuous, time-invariant systems, observability is equivalent to constructability. In the case of discrete-time systems, observability implies constructability, and the two notions are equivalent if the matrix \mathbf{A} of the given discrete system is nonsingular.

3.3 Companion Forms

In this section we discuss two types of special forms of a given linear dynamical system model represented by the triple $(\mathbf{A}, \mathbf{B}, \mathbf{C})$. These forms are obtained by a change of state variables and reveal structural properties of the model. We will use the above-mentioned forms to construct algorithms for pole placement. The pole placement algorithms will then be used to design state-feedback controllers, as well as state estimators. We begin by introducing the controller companion form, or controller form for short.

3.3.1 Controller Form

We begin by considering a single-input, continuous-time model represented by

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{b}u(t),$$

where the pair (\mathbf{A}, \mathbf{b}) is assumed to be reachable, or, equivalently, controllable. This means that

$$\text{rank}[\mathbf{b} \quad \mathbf{A}\mathbf{b} \quad \cdots \quad \mathbf{A}^{n-1}\mathbf{b}] = n.$$

We select the last row of the inverse of the controllability matrix. Let \mathbf{q}_1 be that row. Next, we form the matrix

$$\mathbf{T} = \begin{bmatrix} \mathbf{q}_1 \\ \mathbf{q}_1\mathbf{A} \\ \vdots \\ \mathbf{q}_1\mathbf{A}^{n-1} \end{bmatrix}.$$

Note that \mathbf{T} is invertible. Indeed, because

$$\mathbf{T}[\mathbf{b} \quad \mathbf{A}\mathbf{b} \quad \cdots \quad \mathbf{A}^{n-1}\mathbf{b}] = \begin{bmatrix} 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & \cdots & 1 & x \\ \vdots & \vdots & & \vdots & \vdots \\ 1 & x & \cdots & x & x \end{bmatrix},$$

this implies that \mathbf{T} must be nonsingular. In the above, we used the symbol x to denote “don’t care,” unspecified scalars in our present discussion. Consider now the state variable transformation, $\tilde{\mathbf{x}} = \mathbf{T}\mathbf{x}$. Then, in the new coordinates, the system model is

$$\begin{aligned} \dot{\tilde{\mathbf{x}}}(t) &= \mathbf{T}\mathbf{A}\mathbf{T}^{-1}\tilde{\mathbf{x}}(t) + \mathbf{T}\mathbf{b}u(t) \\ &= \tilde{\mathbf{A}}\tilde{\mathbf{x}}(t) + \tilde{\mathbf{b}}u(t). \end{aligned}$$

The matrices $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{b}}$ have particular structures that we now discuss. We first analyze the structure of $\tilde{\mathbf{b}}$. Because \mathbf{q}_1 is the last row of the inverse of the controllability matrix, we have

$$\mathbf{q}_1\mathbf{b} = \mathbf{q}_1\mathbf{A}\mathbf{b} = \cdots = \mathbf{q}_1\mathbf{A}^{n-2}\mathbf{b} = 0$$

and

$$\mathbf{q}_1\mathbf{A}^{n-1}\mathbf{b} = 1.$$

Hence,

$$\tilde{\mathbf{b}} = \mathbf{T}\mathbf{b} = \begin{bmatrix} \mathbf{q}_1\mathbf{b} \\ \mathbf{q}_1\mathbf{A}\mathbf{b} \\ \vdots \\ \mathbf{q}_1\mathbf{A}^{n-2}\mathbf{b} \\ \mathbf{q}_1\mathbf{A}^{n-1}\mathbf{b} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}.$$

The structure of $\tilde{\mathbf{A}}$ is revealed by considering the relation $\mathbf{T}\mathbf{A}\mathbf{T}^{-1} = \tilde{\mathbf{A}}$ represented as

$$\mathbf{T}\mathbf{A} = \tilde{\mathbf{A}}\mathbf{T}.$$

The left-hand side of the above matrix equation is

$$T\mathbf{A} = \begin{bmatrix} q_1\mathbf{A} \\ q_1\mathbf{A}^2 \\ \vdots \\ q_1\mathbf{A}^{n-1} \\ q_1\mathbf{A}^n \end{bmatrix}.$$

By the Cayley–Hamilton theorem, we have

$$\mathbf{A}^n = -a_0\mathbf{I}_n - a_1\mathbf{A} - \cdots - a_{n-1}\mathbf{A}^{n-1},$$

and hence

$$q_1\mathbf{A}^n = -a_0q_1 - a_1q_1\mathbf{A} - \cdots - a_{n-1}q_1\mathbf{A}^{n-1}.$$

Comparing both sides of the equation $T\mathbf{A} = \tilde{\mathbf{A}}T$ and taking into account the Cayley–Hamilton theorem, we get

$$\tilde{\mathbf{A}} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & & & & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ -a_0 & -a_1 & -a_2 & \cdots & -a_{n-2} & -a_{n-1} \end{bmatrix}.$$

Note that the coefficients of the characteristic polynomial of \mathbf{A} are immediately apparent by inspecting the last row of $\tilde{\mathbf{A}}$. There are other benefits of the pair (\mathbf{A}, \mathbf{b}) being given in its controller form that we discuss in Section 3.4. We now illustrate the above algorithm for reducing a single-input system into the controller form with a numerical example.

◆ Example 3.6

Consider the following dynamical system model:

$$\dot{\mathbf{x}} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -2 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} u.$$

The controllability matrix of the pair (\mathbf{A}, \mathbf{b}) is

$$[\mathbf{b} \quad \mathbf{A}\mathbf{b} \quad \mathbf{A}^2\mathbf{b} \quad \mathbf{A}^3\mathbf{b}] = \begin{bmatrix} 0 & 1 & 0 & -2 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & -2 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}.$$

The pair (\mathbf{A}, \mathbf{b}) is controllable because the controllability matrix is nonsingular. We find the last row of the inverse of the controllability matrix to be

$\mathbf{q}_1 = [0 \ 1 \ 0 \ 0]$. Hence the transformation matrix is

$$T = \begin{bmatrix} \mathbf{q}_1 \\ \mathbf{q}_1 \mathbf{A} \\ \mathbf{q}_1 \mathbf{A}^2 \\ \mathbf{q}_1 \mathbf{A}^3 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{bmatrix}.$$

The matrices \mathbf{A} and \mathbf{b} in the new coordinates have the form

$$T \mathbf{A} T^{-1} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -2 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -1 & 0 & -3 & 0 \end{bmatrix}$$

and

$$T \mathbf{b} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}.$$

We now generalize the above method of transforming the given pair (\mathbf{A}, \mathbf{B}) into the controller form for multi-input systems. We consider a controllable pair (\mathbf{A}, \mathbf{B}) , where $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{B} \in \mathbb{R}^{n \times m}$, $m \leq n$, and $\text{rank } \mathbf{B} = m$. The last assumption we made is just for convenience to simplify further manipulations. This assumption implies that all inputs are mutually independent, which is usually the case in practical applications. We first form the controllability matrix of the pair (\mathbf{A}, \mathbf{B}) and represent it as

$$[\mathbf{b}_1 \ \mathbf{b}_2 \ \cdots \ \mathbf{b}_m \ \mathbf{A}\mathbf{b}_1 \ \mathbf{A}\mathbf{b}_2 \ \cdots \ \mathbf{A}\mathbf{b}_m \ \cdots \ \mathbf{A}^{n-1}\mathbf{b}_1 \ \cdots \ \mathbf{A}^{n-1}\mathbf{b}_m], \quad (3.16)$$

where \mathbf{b}_i , $i = 1, 2, \dots, m$, is the i th column of \mathbf{B} . The above controllability matrix is of full rank because we assumed that the pair (\mathbf{A}, \mathbf{B}) is controllable. Thus, we can select n linearly independent columns from (3.16). We select these columns proceeding from left to right. Then, we rearrange the selected columns to form a $n \times n$ nonsingular matrix \mathbf{L} of the form

$$\mathbf{L} = [\mathbf{b}_1 \ \mathbf{A}\mathbf{b}_1 \ \cdots \ \mathbf{A}^{d_1-1}\mathbf{b}_1 \ \mathbf{b}_2 \ \cdots \ \mathbf{A}^{d_2-1}\mathbf{b}_2 \ \cdots \ \mathbf{b}_m \ \cdots \ \mathbf{A}^{d_m-1}\mathbf{b}_m].$$

We call the m integers d_i the *controllability indices* of the pair (\mathbf{A}, \mathbf{B}) . Note that

$$\sum_{i=1}^m d_i = n.$$

We call the integer

$$d = \max\{d_i : i = 1, 2, \dots, m\}$$

the *controllability index* of the pair (\mathbf{A}, \mathbf{B}) . Observe that all m columns of \mathbf{B} are present in \mathbf{L} because we assumed that \mathbf{B} is of full rank, and observe the way we selected mutually independent columns from the controllability matrix (3.16) to construct the matrix \mathbf{L} . Furthermore, if the

column $A^k \mathbf{b}_j$ is present in \mathbf{L} , then so is $A^{k-1} \mathbf{b}_j$. To proceed further, we let

$$\sigma_k = \sum_{i=1}^k d_i \quad \text{for } k = 1, 2, \dots, m.$$

Note that $\sigma_1 = d_1$, $\sigma_2 = d_1 + d_2$, \dots , $\sigma_m = \sum_{i=1}^m d_i = n$. We now select m rows from \mathbf{L}^{-1} . We denote the selected rows by \mathbf{q}_k , where \mathbf{q}_k is the σ_k th row of \mathbf{L}^{-1} . Next, we form the $n \times n$ matrix

$$\mathbf{T} = \begin{bmatrix} \mathbf{q}_1 \\ \mathbf{q}_1 \mathbf{A} \\ \vdots \\ \mathbf{q}_1 \mathbf{A}^{d_1-1} \\ \mathbf{q}_2 \\ \vdots \\ \mathbf{q}_2 \mathbf{A}^{d_2-1} \\ \vdots \\ \mathbf{q}_m \\ \vdots \\ \mathbf{q}_m \mathbf{A}^{d_m-1} \end{bmatrix}.$$

We can show that \mathbf{T} is nonsingular by computing the product $\mathbf{T}\mathbf{L}$ and noticing that $\det(\mathbf{T}\mathbf{L}) = \pm 1$. We define a state-variable transformation $\tilde{\mathbf{x}} = \mathbf{T}\mathbf{x}$; and then by using similar arguments as in the single-input case, we arrive at the special structure of the pair

$$(\tilde{\mathbf{A}}, \tilde{\mathbf{B}}) = (\mathbf{T}\mathbf{A}\mathbf{T}^{-1}, \mathbf{T}\mathbf{B}),$$

where

$$\tilde{\mathbf{A}} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 & 0 & \cdots & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 & 0 & \cdots & 0 & \cdots & 0 & 0 & \cdots & 0 \\ & & & \ddots & & \vdots & & & \vdots & & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & 1 & 0 & 0 & \cdots & 0 & \cdots & 0 & 0 & \cdots & 0 \\ x & x & x & \cdots & x & x & x & \cdots & x & \cdots & x & x & \cdots & x \\ \hline 0 & 0 & 0 & \cdots & 0 & 0 & 1 & \cdots & 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & & & & \vdots & \vdots & & \ddots & \vdots & & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 0 & 0 & \cdots & 1 & \cdots & 0 & 0 & \cdots & 0 \\ x & x & x & \cdots & x & x & x & \cdots & x & \cdots & x & x & \cdots & x \\ \hline \vdots & & & & \vdots & \vdots & & \ddots & \vdots & & \vdots & & \vdots \\ \hline 0 & 0 & 0 & \cdots & 0 & 0 & 0 & \cdots & 0 & \cdots & 0 & 1 & \cdots & 0 \\ \vdots & & & & \vdots & \vdots & & & \vdots & & \vdots & & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 0 & 0 & \cdots & 0 & \cdots & 0 & 0 & \cdots & 1 \\ x & x & x & \cdots & x & x & x & \cdots & x & \cdots & x & x & \cdots & x \end{bmatrix}.$$

Note that the m diagonal blocks of $\tilde{\mathbf{A}}$ have the structure of the controller form as in the single-input case. The dimensions of the diagonal blocks are $d_i \times d_i$, $i = 1, 2, \dots, m$. The off-diagonal blocks consist of zero elements except for their respective last rows. Thus, all information about $\tilde{\mathbf{A}}$, as well as \mathbf{A} , can be derived from knowledge of the controllability indices, as well as the m nonzero rows containing the elements marked using x , of the matrix $\tilde{\mathbf{A}}$. The matrix $\tilde{\mathbf{B}}$ has the form

$$\tilde{\mathbf{B}} = \begin{bmatrix} 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ 1 & x & x & x & x \\ \hline 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ 0 & 1 & x & x & x \\ \hline \vdots & & & & \vdots \\ \hline 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}.$$

The nonzero rows of $\tilde{\mathbf{B}}$ correspond to the rows of $\tilde{\mathbf{A}}$ containing the elements marked using x . We now illustrate the above algorithm for reducing a multi-input system model into the controller form with a numerical example.

◆ Example 3.7

Consider the pair (\mathbf{A}, \mathbf{B}) , where

$$\mathbf{A} = \begin{bmatrix} -1 & 1 & 0 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 4 & 0 \\ 1 & 1 & -1 & 2 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 1 & 1 \\ 2 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

We transform the above pair into its controller companion form. Following the algorithm, we find

$$\text{rank}[\mathbf{b}_1 \quad \mathbf{b}_2 \quad \mathbf{A}\mathbf{b}_1 \quad \mathbf{A}\mathbf{b}_2] = 4.$$

Hence $d_1 = d_2 = 2$. We next form the matrix \mathbf{L} ;

$$\mathbf{L} = [\mathbf{b}_1 \quad \mathbf{A}\mathbf{b}_1 \quad \mathbf{b}_2 \quad \mathbf{A}\mathbf{b}_2] = \begin{bmatrix} 1 & 1 & 1 & 0 \\ 2 & 3 & 0 & 1 \\ 1 & 4 & 0 & 0 \\ 0 & 2 & 1 & 3 \end{bmatrix}$$

and

$$\mathbf{L}^{-1} = \begin{bmatrix} 0.20 & 0.60 & 0.40 & -0.20 \\ -0.05 & -0.15 & 0.35 & 0.05 \\ 0.85 & -0.45 & 0.05 & 0.15 \\ -0.25 & 0.25 & -0.25 & 0.25 \end{bmatrix}.$$

Hence

$$\mathbf{q}_1 = [-0.05 \quad -0.15 \quad 0.35 \quad 0.05]$$

and

$$\mathbf{q}_2 = [-0.25 \quad 0.25 \quad -0.25 \quad 0.25].$$

The transformation matrix is

$$\mathbf{T} = \begin{bmatrix} \mathbf{q}_1 \\ \mathbf{q}_1 \mathbf{A} \\ \mathbf{q}_2 \\ \mathbf{q}_2 \mathbf{A} \end{bmatrix} = \begin{bmatrix} -0.05 & -0.15 & 0.35 & 0.05 \\ -0.05 & -0.15 & 1.35 & 0.05 \\ -0.25 & 0.25 & -0.25 & 0.25 \\ 0.75 & 0.25 & -1.25 & 0.25 \end{bmatrix}.$$

The matrices \mathbf{A} and \mathbf{B} in the new coordinates have the form

$$\mathbf{T} \mathbf{A} \mathbf{T}^{-1} = \left[\begin{array}{cc|cc} 0 & 1 & 0 & 0 \\ -4 & 5 & 0 & 0 \\ \hline 0 & 0 & 0 & 1 \\ 3 & -3 & 4 & 1 \end{array} \right] \quad \text{and} \quad \mathbf{T} \mathbf{B} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

3.3.2 Observer Form

In this subsection, we use duality and the algorithms for reducing reachable systems into controller forms to derive analogous results for observable systems. We consider an observable dynamical system model

$$\begin{aligned} \dot{\mathbf{x}}(t) &= \mathbf{A}\mathbf{x}(t), \\ \mathbf{y}(t) &= \mathbf{C}\mathbf{x}(t), \end{aligned} \tag{3.17}$$

or, equivalently, the observable pair (\mathbf{A}, \mathbf{C}) , where $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{C} \in \mathbb{R}^{p \times n}$, $p \leq n$, and $\text{rank } \mathbf{C} = p$. We omitted the input variable \mathbf{u} for convenience because it plays no role in the subsequent discussion. We next consider the *dual system* model associated with the system (3.17),

$$\dot{\mathbf{z}}(t) = \mathbf{A}^T \mathbf{z}(t) + \mathbf{C}^T \mathbf{v}(t).$$

Because the pair (\mathbf{A}, \mathbf{C}) is observable, the dual pair $(\mathbf{A}^T, \mathbf{C}^T)$ is reachable. Therefore, the pair $(\mathbf{A}^T, \mathbf{C}^T)$ can be transformed into the controller form using the algorithm from the previous subsection. We then take the dual of the result to get

$$\begin{aligned} \dot{\hat{\mathbf{x}}}(t) &= \hat{\mathbf{A}}\hat{\mathbf{x}}, \\ \mathbf{y}(t) &= \hat{\mathbf{C}}\hat{\mathbf{x}}, \end{aligned}$$

where

$$\hat{\mathbf{A}} = \begin{bmatrix} \hat{\mathbf{A}}_{11} & \hat{\mathbf{A}}_{12} & \cdots & \hat{\mathbf{A}}_{1p} \\ \hat{\mathbf{A}}_{21} & \hat{\mathbf{A}}_{22} & \cdots & \hat{\mathbf{A}}_{2p} \\ \vdots & & & \vdots \\ \hat{\mathbf{A}}_{p1} & \hat{\mathbf{A}}_{p2} & \cdots & \hat{\mathbf{A}}_{pp} \end{bmatrix}.$$

Each diagonal submatrix $\hat{\mathbf{A}}_{ii} \in \mathbb{R}^{\hat{d}_i \times \hat{d}_i}$ of $\hat{\mathbf{A}}$ has the form

$$\hat{\mathbf{A}}_{ii} = \begin{bmatrix} 0 & 0 & \cdots & 0 & x \\ 1 & 0 & \cdots & 0 & x \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & \cdots & 1 & x \end{bmatrix},$$

which is just the transpose of the controller form for the single-input case. The integers \hat{d}_i , $i = 1, 2, \dots, p$ are the controllability indices of the controllable pair $(\mathbf{A}^T, \mathbf{C}^T)$. The off-diagonal blocks, $j \neq i$, have the form

$$\hat{\mathbf{A}}_{ji} = \begin{bmatrix} 0 & 0 & \cdots & 0 & x \\ 0 & 0 & \cdots & 0 & x \\ \vdots & & & & \vdots \\ 0 & 0 & \cdots & 0 & x \end{bmatrix};$$

that is, they consist of zeros except for their last columns. The matrix $\hat{\mathbf{C}}$ consists of p blocks and has the form

$$\hat{\mathbf{C}} = \begin{bmatrix} 0 & 0 & \cdots & 1 & \left| \begin{array}{c} 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 1 \\ \vdots & & & \vdots \end{array} \right| \cdots & \left| \begin{array}{c} 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \\ \vdots & & & \vdots \end{array} \right| \\ 0 & 0 & \cdots & x & \left| \begin{array}{c} 0 & 0 & \cdots & 1 \\ 0 & 0 & \cdots & x \\ \vdots & & & \vdots \end{array} \right| \cdots & \left| \begin{array}{c} 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \\ \vdots & & & \vdots \end{array} \right| \\ \vdots & & & \vdots & \left| \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right| & \cdots & \left| \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right| \\ 0 & 0 & \cdots & x & \left| \begin{array}{c} 0 & 0 & \cdots & x \\ 0 & 0 & \cdots & x \\ \vdots & & & \vdots \end{array} \right| \cdots & \left| \begin{array}{c} 0 & 0 & \cdots & 1 \\ 0 & 0 & \cdots & 1 \\ \vdots & & & \vdots \end{array} \right| \end{bmatrix}.$$

The p controllability indices of the pair $(\mathbf{A}^T, \mathbf{C}^T)$ are called the *observability indices* of the pair (\mathbf{A}, \mathbf{C}) . Note that

$$\sum_{i=1}^p \hat{d}_i = n.$$

Similarly as in the controller form, we define the *observability index* of the pair (\mathbf{A}, \mathbf{C}) to be the integer

$$\hat{d} = \max\{\hat{d}_i : i = 1, 2, \dots, p\}.$$

We now illustrate the above algorithm for transforming an observable pair into the observer form with a numerical example.

◆ Example 3.8

Given the observable pair (\mathbf{A}, \mathbf{C}) ,

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 2 & 0 \\ 0 & 1 & 3 & 0 \\ 0 & 0 & -21 & 5 \end{bmatrix}, \quad \mathbf{C} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

we transform the pair $(\mathbf{A}^T, \mathbf{C}^T)$ into the controller form. Thus, we first form the controllability matrix for $(\mathbf{A}^T, \mathbf{C}^T)$; then proceeding from left to right, we select the first $n = 4$ linearly independent columns. They are

$$\mathbf{c}_1^T, \mathbf{c}_2^T, \mathbf{A}^T \mathbf{c}_1^T, (\mathbf{A}^T)^2 \mathbf{c}_1^T.$$

Therefore, the observability indices are: $\hat{d}_1 = 3$ and $\hat{d}_2 = 1$. Hence,

$$\mathbf{L} = [\mathbf{c}_1^T \quad \mathbf{A}^T \mathbf{c}_1^T \quad (\mathbf{A}^T)^2 \mathbf{c}_1^T \quad \mathbf{c}_2^T] = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 3 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

and

$$\mathbf{L}^{-1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & -3 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

The two vectors that we need to form the desired transformation are

$$\mathbf{q}_1 = [0 \quad 1 \quad 0 \quad 0] \quad \text{and} \quad \mathbf{q}_2 = [0 \quad 0 \quad 0 \quad 1].$$

The transformation matrix is

$$\hat{\mathbf{T}} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & -0 & 1 & 0 \\ 1 & 2 & 3 & -21 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

The observer form of the pair (\mathbf{A}, \mathbf{C}) is given by the matrices

$$\hat{\mathbf{A}} = \hat{\mathbf{T}}^{-T} \mathbf{A} \hat{\mathbf{T}}^T = \left[\begin{array}{ccc|c} 0 & 0 & 1 & 0 \\ 1 & 0 & 2 & 0 \\ 0 & 1 & 3 & 0 \\ \hline 0 & 0 & -105 & 5 \end{array} \right] \quad \text{and} \quad \hat{\mathbf{C}} = \mathbf{C} \hat{\mathbf{T}}^T = \left[\begin{array}{ccc|c} 0 & 0 & 1 & 0 \\ 0 & 0 & -21 & 1 \end{array} \right].$$

3.4 Linear State-Feedback Control

In the following sections we discuss methods for a feedback controller design for linear systems. The design process involves three steps. The first step is the design of a state-feedback control law. In this step we assume that we have all states available to us for feedback purposes. In general, this may not be the case. The assumption that all the states are available allows us to proceed with the state-feedback control law design. We then proceed with the second step; this involves the design of an estimator, also referred to as an observer of the state vector. The estimator should use only the input and output of the plant to be controlled to generate an estimate of the plant state. The last step is to combine the first two steps in that the control law, designed in the first step, uses the state estimate instead of the true state vector. The result of this step is a combined controller–estimator compensator. We first discuss the problem of designing a linear state-feedback control law. In the following sections we present methods for constructing estimators and combined controller–estimator compensators.

The linear state-feedback control law, for a system modeled by $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$, is the feedback of a linear combination of all the state variables and has the form

$$\mathbf{u} = -\mathbf{K}\mathbf{x},$$

where $\mathbf{K} \in \mathbb{R}^{m \times n}$ is a constant matrix. The closed-loop system then is

$$\dot{\mathbf{x}} = (\mathbf{A} - \mathbf{B}\mathbf{K})\mathbf{x}.$$

The poles of the closed-loop system are the roots of the characteristic equation

$$\det(s\mathbf{I}_n - \mathbf{A} + \mathbf{B}\mathbf{K}) = 0.$$

The linear state-feedback control law design consists of selecting the gains

$$k_{ij}, \quad i = 1, 2, \dots, m, \quad j = 1, 2, \dots, n,$$

so that the roots of the closed-loop characteristic equation

$$\det(s\mathbf{I}_n - \mathbf{A} + \mathbf{B}\mathbf{K}) = 0$$

are in desirable locations in the complex plane. We assume that a designer has made the selection of the desired poles of the closed-loop system, and they are

$$s_1, s_2, \dots, s_n.$$

The desired closed-loop poles can be real or complex. If they are complex, then they must come in complex conjugate pairs. This is because we use only real gains k_{ij} . Having selected the desired closed-loop poles, we can form the desired closed-loop characteristic polynomial,

$$\begin{aligned} \alpha_c(s) &= (s - s_1)(s - s_2) \cdots (s - s_n) \\ &= s^n + \alpha_{n-1}s^{n-1} + \cdots + \alpha_1s + \alpha_0. \end{aligned}$$

Our goal is to select a feedback matrix \mathbf{K} such that

$$\det(s\mathbf{I}_n - \mathbf{A} + \mathbf{B}\mathbf{K}) = s^n + \alpha_{n-1}s^{n-1} + \cdots + \alpha_1s + \alpha_0.$$

The above problem is also referred to as the *pole placement* problem. We first discuss the pole placement problem for the single-input plants. In this case, $\mathbf{K} = \mathbf{k} \in \mathbb{R}^{1 \times n}$. The solution to the

problem is easily obtained if the pair (A, b) is already in the controller companion form. In such a case we have

$$A - bk = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & & & & & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ -a_0 - k_1 & -a_1 - k_2 & -a_2 - k_3 & \cdots & -a_{n-2} - k_{n-1} & -a_{n-1} - k_n \end{bmatrix}.$$

Hence, the desired gains are

$$\begin{aligned} k_1 &= \alpha_0 - a_0, \\ k_2 &= \alpha_1 - a_1, \\ &\vdots \\ k_n &= \alpha_{n-1} - a_{n-1}. \end{aligned}$$

If the pair (A, b) is not in the controller companion form, we first transform it into the companion form and then compute the gain vector \tilde{k} such that

$$\det(sI_n - \tilde{A} + \tilde{b}\tilde{k}) = s^n + \alpha_{n-1}s^{n-1} + \cdots + \alpha_1s + \alpha_0.$$

Thus,

$$\tilde{k} = [\alpha_0 - a_0 \quad \alpha_1 - a_1 \quad \cdots \quad \alpha_{n-1} - a_{n-1}].$$

Then,

$$k = \tilde{k}T,$$

where T is the transformation that brings the pair (A, b) into the controller companion form.

We can represent the above formula for the gain matrix in an alternative way. For this, note that

$$\begin{aligned} \tilde{k}T &= [\alpha_0 - a_0 \quad \alpha_1 - a_1 \quad \cdots \quad \alpha_{n-1} - a_{n-1}] \begin{bmatrix} q_1 \\ q_1 A \\ \vdots \\ q_1 A^{n-1} \end{bmatrix} \\ &= q_1(\alpha_0 I_n + \alpha_1 A + \cdots + \alpha_{n-1} A^{n-1}) - q_1(a_0 I_n + a_1 A + \cdots + a_{n-1} A^{n-1}). \end{aligned}$$

By the Cayley–Hamilton theorem, we have

$$A^n = -(a_0 I_n + a_1 A + \cdots + a_{n-1} A^{n-1}).$$

Hence,

$$k = q_1 \alpha_c(A).$$

The above expression for the gain row vector was proposed by Ackermann in 1972 and is now referred to as *Ackermann's formula* for pole placement.

◆ Example 3.9

For the linear dynamical system

$$\dot{\mathbf{x}} = \begin{bmatrix} 1 & -1 \\ 1 & -2 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 2 \\ 1 \end{bmatrix} u,$$

we will use Ackermann's formula to design a state-feedback controller, $u = -\mathbf{k}\mathbf{x}$, such that the closed-loop poles are located at $\{-1, -2\}$.

To use Ackermann's formula, we first form the controllability matrix of the pair (\mathbf{A}, \mathbf{b}) and then find the last row of its inverse, denoted \mathbf{q}_1 . The controllability matrix is

$$[\mathbf{b} \quad \mathbf{A}\mathbf{b}] = \begin{bmatrix} 2 & 1 \\ 1 & 0 \end{bmatrix}.$$

The inverse of the above matrix is

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

Hence, $\mathbf{q}_1 = [1 \quad -2]$. The desired closed-loop characteristic polynomial is

$$\alpha_c(s) = (s + 1)(s + 2) = s^2 + 3s + 2.$$

Therefore,

$$\begin{aligned} \mathbf{k} &= \mathbf{q}_1 \alpha_c(\mathbf{A}) \\ &= \mathbf{q}_1 (\mathbf{A}^2 + 3\mathbf{A} + 2\mathbf{I}_2) \\ &= \mathbf{q}_1 \left(\begin{bmatrix} 1 & -1 \\ 1 & -2 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 1 & -2 \end{bmatrix} + 3 \begin{bmatrix} 1 & -1 \\ 1 & -2 \end{bmatrix} + 2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) \\ &= \mathbf{q}_1 \left(\begin{bmatrix} 0 & 1 \\ -1 & 3 \end{bmatrix} + \begin{bmatrix} 3 & -3 \\ 3 & -6 \end{bmatrix} + \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \right) \\ &= [1 \quad -2] \begin{bmatrix} 5 & -2 \\ 2 & -1 \end{bmatrix} \\ &= [1 \quad 0] \end{aligned}$$

We now present a pole placement algorithm for multi-input systems. If the pair (\mathbf{A}, \mathbf{B}) is already in the controller companion form, then we may proceed as follows. First, represent the

matrix \mathbf{B} as

$$\mathbf{B} = \begin{bmatrix} 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ 1 & x & x & x & x \\ \hline 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ 0 & 1 & x & x & x \\ \hline \vdots & & & & \vdots \\ 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ 1 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ 0 & 1 & 0 & 0 & 0 \\ \hline \vdots & & & & \vdots \\ 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & x & \cdots & x & x \\ 0 & 1 & \cdots & x & x \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & \cdots & 1 & x \\ 0 & 0 & \cdots & 0 & 1 \end{bmatrix} = \hat{\mathbf{B}}\mathbf{\Gamma},$$

where the square and nonsingular matrix $\mathbf{\Gamma}$ consists of the nonzero rows of \mathbf{B} . Next, let

$$\hat{\mathbf{K}} = \mathbf{\Gamma}\mathbf{K}.$$

Notice that

$$\hat{\mathbf{B}}\hat{\mathbf{K}} = \begin{bmatrix} 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ \hat{k}_{11} & \hat{k}_{12} & \cdots & \hat{k}_{1n-1} & \hat{k}_{1n} \\ \hline 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ \hat{k}_{21} & \hat{k}_{22} & \cdots & \hat{k}_{2n-1} & \hat{k}_{2n} \\ \hline \vdots & & & & \vdots \\ 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & & & & \vdots \\ \hat{k}_{m1} & \hat{k}_{m2} & \cdots & \hat{k}_{mn-1} & \hat{k}_{mn} \end{bmatrix};$$

that is, the nonzero rows in the product $\hat{\mathbf{B}}\hat{\mathbf{K}}$ match the nonzero rows of the matrix \mathbf{A} in the controller companion form. We can then, for example, select the gains k_{ij} , $i = 1, 2, \dots, m$,

$j = 1, 2, \dots, n$, so that

$$\mathbf{A} - \mathbf{B}\mathbf{K} = \mathbf{A} - \hat{\mathbf{B}}\hat{\mathbf{K}} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & & & & & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ -\alpha_0 & -\alpha_1 & -\alpha_2 & \cdots & -\alpha_{n-2} & -\alpha_{n-1} \end{bmatrix},$$

where

$$\mathbf{K} = \mathbf{\Gamma}^{-1}\hat{\mathbf{K}}.$$

If the pair (\mathbf{A}, \mathbf{B}) is not in the controller companion form, we first transform it into this form and then compute the gain matrix that allocates the closed-loop poles into the desired locations for the pair (\mathbf{A}, \mathbf{B}) in the controller form. The gain matrix \mathbf{K} that allocates the closed-loop poles into the prespecified locations for the pair (\mathbf{A}, \mathbf{B}) in the original coordinates is then given by

$$\mathbf{K} = \mathbf{\Gamma}^{-1}\hat{\mathbf{K}}\mathbf{T},$$

where \mathbf{T} is the transformation matrix that brings the pair (\mathbf{A}, \mathbf{B}) into the controller companion form.

◆ Example 3.10

For the pair

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 2 & 0 \\ 0 & 1 & 3 & 0 \\ 0 & 0 & -21 & 5 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix},$$

we will use its controller companion form to find the matrix $\mathbf{K} \in \mathbb{R}^{2 \times 4}$ so that closed-loop poles are located at

$$-2, -3 + j, -3 - j, -4.$$

We first transform the pair (\mathbf{A}, \mathbf{B}) into its controller companion form. We begin by forming the controllability matrix

$$[\mathbf{b}_1 \quad \mathbf{b}_2 \quad \mathbf{A}\mathbf{b}_1 \quad \mathbf{A}\mathbf{b}_2 \quad \mathbf{A}^2\mathbf{b}_1 \quad \cdots] = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & \cdots \\ 0 & 0 & 1 & 0 & 0 & \cdots \\ 0 & 0 & 0 & 0 & 1 & \cdots \\ 0 & 1 & 0 & 5 & 0 & \cdots \end{bmatrix}.$$

We then select, proceeding from left to right, the first four linearly independent columns from the above controllability matrix. We obtain

$$[\mathbf{b}_1 \quad \mathbf{b}_2 \quad \mathbf{A}\mathbf{b}_1 \quad \mathbf{A}^2\mathbf{b}_1].$$

Thus, the controllability indices are $d_1 = 3$ and $d_2 = 1$. We rearrange the above columns and form the matrix \mathbf{L} of the form

$$\mathbf{L} = [\mathbf{b}_1 \quad \mathbf{A}\mathbf{b}_1 \quad \mathbf{A}^2\mathbf{b}_1 \quad \mathbf{b}_2] = \mathbf{I}_4 = \mathbf{L}^{-1}.$$

The two row vectors that we need in constructing the transformation are

$$\mathbf{q}_1 = [0 \quad 0 \quad 1 \quad 0] \quad \text{and} \quad \mathbf{q}_2 = [0 \quad 0 \quad 0 \quad 1].$$

The transformation matrix is

$$\mathbf{T} = \begin{bmatrix} \mathbf{q}_1 \\ \mathbf{q}_1 \mathbf{A} \\ \mathbf{q}_1 \mathbf{A}^2 \\ \mathbf{q}_2 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 3 & 0 \\ 1 & 3 & 11 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

The pair (\mathbf{A}, \mathbf{B}) in the new coordinates has the form

$$\tilde{\mathbf{A}} = \mathbf{T} \mathbf{A} \mathbf{T}^{-1} = \left[\begin{array}{ccc|c} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 2 & 3 & 0 \\ \hline -21 & 0 & 0 & 5 \end{array} \right] \quad \text{and} \quad \tilde{\mathbf{B}} = \mathbf{T} \mathbf{B} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

The desired closed-loop characteristic polynomial is

$$\alpha_c(s) = (s+2)(s+3-j)(s+3+j)(s+4) = s^4 + 12s^3 + 54s^2 + 108s + 80.$$

A possible choice of the gain matrix $\tilde{\mathbf{K}}$, from among many, that does the job is $\tilde{\mathbf{K}}$ such that

$$\tilde{\mathbf{A}} - \tilde{\mathbf{B}} \tilde{\mathbf{K}} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -80 & -108 & -54 & -12 \end{bmatrix}.$$

Note that $\mathbf{\Gamma} = \mathbf{I}_2$. Therefore,

$$\tilde{\mathbf{K}} = \begin{bmatrix} 1 & 2 & 3 & -1 \\ 59 & 108 & 54 & 17 \end{bmatrix},$$

and thus

$$\mathbf{K} = \tilde{\mathbf{K}} \mathbf{T} = \begin{bmatrix} 3 & 11 & 40 & -1 \\ 54 & 270 & 977 & 17 \end{bmatrix}.$$

The above algorithm for pole placement for multi-input systems is more of theoretical value rather than practical. This algorithm is not well-suited for numerical implementations because the transformation of a given pair (\mathbf{A}, \mathbf{B}) into the controller companion form suffers from poor numerical properties. For a discussion of pole placement numerically stable algorithms for multi-input systems, we refer the reader to Jamshidi, Tarokh, and Shafai [141].

Unlike in the single-input case, a solution to the pole placement problem for a multi-input system is not unique. Therefore, one can use the remaining degrees of freedom to achieve secondary goals. In Subsection 5.3.3 we discuss a method for constructing a linear state-feedback control law that allocates the closed-loop poles into prespecified locations and at the same time minimizes a quadratic performance index.

Using the results preceding the above discussion, we can state and prove a fundamental theorem of linear systems.

Theorem 3.7 The pole placement problem is solvable for all choices of n desired closed-loop poles, symmetric with respect to the real axis, if and only if the given pair (\mathbf{A}, \mathbf{B}) is reachable.

Proof The proof of necessity (\Leftarrow) follows immediately from the discussion preceding Example 3.10. Indeed, the pair (\mathbf{A}, \mathbf{B}) is reachable, if and only if it can be transformed into the controller canonical form. Once the transformation is performed, we can solve the pole placement problem for the given set of desired closed-loop poles, symmetric with respect to the real axis. We then transform the pair (\mathbf{A}, \mathbf{B}) and the gain matrix back into the original coordinates. Because the similarity transformation affects neither controllability nor the characteristic polynomial, the eigenvalues of $\mathbf{A} - \mathbf{B}\mathbf{K}$ are precisely the desired prespecified closed-loop poles. This completes the proof of the necessity part.

We use a proof by contraposition to prove the sufficiency part (\Rightarrow). Assume that the pair (\mathbf{A}, \mathbf{B}) is nonreachable. By Theorem 3.4, there is a similarity transformation $\mathbf{z} = \mathbf{T}\mathbf{x}$ such that the pair (\mathbf{A}, \mathbf{B}) in the new coordinates has the form

$$\tilde{\mathbf{A}} = \mathbf{T}\mathbf{A}\mathbf{T}^{-1} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{O} & \mathbf{A}_4 \end{bmatrix}, \quad \tilde{\mathbf{B}} = \mathbf{T}\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{O} \end{bmatrix},$$

where the pair $(\mathbf{A}_1, \mathbf{B}_1)$ is reachable, $\mathbf{A}_1 \in \mathbb{R}^{r \times r}$, and $\mathbf{B}_1 \in \mathbb{R}^{r \times m}$. Let

$$\mathbf{u} = -\tilde{\mathbf{K}}\mathbf{z} = -[\mathbf{K}_1 \quad \mathbf{K}_2]\mathbf{z},$$

where $\mathbf{K}_1 \in \mathbb{R}^{m \times r}$ and $\mathbf{K}_2 \in \mathbb{R}^{m \times (n-r)}$. Then,

$$\tilde{\mathbf{A}} - \tilde{\mathbf{B}}\tilde{\mathbf{K}} = \begin{bmatrix} \mathbf{A}_1 - \mathbf{B}_1\mathbf{K}_1 & \mathbf{A}_2 - \mathbf{B}_1\mathbf{K}_2 \\ \mathbf{O} & \mathbf{A}_4 \end{bmatrix}. \quad (3.18)$$

It follows from (3.18) that the nonreachable portion of the system is not affected by the state feedback; that is, the eigenvalues of \mathbf{A}_4 cannot be allocated. Hence, the pole placement problem cannot be solved if the pair $(\tilde{\mathbf{A}}, \tilde{\mathbf{B}})$, or equivalently the pair (\mathbf{A}, \mathbf{B}) , is nonreachable, which completes the proof of the sufficiency part.

For an alternative proof of this important result, the reader may consult Wonham [304].

3.5 State Estimators

In the previous section, we presented algorithms for solving the pole placement problem using the linear state-feedback control law. To implement such a control law, one needs availability

of all the state variables. Often, this requirement is not met, either because measuring all the state variables would require excessive number of sensors or because the state variables are not accessible for direct measurement. Instead, only a subset of state variables or their combination may be available. To implement the control law constructed in the previous section, we will use estimates of state variables rather than the true states. We now discuss the problem of constructing state estimators. Much of the literature refers to state estimators as “observers.” However, we agree with Franklin and Powell [89, p. 139] that “estimator” is much more descriptive in its function because “observer” implies a direct measurement. We present two types of state estimators. First, full-order state estimators will be discussed and then a reduced-order state estimator will be introduced.

3.5.1 Full-Order Estimator

We now discuss the problem of constructing a full-order state estimator for dynamical systems modeled by

$$\begin{aligned}\dot{\mathbf{x}}(t) &= \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t), \\ \mathbf{y}(t) &= \mathbf{C}\mathbf{x}(t),\end{aligned}$$

where $\mathbf{C} \in \mathbb{R}^{p \times n}$, $p \leq n$. We assume that the pair (\mathbf{A}, \mathbf{C}) is observable. Our goal is to construct a dynamical system that will estimate the state vector \mathbf{x} based on the plant input \mathbf{u} and the plant output \mathbf{y} . One could consider constructing a model of the plant dynamics and connecting the resulting dynamical system, referred to as an *open-loop estimator*, as in Figure 3.4. This open-loop estimator is described by

$$\dot{\tilde{\mathbf{x}}}(t) = \mathbf{A}\tilde{\mathbf{x}}(t) + \mathbf{B}\mathbf{u}(t),$$

where $\tilde{\mathbf{x}}(t)$ is the estimate of $\mathbf{x}(t)$. Let

$$\mathbf{e}(t) = \tilde{\mathbf{x}}(t) - \mathbf{x}(t)$$

be the estimation error. Then, the dynamics of the estimation error are described by

$$\dot{\mathbf{e}}(t) = \dot{\tilde{\mathbf{x}}}(t) - \dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{e}(t),$$

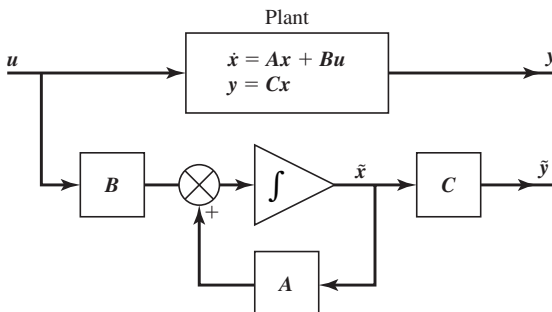


Figure 3.4 Open-loop estimator.

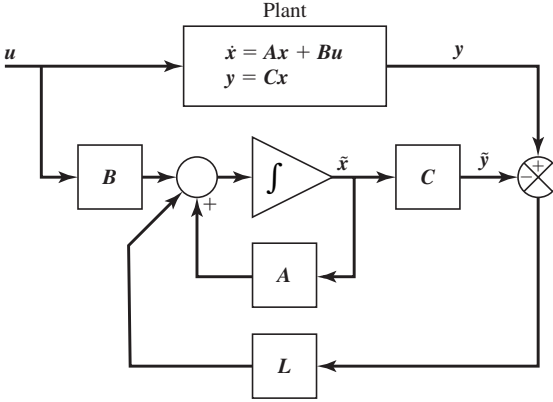


Figure 3.5 Closed-loop estimator.

with the initial estimation error

$$e(0) = \tilde{x}(0) - x(0).$$

If the eigenvalues of the matrix A are in the open left-hand plane, then the error converges to zero. However, we have no control over the convergence rate. Furthermore, the matrix A does not have to have all its eigenvalues in the open left-hand plane. Thus, the open-loop estimator is impractical. We modify this estimator by adding a feedback term to it. The resulting structure, depicted in Figure 3.5, is called the *closed-loop estimator*, the *Luenberger observer*, or the asymptotic full-order estimator. The dynamics of the closed-loop estimator are described by

$$\dot{\tilde{x}}(t) = A\tilde{x}(t) + Bu(t) + L(y(t) - C\tilde{x}(t)),$$

and the dynamics of the estimation error are governed by

$$\begin{aligned} \dot{e}(t) &= \dot{\tilde{x}}(t) - \dot{x}(t) \\ &= (A - LC)e(t), \quad e(0) = \tilde{x}(0) - x(0). \end{aligned}$$

The pair (A, C) is observable if and only if the dual pair (A^T, C^T) is reachable. By assumption, the pair (A, C) is observable, and therefore the pair (A^T, C^T) is reachable. By Theorem 3.7, we can solve the pole placement problem for the dual pair (A^T, C^T) ; that is, for any set of prespecified n complex numbers, symmetric with respect to the real axis, there is a matrix, call it L^T , such that the eigenvalues of $A^T - C^T L^T$ and hence of $A - LC$ are in the prespecified locations. It follows from the above that if the pair (A, C) is observable, then in addition to forcing the estimation error to converge to zero, we can also control its rate of convergence by appropriately selecting the eigenvalues of the matrix $A - LC$. We also note that the selection of the closed-loop estimator gain matrix L can be approached in exactly the same fashion as the construction of the gain matrix K in the linear state-feedback control law design. We illustrate the above point with a numerical example.

◆ Example 3.11

For the given observable pair

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 2 & 0 \\ 0 & 1 & 3 & 0 \\ 0 & 0 & -21 & 5 \end{bmatrix}, \quad \mathbf{C} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

construct the matrix $\mathbf{L} \in \mathbb{R}^{4 \times 2}$ so that the eigenvalues of $\mathbf{A} - \mathbf{L}\mathbf{C}$ are located at

$$\{-2, -3 + j, -3 - j, -4\}.$$

In our construction of the gain matrix \mathbf{L} we use the observer companion form of the pair (\mathbf{A}, \mathbf{C}) that we computed in Example 3.8. Our goal then is to construct a matrix \mathbf{L} so that the characteristic polynomial of the matrix $\mathbf{A} - \mathbf{L}\mathbf{C}$ is

$$\det(s\mathbf{I}_4 - \mathbf{A} + \mathbf{L}\mathbf{C}) = s^4 + 12s^3 + 54s^2 + 108s + 80.$$

Using the observer form of the pair (\mathbf{A}, \mathbf{C}) , we select $\hat{\mathbf{L}}$ so that

$$\hat{\mathbf{A}} - \hat{\mathbf{L}}\hat{\mathbf{C}} = \begin{bmatrix} 0 & 0 & 0 & -80 \\ 1 & 0 & 0 & -108 \\ 0 & 1 & 0 & -54 \\ 0 & 0 & 1 & -12 \end{bmatrix}.$$

We have

$$\hat{\mathbf{L}} = \begin{bmatrix} 1681 & 80 \\ 2270 & 108 \\ 1137 & 54 \\ 251 & 17 \end{bmatrix},$$

and hence

$$\mathbf{L} = \hat{\mathbf{T}}^T \hat{\mathbf{L}} = \begin{bmatrix} 1137 & 54 \\ 3955 & 188 \\ 5681 & 270 \\ -23,626 & -1117 \end{bmatrix}.$$

We note that there are other possible gain matrices \mathbf{L} that could be used to allocate the eigenvalues of $\mathbf{A} - \mathbf{L}\mathbf{C}$ into desired locations.

We now present an alternative approach to the closed-loop estimator design. This approach will then be used to construct lower-order estimators that we discuss in the subsequent section. We assume that the dynamics of an estimator are given by

$$\dot{\mathbf{z}}(t) = \mathbf{D}\mathbf{z}(t) + \mathbf{E}u(t) + \mathbf{G}y(t), \quad (3.19)$$

where $\mathbf{z}(t) \in \mathbb{R}^n$. To begin with the development, suppose that we would like to obtain

$$\mathbf{z} = \mathbf{T}\mathbf{x},$$

where $\mathbf{T} \in \mathbb{R}^{n \times n}$ is invertible. Premultiplying $\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$ by \mathbf{T} yields

$$\mathbf{T}\dot{\mathbf{x}}(t) = \mathbf{T}\mathbf{A}\mathbf{x}(t) + \mathbf{T}\mathbf{B}\mathbf{u}(t). \quad (3.20)$$

Substituting into (3.19) the relation $\mathbf{y}(t) = \mathbf{C}\mathbf{x}(t)$ gives

$$\dot{\mathbf{z}}(t) = \mathbf{D}\mathbf{z}(t) + \mathbf{E}\mathbf{u}(t) + \mathbf{G}\mathbf{C}\mathbf{x}(t). \quad (3.21)$$

Subtracting (3.21) from (3.20) and taking into account that $\mathbf{z} = \mathbf{T}\mathbf{x}$ and $\dot{\mathbf{z}} = \mathbf{T}\dot{\mathbf{x}}$ yields

$$\mathbf{0} = (\mathbf{T}\mathbf{A} - \mathbf{D}\mathbf{T} - \mathbf{G}\mathbf{C})\mathbf{x}(t) + (\mathbf{T}\mathbf{B} - \mathbf{E})\mathbf{u}(t)$$

for all t and for arbitrary input \mathbf{u} . Hence, we must have

$$\begin{aligned} \mathbf{T}\mathbf{A} - \mathbf{D}\mathbf{T} &= \mathbf{G}\mathbf{C}, \\ \mathbf{E} &= \mathbf{T}\mathbf{B}. \end{aligned}$$

If $\mathbf{z} \neq \mathbf{T}\mathbf{x}$ but the two above equations hold, then subtracting (3.20) from (3.21) and taking into account the above relations yields

$$\begin{aligned} \frac{d}{dt}(\mathbf{z} - \mathbf{T}\mathbf{x}) &= \mathbf{D}\mathbf{z} - \mathbf{T}\mathbf{A}\mathbf{x} + \mathbf{G}\mathbf{C}\mathbf{x} \\ &= \mathbf{D}\mathbf{z} - \mathbf{D}\mathbf{T}\mathbf{x} \\ &= \mathbf{D}(\mathbf{z} - \mathbf{T}\mathbf{x}). \end{aligned}$$

If the eigenvalues of \mathbf{D} all have negative real parts, then

$$\mathbf{z} - \mathbf{T}\mathbf{x} \rightarrow \mathbf{0} \quad \text{as } t \rightarrow \infty.$$

Note that if we try $\mathbf{T} = \mathbf{I}_n$, then we have the previously analyzed case, where $\mathbf{G} = \mathbf{L}$, and

$$\begin{aligned} \mathbf{E} &= \mathbf{B}, \\ \mathbf{A} - \mathbf{L}\mathbf{C} &= \mathbf{D}. \end{aligned}$$

The resulting estimator, when $\mathbf{T} = \mathbf{I}_n$, was called the identity observer by its developer, Luenberger [188, 189]. Given the desired eigenvalues, symmetric with respect to the real axis, the equation $\mathbf{A} - \mathbf{L}\mathbf{C} = \mathbf{D}$ can be solved for \mathbf{L} so that the resulting \mathbf{D} has its eigenvalues in the desired locations, if and only if the pair (\mathbf{A}, \mathbf{C}) is observable. However, it is not true that given an arbitrary \mathbf{D} we can solve for \mathbf{L} so that $\mathbf{A} - \mathbf{L}\mathbf{C} = \mathbf{D}$ even if the pair (\mathbf{A}, \mathbf{C}) is observable. We will use the results of the above discussion to construct estimators of reduced order in the following section.

3.5.2 Reduced-Order Estimator

If the output \mathbf{y} is a p vector with $p < n$, we would like to use the p outputs to determine p of the x_i 's and design an estimator of order $n - p$ to estimate the rest of the state variables. Specifically, if $\text{rank } \mathbf{C} = p$, then we can use the equation $\mathbf{y} = \mathbf{C}\mathbf{x}$ to solve for p of the x_i 's in terms of the y_i 's and the remaining $n - p$ state variables x_k 's. Then, we construct an estimator

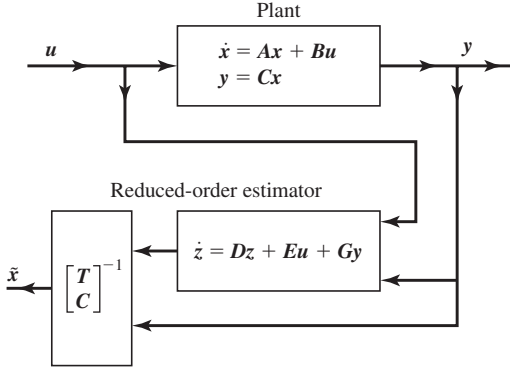


Figure 3.6 Reduced-order estimator.

of order $n - p$. To recover all states, we should like to solve for x the equations

$$\begin{aligned} z &= Tx, \\ y &= Cx, \end{aligned}$$

where now $T \in \mathbb{R}^{(n-p) \times n}$. This can be done if and only if

$$\text{rank} \begin{bmatrix} T \\ C \end{bmatrix} = n. \quad (3.22)$$

Given $D \in \mathbb{R}^{(n-p) \times (n-p)}$, $G \in \mathbb{R}^{(n-p) \times p}$, and the triple (A, B, C) , we can use the Kronecker product to show that the equation

$$TA - DT = GC$$

can be uniquely solved for T if and only if no eigenvalue of D is also an eigenvalue of A . Indeed, the above equation is a special case of the matrix equation (A.54) that we analyze in Section A.5. Note that when choosing the matrices D and G , we must make sure that the condition (3.22) is satisfied. The structure of the resulting reduced-order estimator is shown in Figure 3.6. The last step in the control design process is to combine the control law that we constructed in the previous section with an estimator.

3.6 Combined Controller–Estimator Compensator

We start by combining a control law with a full-order estimator. We write the equations that govern the behavior of the dynamical system consisting of the plant model and the full-order estimator. The resulting system's order is $2n$, and the equations that describe the system are

$$\left. \begin{aligned} \begin{bmatrix} \dot{x}(t) \\ \dot{\tilde{x}}(t) \end{bmatrix} &= \begin{bmatrix} A & O \\ LC & A - LC \end{bmatrix} \begin{bmatrix} x(t) \\ \tilde{x}(t) \end{bmatrix} + \begin{bmatrix} B \\ B \end{bmatrix} u(t) \\ y(t) &= [C \quad O] \begin{bmatrix} x(t) \\ \tilde{x}(t) \end{bmatrix} \end{aligned} \right\} \quad (3.23)$$

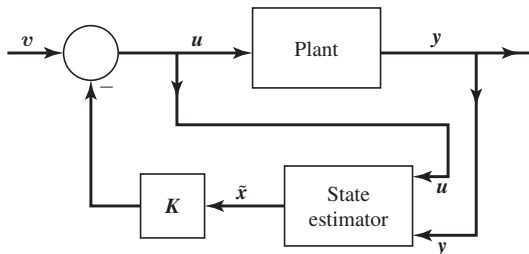


Figure 3.7 Closed-loop system driven by the combined controller–estimator compensator.

Suppose that we now apply the control law

$$u(t) = -K\tilde{x}(t) + v(t) \quad (3.24)$$

instead of the actual state-feedback control law. In the above, the vector v denotes an external input signal. A schematic of a closed-loop system driven by the combined controller-estimator compensator is shown in Figure 3.7. A model of the closed-loop system is obtained by combining equations (3.23) and (3.24):

$$\begin{aligned} \begin{bmatrix} \dot{x}(t) \\ \dot{\tilde{x}}(t) \end{bmatrix} &= \begin{bmatrix} A & -BK \\ LC & A - LC - BK \end{bmatrix} \begin{bmatrix} x(t) \\ \tilde{x}(t) \end{bmatrix} + \begin{bmatrix} B \\ B \end{bmatrix} v(t), \\ y(t) &= [C \quad O] \begin{bmatrix} x(t) \\ \tilde{x}(t) \end{bmatrix}. \end{aligned}$$

To analyze the above closed-loop system, it is convenient to perform a change of state variables using the following transformation:

$$\begin{bmatrix} x \\ \tilde{x} - x \end{bmatrix} = \begin{bmatrix} I_n & O \\ -I_n & I_n \end{bmatrix} \begin{bmatrix} x \\ \tilde{x} \end{bmatrix}.$$

Note that

$$\begin{bmatrix} I_n & O \\ -I_n & I_n \end{bmatrix}^{-1} = \begin{bmatrix} I_n & O \\ I_n & I_n \end{bmatrix}.$$

In the new coordinates the equations describing the closed-loop system are

$$\begin{aligned} \begin{bmatrix} \dot{x}(t) \\ \dot{\tilde{x}}(t) - \dot{x}(t) \end{bmatrix} &= \begin{bmatrix} A - BK & -BK \\ O & A - LC \end{bmatrix} \begin{bmatrix} x(t) \\ \tilde{x}(t) - x(t) \end{bmatrix} + \begin{bmatrix} B \\ O \end{bmatrix} v(t), \\ y(t) &= [C \quad O] \begin{bmatrix} x(t) \\ \tilde{x}(t) - x(t) \end{bmatrix}. \end{aligned}$$

Note that in the above system the subsystem corresponding to the error component $e(t) = \tilde{x}(t) - x(t)$ is uncontrollable. Furthermore, the $2n$ poles of the closed-loop system are equal to the individual eigenvalues of both $A - BK$ and $A - LC$, which means that the design of the control law is separated from the construction of the estimator. This is what is known as the

separation principle. The closed-loop transfer function relating $Y(s)$ and $V(s)$ is

$$\begin{aligned} Y(s) &= [C \quad O] \begin{bmatrix} sI_n - A + BK & BK \\ O & sI_n - A + LC \end{bmatrix}^{-1} \begin{bmatrix} B \\ O \end{bmatrix} V(s) \\ &= C(sI_n - A + BK)^{-1}BV(s). \end{aligned}$$

The above expression is identical to that for the closed-loop system if we applied the state-feedback control law $u(t) = -Kx(t) + v(t)$. Thus, the combined estimator-controller compensator yields the same closed-loop transfer function as the actual state-feedback control law.

We now briefly discuss the issue of the estimator pole selection. Franklin, Powell, and Emami-Naeini [90, p. 552] recommend that the real parts of the estimator poles—that is, the real parts of the eigenvalues of the matrix $A - LC$ —be a factor of 2 to 6 times deeper in the open left-half plane than the real parts of the controller poles, which are the eigenvalues of the matrix $A - BK$. Such a choice ensures a faster decay of the estimator error $e(t) = \tilde{x}(t) - x(t)$ compared with the desired controller dynamics. This, in turn, causes the controller poles to dominate the closed-loop system response. Because the estimator poles represent a measure of the speed with which the error $e(t) = \tilde{x}(t) - x(t)$ decays to zero, one would tend to assign estimator poles deep in the left-hand plane. However, fast decay requires large gains that may lead to saturation of some signals and unpredictable nonlinear effects. If the estimator poles were slower than the controller poles, the closed-loop system response would be dominated by the estimator, which is undesirable. As is usual in engineering practice, the term *compromise* can be used to describe the process of constructing the final compensator structure.

We now analyze the combined estimator–controller compensator that uses the reduced-order estimator. The closed-loop system driven by this type of compensator also enjoys the separation property as in the case of the full-order estimator. To prove the above statement, we begin with writing the equations describing the closed-loop system:

$$\begin{aligned} \dot{x}(t) &= Ax(t) + Bu(t), \\ \dot{z}(t) &= Dz(t) + TBu(t) + GCx(t), \\ y(t) &= Cx(t), \\ u(t) &= -K\tilde{x}(t) + v(t), \end{aligned}$$

where

$$\tilde{x} = \begin{bmatrix} T \\ C \end{bmatrix}^{-1} \begin{bmatrix} z \\ y \end{bmatrix}.$$

Let

$$\begin{bmatrix} T \\ C \end{bmatrix}^{-1} = [M \quad N],$$

where $M \in \mathbb{R}^{n \times (n-p)}$ and $N \in \mathbb{R}^{n \times p}$. Note that

$$MT + NC = I_n.$$

Hence, the control input can be represented as

$$\begin{aligned} u(t) &= -K \begin{bmatrix} M & N \end{bmatrix} \begin{bmatrix} z(t) \\ y(t) \end{bmatrix} + v(t) \\ &= -KMz(t) - KNy(t) + v(t). \end{aligned}$$

Using the above equations and noting that $y(t) = Cx(t)$, we represent equations that describe the closed-loop system as

$$\begin{aligned} \begin{bmatrix} \dot{x}(t) \\ \dot{z}(t) \end{bmatrix} &= \begin{bmatrix} A & O \\ GC & D \end{bmatrix} \begin{bmatrix} x(t) \\ z(t) \end{bmatrix} + \begin{bmatrix} B \\ TB \end{bmatrix} (-KMz(t) - KNy(t) + v(t)) \\ &= \begin{bmatrix} A - BKN C & -BKM \\ GC - TBKN C & D - TBKM \end{bmatrix} \begin{bmatrix} x(t) \\ z(t) \end{bmatrix} + \begin{bmatrix} B \\ TB \end{bmatrix} v(t), \\ y(t) &= \begin{bmatrix} C & O \end{bmatrix} \begin{bmatrix} x(t) \\ z(t) \end{bmatrix}. \end{aligned}$$

To proceed further, we perform a transformation of the state variable of the above system:

$$\begin{bmatrix} x \\ z - Tx \end{bmatrix} = \begin{bmatrix} I_n & O \\ -T & I_{n-p} \end{bmatrix} \begin{bmatrix} x \\ z \end{bmatrix}.$$

Note that

$$\begin{bmatrix} I_n & O \\ -T & I_{n-p} \end{bmatrix}^{-1} = \begin{bmatrix} I_n & O \\ T & I_{n-p} \end{bmatrix}.$$

In the new coordinates the equations describing the closed-loop system become

$$\begin{aligned} \begin{bmatrix} \dot{x}(t) \\ \dot{z}(t) - T\dot{x}(t) \end{bmatrix} &= \begin{bmatrix} A - BKN C - BKM T & -BKM \\ -TA + GC + DT & D \end{bmatrix} \begin{bmatrix} x(t) \\ z(t) - Tx(t) \end{bmatrix} + \begin{bmatrix} B \\ O \end{bmatrix} v(t) \\ y(t) &= \begin{bmatrix} C & O \end{bmatrix} \begin{bmatrix} x(t) \\ z(t) - Tx(t) \end{bmatrix}. \end{aligned}$$

Taking into account the equations

$$TA - DT - GC = 0 \quad \text{and} \quad MT = I_n - NC,$$

we obtain

$$\begin{aligned} \begin{bmatrix} \dot{x}(t) \\ \dot{z}(t) - T\dot{x}(t) \end{bmatrix} &= \begin{bmatrix} A - BK & -BKM \\ O & D \end{bmatrix} \begin{bmatrix} x(t) \\ z(t) - Tx(t) \end{bmatrix} + \begin{bmatrix} B \\ O \end{bmatrix} v(t) \\ y(t) &= \begin{bmatrix} C & O \end{bmatrix} \begin{bmatrix} x(t) \\ z(t) - Tx(t) \end{bmatrix}. \end{aligned}$$

It is now apparent that the separation principle also holds if we use the reduced-order estimator in constructing the compensator; that is, the selection of the gain matrix K that assigns closed-loop poles is independent, separate from the selection of the estimator matrix D that assigns the estimator poles. Furthermore, the transfer function matrix relating $Y(s)$ and $V(s)$ in the

closed-loop system with the reduced-order estimator is exactly the same as the one that would be obtained using the actual state-feedback control law $\mathbf{u}(t) = -\mathbf{K}\mathbf{x}(t) + \mathbf{v}(t)$.

◆ Example 3.12

Consider a schematic of an armature-controlled DC motor system, shown in Figure 3.8, where the system parameters are: $R_a = 5 \Omega$, $L_a = 200 \text{ mH}$, $K_b = 0.1 \text{ V/rad/sec}$, $K_j = 0.1 \text{ Nm/A}$, the gear ratio $N_1/N_2 = 1/50$. The armature inertia is $I_{\text{armature}} = 2 \times 10^{-3} \text{ kg} \cdot \text{m}^2$. The load of 10 kg is located at an effective radius of 0.2 m . The gear inertia and friction are negligible.

We first construct the state-space model of the DC motor system. We start with writing the equation relating torque to angular acceleration:

$$T_m = I_{eq} \ddot{\theta}, \quad (3.25)$$

where

$$\begin{aligned} I_{eq} &= I_{\text{armature}} + (N_1/N_2)^2 I_l \\ &= 2 \times 10^{-3} \text{ kg} \cdot \text{m}^2 + (1/50)^2 \times 10 \times (0.2)^2 \text{ kg} \cdot \text{m}^2 \\ &= 2.16 \times 10^{-3} \text{ kg} \cdot \text{m}^2. \end{aligned}$$

Kirchhoff's voltage law applied to the armature circuit gives

$$R_a i_a + L_a \frac{di_a}{dt} + e_b = e_a,$$

where $e_b = K_b \frac{d\theta}{dt}$. The equation for the developed torque is

$$T_m = K_j i_a. \quad (3.26)$$

Combining (3.25) and (3.26) gives

$$\ddot{\theta} = \frac{K_j}{I_{eq}} i_a.$$

Let $x_1 = i_a$, $x_2 = \theta$, $x_3 = \dot{\theta} = \omega$, $u = e_a$, and $y = \theta$. Taking into account the definition of the state variables, we represent the above modeling equations in

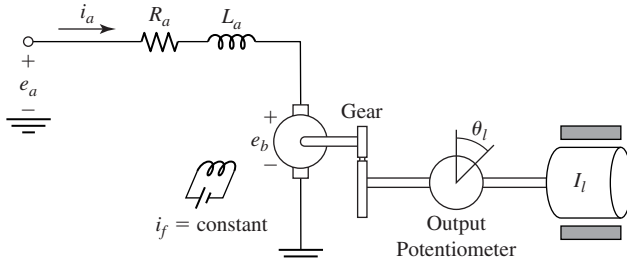


Figure 3.8 Schematic of an armature-controlled DC motor system of Example 3.12.

state-space format:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} -\frac{R_a}{L_a} & 0 & -\frac{K_b}{L_a} \\ 0 & 0 & 1 \\ \frac{K_i}{T_{eq}} & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} \frac{1}{L_a} \\ 0 \\ 0 \end{bmatrix} u$$

$$y = \begin{bmatrix} 0 & \frac{N_1}{N_2} & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}.$$

Substituting the given parameter values, we get

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} -25 & 0 & -0.5 \\ 0 & 0 & 1 \\ 46.296 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 5 \\ 0 \\ 0 \end{bmatrix} u = \mathbf{Ax} + \mathbf{bu},$$

$$y = [0 \quad 0.02 \quad 0] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \mathbf{cx}.$$

Our next step is to design a state-feedback controller $u = -\mathbf{kx} + v$, such that the eigenvalues of $\mathbf{A} - \mathbf{bk}$ are

$$-1 + j, \quad -1 - j, \quad -10.$$

To allocate the poles into the desired locations, we first transform the system model into the controller companion form. We first form the controllability matrix:

$$[\mathbf{b} \quad \mathbf{Ab} \quad \mathbf{A}^2\mathbf{b}] = \begin{bmatrix} 5 & -125 & 3009.3 \\ 0 & 0 & 231.5 \\ 0 & 231.5 & -5787 \end{bmatrix}.$$

The system is controllable. The last row of the inverse of the controllability matrix is

$$\mathbf{q}_1 = [0 \quad 0.0043 \quad 0],$$

and hence the transformation matrix we are seeking has the form

$$\mathbf{T} = \begin{bmatrix} \mathbf{q}_1 \\ \mathbf{q}_1 \mathbf{A} \\ \mathbf{q}_1 \mathbf{A}^2 \end{bmatrix} = \begin{bmatrix} 0 & 0.0043 & 0 \\ 0 & 0 & 0.0043 \\ 0.2 & 0 & 0 \end{bmatrix}.$$

The system's matrices in the new coordinates are

$$\tilde{\mathbf{A}} = \mathbf{TAT}^{-1} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -23.148 & -25 \end{bmatrix},$$

$$\tilde{\mathbf{b}} = \mathbf{Tb} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

and

$$\tilde{\mathbf{c}} = \mathbf{cT}^{-1} = [4.63 \quad 0 \quad 0].$$

The desired closed-loop characteristic polynomial $\alpha_c(s)$ is

$$\alpha_s(s) = (s + 1 - j)(s + 1 + j)(s + 10) = s^3 + 12s^2 + 22s + 20.$$

We now find $\tilde{\mathbf{k}}$ so that

$$\tilde{\mathbf{A}} - \tilde{\mathbf{b}}\tilde{\mathbf{k}} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -20 & -22 & -12 \end{bmatrix}.$$

We have

$$\tilde{\mathbf{k}} = [20 \quad -1.148 \quad -13]$$

and

$$\mathbf{k} = \tilde{\mathbf{k}}\mathbf{T} = [-2.6 \quad 0.0864 \quad -0.005].$$

Hence,

$$\mathbf{A} - \mathbf{b}\mathbf{k} = \begin{bmatrix} -12 & -0.432 & -0.4752 \\ 0 & 0 & 1 \\ 46.294 & 0 & 0 \end{bmatrix}.$$

The transfer function of the closed-loop system is

$$\frac{\mathcal{L}(y(t))}{\mathcal{L}(v(t))} = \frac{Y(s)}{V(s)} = \mathbf{c}(s\mathbf{I} - \mathbf{A} + \mathbf{b}\mathbf{k})^{-1}\mathbf{b} = \frac{4.6296}{s^3 + 12s^2 + 22s + 20}.$$

The unit-step response of the closed-loop system is shown in Figure 3.9.

Next, we design a full-order state estimator placing the estimator poles at $\{-4, -5 + 2j, -5 - 2j\}$ and then synthesize the combined controller-estimator

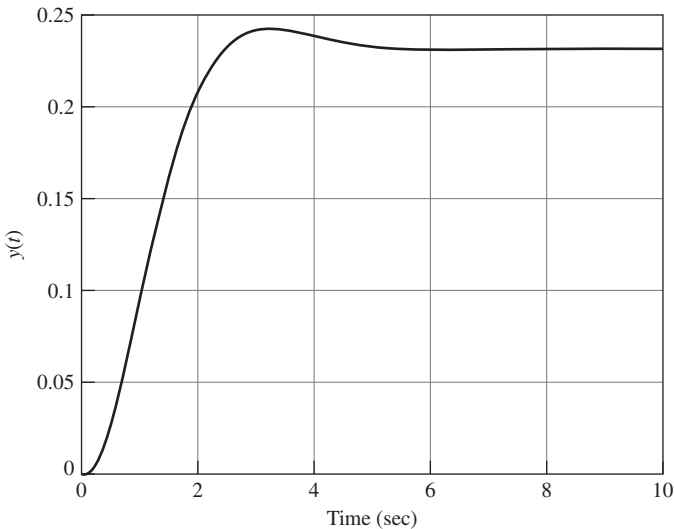


Figure 3.9 Step response of the closed-loop DC motor system with state feedback.

compensator. To compute the estimator gain vector \mathbf{I} , we first transform the pair (\mathbf{A}, \mathbf{c}) into the observer companion form. This is equivalent to transforming the pair $(\mathbf{A}^T, \mathbf{c}^T)$ into the controller companion form. The transformation matrix that brings the pair (\mathbf{A}, \mathbf{c}) into its observer companion form is

$$\hat{\mathbf{Q}} = \hat{\mathbf{T}}^T = \begin{bmatrix} 1.1 & -27 & 650 \\ 0 & 0 & 50 \\ 0 & 50 & -1250 \end{bmatrix}.$$

The matrices \mathbf{A} and \mathbf{c} in the new coordinates have the form

$$\hat{\mathbf{A}} = \hat{\mathbf{Q}}^{-1} \mathbf{A} \hat{\mathbf{Q}} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & -23.1481 \\ 0 & 1 & -25 \end{bmatrix} \quad \text{and} \quad \hat{\mathbf{c}} = \mathbf{c} \hat{\mathbf{Q}} = [0 \quad 0 \quad 1].$$

The desired characteristic polynomial of the estimator is

$$\det(\mathbf{I}_3 - \mathbf{A} + \mathbf{I}\mathbf{c}) = (s + 4)(s + 5 - 2j)(s + 5 + 2j) = s^3 + 14s^2 + 69s + 116.$$

We now find $\hat{\mathbf{I}}$ so that

$$\hat{\mathbf{A}} - \hat{\mathbf{I}}\hat{\mathbf{c}} = \begin{bmatrix} 0 & 0 & -116 \\ 1 & 0 & -69 \\ 0 & 1 & -14 \end{bmatrix}.$$

We have

$$\hat{\mathbf{I}} = \begin{bmatrix} 116 \\ 45.8519 \\ -11 \end{bmatrix} \quad \text{and therefore} \quad \mathbf{I} = \hat{\mathbf{Q}}\hat{\mathbf{I}} = \begin{bmatrix} -8263 \\ -550 \\ 16043 \end{bmatrix}.$$

The dynamics of the full-order state estimator are described by

$$\dot{\tilde{\mathbf{x}}} = (\mathbf{A} - \mathbf{I}\mathbf{c})\tilde{\mathbf{x}} + \mathbf{b}u + \mathbf{I}y;$$

that is,

$$\dot{\tilde{\mathbf{x}}} = \begin{bmatrix} -25 & 165.2544 & -0.5 \\ 0 & 11 & 1 \\ 46.2963 & -320.8519 & 0 \end{bmatrix} \tilde{\mathbf{x}} + \begin{bmatrix} 5 \\ 0 \\ 0 \end{bmatrix} u + \begin{bmatrix} -8263 \\ -550 \\ 16043 \end{bmatrix} y.$$

We connect the full-order estimator to the DC motor system, thus obtaining a closed-loop system with the combined controller–full-order-estimator compensator in the loop as in Figure 3.7. In Figure 3.10 we show a plot of x_1 and its estimate for the closed-loop system with the estimator in the loop, where $v = 0$, $\mathbf{x}(0) = [1 \ 0.2 \ -0.1]^T$, and $\tilde{\mathbf{x}}(0) = \mathbf{0}$.

Our last part of this example is to design the reduced-order state estimator. We choose the reduced-order estimator poles at $\{-5 + 2j, -5 - 2j\}$. The reduced-order estimator dynamics are given by

$$\dot{\mathbf{z}} = \mathbf{D}\mathbf{z} + \mathbf{E}u + \mathbf{G}y.$$

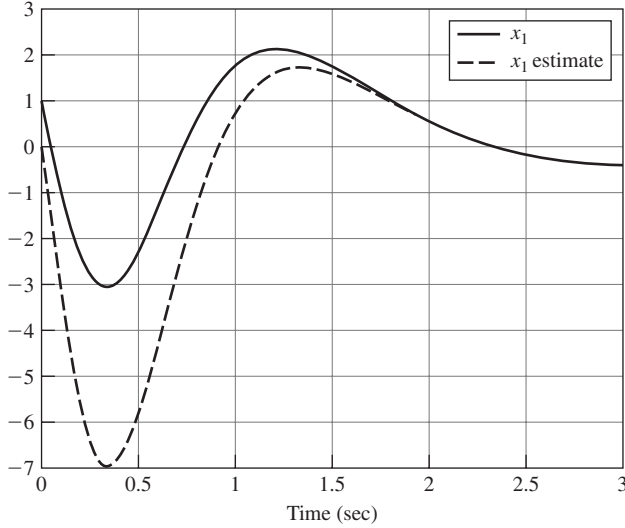


Figure 3.10 A plot of x_1 and its estimate versus time for the closed-loop system with the full-order estimator in the loop.

Let

$$\mathbf{G} = \begin{bmatrix} 100 \\ 100 \end{bmatrix}.$$

The estimator's poles are supposed to be $-5 + 2j$ and $-5 - 2j$. Thus, the desired characteristic polynomial of \mathbf{D} should be

$$\det(s\mathbf{I}_2 - \mathbf{D}) = s^2 + 10s + 29.$$

Let

$$\mathbf{D} = \begin{bmatrix} 0 & 1 \\ -29 & -10 \end{bmatrix}.$$

We now solve the matrix equation,

$$-\mathbf{D}\mathbf{T} + \mathbf{T}\mathbf{A} = \mathbf{G}\mathbf{c},$$

for \mathbf{T} . We first convert the above equation in the format similar to that given by (A.55) in Section A.5 and then solve the resulting equation. Employing the following MATLAB's commands,

```
Gc=G*c;
t=(-kron(eye(3),D)+kron(A',eye(2)))\Gc(:);
T=[t(1:2) t(3:4) t(5:6)]
```

we obtain

$$\mathbf{T} = \begin{bmatrix} -0.5493 & 0.7586 & -0.2507 \\ 2.1271 & -2.0000 & 1.0332 \end{bmatrix}.$$

Note that

$$\det \begin{bmatrix} \mathbf{T} \\ \mathbf{c} \end{bmatrix} = 6.8691 \times 10^{-4} \neq 0$$

and

$$\begin{bmatrix} \mathbf{T} \\ \mathbf{c} \end{bmatrix}^{-1} = \begin{bmatrix} -30.0840 & -7.2980 & 411.3200 \\ 0 & 0 & 50 \\ 61.9338 & 15.9921 & -750 \end{bmatrix}.$$

We also have

$$\mathbf{E} = \mathbf{T}\mathbf{b} = \begin{bmatrix} -2.7463 \\ 10.6357 \end{bmatrix}.$$

Hence, the dynamics of the reduced-order estimator are described by

$$\dot{\mathbf{z}} = \begin{bmatrix} 0 & 1 \\ -29 & -10 \end{bmatrix} \mathbf{z} + \begin{bmatrix} -2.7463 \\ 10.6357 \end{bmatrix} u + \begin{bmatrix} 100 \\ 100 \end{bmatrix} y.$$

An estimate of the entire state can be obtained from

$$\hat{\mathbf{x}} = \begin{bmatrix} \mathbf{T} \\ \mathbf{c} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{z} \\ y \end{bmatrix} = \begin{bmatrix} -30.0840 & -7.2980 & 411.3200 \\ 0 & 0 & 50 \\ 61.9338 & 15.9921 & -750 \end{bmatrix} \begin{bmatrix} \mathbf{z} \\ y \end{bmatrix}.$$

A plot of x_1 and its estimate versus time for the closed-loop system with the combined controller–reduced-order-estimator compensator is depicted in Figure 3.11, where $v = 0$, $\mathbf{x}(0) = [1 \ 0.2 \ -0.1]^T$, and $\dot{\mathbf{z}}(0) = \mathbf{0}$. Of course, for zero initial

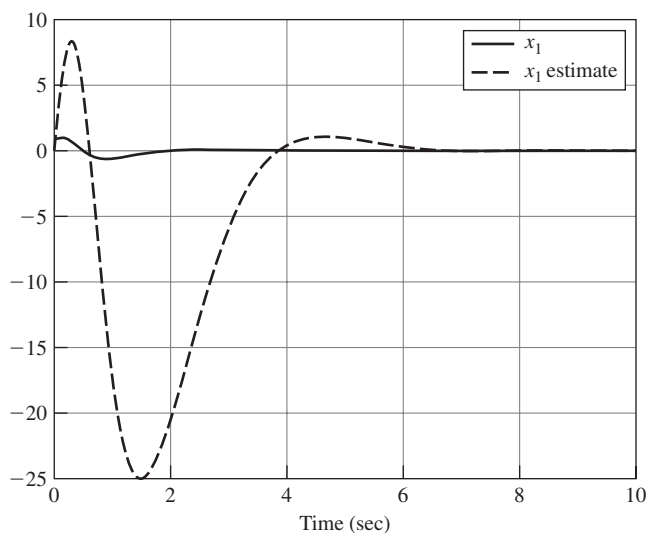


Figure 3.11 A plot of x_1 and its estimate versus time for the closed-loop system with the reduced-order estimator in the loop.

conditions, the step responses of the closed-loop systems with the full-order as well as with the reduced-order estimators in the loop are the same as the step response of the closed-loop system with the state-feedback controller $u = -\mathbf{k}\mathbf{x} + v$. This follows from our analyses of the combined controller–estimator compensators.

Notes

For further reading on the subject of linear system theory, see Brockett [34], Wolovich [301], Barnett [20], Kailath [146], Skelton [261], DeCarlo [59], Rugh [245], Antsaklis and Michel [12], or Chen [48]. For more information on canonical forms, see the papers by Popov [236] and Denham [62].

One may argue that Maxwell’s [197] analysis of Watt’s governor constitutes the first work in linear system theory: the physical object, its linear model, and a method of analysis. For example, MacFarlane [192, p. 251] states that “Maxwell’s fundamentally important contribution lay in recognizing that the behavior of an automatic feedback control system in the vicinity of an equilibrium condition could be approximated by a linear differential equation, and hence the stability of the control system could be discussed in terms of the location of the roots of an associated algebraic equation.” On the other hand, L. Markus in Elworthy et al. [76, p. xxvii] makes the following statement: “But certainly Watt’s application of the centrifugal governor for regulating steam engines (1775) can be taken as the start of modern control engineering.”

In Kalman’s [147] well-known 1963 paper, there is a theorem which states that the state space, \mathcal{X} , of a linear system model may be decomposed as a direct sum:

$$\mathcal{X} = \mathcal{X}_A \oplus \mathcal{X}_B \oplus \mathcal{X}_C \oplus \mathcal{X}_D,$$

where

\mathcal{X}_A = states that are reachable but not observable,

\mathcal{X}_B = states that are reachable and observable,

\mathcal{X}_C = states that are neither reachable nor observable,

\mathcal{X}_D = states that are not reachable but observable.

This decomposition became, to a certain extent, an icon of linear systems. For example, the logo on the covers of the old issues of the *International Journal of Systems Science* was the above decomposition. Many authors included Kalman’s original decomposition algorithm in their textbooks on linear systems. Yet almost twenty years after Kalman’s publication of the decomposition theorem, Boley [29] generated an example that revealed some shortcomings of Kalman’s original algorithm. In his reply, Kalman [148] revised his decomposition algorithm and fixed his original proof of the decomposition theorem.

EXERCISES

- 3.1** (Adapted from Barnett [20, p. 94]) Consider the system, shown in Figure 3.12, which is composed of two masses moving on a smooth—that is, frictionless—horizontal plane.

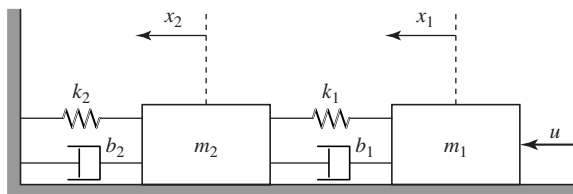


Figure 3.12 A mechanical system for Exercise 3.1.

The springs are linear; that is, they obey Hooke's law:

$$\text{force} = k \times \text{extension},$$

where k is a constant called the force constant of the spring. In addition to the springs, the masses are joined by linear dampers. In a linear damper, the force is proportional to the velocity of one end of the damper relative to the other end. The proportionality constant is called the damping coefficient. In this exercise, damping coefficients are b_1 and b_2 , respectively. A force u is applied to the right end mass. Obtain the system's state-space model in the form $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b}u$. Determine under what condition, in terms of k_1 , the system model is controllable for $m_1 = m_2 = 1$ kg, $b_1 = b_2 = 1$ Nsec/m, and $k_2 = 1/4$ N/m.

- 3.2** In the circuit shown in Figure 3.13, take the state variables to be the voltage across the capacitor and the current through the inductor. The source voltage v_s is the control variable, and the current i is the system output. Obtain a state-space model of this circuit. Obtain the condition, in terms of R_1 , R_2 , C , and L , under which the model is uncontrollable and unobservable. Find the transfer function of the circuit model.

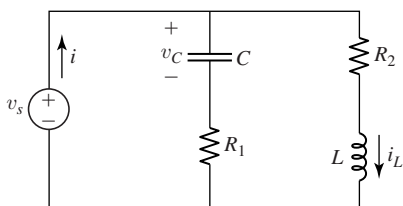


Figure 3.13 Electric circuit for Exercise 3.2.

- 3.3** In the circuit shown in Figure 3.14, take the state variables to be the voltage across the capacitor C_1 , the current through the inductor L_1 , and the current through the inductor L_2 . The source current i_s is the control variable, and the current through the inductor L_2 is the system output.

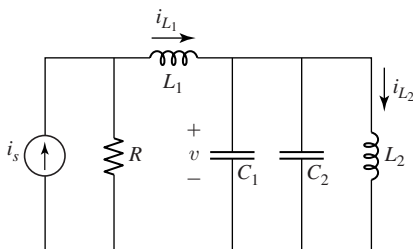


Figure 3.14 Electric circuit for Exercise 3.3.

L_2 is the system output. Obtain a state-space model of this circuit. Analyze the circuit model controllability and observability properties. Find the transfer function of the circuit model.

3.4 For the system

$$\dot{x} = x + u,$$

find the set of points that can be steered to the origin $x = 0$ by an admissible control $|u| \leq 1$.

3.5 (a) Find e^{At} for

$$A = \begin{bmatrix} 0 & 1 \\ -ab & -a-b \end{bmatrix}.$$

(b) Given a dynamical system modeled by the differential equation

$$\ddot{x} + 2\dot{x} + x = \dot{r} + r,$$

where $r = r(t) = \sin t$, and the initial conditions are $x(0) = 1$ and $\dot{x}(0) = 0$, choose the state variables $x_1 = x$ and $x_2 = \dot{x}_1$. Represent the differential equation in state-space format. Solve for $x(t)$ for $t \geq 0$.

In both cases, verify your calculations using MATLAB's Symbolic Math Toolbox.

3.6 Find sufficient conditions, in terms of a , b , c , and d , for reachability of the following pairs:

(a)

$$A = \begin{bmatrix} -1 & 1 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & -2 & 1 \\ 0 & 0 & 0 & 0 & -2 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ a & b \\ 0 & 0 \\ c & d \end{bmatrix};$$

(b)

$$A = \begin{bmatrix} -1 & 1 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 \\ 0 & 0 & 0 & 0 & -1 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ a & b \\ 0 & 0 \\ c & d \end{bmatrix}.$$

3.7 A system composed of two subsystems connected in cascade is depicted in the block diagram shown in Figure 3.15. Take the state variables of the system to be $x_1 = y_1 - u_1$ and $x_2 = y_2$. Assume that the system input is $u = u_1$ and the system output is $y = y_2$. Obtain the state-space model of the interconnected system. When is the system uncontrollable? Find the transfer function of the uncontrollable system.

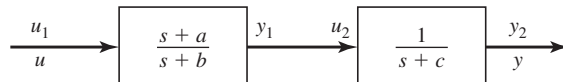


Figure 3.15 A block diagram for Exercise 3.7.

- 3.8** Derive a state-space description of a system represented by the block diagram shown in Figure 3.16. Construct a state-variable transformation $z = T\mathbf{x}$ such that the matrix TAT^{-1} has a diagonal form. Find the eigenvalue of TAT^{-1} corresponding to the uncontrollable state, and find the eigenvalue corresponding to the nonobservable state. Find the transfer function of the system.

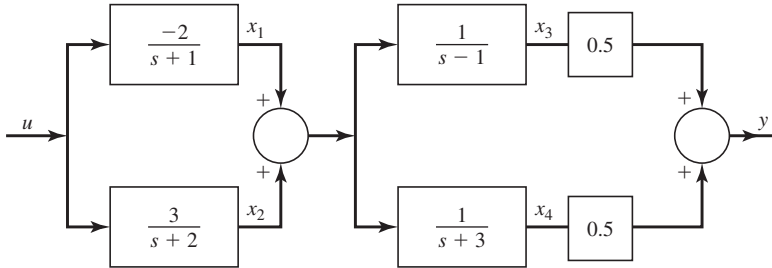


Figure 3.16 A block diagram for Exercise 3.8.

- 3.9** (Adapted from Barnett [20, p. 95]) A simple representation of a vehicle suspension system is shown in Figure 3.17. Assume that the forces produced by the springs and dampers act at the end points P and Q and that x_1 and x_2 are the displacements of these points from equilibrium. The force exerted by a damper is proportional to the velocity of one end of the damper relative to the other end, and the springs obey Hooke's law. Assume that the mass of the platform can be neglected so that the spring motions can be regarded as independent. The control force $4u$ is applied one-quarter of the way from one end of the platform. Derive the two-dimensional state-space model of this simple representation of a vehicle suspension system. Next, if the ends of the platform are each given an initial displacement of 10 units, find a control law that returns the platform to equilibrium at $t = 1$ sec.

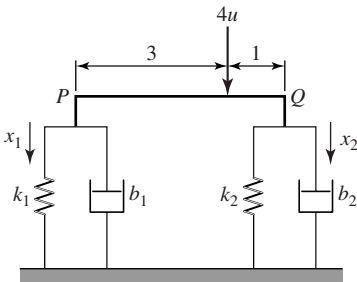


Figure 3.17 A simple representation of a vehicle suspension system of Exercise 3.9.

- 3.10** For the dynamical system model, $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b}u$, where

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \end{bmatrix},$$

design a state-feedback controller that stabilizes the system about the operating state corresponding to the constant input $u = 2$. The desired closed-loop poles are $-2, -3$.

- 3.11** Consider the vehicle model given in Subsection 1.7.2. Let the input signal be $u = \delta_f$ (the front-wheel steer angle) and $\delta_r = 0$. Compute the transfer function $C(s)/U(s)$. Plot the vehicle response Y versus X for X changing from 0 m to 100 m, when the input δ_f , in the vehicle coordinates, has the form $\delta_f(t) = 0.0872 h(t)$ rad, where $h(t)$ is the unit step function.

Suppose now that the plot of δ_f versus time, in the car coordinates, has the form shown in Figure 3.18. Note that the values of δ_f are given in deg. The value of δ_f changes from -2.5 deg to 2.5 deg. In your simulations you will need to convert deg to rad. Plot the vehicle response Y versus X for X changing from 0 m to 100 m.

Construct the state-feedback controller $u = -kx + r$ such that $\text{eig}(A - bk) = \{-2, -3 + j, -3 - j, -4\}$. Plot u and r versus time, and Y versus X , when $r = 0.0872 h(t)$ rad.

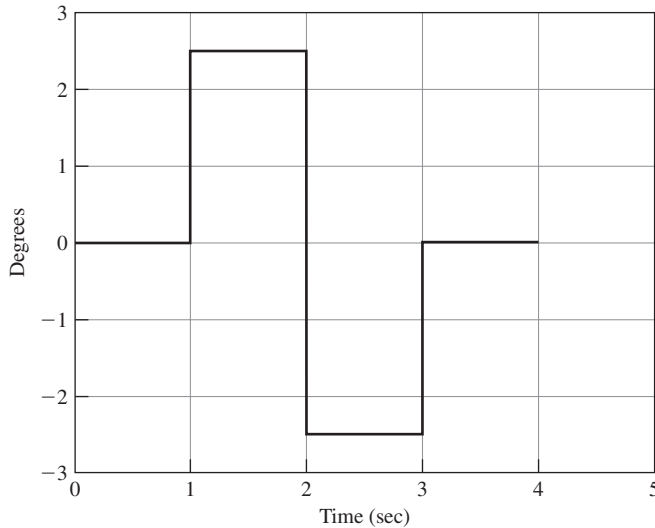


Figure 3.18 Time history of the front steer angle in the car coordinates of Exercise 3.11.

- 3.12** Show that the controllability property is invariant under the following set of transformations:
- Nonsingular coordinate transformation in the state space
 - Nonsingular transformation of the inputs
 - State feedback
 - Output feedback
- 3.13** Given the following model of a dynamical system,

$$\dot{\mathbf{x}}(t) = \begin{bmatrix} 0 & 2 \\ 1 & 0 \end{bmatrix} \mathbf{x}(t) + \begin{bmatrix} 1 \\ 0 \end{bmatrix} u(t),$$

$$y(t) = [0 \quad 1] \mathbf{x}(t),$$

do the following:

- Design a full-order asymptotic state estimator for the above system with the estimator poles located at -2 and -3 . Write the equations of the estimator dynamics.

Note that the pair (A, c) is already in the observer form, which should facilitate your estimator design.

- (b) Denote the estimator state vector by z . Let $u = -[2 \ 3]z + v$. Find the transfer function of the closed-loop system $Y(s)/V(s)$.

3.14 Let

$$\begin{aligned}\dot{x} &= Ax + bu = \begin{bmatrix} -1 & 1 \\ 0 & 2 \end{bmatrix} x + \begin{bmatrix} 1 \\ 1 \end{bmatrix} u, \\ y &= cx = [1 \ 0]x.\end{aligned}$$

- (a) Find the state transformation that transforms the pair (A, c) into the observer companion form.
 (b) Find the representation of the dynamical system in the new coordinates.

3.15 For the following model of a dynamical system,

$$\begin{aligned}\dot{x} &= \begin{bmatrix} 0 & 1 \\ 2 & 0 \end{bmatrix} x + \begin{bmatrix} 1 \\ 0 \end{bmatrix} u, \\ y &= [1 \ 0]x + 2u,\end{aligned}$$

do the following:

- (a) Design a full-order asymptotic state estimator with the estimator poles located at -3 and -4 . Write the equations of the estimator dynamics.
 (b) Denote the estimator state by z . Let $u = -[3 \ 2]z + r$. Find the transfer function, $Y(s)/R(s)$, of the closed-loop system.

3.16 Given a dynamical system modeled by

$$\begin{aligned}\dot{x} &= Ax + bu, \\ y &= cx,\end{aligned}$$

where

$$A = \begin{bmatrix} 1 & -3 \\ 4 & 2 \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad c = [2 \ 3],$$

do the following:

- (a) Design a reduced-order state estimator

$$\dot{z} = dz + eu + gy,$$

where $d = -5$ and $g = 1$.

- (b) Write an expression for the estimated states.

3.17 For the plant

$$\begin{aligned}\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} &= \begin{bmatrix} 1 & -1 \\ 1 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 2 \\ 1 \end{bmatrix} u \\ y &= [0 \ 1] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix},\end{aligned}$$

do the following:

- Design the reduced-order state estimator for $G = [1]$ and $D = [-2]$.
- Write the equations of the estimated state in terms of the plant's output and the estimator state.

3.18 Consider a dynamical system that consists of a cart with an inverted pendulum (point mass on a massless shaft) attached to it as depicted in Figure 3.19.

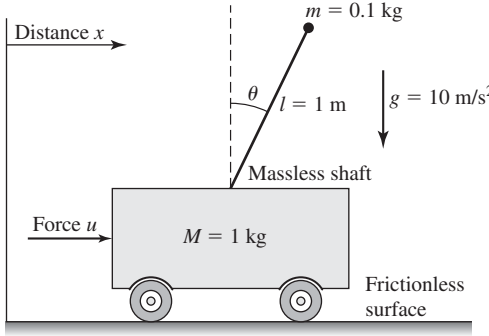


Figure 3.19 An inverted pendulum moving on a cart of Exercise 3.18.

- Obtain a nonlinear state-space model of the system using the state variables

$$x_1 = x, \quad x_2 = \dot{x}, \quad x_3 = \theta, \quad x_4 = \dot{\theta}.$$

Designate x_1 as the system output.

- Linearize the model about the operating point $\mathbf{x} = \mathbf{0}$, $u = 0$.
- Design a state feedback control law $u = -\mathbf{k}\mathbf{x}$ such that the closed-loop poles are located at

$$-1, \quad -2, \quad -1 \pm j$$

- Write a simple MATLAB script that animates the motion of the linearized closed-loop model of the inverted pendulum on the cart along with animation of the displacement x and the angle θ data. The basic handle graphics of MATLAB can be used to write the script. A script contained in Example 4 in Pratap [238, p. 194] can be helpful here. Alternatively, you can employ SIMULINK.

3.19 Consider the nonlinear model of the one-link robot manipulator of Example 5.13. In the derivation, assume that $L_a = 100 \text{ mH}$ and $g = 9.8 \text{ m/s}^2$. Represent the obtained third-order model in the form $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, u)$, where $x_1 = \theta_p$, $x_2 = \dot{\theta}_p$, and $x_3 = i_a$.

- Find the constant value of $u = u_e$ and x_{3e} so that the state $\mathbf{x}_e = [30\pi/180 \ 0 \ x_{3e}]^T$ becomes the equilibrium state; that is, the obtained \mathbf{x}_e and u_e satisfy the equation

$$\mathbf{0} = \mathbf{f}(\mathbf{x}_e, u_e).$$

- Linearize the nonlinear model, $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, u)$, about the operating point $[\mathbf{x}_e^T \ u_e]^T$ to obtain

$$\frac{d}{dt}\Delta\mathbf{x} = \mathbf{A}\Delta\mathbf{x} + \mathbf{b}\Delta u,$$

where

$$\Delta \mathbf{x} = \mathbf{x} - \mathbf{x}_e, \quad \Delta u = u - u_e,$$

and

$$\mathbf{A} = \frac{\partial \mathbf{f}}{\partial \mathbf{x}}(\mathbf{x}_e, u_e) \quad \text{and} \quad \mathbf{b} = \frac{\partial \mathbf{f}}{\partial u}(\mathbf{x}_e, u_e).$$

- (c) Design a state feedback control law $\Delta u = -\mathbf{k} \Delta \mathbf{x}$ such that the closed-loop poles are located at

$$-2, \quad -2 \pm j$$

- (d) Implement the obtained control law on the nonlinear model $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, u)$. Write a MATLAB script that animates the dynamical behavior of the nonlinear closed-loop system along with animation of the angle θ_p and the input u data.
- (e) Design a full-order state estimator for the linearized model. The estimator's poles should be located at

$$-5, \quad -5 \pm j$$

- (f) Implement the combined controller–estimator compensator on the nonlinear model and write a MATLAB script that animates the dynamical behavior of the nonlinear closed-loop system along with animation of the angle θ_p and the input u data.



CHAPTER 4

Stability

Examine each question in terms of what is ethically and aesthetically right, as well as what is economically expedient. A thing is right when it tends to preserve the integrity, stability, and beauty of the biotic community. It is wrong when it tends otherwise.

—Aldo Leopold

4.1 Informal Introduction to Stability

In general, it can be difficult to obtain an explicit solution of nonlinear differential equations modeling a dynamical system. However, often we do not need to know the explicit solution of the modeling equations. Instead, we are interested in the behavior of the solution as time tends to infinity. Furthermore, in many models of physical systems a “small” change in the initial conditions results in a “small” change in the solution. Specifically, if a system is perturbed from its rest, or equilibrium position, then it starts to move. We can roughly say that the equilibrium position is stable if the system does not go far from this position for small initial perturbations. This concept may be illustrated, as in Willems [300], Brogan [35], or Barnett [20], by means of a ball resting in equilibrium on a sheet of metal bent into various shapes with cross sections depicted in Figure 4.1. This figure illustrates four possible scenarios. If we neglect frictional forces, a small perturbation from the equilibrium leads to the following:

1. Oscillatory motion about equilibrium irrespective of the magnitude of the initial perturbation. In this case the equilibrium position is globally stable.
2. The new position is also an equilibrium, and the ball will remain in it. Here the equilibrium position is stable.
3. The ball moving away without returning to the equilibrium. This equilibrium position is unstable.
4. Oscillatory motion about the equilibrium, unless the initial perturbation is so large that the ball is forced to oscillate about the new equilibrium position. This equilibrium position is locally stable.

If friction is taken into account, then the oscillatory motion steadily decreases until the ball returns to the equilibrium positions as in items 1 and 4 above. We refer to such equilibrium positions as being asymptotically stable. In our discussion of stability concepts, we implicitly

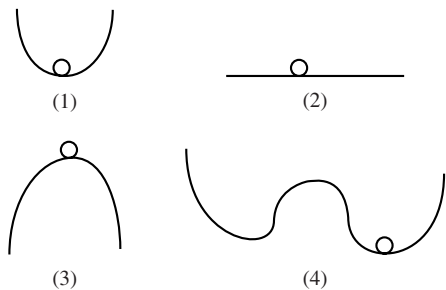


Figure 4.1 Illustration of stable and unstable equilibrium positions.

assume, as in item 4, that there is some “small” region for allowed initial perturbations for which subsequent motions converge to the equilibrium position. If however, as in item 1, any initial perturbation yields a motion that converges to the equilibrium position, then we have global asymptotic stability, which is also referred to as asymptotic stability in the large.

The above informal discussion should serve as a preparation for an in-depth analysis of stability of dynamical systems. We begin our analysis by formally defining the equilibrium state. We then define stability, asymptotic stability, and other relevant notions.

4.2 Basic Definitions of Stability

We consider a class of dynamical systems modeled by the equation:

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t)), \quad \mathbf{x}(t_0) = \mathbf{x}_0, \quad (4.1)$$

where $\mathbf{x}(t) \in \mathbb{R}^n$ is the state vector and $\mathbf{f} : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a vector-valued function with the components

$$f_i(t, x_1, x_2, \dots, x_n) : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}, \quad i = 1, 2, \dots, n.$$

We assume that f_i 's are continuous and have continuous first partial derivatives. Under these assumptions, the solution to $\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t))$, $\mathbf{x}(t_0) = \mathbf{x}_0$, exists and is unique. We denote this solution as $\mathbf{x}(t) = \mathbf{x}(t; t_0, \mathbf{x}_0)$. If the f_i 's do not depend explicitly on t , then the system represented by $\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t)) = \mathbf{f}(\mathbf{x}(t))$ is called *autonomous*; otherwise it is called *nonautonomous*. An *equilibrium point* or state is a constant vector, say \mathbf{x}_e , such that

$$\mathbf{f}(t, \mathbf{x}_e) = \mathbf{0} \quad \text{for all } t.$$

Observe that an equilibrium state \mathbf{x}_e is a constant solution to (4.1).

If \mathbf{x}_e is an equilibrium point, we can arrange for this equilibrium point to be transferred to the origin of \mathbb{R}^n by introducing the new variable:

$$\tilde{\mathbf{x}} = \mathbf{x} - \mathbf{x}_e.$$

We assume that this has been done for an equilibrium point under consideration. Therefore, we have

$$\mathbf{f}(t, \mathbf{0}) = \mathbf{0} \quad \text{for all } t.$$

We now introduce basic definitions of stability in the sense of A. M. Lyapunov* (1857–1918). In the following considerations we use the standard Euclidean norm of a vector. Thus, if $\mathbf{x} \in \mathbb{R}^n$, then

$$\|\mathbf{x}\| = \|\mathbf{x}\|_2 = (\mathbf{x}^T \mathbf{x})^{1/2}.$$

Definition 4.1 An equilibrium state \mathbf{x}_e is said to be *stable* if for any given t_0 and any positive scalar ε , there exists a positive scalar $\delta = \delta(t_0, \varepsilon)$ such that if $\|\mathbf{x}(t_0) - \mathbf{x}_e\| < \delta$, then $\|\mathbf{x}(t; t_0, \mathbf{x}_0) - \mathbf{x}_e\| < \varepsilon$ for all $t \geq t_0$.

Friedland [91, p. 43] makes the following comment regarding the above definition of stability: “If this ‘epsilonics’ definition leaves you cold, think of a contest between you, the system designer, and an adversary (nature?). The adversary picks a region in state space of radius ε and challenges you to find another region, of radius δ , such that if the initial state starts out inside your region, it remains inside his region. If you can do this, your system design is stable.” We add that you have to find an appropriate δ for any choice of ε for the equilibrium state to be stable.

Definition 4.2 An equilibrium state \mathbf{x}_e is said to be *convergent*, or *attractive*, if for any given t_0 there is a positive scalar $\delta_1 = \delta_1(t_0)$ such that if $\|\mathbf{x}(t_0) - \mathbf{x}_e\| < \delta_1$, then

$$\lim_{t \rightarrow \infty} \mathbf{x}(t; t_0, \mathbf{x}_0) = \mathbf{x}_e.$$

We can restate Definition 4.2 as follows. An equilibrium state \mathbf{x}_e is convergent if for any t_0 there is a positive scalar $\delta_1 = \delta_1(t_0)$ such that if $\|\mathbf{x}(t_0) - \mathbf{x}_e\| < \delta_1$, then for any given $\varepsilon_1 > 0$ there exists a $T = T(\varepsilon_1, t_0, \mathbf{x}_0)$ such that $\|\mathbf{x}(t; t_0, \mathbf{x}_0) - \mathbf{x}_e\| < \varepsilon_1$ for all $t > t_0 + T$.

Definition 4.3 An equilibrium state \mathbf{x}_e is *asymptotically stable* if it is stable and convergent.

Definition 4.4 An equilibrium state \mathbf{x}_e is said to be *uniformly stable* if the positive scalar $\delta = \delta(t_0, \varepsilon)$ introduced in Definition 4.1 can be taken independent of the initial time t_0 .

◆ Example 4.1

The following example that comes from Rugh [245, p. 100] nicely illustrates the notion of uniform stability. The solution to the differential equation,

$$\dot{x}(t) = (4t \sin t - 2t)x(t), \quad x(t_0) = x_0, \quad (4.2)$$

is

$$x(t) = x_0 \exp(4 \sin t - 4t \cos t - t^2 - 4 \sin t_0 + 4t_0 \cos t_0 + t_0^2). \quad (4.3)$$

One can verify this using MATLAB’s Symbolic Toolbox by typing the command

```
x=dsolve('Dx=(4*t*sin(t)-2*t)*x, x(t_0)=x_0','t')
```

It is easy to show that for fixed t_0 there is a γ such that (4.3) is bounded by $\gamma|x_0|$ for all $t \geq t_0$. This is because the term $(-t^2)$ dominates the exponent as t increases. Therefore, for a given ε , it is enough to choose x_0 so that $\gamma|x_0| < \varepsilon$, or equivalently $\delta < \varepsilon/\gamma$, to show that the equilibrium solution $x_e = 0$ is stable. However, the

*Sometimes transliterated in the literature as Liapunov or Liapunoff.

equilibrium solution $x_e = 0$ is not uniformly stable. To show this, fix initial state x_0 and consider a sequence of initial times $t_0 = 2k\pi$, where $k = 0, 1, \dots$, and the values of the respective solutions at times π units later:

$$x(2k\pi + \pi) = x_0 \exp((4k + 1)\pi(4 - \pi)).$$

One can see that there is no bound on the exponential factor and hence a candidate bound γ increases as k increases, where increasing k corresponds to the increasing initial time.

Definition 4.5 An equilibrium state x_e is said to be uniformly convergent, or uniformly attractive, if the positive scalars $\delta_1 = \delta_1(t_0)$ and $T = T(\varepsilon_1, t_0, x_0)$ introduced in Definition 4.2 are independent of t_0 .

The last two definitions allow us to introduce the notion of uniform asymptotic stability.

Definition 4.6 An equilibrium state x_e is said to be uniformly asymptotically stable if it is uniformly stable and uniformly convergent.

Notice that the solution $x(t) = x_e$ must first be stable in order to qualify to be asymptotically stable. This is to prevent a system trajectory from straying arbitrarily far from the equilibrium state before converging toward it. Note also that the concepts of Lyapunov stability and asymptotic stability are local concepts. They apply to the system's behavior in a small neighborhood about the equilibrium point (Hsu and Meyer [128]).

A special case of uniform asymptotic stability is *uniform exponential stability* that implies uniform stability, and it imposes an additional requirement that all solutions originating “close” to the given equilibrium, $x_e = 0$, approach this equilibrium exponentially as $t \rightarrow \infty$. More rigorously, we say that $x_e = 0$ is uniformly exponentially stable if it is stable and there exist finite positive constants γ and λ such that for any t_0 and x_0 , any solution starting in a neighborhood of 0 satisfies

$$\|x(t)\| \leq \gamma e^{-\lambda(t-t_0)} \|x_0\|.$$

Note that γ is no less than unity, and uniform implies that γ and λ are independent of t_0 (see Figure 4.2).

An equilibrium state is unstable if it is not stable. An ε - δ definition of instability is given next.

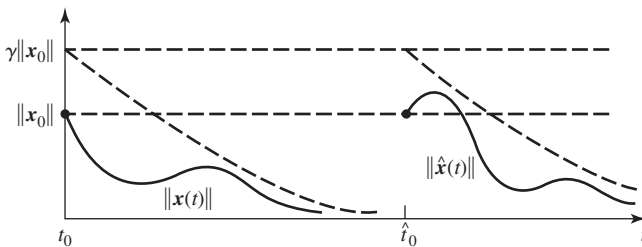


Figure 4.2 Illustration of uniform exponential stability: A decaying-exponential bound is independent of the initial time t_0 .

Definition 4.7 An equilibrium state \mathbf{x}_e is said to be unstable if there is an $\varepsilon > 0$ so that for any $\delta > 0$ there exists $\mathbf{x}(t_0)$ such that if $\|\mathbf{x}(t_0) - \mathbf{x}_e\| < \delta$, then $\|\mathbf{x}(t_1) - \mathbf{x}_e\| \geq \varepsilon$ for some $t_1 > t_0$.

A two-dimensional geometric interpretation of the above definitions is shown in Figure 4.3. We assume that $\mathbf{x}_e = 0$ is the equilibrium state. If the equilibrium state is stable, as in Figure 4.3(a), then given an outer circle of radius ε , there exists an inner circle of radius δ such that trajectories starting within the δ -circle never leave the ε -circle. If the equilibrium state is asymptotically stable, then trajectories tend to the equilibrium state as t tends to ∞ . This situation is depicted in Figure 4.3(b). Finally, if the equilibrium state is unstable, as in Figure 4.3(c), then there exists an ε -circle such that for every δ -circle there exists a trajectory starting within it that leaves the ε -circle at some later time. A geometric interpretation of the above described scenarios in one dimension is given in Figure 4.4.

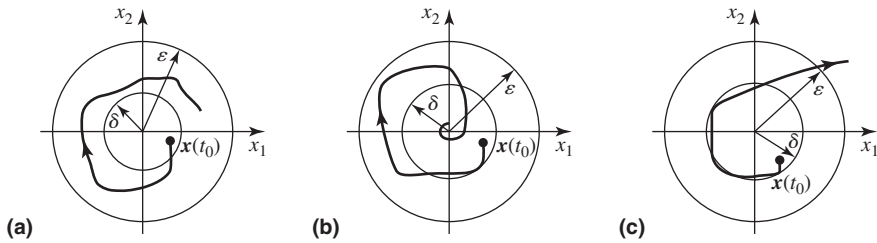


Figure 4.3 Illustration of stability concepts in two dimensions. (a) Stable equilibrium state in the sense of Lyapunov. (b) Asymptotically stable equilibrium state. (c) Unstable equilibrium state.

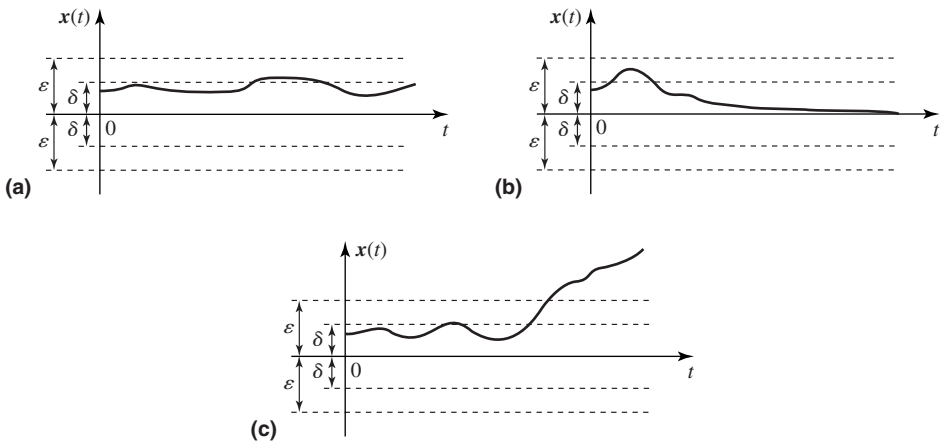


Figure 4.4 Illustration of stability concepts in one dimension. (a) Stable equilibrium state in the sense of Lyapunov. (b) Asymptotically stable equilibrium state. (c) Unstable equilibrium state.

4.3 Stability of Linear Systems

It follows from the discussion in the previous sections that a linear time-invariant dynamical system modeled by the equations

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}, \quad \mathbf{x}(t_0) = \mathbf{x}_0 \quad (4.4)$$

is asymptotically stable if and only if the solution to equation (4.4) decays to zero as $t \rightarrow \infty$ for any initial state \mathbf{x}_0 . We can view the vector $\mathbf{x}(t) \in \mathbb{R}^n$ as defining the coordinates of a point moving in an n -dimensional state space \mathbb{R}^n . In an asymptotically stable system modeled by equation (4.4), this point converges to the origin of \mathbb{R}^n . We can argue that if a trajectory is converging to the origin of the state space, then it should be possible to find a family of nested surfaces described by $V(x_1, x_2, \dots, x_n) = c$, $c \geq 0$, such that monotonically decreasing values of c give surfaces progressively shrinking in on the origin with the limiting surface $V(x_1, x_2, \dots, x_n) = 0$ being the origin $\mathbf{x} = \mathbf{0}$. Furthermore, along every trajectory of (4.4), the parameter c should steadily decrease. We now introduce some notation and definitions that we use in the ensuing discussion.

Let $V(\mathbf{x}) = V(x_1, x_2, \dots, x_n) : \mathbb{R}^n \rightarrow \mathbb{R}$ be a real-valued function and let $S \subseteq \mathbb{R}^n$ be a compact region containing the origin $\mathbf{x} = \mathbf{0}$ in its interior.

Definition 4.8 We say that the function $V = V(\mathbf{x})$ is *positive semidefinite* (p.s.d.) in S , with respect to $\mathbf{x} = \mathbf{0}$, if:

1. V is continuously differentiable, that is, $V \in C^1$.
2. $V(\mathbf{0}) = 0$.
3. $V(\mathbf{x}) \geq 0$ for all $\mathbf{x} \in S$.

Definition 4.9 We say that the function $V = V(\mathbf{x})$ is *positive definite* (p.d.) in S , with respect to $\mathbf{x} = \mathbf{0}$, if:

1. $V \in C^1$.
2. $V(\mathbf{0}) = 0$.
3. $V(\mathbf{x}) > 0$ for all $\mathbf{x} \in S \setminus \{\mathbf{0}\}$.

The symbol $S \setminus \{\mathbf{0}\}$ in item 3 above denotes the set $\{\mathbf{x} : \mathbf{x} \in S, \mathbf{x} \neq \mathbf{0}\}$.

Negative semidefinite (n.s.d.) and *negative definite* (n.d.) functions are defined in a similar fashion by reversing inequality signs in item 3 in both definitions above.

There are many functions that satisfy the above definitions. In the stability analysis of dynamical systems, one class is particularly useful. This is the positive definite quadratic form, a label attached to a nested family of ellipses and their extensions to higher-dimensional spaces—ellipsoids for $n = 3$ and hyperellipsoids for $n \geq 4$. Quadratic forms are reviewed in Section A.4. We consider the quadratic form

$$V(\mathbf{x}) = \mathbf{x}^T \mathbf{P} \mathbf{x},$$

where \mathbf{P} is a real symmetric $n \times n$ matrix. The time derivative of $V(\mathbf{x}(t))$ evaluated on a solution of $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$ is

$$\begin{aligned} \frac{dV(\mathbf{x}(t))}{dt} &= \dot{V}(t) \\ &= \dot{\mathbf{x}}^T(t) \mathbf{P} \mathbf{x}(t) + \mathbf{x}^T(t) \mathbf{P} \dot{\mathbf{x}}(t). \end{aligned}$$

Substituting $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$ and $\dot{\mathbf{x}}^T = \mathbf{x}^T \mathbf{A}^T$ into the above equation yields

$$\begin{aligned}\dot{V} &= \mathbf{x}^T \mathbf{A}^T \mathbf{P} \mathbf{x} + \mathbf{x}^T \mathbf{P} \mathbf{A} \mathbf{x} \\ &= \mathbf{x}^T (\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A}) \mathbf{x}.\end{aligned}$$

Let

$$\mathbf{Q} = -(\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A}).$$

Then,

$$\dot{V} = -\mathbf{x}^T \mathbf{Q} \mathbf{x}. \quad (4.5)$$

We are now ready to state and prove the following theorem concerning the stability of the system $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$, or, equivalently, the matrix $\mathbf{A} \in \mathbb{R}^n$.

Theorem 4.1 (A. M. Lyapunov, 1892) The system $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$, $\mathbf{x}(t_0) = \mathbf{x}_0$, is asymptotically stable if and only if for any given real symmetric positive definite (r.s.p.d.) matrix \mathbf{Q} the solution \mathbf{P} to the continuous Lyapunov matrix equation,

$$\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{Q} \quad (4.6)$$

is also (real symmetric) positive definite.

Proof (\Rightarrow) We prove necessity, for asymptotic stability, by contradiction. We have $-\mathbf{x}^T \mathbf{Q} \mathbf{x} < 0$ for all $\mathbf{x} \neq \mathbf{0}$. We assume that \mathbf{A} is an asymptotically stable matrix and that for some nonzero \mathbf{x}_0 we have $\mathbf{x}_0^T \mathbf{P} \mathbf{x}_0 \leq 0$. From the discussion before the theorem, and in particular from (4.5), it follows that

$$\frac{d}{dt}(\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)) = -\mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) < 0.$$

The above means that $V(\mathbf{x}(t)) = \mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)$ decreases on the trajectories of the system $\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t)$ as $t \rightarrow \infty$. On the other hand, because \mathbf{A} is asymptotically stable, we have

$$\lim_{t \rightarrow \infty} \mathbf{x}(t) = \mathbf{0}.$$

Hence,

$$\lim_{t \rightarrow \infty} \mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t) = \lim_{\mathbf{x}(t) \rightarrow \mathbf{0}} \mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t) = 0.$$

Thus, we arrived at a contradiction. This completes the proof of the first part of the theorem that states that asymptotic stability of \mathbf{A} implies that for any r.s.p.d matrix \mathbf{Q} the solution \mathbf{P} to the equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{Q}$ is r.s.p.d..

(\Leftarrow) We also prove sufficiency by contradiction. We assume that

$$\mathbf{x}^T \mathbf{P} \mathbf{x} > 0 \quad \text{and} \quad -\mathbf{x}^T \mathbf{Q} \mathbf{x} < 0 \quad \text{for all } \mathbf{x} \neq \mathbf{0},$$

and that the matrix \mathbf{A} is not asymptotically stable. This means that the solution to $\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t)$, $\mathbf{x}_0 \neq \mathbf{0}$, which has the form $\mathbf{x}(t) = e^{\mathbf{A}(t-t_0)} \mathbf{x}_0$, is bounded away from the origin $\mathbf{x} = \mathbf{0}$. Thus, there is a constant $b > 0$ such that

$$\mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) > b > 0 \quad \text{for all } t > t_0.$$

Integrating $\frac{d}{dt}(\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)) = -\mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t)$ yields

$$\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t) = \mathbf{x}^T(t_0) \mathbf{P} \mathbf{x}(t_0) + \int_{t_0}^t (-\mathbf{x}^T(s) \mathbf{Q} \mathbf{x}(s)) ds.$$

Taking into account that $-\mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) < -b < 0$ for all $t > t_0$, we obtain for $t > t_0$,

$$\begin{aligned} \mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t) &= \mathbf{x}^T(t_0) \mathbf{P} \mathbf{x}(t_0) + \int_{t_0}^t (-\mathbf{x}^T(s) \mathbf{Q} \mathbf{x}(s)) ds \\ &< \mathbf{x}^T(t_0) \mathbf{P} \mathbf{x}(t_0) - b(t - t_0). \end{aligned}$$

The above implies that for a sufficiently large t we would have $\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t) < 0$, which contradicts the assumption that $\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t) > 0$. This means that \mathbf{A} must be asymptotically stable, and the proof is complete.

◆ Example 4.2

Let

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

and $\mathbf{Q} = \mathbf{I}_2$. We solve the Lyapunov equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{Q}$ for \mathbf{P} . Assume \mathbf{P} to have the form

$$\mathbf{P} = \begin{bmatrix} p_1 & p_2 \\ p_2 & p_3 \end{bmatrix} = \mathbf{P}^T.$$

Then, we have

$$\begin{aligned} \mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} &= \begin{bmatrix} 0 & -1 \\ 1 & -2 \end{bmatrix} \begin{bmatrix} p_1 & p_2 \\ p_2 & p_3 \end{bmatrix} + \begin{bmatrix} p_1 & p_2 \\ p_2 & p_3 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ -1 & -2 \end{bmatrix} \\ &= \begin{bmatrix} -2p_2 & -p_3 + p_1 - 2p_2 \\ -p_3 + p_1 - 2p_2 & 2p_2 - 4p_3 \end{bmatrix} \\ &= -\mathbf{Q} = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix}. \end{aligned}$$

Solving the corresponding linear equations yields

$$\mathbf{P} = \begin{bmatrix} 3/2 & 1/2 \\ 1/2 & 1/2 \end{bmatrix}.$$

The solution matrix \mathbf{P} is positive definite. Hence, the matrix \mathbf{A} must have its eigenvalues in the open left-half plane. Indeed, the eigenvalues of \mathbf{A} are both located at -1 .

Definition 4.10 A positive definite function $V(\mathbf{x})$ whose time derivative evaluated on the solutions of the system $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$ is negative definite (or negative semidefinite) is called a Lyapunov function for this system.

Lyapunov's theorem gives us a necessary and sufficient condition for asymptotic stability of linear time-invariant dynamical systems. To check if a given matrix \mathbf{A} is asymptotically stable or not, it is enough to take an arbitrary r.s.p.d. matrix \mathbf{Q} , solve the equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{Q}$ for real symmetric \mathbf{P} , and check if the obtained \mathbf{P} is positive definite, or not. If $\mathbf{P} > 0$, then \mathbf{A} is asymptotically stable. If \mathbf{P} is not positive definite, then \mathbf{A} cannot be asymptotically stable. Because \mathbf{Q} can be an arbitrary r.s.p.d. matrix, $\mathbf{Q} = \mathbf{I}_n$ (the $n \times n$ identity matrix) is a good choice. It is important to note that there would be no use, as the following example shows, in first choosing a matrix \mathbf{P} and then calculating \mathbf{Q} from the equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{Q}$.

◆ Example 4.3

Let

$$\mathbf{A} = \begin{bmatrix} -1 & 3 \\ 0 & -1 \end{bmatrix},$$

which is asymptotically stable. Let

$$\mathbf{P} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbf{I}_2 > 0.$$

Then,

$$\mathbf{Q} = -(\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A}) = \begin{bmatrix} 2 & -3 \\ -3 & 2 \end{bmatrix},$$

which is indefinite. Thus, unless \mathbf{Q} turns out to be definite, nothing can be inferred about asymptotic stability from the Lyapunov theorem if we start with \mathbf{P} and then calculate \mathbf{Q} .

We now give a geometric interpretation of the Lyapunov theorem. First, note that the function

$$\dot{V} = \frac{dV(\mathbf{x}(t))}{dt}$$

can be expressed, using the chain rule, as

$$\dot{V} = \nabla V^T \dot{\mathbf{x}} = \|\nabla V\| \|\dot{\mathbf{x}}\| \cos \theta,$$

where θ is the angle between the vectors ∇V and $\dot{\mathbf{x}}$. If $\|\dot{\mathbf{x}}\| = 1$, then $\nabla V^T \dot{\mathbf{x}}$ can be viewed as the rate of increase of V at \mathbf{x} in the direction $\dot{\mathbf{x}}$. Thus, $\dot{V} < 0$ if and only if $90^\circ < \theta < 270^\circ$. In such a case the velocity vector $\dot{\mathbf{x}}(t)$ points in the direction of decreasing V . We illustrate the above discussion in Figure 4.5.

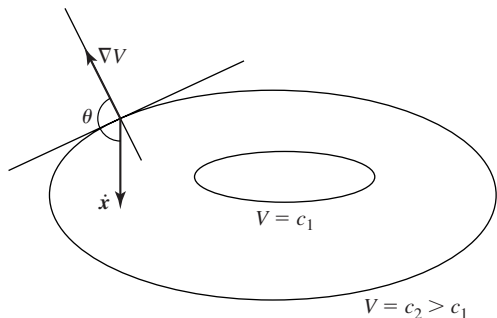


Figure 4.5 Geometric interpretation of the Lyapunov theorem.

We now prove a variant of Theorem 4.1, which is very useful in numerical verification of stability.

Theorem 4.2 A matrix \mathbf{A} is asymptotically stable if and only if for any $\mathbf{Q} = \mathbf{C}^T \mathbf{C}$, $\mathbf{C} \in \mathbb{R}^{p \times n}$, such that the pair (\mathbf{A}, \mathbf{C}) is observable, the solution \mathbf{P} to the Lyapunov matrix equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{C}^T \mathbf{C}$ is positive definite.

Proof We first prove the necessity condition (\Rightarrow) for asymptotic stability. If \mathbf{A} is asymptotically stable—that is, $\Re \lambda_i(\mathbf{A}) < 0$ —then we will show that

$$\mathbf{P} = \int_0^\infty e^{\mathbf{A}^T t} \mathbf{C}^T \mathbf{C} e^{\mathbf{A} t} dt \quad (4.7)$$

is the symmetric positive definite solution to the Lyapunov equation. Indeed, the integrand of (4.7) is a sum of terms of the form

$$t^k e^{\lambda_i t},$$

where λ_i is the i th eigenvalue of \mathbf{A} . Since $\Re \lambda_i(\mathbf{A}) < 0$, the integral in (4.7) exists because the terms in \mathbf{P} do not diverge when $t \rightarrow \infty$. Indeed, if we let $-\alpha_i = \Re \lambda_i$, where $\alpha_i > 0$, then

$$\lim_{t \rightarrow \infty} t^k e^{-\alpha_i t} = \lim_{t \rightarrow \infty} \frac{t^k}{e^{\alpha_i t}},$$

We differentiate the numerator and denominator of the above k times and use L'Hôpital's rule to get

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{t^k}{e^{\alpha_i t}} &= \lim_{t \rightarrow \infty} \frac{k t^{k-1}}{\alpha_i e^{\alpha_i t}} \\ &= \lim_{t \rightarrow \infty} \frac{k(k-1) t^{k-2}}{\alpha_i^2 e^{\alpha_i t}} \\ &\vdots \\ &= \lim_{t \rightarrow \infty} \frac{k!}{\alpha_i^k e^{\alpha_i t}} \\ &= 0. \end{aligned} \quad (4.8)$$

Thus, by (4.8), the formula for \mathbf{P} given by (4.7) is well-defined. Next, note that $\mathbf{P} = \mathbf{P}^T$. The matrix \mathbf{P} is also positive definite because the pair (\mathbf{A}, \mathbf{C}) is observable. Substituting the expression for \mathbf{P} into the Lyapunov equation and performing manipulations yields

$$\begin{aligned} \mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} &= \int_0^\infty \mathbf{A}^T e^{\mathbf{A}^T t} \mathbf{C}^T \mathbf{C} e^{\mathbf{A} t} dt + \int_0^\infty e^{\mathbf{A}^T t} \mathbf{C}^T \mathbf{C} e^{\mathbf{A} t} \mathbf{A} dt \\ &= \int_0^\infty \frac{d}{dt} (e^{\mathbf{A}^T t} \mathbf{C}^T \mathbf{C} e^{\mathbf{A} t}) dt \\ &= e^{\mathbf{A}^T t} \mathbf{C}^T \mathbf{C} e^{\mathbf{A} t} \Big|_0^\infty \\ &= -\mathbf{C}^T \mathbf{C}. \end{aligned}$$

The sufficiency part (\Leftarrow) can be proven using the arguments of the proof of the sufficiency part of the Lyapunov theorem and noting that the observability of the pair (\mathbf{A}, \mathbf{C}) implies that $\mathbf{x}^T(t) \mathbf{C}^T \mathbf{C} \mathbf{x}(t)$ is not identically zero along any nonzero solution of $\dot{\mathbf{x}}(t) = \mathbf{A} \mathbf{x}(t)$.

We now show that if \mathbf{A} is asymptotically stable, then the Lyapunov equation (4.6) has a unique solution for every \mathbf{Q} . Uniqueness can be proven in many ways. We use the Kronecker product in our proof. Using the results of Section A.5, we represent the Lyapunov equation (4.6) as

$$(\mathbf{I}_n \otimes \mathbf{A}^T + \mathbf{A}^T \otimes \mathbf{I}_n) s(\mathbf{P}) = -s(\mathbf{Q}). \quad (4.9)$$

The condition for a solution $s(\mathbf{P})$ of the above equation to be unique is

$$\lambda_i(\mathbf{A}^T) + \lambda_j(\mathbf{A}) \neq 0, \quad i, j = 1, 2, \dots, n.$$

Because $\lambda_i(\mathbf{A}^T) = \lambda_i(\mathbf{A})$, and \mathbf{A} is asymptotically stable by assumption, the condition $\lambda_i(\mathbf{A}^T) + \lambda_j(\mathbf{A}) \neq 0$ is met, and hence the Lyapunov equation has a unique solution.

4.4 Evaluating Quadratic Indices

In this section we apply the Lyapunov theory to evaluate the quadratic indices

$$J_r = \int_0^\infty t^r \mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) dt, \quad r = 0, 1, \dots$$

subject to

$$\dot{\mathbf{x}}(t) = \mathbf{A} \mathbf{x}(t), \quad \mathbf{x}(0) = \mathbf{x}_0,$$

where $\mathbf{Q} = \mathbf{C}^T \mathbf{C}$ such that the pair (\mathbf{A}, \mathbf{C}) is observable, and the matrix \mathbf{A} is asymptotically stable. We can reformulate the above problem as follows. Evaluate the quadratic performance indices

$$J_r = \int_0^\infty t^r \mathbf{y}^T(t) \mathbf{y}(t) dt, \quad r = 0, 1, \dots$$

subject to

$$\begin{aligned}\dot{\mathbf{x}}(t) &= \mathbf{A}\mathbf{x}(t), & \mathbf{x}(0) &= \mathbf{x}_0, \\ \mathbf{y}(t) &= \mathbf{C}\mathbf{x}(t),\end{aligned}$$

where the pair (\mathbf{A}, \mathbf{C}) is observable, and the matrix \mathbf{A} is asymptotically stable.

To evaluate

$$J_0 = \int_0^\infty \mathbf{y}^T(t) \mathbf{y}(t) dt = \int_0^\infty \mathbf{x}^T(t) \mathbf{C}^T \mathbf{C} \mathbf{x}(t) dt = \int_0^\infty \mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) dt,$$

we note that

$$\frac{d}{dt}(\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)) = -\mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t), \quad (4.10)$$

where \mathbf{P} and $\mathbf{Q} = \mathbf{C}^T \mathbf{C}$ satisfy the Lyapunov matrix equation, $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{Q}$. Hence,

$$J_0 = - \int_0^\infty \frac{d}{dt}(\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)) dt. \quad (4.11)$$

Integrating both sides of (4.10) with respect to t yields

$$\begin{aligned}J_0 &= - \int_0^\infty \frac{d}{dt}(\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)) dt \\ &= - (\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)) \Big|_0^\infty \\ &= \lim_{t \rightarrow \infty} (-\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)) + \mathbf{x}^T(0) \mathbf{P} \mathbf{x}(0) \\ &= \mathbf{x}^T(0) \mathbf{P} \mathbf{x}(0)\end{aligned}$$

because, by assumption, \mathbf{A} is asymptotically stable and therefore $\lim_{t \rightarrow \infty} \mathbf{x}(t) = \mathbf{0}$ for all $\mathbf{x}(0)$. Furthermore, $\mathbf{P} = \mathbf{P}^T > 0$, and thus $J_0 > 0$ for all $\mathbf{x}(0) \neq \mathbf{0}$.

Repeating the above arguments leads to a similar expression for J_1 . First observe that

$$\begin{aligned}\frac{d}{dt}(t \mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)) &= \mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t) + t \frac{d}{dt}(\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)) \\ &= \mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t) - t \mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t),\end{aligned}$$

where \mathbf{P} satisfies the Lyapunov matrix equation, $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{Q}$. Hence,

$$\begin{aligned}J_1 &= \int_0^\infty t \mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) dt \\ &= \int_0^\infty \mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t) dt - \int_0^\infty \frac{d}{dt}(t \mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)) dt.\end{aligned} \quad (4.12)$$

Because $\mathbf{P} = \mathbf{P}^T > 0$, we can solve

$$\mathbf{A}^T \mathbf{P}_1 + \mathbf{P}_1 \mathbf{A} = -\mathbf{P}$$

for $\mathbf{P}_1 = \mathbf{P}_1^T > 0$. The above equation can also be represented as

$$\frac{d}{dt}(\mathbf{x}^T(t)\mathbf{P}_1\mathbf{x}(t)) = -\mathbf{x}^T(t)\mathbf{P}\mathbf{x}(t). \quad (4.13)$$

Substituting (4.13) into (4.12) yields

$$\begin{aligned} J_1 &= - \int_0^\infty \frac{d}{dt}(\mathbf{x}^T(t)\mathbf{P}_1\mathbf{x}(t))dt - (\mathbf{x}^T(t)\mathbf{P}\mathbf{x}(t)) \Big|_0^\infty \\ &= \lim_{t \rightarrow \infty} (-\mathbf{x}^T(t)\mathbf{P}_1\mathbf{x}(t)) + \mathbf{x}^T(0)\mathbf{P}_1\mathbf{x}(0) - 0 + 0 \\ &= \mathbf{x}^T(0)\mathbf{P}_1\mathbf{x}(0). \end{aligned}$$

Using similar arguments, we can evaluate J_r , $r > 1$, which involves successive solutions to the Lyapunov matrix equations.

In many control problems a dynamical system model is given in terms of a transfer function rather than in terms of a system of differential equations. We now discuss a method of converting a given rational function into a system of first-order differential equations. After this process is completed, we use the Lyapunov theory to evaluate a performance index of interest. Suppose we are given a rational function of the form

$$E(s) = \frac{b_{n-1}s^{n-1} + \cdots + b_1s + b_0}{s^n + a_{n-1}s^{n-1} + \cdots + a_1s + a_0},$$

where $E(s) = \mathcal{L}(e(t))$ is the Laplace transform of $e(t)$. We now show that $e(t)$ can be generated as the solution of the differential equation

$$\frac{d^n e(t)}{dt^n} + a_{n-1} \frac{d^{n-1} e(t)}{dt^{n-1}} + \cdots + a_1 \frac{de(t)}{dt} + a_0 e(t) = 0, \quad (4.14)$$

subject to a set of initial conditions that we shall derive. We follow the method of Power and Simpson [237, pp. 220–221]. We first take the Laplace transform of the differential equation and then rearrange the resulting terms in an appropriate matrix equation. Recall that

$$\mathcal{L}\left(\frac{d^i e(t)}{dt^i}\right) = s^i E(s) - s^{i-1}e(0) - s^{i-2} \frac{de(0)}{dt} - \cdots - \frac{d^{i-1}e(0)}{dt^{i-1}}. \quad (4.15)$$

Applying the above to (4.14) and comparing corresponding terms with $E(s)$ leads to the equation for the initial conditions:

$$\begin{bmatrix} a_1 & a_2 & \cdots & a_{n-1} & 1 \\ a_2 & a_3 & \cdots & 1 & 0 \\ \vdots & \vdots & & \vdots & \vdots \\ a_{n-1} & 1 & \cdots & 0 & 0 \\ 1 & 0 & \cdots & 0 & 0 \end{bmatrix} \begin{bmatrix} e(0) \\ \frac{de(0)}{dt} \\ \vdots \\ \frac{d^{n-2}e(0)}{dt^{n-2}} \\ \frac{d^{n-1}e(0)}{dt^{n-1}} \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_{n-2} \\ b_{n-1} \end{bmatrix}.$$

The coefficient matrix in the above equation is nonsingular. Thus, there is a unique solution to this equation, that is, we can uniquely determine the initial conditions. We next define the state vector

$$\mathbf{x} = \begin{bmatrix} e & \frac{de}{dt} & \cdots & \frac{d^{n-1}e}{dt^{n-1}} \end{bmatrix}^T,$$

and represent the differential equation (4.14) as

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x},$$

where

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ -a_0 & -a_1 & -a_2 & \cdots & -a_{n-2} & -a_{n-1} \end{bmatrix}.$$

The vector of initial conditions is

$$\mathbf{x}(0) = \begin{bmatrix} e(0) \\ \frac{de(0)}{dt} \\ \vdots \\ \frac{d^{n-2}e(0)}{dt^{n-2}} \\ \frac{d^{n-1}e(0)}{dt^{n-1}} \end{bmatrix} = \begin{bmatrix} a_1 & a_2 & \cdots & a_{n-1} & 1 \\ a_2 & a_3 & \cdots & 1 & 0 \\ \vdots & \vdots & & \vdots & \vdots \\ a_{n-1} & 1 & \cdots & 0 & 0 \\ 1 & 0 & \cdots & 0 & 0 \end{bmatrix}^{-1} \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_{n-2} \\ b_{n-1} \end{bmatrix}.$$

Note that $e = [1 \ 0 \ \cdots \ 0 \ 0]\mathbf{x}$.

The above procedure for converting a rational function description of a process into its equivalent state-space representation shows a relation between the zeros of the given rational function and the initial conditions of the associated state-space realization.

We will now use the above method of constructing a state-space realization of a given rational function in the following example that involves evaluation of a quadratic index.

◆ Example 4.4

Consider a block diagram of a dynamical control system model shown in Figure 4.6. This system comes from Dorf and Bishop [66, p. 257]. The values of the parameters are $K_1 = 0.5$ and $K_1 K_2 K_p = 2.5$. We shall use the Lyapunov theory to select the gain K_3 that minimizes the effect of the disturbance $u(t) = \mathcal{L}^{-1}(U(s))$ in the sense of the integral of the square error (ISE) criterion. In other words, we wish to find K_3 that leads to minimization of the performance index

$$J_0 = \int_0^\infty c^2(t) dt,$$

where $c(t) = \mathcal{L}^{-1}(C(s))$. We assume a unit step disturbance and $R(s) = 0$.

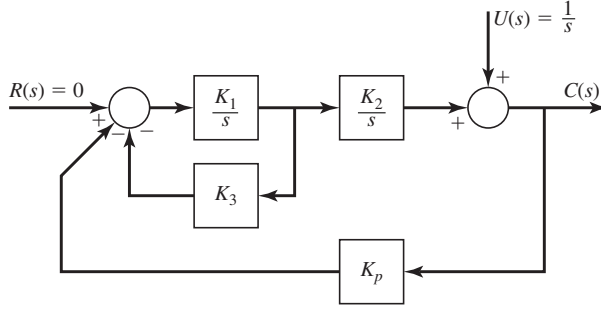


Figure 4.6 Block diagram of the optimized control system model of Example 4.4.

We first find a rational function that models $C(s)$. From the block diagram shown in Figure 4.6, we obtain

$$C(s) = U(s) - \frac{K_2}{s} \frac{K_1}{s + K_1 K_3} K_p C(s).$$

Hence

$$C(s) = \frac{s + 0.5K_3}{s^2 + 0.5K_3s + 2.5}.$$

The next step is to determine the state-space representation corresponding to $C(s)$. We use the method given before this example. We first form the differential equation for $c(t)$:

$$\frac{d^2 c(t)}{dt^2} + 0.5K_3 \frac{dc(t)}{dt} + 2.5c(t) = 0$$

with initial conditions

$$\begin{bmatrix} 0.5K_3 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} c(0) \\ \frac{dc(0)}{dt} \end{bmatrix} = \begin{bmatrix} 0.5K_3 \\ 1 \end{bmatrix}.$$

Solving the above equation yields

$$\begin{bmatrix} c(0) \\ \frac{dc(0)}{dt} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

Let $x_1 = c$, $x_2 = \frac{dc}{dt}$. The differential equation in vector-matrix form corresponding to $C(s)$ is

$$\dot{\mathbf{x}}(t) = \begin{bmatrix} 0 & 1 \\ -2.5 & -0.5K_3 \end{bmatrix} \mathbf{x}(t), \quad \mathbf{x}(0) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$c(t) = \mathbf{c}\mathbf{x}(t),$$

where $\mathbf{c} = [1 \ 0]$. We are now ready to evaluate the performance index J_0 :

$$J_0 = \int_0^\infty c^2(t) dt = \int_0^\infty \mathbf{x}^T(t) \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \mathbf{x}(t) dt = \int_0^\infty \mathbf{x}^T(t) \mathbf{c}^T \mathbf{c} \mathbf{x}(t) dt.$$

The pair (\mathbf{A}, \mathbf{c}) is observable. Therefore, we can solve the Lyapunov equation

$$\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{c}^T \mathbf{c}$$

for \mathbf{P} to obtain

$$\mathbf{P} = \begin{bmatrix} \frac{K_3}{10} + \frac{1}{K_3} & \frac{1}{5} \\ \frac{1}{5} & \frac{2}{5K_3} \end{bmatrix}.$$

Hence

$$J_0 = \mathbf{x}^T(0) \mathbf{P} \mathbf{x}(0) = \frac{K_3}{10} + \frac{1}{K_3}.$$

Differentiating J_0 and equating the result to zero gives

$$\frac{dJ_0}{dK_3} = \frac{1}{10} - \frac{1}{K_3^2} = 0.$$

Selecting

$$K_3 = \sqrt{10}$$

yields a positive definite \mathbf{P} . Note that because

$$\left. \frac{d^2 J_0}{dK_3^2} \right|_{K_3=\sqrt{10}} = \left. \frac{2}{K_3^3} \right|_{K_3=\sqrt{10}} > 0,$$

that is, the second-order sufficiency condition for the minimum is satisfied, we conclude that $K_3 = \sqrt{10} \approx 3.16$ minimizes J_0 .

We now wish to determine the gain K_3 that leads to minimization of the effect of the disturbance $u(t)$ in the sense of the integral of time multiplied by the squared error (ITSE) criterion. That is, we wish to select K_3 to minimize the performance index:

$$J_1 = \int_0^\infty t c^2(t) dt.$$

We first solve the Lyapunov equation,

$$\mathbf{A}^T \mathbf{P}_1 + \mathbf{P}_1 \mathbf{A} = -\mathbf{P},$$

to get

$$\mathbf{P}_1 = \begin{bmatrix} \frac{2}{K_3^2} + \frac{K_3^2}{100} & \frac{1}{5} \left(\frac{1}{K_3} + \frac{K_3}{10} \right) \\ \frac{1}{5} \left(\frac{1}{K_3} + \frac{K_3}{10} \right) & \frac{4}{5K_3^2} + \frac{1}{25} \end{bmatrix}.$$

Hence

$$J_1 = \mathbf{x}^T(0) \mathbf{P}_1 \mathbf{x}(0) = \frac{2}{K_3^2} + \frac{K_3^2}{100}.$$

Differentiating J_1 and equating the result to zero gives

$$\frac{dJ_1}{dK_3} = -\frac{4}{K_3^3} - \frac{K_3}{50} = 0.$$

We select $K_3 = \sqrt[4]{200}$ to keep \mathbf{P}_1 positive definite. For this choice of K_3 the second-order sufficiency condition for the minimum of J_1 is satisfied. Thus, $K_3 \approx 3.76$ minimizes J_1 .

4.5 Discrete-Time Lyapunov Equation

We will now use the Lyapunov approach to the stability analysis of the discrete-time linear time-invariant systems modeled by

$$\mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k), \quad \mathbf{x}(k_0) = \mathbf{x}_0. \quad (4.16)$$

By analogy with the continuous-time case, consider a quadratic positive definite form,

$$V(\mathbf{x}(k)) = \mathbf{x}^T(k) \mathbf{P} \mathbf{x}(k).$$

We evaluate the first difference of V , defined as $\Delta V = V(\mathbf{x}(k+1)) - V(\mathbf{x}(k))$, on the trajectories of (4.16) to obtain

$$\begin{aligned} \Delta V(\mathbf{x}(k)) &= V(\mathbf{x}(k+1)) - V(\mathbf{x}(k)) \\ &= \mathbf{x}^T(k+1) \mathbf{P} \mathbf{x}(k+1) - \mathbf{x}^T(k) \mathbf{P} \mathbf{x}(k) \\ &= \mathbf{x}^T(k) \mathbf{A}^T \mathbf{P} \mathbf{A} \mathbf{x}(k) - \mathbf{x}^T(k) \mathbf{P} \mathbf{x}(k) \\ &= \mathbf{x}^T(k) (\mathbf{A}^T \mathbf{P} \mathbf{A} - \mathbf{P}) \mathbf{x}(k), \end{aligned} \quad (4.17)$$

where we replaced $\mathbf{x}(k+1)$ in the above with $\mathbf{A}\mathbf{x}(k)$. For the solutions of (4.16) to decrease the “energy function” V , at each $k > k_0$, it is necessary and sufficient that $\Delta V(\mathbf{x}(k)) < 0$, or, equivalently, that for some positive definite \mathbf{Q} ,

$$\Delta V(\mathbf{x}(k)) = -\mathbf{x}^T(k) \mathbf{Q} \mathbf{x}(k). \quad (4.18)$$

Combining (4.16) and (4.18), we obtain the so-called *discrete-time Lyapunov equation*,

$$\boxed{\mathbf{A}^T \mathbf{P} \mathbf{A} - \mathbf{P} = -\mathbf{Q}.} \quad (4.19)$$

Theorem 4.3 The matrix \mathbf{A} has its eigenvalues in the open unit disk if and only if for any positive definite \mathbf{Q} the solution \mathbf{P} to (4.19) is positive definite.

Proof (\Rightarrow) We can use the Kronecker product, and in particular relation (A.62), to represent the discrete-time Lyapunov equation as

$$(\mathbf{I}_{n^2} - \mathbf{A}^T \otimes \mathbf{A}^T) \mathbf{s}(\mathbf{P}) = \mathbf{s}(\mathbf{Q}), \quad (4.20)$$

where \otimes denotes the Kronecker product and $\mathbf{s}(\cdot)$ is the stacking operator. It is easy to see that the above equation has a unique solution if and only if

$$1 - \lambda_i \lambda_j \neq 0 \quad \text{for all } i, j = 1, 2, \dots, n,$$

where $\{\lambda_i\}$ is the set of eigenvalues of \mathbf{A} . By assumption, $|\lambda_i| < 1$, which is a sufficient condition for (4.20) and hence (4.19) to have a unique solution. The solution \mathbf{P} can be expressed as

$$\mathbf{P} = \sum_{k=0}^{\infty} (\mathbf{A}^T)^k \mathbf{Q} \mathbf{A}^k. \quad (4.21)$$

The above expression is well-defined because all eigenvalues of \mathbf{A} have magnitude (strictly) less than unity. Substituting (4.21) into (4.19) gives

$$\begin{aligned} \mathbf{A}^T \mathbf{P} \mathbf{A} - \mathbf{P} &= \mathbf{A}^T \left(\sum_{k=0}^{\infty} (\mathbf{A}^T)^k \mathbf{Q} \mathbf{A}^k \right) \mathbf{A} - \sum_{k=0}^{\infty} (\mathbf{A}^T)^k \mathbf{Q} \mathbf{A}^k \\ &= \sum_{k=1}^{\infty} (\mathbf{A}^T)^k \mathbf{Q} \mathbf{A}^k - \sum_{k=1}^{\infty} (\mathbf{A}^T)^k \mathbf{Q} \mathbf{A}^k - \mathbf{Q} \\ &= -\mathbf{Q}. \end{aligned}$$

Thus, indeed (4.21) satisfies the discrete-time Lyapunov equation. Note that if \mathbf{Q} is symmetric, so is \mathbf{P} . If \mathbf{Q} is positive definite, so is \mathbf{P} .

(\Leftarrow) Let λ_i be an eigenvalue of \mathbf{A} and let \mathbf{v}_i be a corresponding eigenvector; that is, $\mathbf{A} \mathbf{v}_i = \lambda_i \mathbf{v}_i$. Observe that λ_i and \mathbf{v}_i may be complex. We denote by \mathbf{v}_i^* the complex conjugate transpose of \mathbf{v}_i . Note that $\mathbf{v}_i^* \mathbf{P} \mathbf{v}_i$ is a scalar. We will show that its complex conjugate equals itself, which means that $\mathbf{v}_i^* \mathbf{P} \mathbf{v}_i$ is real. Indeed,

$$\overline{\mathbf{v}_i^* \mathbf{P} \mathbf{v}_i} = \mathbf{v}_i^* \mathbf{P}^* \mathbf{v}_i = \mathbf{v}_i^* \mathbf{P}^T \mathbf{v}_i = \mathbf{v}_i^* \mathbf{P} \mathbf{v}_i,$$

where we used the fact that the complex conjugate transpose of a real matrix reduces to its transpose. Thus, $\mathbf{v}_i^* \mathbf{P} \mathbf{v}_i$ is real for any complex \mathbf{v}_i .

Premultiplying now (4.19) by \mathbf{v}_i^* and postmultiplying by \mathbf{v}_i gives

$$\mathbf{v}_i^* (\mathbf{A}^T \mathbf{P} \mathbf{A} - \mathbf{P}) \mathbf{v}_i = -\mathbf{v}_i^* \mathbf{Q} \mathbf{v}_i. \quad (4.22)$$

The left-hand side of the above equality evaluates to

$$\begin{aligned} \mathbf{v}_i^* \mathbf{A}^T \mathbf{P} \mathbf{A} \mathbf{v}_i - \mathbf{v}_i^* \mathbf{P} \mathbf{v}_i &= \bar{\lambda}_i \mathbf{v}_i^* \mathbf{P} \mathbf{v}_i \lambda_i - \mathbf{v}_i^* \mathbf{P} \mathbf{v}_i \\ &= (|\lambda_i|^2 - 1) \mathbf{v}_i^* \mathbf{P} \mathbf{v}_i. \end{aligned}$$

Comparing the above with the right-hand side of (4.22) yields

$$(1 - |\lambda_i|^2) \mathbf{v}_i^* \mathbf{P} \mathbf{v}_i = \mathbf{v}_i^* \mathbf{Q} \mathbf{v}_i > 0$$

because \mathbf{Q} is positive definite. Because \mathbf{P} is positive definite, we have to have $|\lambda_i| < 1$. The proof is complete.

We will now concern ourselves with evaluating performance indices on the trajectories of (4.16). Specifically, we wish to evaluate

$$K_r = \sum_{k=0}^{\infty} k^r \mathbf{x}^T(k) \mathbf{Q} \mathbf{x}(k), \quad r = 0, 1, \dots \quad (4.23)$$

where $\mathbf{Q} = \mathbf{Q}^T > 0$ subject to (4.16). We assume that the matrix \mathbf{A} has its eigenvalues in the open unit disk; that is, the matrix \mathbf{A} is *convergent*. We first consider the case when $k = 0$. For this case, K_r becomes

$$K_0 = \sum_{k=0}^{\infty} \mathbf{x}^T(k) \mathbf{Q} \mathbf{x}(k). \quad (4.24)$$

Because the matrix \mathbf{A} is convergent, the solution \mathbf{P} to (4.19) is positive definite. Using (4.17), (4.18), and (4.16), we write

$$\begin{aligned} \mathbf{x}^T(k) \mathbf{Q} \mathbf{x}(k) &= -\mathbf{x}^T(k) (\mathbf{A}^T \mathbf{P} \mathbf{A} - \mathbf{P}) \mathbf{x}(k) \\ &= \mathbf{x}^T(k) \mathbf{P} \mathbf{x}(k) - \mathbf{x}^T(k+1) \mathbf{P} \mathbf{x}(k+1). \end{aligned} \quad (4.25)$$

Summing both sides of the above from $k = 0$ to $k = \infty$ and using the fact that $\mathbf{x}(\infty) = \mathbf{0}$ because the matrix \mathbf{A} is convergent, we obtain

$$K_0 = \mathbf{x}^T(0) \mathbf{P} \mathbf{x}(0). \quad (4.26)$$

We will now evaluate

$$K_1 = \sum_{k=0}^{\infty} k \mathbf{x}^T(k) \mathbf{Q} \mathbf{x}(k),$$

subject to (4.16). To proceed, we multiply both sides of (4.25) by k to get

$$k \mathbf{x}^T(k) \mathbf{Q} \mathbf{x}(k) = k \mathbf{x}^T(k) \mathbf{P} \mathbf{x}(k) - k \mathbf{x}^T(k+1) \mathbf{P} \mathbf{x}(k+1).$$

Summing the above from $k = 0$ to $k = \infty$ gives

$$\begin{aligned} K_1 &= \mathbf{x}^T(1) \mathbf{P} \mathbf{x}(1) - \mathbf{x}^T(2) \mathbf{P} \mathbf{x}(2) + 2\mathbf{x}^T(2) \mathbf{P} \mathbf{x}(2) - 2\mathbf{x}^T(3) \mathbf{P} \mathbf{x}(3) + \dots \\ &= \sum_{k=0}^{\infty} \mathbf{x}^T(k) \mathbf{P} \mathbf{x}(k) - \mathbf{x}^T(0) \mathbf{P} \mathbf{x}(0). \end{aligned} \quad (4.27)$$

Using the method we employed to find the value of K_0 , we now evaluate the first term on the right-hand side of (4.27) to get

$$\sum_{k=0}^{\infty} \mathbf{x}^T(k) \mathbf{P} \mathbf{x}(k) = \mathbf{x}^T(0) \mathbf{P}_1 \mathbf{x}(0),$$

where \mathbf{P}_1 satisfies

$$\mathbf{A}^T \mathbf{P}_1 \mathbf{A} - \mathbf{P}_1 = -\mathbf{P}.$$

Taking the above into account, we represent (4.27) as

$$K_1 = \mathbf{x}^T(0) (\mathbf{P}_1 - \mathbf{P}) \mathbf{x}(0).$$

Proceeding in a similar manner, we can find expressions for K_r , where $r > 1$.

4.6 Constructing Robust Linear Controllers

In this section we use the Lyapunov theorem to construct robust state-feedback controllers. One of the basic issues in the control of dynamical systems is the effect of uncertainties, or neglected nonlinearities, on the stability of the closed-loop system. A controller is viewed as robust if it maintains stability for all uncertainties in an expected range. We consider a class of dynamical systems modeled by

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) + \mathbf{f}(t, \mathbf{x}(t)) + \mathbf{B}\mathbf{h}(t, \mathbf{x}(t), \mathbf{u}(t)), \quad (4.28)$$

where $\mathbf{x}(t) \in \mathbb{R}^n$, $\mathbf{u}(t) \in \mathbb{R}^m$, and the functions \mathbf{h} and \mathbf{f} model uncertainties, or nonlinearities, in the system. We refer to \mathbf{h} as the *matched uncertainty* because it affects the system dynamics via the input matrix \mathbf{B} in the same fashion as the input \mathbf{u} does. In other words, the uncertainty \mathbf{h} matches the system input \mathbf{u} . The vector \mathbf{f} models the *unmatched uncertainty*. We assume that the uncertain elements \mathbf{h} and \mathbf{f} satisfy the following norm bounds:

1. $\|\mathbf{h}(t, \mathbf{x}, \mathbf{u})\| \leq \gamma_h \|\mathbf{u}\| + \alpha_h \|\mathbf{x}\|$,
2. $\|\mathbf{f}(t, \mathbf{x})\| \leq \alpha_f \|\mathbf{x}\|$,

where γ_h , α_h , and α_f are known nonnegative constants. We further assume that the matrix \mathbf{A} is asymptotically stable. If this is not the case, we apply a preliminary state-feedback controller $\mathbf{u} = -\mathbf{K}\mathbf{x}$ such that $\mathbf{A} - \mathbf{B}\mathbf{K}$ is asymptotically stable. Such a feedback exists provided that the pair (\mathbf{A}, \mathbf{B}) is stabilizable. Our goal is to construct a linear state-feedback controller that would make the closed-loop system asymptotically stable for arbitrary \mathbf{f} and \mathbf{h} that satisfy the above norm bounds.

Theorem 4.4 Suppose that \mathbf{A} is asymptotically stable and that $\mathbf{P} = \mathbf{P}^T > 0$ is the solution to the Lyapunov matrix equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -2\mathbf{Q}$ for some $\mathbf{Q} = \mathbf{Q}^T > 0$. Suppose also that

$$\alpha_f < \frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{P})} \quad \text{and} \quad \gamma_h < 1.$$

Then, the state-feedback controller $\mathbf{u} = -\gamma \mathbf{B}^T \mathbf{P} \mathbf{x}$, where

$$\gamma > \frac{\alpha_h^2}{4(\lambda_{\min}(\mathbf{Q}) - \alpha_f \lambda_{\max}(\mathbf{P}))(1 - \gamma_h)}$$

stabilizes the uncertain system model (4.28) for arbitrary \mathbf{f} and \mathbf{h} that satisfy the norm bounds.

Proof The time derivative of the positive definite function $V = \mathbf{x}^T \mathbf{P} \mathbf{x}$ evaluated on any trajectory of the closed-loop system is

$$\begin{aligned} \dot{V} &= 2\mathbf{x}^T \mathbf{P} \dot{\mathbf{x}} \\ &= -2\mathbf{x}^T \mathbf{Q} \mathbf{x} - 2\gamma \mathbf{x}^T \mathbf{P} \mathbf{B} \mathbf{B}^T \mathbf{P} \mathbf{x} + 2\mathbf{x}^T \mathbf{P} \mathbf{f} + 2\mathbf{x}^T \mathbf{P} \mathbf{B} \mathbf{h}. \end{aligned}$$

Our goal is to determine $\tilde{\gamma} > 0$ such that if $\gamma > \tilde{\gamma}$, then $\dot{V} < 0$, which in turn implies the asymptotic stability of the closed-loop system. To proceed further, recall that if \mathbf{Q} is a symmetric matrix, then

$$\lambda_{\min}(\mathbf{Q}) \|\mathbf{x}\|^2 \leq \mathbf{x}^T \mathbf{Q} \mathbf{x} \leq \lambda_{\max}(\mathbf{Q}) \|\mathbf{x}\|^2,$$

and therefore

$$-\mathbf{x}^T \mathbf{Q} \mathbf{x} \leq -\lambda_{\min}(\mathbf{Q}) \|\mathbf{x}\|^2.$$

Furthermore,

$$\mathbf{x}^T \mathbf{P} \mathbf{B} \mathbf{B}^T \mathbf{P} \mathbf{x} = \mathbf{x}^T \mathbf{P} \mathbf{B} (\mathbf{x}^T \mathbf{P} \mathbf{B})^T = \|\mathbf{x}^T \mathbf{P} \mathbf{B}\|^2.$$

For a symmetric positive definite matrix \mathbf{P} , its induced 2-norm is

$$\|\mathbf{P}\| = \lambda_{\max}(\mathbf{P}).$$

Taking the above relations into account, using the norm bounds on the uncertain elements, and performing some manipulations yields

$$\begin{aligned} \dot{V} &\leq -2\lambda_{\min}(\mathbf{Q}) \|\mathbf{x}\|^2 - 2\gamma \|\mathbf{x}^T \mathbf{P} \mathbf{B}\|^2 + 2\alpha_f \lambda_{\max}(\mathbf{P}) \|\mathbf{x}\|^2 \\ &\quad + 2\|\mathbf{x}^T \mathbf{P} \mathbf{B}\|(\alpha_h \|\mathbf{x}\| + \gamma_h \gamma \|\mathbf{x}^T \mathbf{P} \mathbf{B}\|). \end{aligned} \quad (4.29)$$

Let

$$\tilde{\mathbf{Q}} = \begin{bmatrix} \lambda_{\min}(\mathbf{Q}) - \alpha_f \lambda_{\max}(\mathbf{P}) & -\frac{\alpha_h}{2} \\ -\frac{\alpha_h}{2} & \gamma(1 - \gamma_h) \end{bmatrix} \in \mathbb{R}^{2 \times 2}. \quad (4.30)$$

Then, we can represent (4.29) as

$$\dot{V} \leq -2 \begin{bmatrix} \|\mathbf{x}\| & \|\mathbf{x}^T \mathbf{P} \mathbf{B}\| \end{bmatrix} \tilde{\mathbf{Q}} \begin{bmatrix} \|\mathbf{x}\| \\ \|\mathbf{x}^T \mathbf{P} \mathbf{B}\| \end{bmatrix}. \quad (4.31)$$

By Theorem A.6 on page 634, for \dot{V} to be negative it is enough that the leading principal minors of the matrix $\tilde{\mathbf{Q}}$ given by (4.30) be positive. The first-order leading principal minor is positive by assumption. For the second-order leading principal minor to be positive it is sufficient that $\gamma_h < 1$ and

$$\gamma > \frac{\alpha_h^2}{4(\lambda_{\min}(\mathbf{Q}) - \alpha_f \lambda_{\max}(\mathbf{P}))(1 - \gamma_h)}.$$

The proof is now complete.

Note that if $\gamma_h = 0$, then the closed-loop system is asymptotically stable if

$$\alpha_f < \frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{P})} \quad \text{and} \quad \gamma > \frac{\alpha_h^2}{4(\lambda_{\min}(\mathbf{Q}) - \alpha_f \lambda_{\max}(\mathbf{P}))}.$$

◆ Example 4.5

Consider a dynamical system model

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b}(u + h(\mathbf{x}, u)) = \begin{bmatrix} -2 & 0 \\ 0 & -1 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} (u + h(\mathbf{x}, u)).$$

Let \mathbf{P} be the solution of the matrix Lyapunov equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -2\mathbf{I}_2$. We will design a linear state-feedback controller $u = -\gamma \mathbf{b}^T \mathbf{P} \mathbf{x}$ so that the origin of the closed-loop system is uniformly asymptotically stable for any $h(\mathbf{x}, u)$ such that

$$|h(\mathbf{x}, u)| \leq 4\|\mathbf{x}\| + \frac{1}{2}|u|. \quad (4.32)$$

In particular, we will find a bound $\tilde{\gamma} > 0$ so that for any $\gamma > \tilde{\gamma}$ the closed-loop system is uniformly asymptotically stable in the face of uncertainties satisfying the norm bound (4.32).

Solving the Lyapunov equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -2\mathbf{I}_2$ yields

$$\mathbf{P} = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix}.$$

Hence,

$$u = -\gamma \mathbf{b}^T \mathbf{P} \mathbf{x} = -\gamma \left(\frac{1}{2}x_1 + x_2 \right).$$

We obtain a bound on γ by evaluating $\frac{d}{dt}(\mathbf{x}^T \mathbf{P} \mathbf{x})$ on the trajectories of the closed-loop system. First, we evaluate

$$\begin{aligned} \frac{d}{dt}(\mathbf{x}^T \mathbf{P} \mathbf{x}) &= 2\mathbf{x}^T \mathbf{P} \dot{\mathbf{x}} \\ &= 2\mathbf{x}^T \mathbf{P} (\mathbf{A}\mathbf{x} + \mathbf{b}(u + h(\mathbf{x}, u))) \\ &= -2\mathbf{x}^T \mathbf{x} - 2\gamma \mathbf{x}^T \mathbf{P} \mathbf{b} \mathbf{b}^T \mathbf{P} \mathbf{x} + 2\mathbf{x}^T \mathbf{P} \mathbf{b} h. \end{aligned}$$

Taking the norms and performing manipulations yields

$$\begin{aligned} \frac{d}{dt}(\mathbf{x}^T \mathbf{P} \mathbf{x}) &\leq -2\mathbf{x}^T \mathbf{x} - 2\gamma |\mathbf{x}^T \mathbf{P} \mathbf{b}|^2 + 2|\mathbf{x}^T \mathbf{P} \mathbf{b}|(4\|\mathbf{x}\| + \frac{1}{2}\gamma |\mathbf{x}^T \mathbf{P} \mathbf{b}|) \\ &= -2\|\mathbf{x}\|^2 - \gamma |\mathbf{x}^T \mathbf{P} \mathbf{b}|^2 + 8|\mathbf{x}^T \mathbf{P} \mathbf{b}|\|\mathbf{x}\| \\ &= -2 \begin{bmatrix} \|\mathbf{x}\| & |\mathbf{x}^T \mathbf{P} \mathbf{b}| \end{bmatrix} \begin{bmatrix} 1 & -2 \\ -2 & \gamma/2 \end{bmatrix} \begin{bmatrix} \|\mathbf{x}\| \\ |\mathbf{x}^T \mathbf{P} \mathbf{b}| \end{bmatrix}. \end{aligned}$$

For $\frac{d}{dt}(\mathbf{x}^T \mathbf{P} \mathbf{x})$ to be negative definite, it is enough that the matrix

$$\begin{bmatrix} 1 & -2 \\ -2 & \gamma/2 \end{bmatrix}$$

be positive definite. This is the case whenever

$$\gamma > \tilde{\gamma} = 8.$$

4.7 Hurwitz and Routh Stability Criteria

In this section we present two tests for determining if the zeros of an n th-degree polynomial are located in the open left-half complex plane. In deriving these tests, we will use the Lyapunov matrix equation given by (4.6). We consider the n th-degree polynomial

$$p(s) = a_n s^n + a_{n-1} s^{n-1} + \cdots + a_1 s + a_0 \quad (4.33)$$

with real coefficients. Associated with the polynomial $p(s)$ is the $n \times n$ real Hurwitz matrix,

$$\mathbf{H} = \begin{bmatrix} a_{n-1} & a_{n-3} & a_{n-5} & a_{n-7} & \cdots \\ a_n & a_{n-2} & a_{n-4} & a_{n-6} & \cdots \\ 0 & a_{n-1} & a_{n-3} & a_{n-5} & \cdots \\ 0 & a_n & a_{n-2} & a_{n-4} & \cdots \\ \vdots & & & \vdots & \ddots \end{bmatrix}. \quad (4.34)$$

In our further discussion we assume, without loss of generality, that $a_n = 1$. The leading principal minors of the Hurwitz matrix are

$$\begin{aligned} H_1 &= a_{n-1}, \quad H_2 = \det \begin{bmatrix} a_{n-1} & a_{n-3} \\ 1 & a_{n-2} \end{bmatrix}, \\ H_3 &= \det \begin{bmatrix} a_{n-1} & a_{n-3} & a_{n-5} \\ 1 & a_{n-2} & a_{n-4} \\ 0 & a_{n-1} & a_{n-3} \end{bmatrix}, \dots, H_n = \det \mathbf{H}. \end{aligned} \quad (4.35)$$

We next define

$$b_1 = H_1, \quad b_2 = \frac{H_2}{H_1}, \quad b_3 = \frac{H_3}{H_2 H_1},$$

and

$$b_k = \frac{H_{k-3} H_k}{H_{k-2} H_{k-1}} \quad \text{for } k = 4, 5, \dots, n.$$

Thus, for example, $b_4 = \frac{H_1 H_4}{H_2 H_3}$. Let

$$\mathbf{B} = \begin{bmatrix} 0 & 1 & 0 & 0 & \cdots & 0 & 0 & 0 \\ -b_n & 0 & 1 & 0 & \cdots & 0 & 0 & 0 \\ 0 & -b_{n-1} & 0 & 1 & \cdots & 0 & 0 & 0 \\ \vdots & & & & & & & \vdots \\ 0 & 0 & 0 & 0 & \cdots & -b_3 & 0 & 1 \\ 0 & 0 & 0 & 0 & \cdots & 0 & -b_2 & -b_1 \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

For example, for $n = 4$, we have

$$\mathbf{B} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ -b_4 & 0 & 1 & 0 \\ 0 & -b_3 & 0 & 1 \\ 0 & 0 & -b_2 & -b_1 \end{bmatrix}.$$

Let

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & & & & & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ -a_0 & -a_1 & -a_2 & \cdots & -a_{n-2} & -a_{n-1} \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

Lemma 4.1 The matrices \mathbf{B} and \mathbf{A} are similar; that is, there is a nonsingular matrix $\mathbf{T} \in \mathbb{R}^{n \times n}$ such that

$$\mathbf{B} = \mathbf{T}^{-1} \mathbf{A} \mathbf{T}. \quad (4.36)$$

Proof The similarity matrix is lower triangular and has the form

$$\mathbf{T} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & \cdots & 0 \\ t_{31} & 0 & 1 & 0 & 0 & 0 & \cdots & 0 \\ 0 & t_{42} & 0 & 1 & 0 & 0 & \cdots & 0 \\ t_{51} & 0 & t_{53} & 0 & 1 & 0 & \cdots & 0 \\ 0 & t_{62} & 0 & t_{64} & 0 & 1 & \cdots & 0 \\ \vdots & & & & & & & \vdots \end{bmatrix}. \quad (4.37)$$

To determine the nonfixed coefficients t_{ij} of the matrix \mathbf{T} , we represent (4.36) as

$$\mathbf{T} \mathbf{B} = \mathbf{A} \mathbf{T}, \quad (4.38)$$

and compare both sides of (4.38) to determine the coefficients t_{ij} . We obtain, for example, $t_{31} = -b_n$, $t_{42} = -b_n - b_{n-1}$, and so on.

Let \mathbf{P} be an $n \times n$ diagonal matrix defined as

$$\mathbf{P} = \text{diag}[b_1 b_2 \times \cdots \times b_n \quad b_1 b_2 \times \cdots \times b_{n-1} \quad \cdots \quad b_1 b_2 \quad b_1], \quad (4.39)$$

and \mathbf{Q} be an $n \times n$ diagonal matrix defined as

$$\mathbf{Q} = \text{diag}[0 \quad 0 \quad \cdots \quad 0 \quad 2b_1^2].$$

We represent the matrix \mathbf{Q} as $\mathbf{Q} = \mathbf{c}^T \mathbf{c}$, where

$$\mathbf{c} = [0 \quad 0 \quad \cdots \quad 0 \quad \sqrt{2}b_1].$$

Note that the pair (\mathbf{B}, \mathbf{c}) is observable. The matrices \mathbf{B} , \mathbf{P} , and \mathbf{Q} satisfy the Lyapunov matrix equation,

$$\mathbf{B}^T \mathbf{P} + \mathbf{P} \mathbf{B} = -\mathbf{Q}. \quad (4.40)$$

By Theorem 4.2, on page 158, applied to (4.40), the matrix \mathbf{B} is asymptotically stable if and only if the symmetric matrix \mathbf{P} is positive definite. By Theorem A.6, the matrix \mathbf{P} is positive definite if and only if its leading principal minors are positive. Using the above results, we can prove the following theorem:

Theorem 4.5 (Hurwitz) The real polynomial $p(s)$ of degree n with positive coefficients a_i has its zeros in the open left-half plane if and only if the leading principal minors of the Hurwitz matrix are positive.

Proof The polynomial $p(s)$ has its zeros in the open left plane if and only if the matrix \mathbf{B} is asymptotically stable. This in turn is true if and only if the diagonal elements of \mathbf{P} , given by (4.39), are positive. We have

$$b_1 = H_1, \quad b_1 b_2 = H_2,$$

and for $k = 3, \dots, n$ we obtain

$$b_1 b_2 \times \dots \times b_k = \frac{H_k}{H_{k-2}}.$$

Thus, the diagonal elements of \mathbf{P} are all positive, if and only if the leading principal minors of the Hurwitz matrix are positive, which completes the proof.

We will now prove the Routh necessary and sufficient condition for the real polynomial (4.33) to have its zeros in the open left-half plane. We first use the polynomial coefficients to form the following $(n + 1)$ -row Routh array:

$$\begin{array}{rcl}
 s^n & \text{row} & \left| \begin{array}{cccc} a_n & a_{n-2} & a_{n-4} & a_{n-6} & \cdots \end{array} \right. \\
 s^{n-1} & \text{row} & \left| \begin{array}{cccc} a_{n-1} & a_{n-3} & a_{n-5} & a_{n-7} & \cdots \end{array} \right. \\
 s^{n-2} & \text{row} & \left| \begin{array}{cccc} r_{21} & r_{22} & r_{23} & r_{24} & \cdots \end{array} \right. \\
 s^{n-3} & \text{row} & \left| \begin{array}{cccc} r_{31} & r_{32} & r_{33} & r_{34} & \cdots \end{array} \right. \\
 \vdots & & \left| \begin{array}{cccc} \vdots & & & \vdots & \end{array} \right. \\
 s^0 & \text{row} & \left| \begin{array}{cccc} r_{n1} & & & & \end{array} \right.
 \end{array} \tag{4.41}$$

The elements that comprise the s^n and s^{n-1} rows are appropriately ordered coefficients of $p(s)$. The elements of the s^{n-2} row are formed from its preceding two rows as follows:

$$\begin{aligned}
 r_{21} &= -\frac{1}{a_{n-1}} \det \begin{bmatrix} a_n & a_{n-2} \\ a_{n-1} & a_{n-3} \end{bmatrix} = \frac{a_{n-1}a_{n-2} - a_n a_{n-3}}{a_{n-1}}, \\
 r_{22} &= -\frac{1}{a_{n-1}} \det \begin{bmatrix} a_n & a_{n-4} \\ a_{n-1} & a_{n-5} \end{bmatrix} = \frac{a_{n-1}a_{n-4} - a_n a_{n-5}}{a_{n-1}}, \\
 r_{23} &= -\frac{1}{a_{n-1}} \det \begin{bmatrix} a_n & a_{n-6} \\ a_{n-1} & a_{n-7} \end{bmatrix} = \frac{a_{n-1}a_{n-6} - a_n a_{n-7}}{a_{n-1}}, \\
 \vdots & \quad \quad \quad \vdots
 \end{aligned}$$

The elements of the s^{n-3} row of the Routh array are generated using its preceding two rows:

$$\begin{aligned} r_{31} &= -\frac{1}{r_{21}} \det \begin{bmatrix} a_{n-1} & a_{n-3} \\ r_{21} & r_{22} \end{bmatrix} = \frac{r_{21}a_{n-3} - r_{22}a_{n-1}}{r_{21}}, \\ r_{32} &= -\frac{1}{r_{21}} \det \begin{bmatrix} a_{n-1} & a_{n-5} \\ r_{21} & r_{23} \end{bmatrix} = \frac{r_{21}a_{n-5} - r_{23}a_{n-1}}{r_{21}}, \\ r_{33} &= -\frac{1}{r_{21}} \det \begin{bmatrix} a_{n-1} & a_{n-7} \\ r_{21} & r_{24} \end{bmatrix} = \frac{r_{21}a_{n-7} - r_{24}a_{n-1}}{r_{21}}, \\ &\vdots \quad \quad \quad \vdots \end{aligned}$$

Subsequent rows of the array are formed from their preceding rows in an analogous fashion. Undefined positions can be filled with zeros. The last two rows of the array will have one term each. We will now formulate the stability test that in some texts (for example, in Rohrs et al. [243] or Wolovich [302]) is also called the Routh–Hurwitz criterion.

Theorem 4.6 (Routh) The zeros of the polynomial $p(s) = a_n s^n + a_{n-1} s^{n-1} + \cdots + a_1 s + a_0$, whose coefficients are all nonzero and have the same sign, are all in the open left-half plane if and only if the elements of the first column of the associated Routh array all have the same sign and are nonzero.

Proof Without loss of generality, we assume that the polynomial coefficients are all positive. If they are negative, we can multiply both sides of the equation $p(s) = 0$ by -1 , thus making the coefficients positive. We then perform elementary row operations on the associated Hurwitz matrix, \mathbf{H} . These operations are equivalent to premultiplying the Hurwitz matrix by lower triangular matrices with ones on the diagonal. For example,

$$\begin{bmatrix} 1 & 0 & 0 & 0 & \cdots \\ -a_n/a_{n-1} & 1 & 0 & 0 & \cdots \\ 0 & 0 & 1 & 0 & \cdots \\ \vdots & & & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} a_{n-1} & a_{n-3} & a_{n-5} & a_{n-7} \\ a_n & a_{n-2} & a_{n-4} & a_{n-6} \\ 0 & a_{n-1} & a_{n-3} & a_{n-5} \\ 0 & a_n & a_{n-2} & a_{n-4} \\ \vdots & & \vdots & \ddots \end{bmatrix} = \begin{bmatrix} a_{n-1} & a_{n-3} & a_{n-5} & a_{n-7} \\ 0 & r_{21} & r_{22} & r_{23} \\ 0 & a_{n-1} & a_{n-3} & a_{n-5} \\ 0 & a_n & a_{n-2} & a_{n-4} \\ \vdots & & \vdots & \ddots \end{bmatrix}.$$

Let \mathbf{U} be the product of the matrices corresponding to the above described elementary operations on \mathbf{H} . Then,

$$\mathbf{U} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ u_{12} & 1 & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ u_{n1} & u_{n2} & \cdots & 1 \end{bmatrix},$$

where, for example, $u_{12} = -a_n/a_{n-1}$. Premultiplying the Hurwitz matrix (4.34) by the matrix \mathbf{U} given above yields

$$\mathbf{UH} = \begin{bmatrix} a_{n-1} & a_{n-3} & a_{n-5} & \cdots & \cdots \\ 0 & r_{21} & r_{22} & \cdots & r_{2,n-1} \\ 0 & 0 & r_{31} & \cdots & r_{3,n-2} \\ \vdots & & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & r_{n1} \end{bmatrix}. \quad (4.42)$$

Because the matrix \mathbf{U} is lower triangular with the ones on the diagonal, the leading principal minors of \mathbf{H} and \mathbf{UH} are the same, where

$$H_1 = a_{n-1}, \quad H_2 = a_{n-1}r_{21}, \dots, \quad H_i = a_{n-1}r_{21} \times \cdots \times r_{i1}, \quad i = 1, 2, \dots, n.$$

Hence,

$$a_{n-1} = H_1 \quad \text{and} \quad r_{i1} = \frac{H_i}{H_{i-1}} \quad \text{for } i = 2, \dots, n.$$

By the Hurwitz theorem, the zeros of the polynomial $p(s)$ are in the open left-half plane if and only if the coefficients r_{i1} , $i = 1, 2, \dots, n$, which are the elements of the first column of the Routh array, are all positive. The proof is complete.

It is easy to verify that a necessary and sufficient condition for the first- and second-order polynomials to have their zeros in the open left-half plane is that all their coefficients be nonzero and have the same sign. However, this condition is only necessary and not sufficient for third- and higher-order polynomials to have all their zeros in the open left-half plane.

It is possible to extend the above Routh stability criterion so that it can be used to find the number of polynomial zeros with positive real parts. The number of zeros of $p(s)$ with positive real parts is equal to the number of sign changes in the first column of the Routh array. For a proof of this statement using the law of inertia of quadratic forms, see, for example, Kaczorek [143, Section 8.4.3].

◆ Example 4.6

We will use the Routh stability criterion to analyze stability of the linearized model of the flyball governor. The governor's linearized model is given by (2.26) on page 77. The model's characteristic equation is

$$\det \begin{bmatrix} s & -1 & 0 \\ \frac{g \sin^2 x_{1e}}{\cos x_{1e}} & s + \frac{b}{m} & -\frac{2g \sin x_{1e}}{x_{3e}} \\ \frac{\kappa}{l} \sin x_{1e} & 0 & s \end{bmatrix} = s^3 + \frac{b}{m}s^2 + \frac{g \sin^2 x_{1e}}{\cos x_{1e}}s + \frac{2g\kappa \sin^2 x_{1e}}{l x_{3e}} = s^3 + a_2s^2 + a_1s + a_0. \quad (4.43)$$

The coefficients of the above third-order polynomial are all positive. Therefore, the necessary condition for stability is satisfied. We next construct the Routh array for (4.43) to get

$$\begin{array}{c|cc} s^3 & 1 & a_1 \\ s^2 & a_2 & a_0 \\ s^1 & \frac{1}{a_2}(a_2 a_1 - a_0) & \\ s^0 & a_0 & \end{array} \quad (4.44)$$

Because the characteristic polynomial coefficients are all positive, the linearized model is asymptotically stable if and only if

$$a_2 a_1 - a_0 = \frac{b}{m} \frac{g \sin^2 x_{1e}}{\cos x_{1e}} - \frac{2g\kappa \sin^2 x_{1e}}{I x_{3e}} > 0. \quad (4.45)$$

Let

$$\nu = \frac{x_{3e}}{2\kappa \cos x_{1e}}.$$

Using the above notation, we express (4.45) as

$$\frac{bl}{m} \nu > 1. \quad (4.46)$$

It follows from condition (4.46) that:

1. An increase of the mass m of the weights of the governor has a deteriorating effect on the governor's stability.
2. A decrease of the friction coefficient b has also a harmful effect on stability.
3. A decrease of the moment of inertia I of the flywheel harms stability.
4. A decrease of the coefficient ν harms stability.

The above deductions are attributed by Pontryagin [235, p. 219] to the Russian engineer Vyshnegradskiy. In 1876, Vyshnegradskiy used these observations to propose design rules for governors to regulate engines without engine speeds fluctuating.

4.8 Stability of Nonlinear Systems

In this section we will be concerned with stability analysis of nonlinear time-varying dynamical system models described by (4.1)—that is, models of the form

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t)), \quad \mathbf{x}(t_0) = \mathbf{x}_0,$$

where $\mathbf{x}(t) \in \mathbb{R}^n$ is the state vector and \mathbf{f} is a vector-valued function.

To proceed, we need the following lemma. Let $B_q(\mathbf{0})$ denote an open ball with the center at $\mathbf{0}$ and radius q , and let $\mathcal{S} = \bar{B}_R(\mathbf{0}) \subset \mathbb{R}^n$ be a closed ball, a compact set, that contains the origin $\mathbf{x} = \mathbf{0}$ in its interior.

Lemma 4.2 For a continuously differentiable function $V : \mathbb{R}^n \rightarrow \mathbb{R}$, where $V(\mathbf{0}) = 0$, we have $V(\mathbf{x}) > 0$ for all $\mathbf{x} \in \mathcal{S} \setminus \{\mathbf{0}\}$, if and only if there is a continuous, strictly increasing function $\psi : \mathbb{R} \rightarrow \mathbb{R}$ such that $\psi(0) = 0$, $\psi(r) > 0$ for $r > 0$, and $V(\mathbf{x}) \geq \psi(r)$, where $r = \|\mathbf{x}\|$ and $r \leq R$.

Proof Sufficiency (\Leftarrow) is clear.

To prove necessity (\Rightarrow), we define, following Willems [300, p. 21], the function $\kappa : \mathbb{R} \rightarrow \mathbb{R}$,

$$\kappa(y) = \min_{\substack{\mathbf{x} \in \mathcal{S} \\ \|\mathbf{x}\| = y}} V(\mathbf{x}). \quad (4.47)$$

Note that the function κ is well-defined because the set

$$\{\mathbf{x} : \|\mathbf{x}\| = y\} \cap \mathcal{S} \quad (4.48)$$

is compact and V is continuous. Therefore, by the Weierstrass theorem on page 650, the function V achieves its infimum over the set (4.48).

We will now show that κ is continuous on $\mathcal{S} = B_R(\mathbf{0})$. Then, we will construct an appropriate function ψ satisfying the hypotheses of the lemma. To show that κ is continuous on \mathcal{S} , we will use the fact that V is uniformly continuous in \mathcal{S} . To this end let $\mathbf{x}^{(1)} \in \mathcal{S}$ be such that $\|\mathbf{x}^{(1)}\| = y$ and $\kappa(y) = V(\mathbf{x}^{(1)})$. Next, let Δy be an increment of y and let $\mathbf{x}^{(2)} \in \mathcal{S}$ be a vector such that $y + \Delta y = \|\mathbf{x}^{(2)}\|$. For example, we can take

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \frac{\mathbf{x}^{(1)}}{\|\mathbf{x}^{(1)}\|} \Delta y.$$

Observe that

$$V(\mathbf{x}^{(2)}) \geq \kappa(y + \Delta y).$$

Choose now a vector $\mathbf{x}^{(3)}$ such that

$$\|\mathbf{x}^{(3)}\| = y + \Delta y \quad \text{and} \quad V(\mathbf{x}^{(3)}) = \kappa(y + \Delta y).$$

Note that for any vector, say $\mathbf{x}^{(4)}$, such that $\|\mathbf{x}^{(4)}\| = y$ we have $V(\mathbf{x}^{(4)}) \geq \kappa(y)$ (see Figure 4.7). We can take

$$\mathbf{x}^{(4)} = \mathbf{x}^{(3)} - \frac{\mathbf{x}^{(3)}}{\|\mathbf{x}^{(3)}\|} \Delta y.$$

By the uniform continuity of V , for any $\epsilon > 0$ there exists a $\delta > 0$ such that if

$$\|\mathbf{x}^{(1)} - \mathbf{x}^{(2)}\| < \delta \quad \text{and} \quad \|\mathbf{x}^{(3)} - \mathbf{x}^{(4)}\| < \delta,$$

then

$$|V(\mathbf{x}^{(1)}) - V(\mathbf{x}^{(2)})| < \epsilon \quad \text{and} \quad |V(\mathbf{x}^{(3)}) - V(\mathbf{x}^{(4)})| < \epsilon.$$

From $|V(\mathbf{x}^{(1)}) - V(\mathbf{x}^{(2)})| < \epsilon$, we obtain

$$-\epsilon < V(\mathbf{x}^{(1)}) - V(\mathbf{x}^{(2)}) < \epsilon.$$

Thus

$$V(\mathbf{x}^{(1)}) > V(\mathbf{x}^{(2)}) - \epsilon.$$

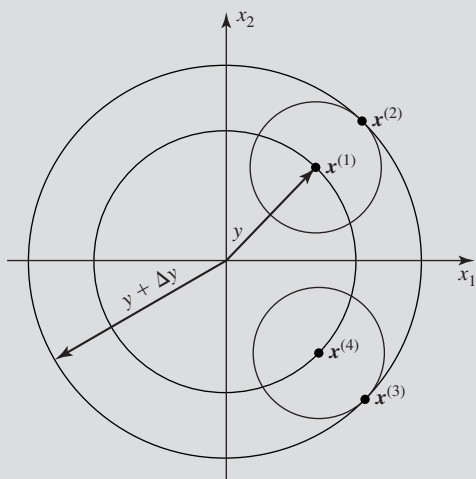


Figure 4.7 A two-dimensional illustration of the proof of the continuity of the function κ given by (4.47).

From $|V(\mathbf{x}^{(3)}) - V(\mathbf{x}^{(4)})| < \epsilon$, we get

$$V(\mathbf{x}^{(3)}) > V(\mathbf{x}^{(4)}) - \epsilon.$$

From the definition of the function κ it follows that

$$V(\mathbf{x}^{(2)}) \geq V(\mathbf{x}^{(3)}) \quad \text{and} \quad V(\mathbf{x}^{(4)}) \geq V(\mathbf{x}^{(1)}).$$

Because $V(\mathbf{x}^{(1)}) > V(\mathbf{x}^{(2)}) - \epsilon$, we obtain

$$V(\mathbf{x}^{(1)}) > V(\mathbf{x}^{(3)}) - \epsilon$$

and therefore

$$-\epsilon < V(\mathbf{x}^{(1)}) - V(\mathbf{x}^{(3)}). \quad (4.49)$$

Because $V(\mathbf{x}^{(3)}) > V(\mathbf{x}^{(4)}) - \epsilon$, we have

$$V(\mathbf{x}^{(3)}) > V(\mathbf{x}^{(1)}) - \epsilon,$$

and thus

$$V(\mathbf{x}^{(1)}) - V(\mathbf{x}^{(3)}) < \epsilon. \quad (4.50)$$

Combining (4.49) and (4.50), we get

$$|V(\mathbf{x}^{(1)}) - V(\mathbf{x}^{(3)})| < \epsilon,$$

or, equivalently,

$$|\kappa(y) - \kappa(y + \Delta y)| < \epsilon$$

if

$$\Delta y < \delta,$$

which proves the continuity of κ .

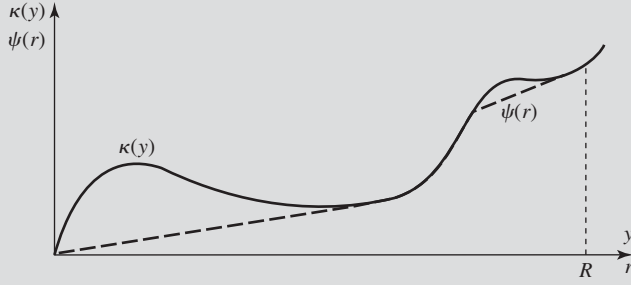


Figure 4.8 Constructing the function ψ given by (4.51).

We are now ready to construct a function ψ that satisfies the hypotheses of the lemma. This function comes from Willems [300, p. 22] and has the form

$$\psi(r) = r \min_{r \leq y \leq R} \frac{\kappa(y)}{y}. \quad (4.51)$$

Note that ψ is continuous, $\psi(0) = 0$, and it is positive and strictly increasing; see Figure 4.8 for a graphical illustration of the way ψ is constructed from κ . The proof of the lemma is now complete.

We can characterize a positive definite function V in another way as in the following lemma.

Lemma 4.3 A continuously differentiable function $V: \mathbb{R}^n \rightarrow \mathbb{R}$ is positive definite in $\mathcal{S} = \bar{B}_R(\mathbf{0}) \subset \mathbb{R}^n$, with respect to $\mathbf{x} = \mathbf{0}$, if and only if there exists a continuous nondecreasing function $\alpha: \mathbb{R} \rightarrow \mathbb{R}$, such that $\alpha(0) = 0$, $\alpha(r) > 0$ for $r > 0$, and $V(\mathbf{x}) \geq \alpha(r)$, where $r = \|\mathbf{x}\|$ and $r \leq R$.

A possible choice for α is

$$\alpha(r) = \min_{r \leq \|\mathbf{x}\| \leq R} V(\mathbf{x}).$$

Note that the function α is nondecreasing because as r increases, the minimum is taken over a smaller region. Furthermore, $\alpha(r) > 0$ whenever $r > 0$ because $V(\mathbf{x}) > 0$ for $\mathbf{x} \in \mathcal{S} \setminus \{\mathbf{0}\}$. The function α is continuous because V is continuous.

The main difference between the above two characterizations of positive definite functions is that in Lemma 4.2 the function ψ is strictly increasing rather than just nondecreasing. In our further discussion, we will exploit extensively the existence of such a function ψ associated with a positive definite function V .

So far we only considered positive definite, or negative definite functions, that are time-invariant. We now extend our definition of positive definiteness to time-varying functions of the form $W(t, \mathbf{x})$. We first define a time-varying positive semidefinite function. As before, we let \mathcal{S} to denote a compact subset of \mathbb{R}^n containing the origin $\mathbf{x} = \mathbf{0}$ in its interior.

Definition 4.11 A real-valued function $W : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}$ is positive semidefinite, with respect to $\mathbf{x} = \mathbf{0}$, in a region $\mathcal{S} \subseteq \mathbb{R}^n$ if for all t and all \mathbf{x} in \mathcal{S} :

1. W is continuous and has continuous partial derivatives with respect to all of its arguments; that is, $W \in \mathcal{C}^1$.
2. $W(t, \mathbf{0}) = 0$.
3. $W(t, \mathbf{x}) \geq 0$.

Definition 4.12 A real-valued function $W : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}$ is positive definite, with respect to $\mathbf{x} = \mathbf{0}$, in a region $\mathcal{S} \subseteq \mathbb{R}^n$ if for all t and all \mathbf{x} in \mathcal{S} :

1. $W \in \mathcal{C}^1$.
2. $W(t, \mathbf{0}) = 0$.
3. There exists a positive definite time-invariant function V such that $W(t, \mathbf{x}) \geq V(\mathbf{x})$.

We next introduce a notion of a *decreasing* function.

Definition 4.13 A time-varying function $W : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}$ is decreasing in a region \mathcal{S} if a positive definite function \tilde{V} exists such that for all \mathbf{x} in \mathcal{S} and all t :

$$|W(t, \mathbf{x})| \leq \tilde{V}(\mathbf{x});$$

that is, W tends to zero uniformly with respect to t as $\|\mathbf{x}\| \rightarrow 0$.

Now consider an equilibrium state of the dynamical system modeled by (4.1)—that is, a constant state \mathbf{x}_e such that $\mathbf{f}(t, \mathbf{x}_e) = \mathbf{0}$ for all t . Let $W(t, \mathbf{x})$ be a time-dependent real-valued function whose time derivative along the solution $\mathbf{x}(t) = \mathbf{x}(t; t_0, \mathbf{x}_0)$ of the given differential equation is

$$\frac{dW(t, \mathbf{x})}{dt} = \dot{W}(t, \mathbf{x}) = \frac{\partial W}{\partial t} + \nabla W^T \dot{\mathbf{x}}.$$

Then, we have the following theorem concerning the stability of the equilibrium point in the sense of Lyapunov.

Theorem 4.7 If in some neighborhood \mathcal{S} of the equilibrium state \mathbf{x}_e there is a positive definite function W , with respect to \mathbf{x}_e , such that its time derivative along any solution of $\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t))$ is negative semidefinite in \mathcal{S} , then the equilibrium state \mathbf{x}_e is stable. If W is also decreasing, then the equilibrium state is uniformly stable.

Proof Without loss of generality, we assume that $\mathbf{x}_e = \mathbf{0}$ and follow the arguments of Willems [300, pp. 29, 30]. Because W is positive definite, by Definition 4.12 and Lemma 4.2 there is a strictly increasing function $\psi_1(r)$, $\psi(0) = 0$, such that

$$W(t, \mathbf{x}) \geq \psi_1(\|\mathbf{x}\|),$$

where $\|\mathbf{x}\| \leq R$, and R is the radius of the largest ball contained in \mathcal{S} . Let

$$\psi_2(t_0, r) = \max_{\|\mathbf{x}\| \leq r} W(t_0, \mathbf{x}).$$

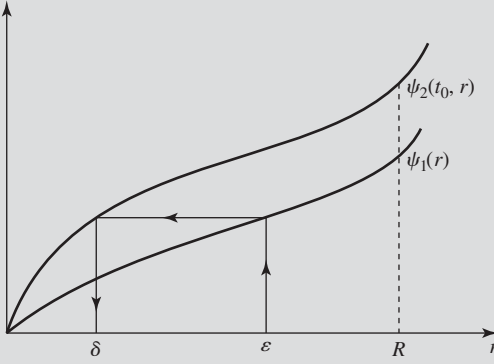


Figure 4.9 Determining a δ given an ϵ .

Then

$$\psi_1(r) \leq \psi_2(t_0, r).$$

Thus, for any $\epsilon > 0$ there is a positive δ that depends on ϵ and t_0 , such that

$$\psi_2(t_0, \delta) < \min\{\psi_1(\epsilon), \psi_1(R)\}.$$

Figure 4.9 illustrates how to determine a δ given an ϵ . By assumption, $dW(t, \mathbf{x})/dt \leq 0$ along any solution of the system. Hence $W(t, \mathbf{x})$ is a nonincreasing function and therefore

$$W(t, \mathbf{x}(t; t_0, \mathbf{x}_0)) \leq W(t_0, \mathbf{x}_0) \quad \text{for all } t \geq t_0.$$

If $\|\mathbf{x}_0\| < \delta$, then for $t \geq t_0$ we have

$$\psi_1(\epsilon) > \psi_2(t_0, \delta) \geq W(t_0, \mathbf{x}_0) \geq W(t, \mathbf{x}(t; t_0, \mathbf{x}_0)) \geq \psi_1(\|\mathbf{x}(t; t_0, \mathbf{x}_0)\|).$$

Therefore,

$$\psi_1(\epsilon) > \psi_1(\|\mathbf{x}(t; t_0, \mathbf{x}_0)\|).$$

Because ψ_1 is strictly increasing, we have to have

$$\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| < \epsilon$$

for all $t \geq t_0$. This means that the null solution is (locally) stable, which completes the proof of the first part of the theorem.

If W is also decrescent, then there is a nondecreasing function $\tilde{\psi}_2 : \mathbb{R} \rightarrow \mathbb{R}$, $\tilde{\psi}_2(0) = 0$, and

$$W(t, \mathbf{x}) \leq \tilde{\psi}_2(\|\mathbf{x}\|)$$

for all \mathbf{x} such that $\|\mathbf{x}\| \leq R$ and all t . Indeed, because $W(t, \mathbf{x})$ is positive definite and decrescent, there exists a positive definite $V(\mathbf{x})$ such that $W(t, \mathbf{x}) \leq V(\mathbf{x})$. Then, we can take $\tilde{\psi}_2$ defined as

$$\tilde{\psi}_2(r) = \max_{\|\mathbf{x}\| \leq r} V(\mathbf{x}).$$

Let

$$\tilde{\psi}_2(\delta) < \min\{\psi_1(\epsilon), \psi_1(R)\}.$$

Note that δ depends now only on ϵ , and not on t_0 as before. To show the uniform stability of the null solution, we repeat the argument above with $\tilde{\psi}_2$ replacing ψ_2 . This completes the proof of the theorem.

The next theorem gives us a sufficiency condition for an equilibrium point to be uniformly asymptotically stable.

Theorem 4.8 If in some neighborhood \mathcal{S} of the given equilibrium point \mathbf{x}_e there is a positive definite, with respect to \mathbf{x}_e , and decrescent function W such that its time derivative along any solution of the system model (4.1) is negative definite in \mathcal{S} , then the equilibrium state is uniformly asymptotically stable.

Proof As before, without loss of generality, we assume that $\mathbf{x}_e = \mathbf{0}$ and then follow the arguments of Willems [300, pp. 30, 31]. The previous theorem ensures uniform stability of the equilibrium state \mathbf{x}_e . Thus, only the uniform convergence remains to be shown. From the assumptions of the theorem it follows that there exist strictly increasing functions ψ_1 and ψ_3 , along with a nondecreasing function ψ_2 such that

$$\begin{aligned}\psi_1(0) &= \psi_2(0) = \psi_3(0), \\ \psi_1(\|\mathbf{x}\|) &\leq W(t, \mathbf{x}) \leq \psi_2(\|\mathbf{x}\|), \\ \psi_3(\|\mathbf{x}\|) &\leq -\dot{W}(t, \mathbf{x}),\end{aligned}$$

for all t and all \mathbf{x} such that $\|\mathbf{x}\| \leq R$, where R is the radius of the largest ball contained in \mathcal{S} . First select δ_1 so that

$$\psi_2(\delta_1) < \psi_1(R), \quad (4.52)$$

Next, given an $\epsilon_1 > 0$, we define positive constants δ_2 and T as follows:

$$\psi_2(\delta_2) < \min\{\psi_1(\epsilon_1), \psi_2(\delta_1)\}, \quad (4.53)$$

$$T = \psi_1(R)/\psi_3(\delta_2). \quad (4.54)$$

Note that δ_1 and T are independent of t_0 .

We shall prove the theorem in two stages. First, we shall show that for some $t_2 \in [t_0, t_0 + T]$,

$$\|\mathbf{x}(t_2; t_0, \mathbf{x}_0)\| < \delta_2$$

whenever $\|\mathbf{x}_0\| < \delta_1$. Then, we shall show that once $\|\mathbf{x}(t_2; t_0, \mathbf{x}_0)\| < \delta_2$, then

$$\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| < \epsilon_1 \quad \text{for all } t \geq t_2.$$

We begin with proving the claim that for some $t_2 \in [t_0, t_0 + T]$,

$$\|\mathbf{x}(t_2; t_0, \mathbf{x}_0)\| < \delta_2$$

whenever $\|\mathbf{x}_0\| < \delta_1$. We prove this claim by contradiction. We assume that for all $t \in [t_0, t_0 + T]$,

$$\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| \geq \delta_2 \quad \text{if } \|\mathbf{x}_0\| < \delta_1. \quad (4.55)$$

Using (4.55), (4.54), and (4.52), we obtain

$$\begin{aligned}
 0 < \psi_1(\delta_2) &\leq \psi_1(\|\mathbf{x}(t; t_0, \mathbf{x}_0)\|) \\
 &\leq W(t_0 + T, \mathbf{x}(t_0 + T; t_0, \mathbf{x}_0)) \\
 &= W(t_0, \mathbf{x}_0) + \int_{t_0}^{t_0+T} \dot{W}[t, \mathbf{x}(t; t_0, \mathbf{x}_0)] dt \\
 &\leq W(t_0, \mathbf{x}_0) - \int_{t_0}^{t_0+T} \psi_3(\|\mathbf{x}(t; t_0, \mathbf{x}_0)\|) dt \\
 &\leq W(t_0, \mathbf{x}_0) - T\psi_3(\delta_2) \\
 &\leq \psi_2(\|\mathbf{x}_0\|) - T\psi_3(\delta_2) \\
 &< \psi_2(\delta_1) - \psi_1(R) \\
 &< 0,
 \end{aligned}$$

which is a contradiction, and thus we proved the first part of the theorem. We now proceed with the second part. It follows from the first part that if $\|\mathbf{x}_0\| < \delta_1$, then for some $t_2 \in [t_0, t_0 + T]$ we have $\|\mathbf{x}(t_2; t_0, \mathbf{x}_0)\| < \delta_2$. Hence, for $t \geq t_2$,

$$\psi_1(\|\mathbf{x}(t; t_0, \mathbf{x}_0)\|) \leq W(t, \mathbf{x}(t; t_0, \mathbf{x}_0)) \leq W(t_2, \mathbf{x}(t_2; t_0, \mathbf{x}_0)) < \psi_2(\delta_2) < \psi_1(\epsilon_1).$$

Therefore,

$$\psi_1(\|\mathbf{x}(t; t_0, \mathbf{x}_0)\|) < \psi_1(\epsilon_1),$$

and because ψ_1 is strictly increasing we obtain

$$\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| < \epsilon_1$$

for all $t \geq t_2$, which proves the uniform convergence. The proof of the theorem is thus complete.

◆ Example 4.7

Given a dynamical system model

$$\dot{\mathbf{x}}(t) = \begin{bmatrix} -0.75 & 0.3536 \\ 0.3536 & -0.5 \end{bmatrix} \mathbf{x}(t), \quad \mathbf{x}(0) = \mathbf{x}_0.$$

The solution $\mathbf{P} = \mathbf{P}^T$ to the Lyapunov matrix equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{Q}$ where $\mathbf{Q} = 2\mathbf{I}_2$, is

$$\mathbf{P} = \begin{bmatrix} 2 & \sqrt{2} \\ \sqrt{2} & 3 \end{bmatrix}.$$

The eigenvalues of \mathbf{P} are

$$\lambda_1 = \lambda_{\min}(\mathbf{P}) = 1, \quad \lambda_2 = \lambda_{\max}(\mathbf{P}) = 4.$$

In general, for a symmetric matrix \mathbf{P} ,

$$\lambda_{\min}(\mathbf{P}) \|\mathbf{x}\|^2 \leq \mathbf{x}^T \mathbf{P} \mathbf{x} \leq \lambda_{\max}(\mathbf{P}) \|\mathbf{x}\|^2.$$

Therefore, we can take

$$\psi_1(\|\mathbf{x}\|) = \lambda_{\min}(\mathbf{P}) \|\mathbf{x}\|^2, \quad \psi_2(\|\mathbf{x}\|) = \lambda_{\max}(\mathbf{P}) \|\mathbf{x}\|^2, \quad \text{and} \quad \psi_3(\|\mathbf{x}\|) = \lambda_{\min}(\mathbf{Q}) \|\mathbf{x}\|^2.$$

In our example, $\lambda_{\min}(\mathbf{Q}) = 2$.

- Given $\varepsilon = 4$. We will find $\delta > 0$ such that if $\|\mathbf{x}_0\| \leq \delta$, then $\|\mathbf{x}(t)\| \leq \varepsilon$ for $t \geq 0$. We obtain $\delta = \delta(\varepsilon)$ by solving the equation,

$$\psi_1(\varepsilon) = \psi_2(\delta),$$

to get

$$\delta = \sqrt{\frac{\lambda_{\min}(\mathbf{P})}{\lambda_{\max}(\mathbf{P})}} \varepsilon = 2.$$

- Given $\mathbf{x}_0 = [8 \ 6]^T$ and $\varepsilon_1 = 2$. We will find T such that $\|\mathbf{x}(t)\| \leq \varepsilon_1$ for all $t \geq T$. We have $\|\mathbf{x}_0\| = 10$. Let $\delta_1 = 10$. We then compute R , δ_2 , and T from the equations

$$\begin{aligned} \psi_2(\delta_1) &= \psi_1(R), \\ \psi_2(\delta_2) &= \min(\psi_1(\varepsilon_1), \psi_2(\delta_1)), \\ T &= \frac{\psi_1(R)}{\psi_3(\delta_2)}. \end{aligned}$$

Using the first equation, we find

$$R = \sqrt{\frac{\lambda_{\max}(\mathbf{P})}{\lambda_{\min}(\mathbf{P})}} \delta_1 = 20.$$

We then solve for $\delta_2 \leq \delta_1$ to get

$$\delta_2 = \sqrt{\frac{\lambda_{\min}(\mathbf{P})}{\lambda_{\max}(\mathbf{P})}} \varepsilon_1 = 1.$$

Hence,

$$T = \frac{\lambda_{\min}(\mathbf{P}) R^2}{\lambda_{\min}(\mathbf{Q}) \delta_2^2} = \frac{\lambda_{\min}(\mathbf{P}) \frac{\lambda_{\max}(\mathbf{P})}{\lambda_{\min}(\mathbf{P})} \delta_1^2}{\lambda_{\min}(\mathbf{Q}) \frac{\lambda_{\min}(\mathbf{P})}{\lambda_{\max}(\mathbf{P})} \varepsilon_1^2} = \frac{\lambda_{\max}^2(\mathbf{P}) \delta_1^2}{\lambda_{\min}(\mathbf{Q}) \lambda_{\min}(\mathbf{P}) \varepsilon_1^2} = 200.$$

Proposition 4.1 If in some neighborhood \mathcal{S} of $\mathbf{x}_e = \mathbf{0}$ the positive definite and decrescent function W satisfies

$$c_1 \|\mathbf{x}\|^2 \leq W(t, \mathbf{x}) \leq c_2 \|\mathbf{x}\|^2$$

for some positive constants c_1 and c_2 and its time derivative along any solution of the system (4.1) is negative definite in \mathcal{S} and such that for some positive constant c_3 ,

$$\dot{W}(t, \mathbf{x}) \leq -c_3 \|\mathbf{x}\|^2,$$

then the equilibrium state $\mathbf{x}_e = \mathbf{0}$ is uniformly exponentially stable.

Proof We have

$$\begin{aligned} \dot{W}(t, \mathbf{x}) &\leq -c_3 \|\mathbf{x}\|^2 \\ &= -\frac{c_3}{c_2} (c_2 \|\mathbf{x}\|^2) \\ &\leq -\frac{c_3}{c_2} W(t, \mathbf{x}). \end{aligned}$$

Applying to the above Theorem A.20, which is called the comparison theorem, we obtain

$$W(t, \mathbf{x}) \leq W(t_0, \mathbf{x}_0) \exp\left(-\frac{c_3}{c_2}(t - t_0)\right).$$

Hence,

$$\begin{aligned} \|\mathbf{x}\| &\leq \sqrt{\frac{W(t, \mathbf{x})}{c_1}} \\ &\leq \sqrt{\frac{W(t_0, \mathbf{x}_0) \exp\left(-\frac{c_3}{c_2}(t - t_0)\right)}{c_1}} \\ &\leq \left(\frac{c_2}{c_1}\right)^{\frac{1}{2}} \|\mathbf{x}(t_0)\| \exp\left(-\frac{c_3}{2c_2}(t - t_0)\right), \end{aligned}$$

which means that the equilibrium state $\mathbf{x}_e = \mathbf{0}$ is uniformly exponentially stable.

The above theorems were concerned with local stability. We now proceed to prove sufficiency conditions for global stability. For this we need the following definitions. We first define a *radially unbounded* or *coercive* function.

Definition 4.14 A time-varying function $W : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}$ is radially unbounded, or coercive, if as $\|\mathbf{x}\| \rightarrow \infty$ the function W tends to infinity uniformly with respect to t ; that is, for any given positive real number M , there is a positive real number N so that $W(t, \mathbf{x}) > M$ for all t and all \mathbf{x} whenever $\|\mathbf{x}\| > N$.

The next two definitions are concerned with *Lagrange stable equilibrium states*.

Definition 4.15 The equilibrium state \mathbf{x}_e is called Lagrange stable, or bounded, if for every real number $d > 0$ there exists a positive real number $b = b(t_0, d)$ such that if

$$\|\mathbf{x}_0 - \mathbf{x}_e\| < d,$$

then

$$\|\mathbf{x}(t; t_0, \mathbf{x}_0) - \mathbf{x}_e\| < b$$

for all $t \geq t_0$.

Definition 4.16 The equilibrium state \mathbf{x}_e is uniformly bounded if the number b in Definition 4.15 can be taken independent of the initial time t_0 .

We now define *uniformly asymptotically stable in the large* equilibrium state.

Definition 4.17 If the equilibrium state \mathbf{x}_e is uniformly stable and uniformly bounded and all trajectories are uniformly convergent to \mathbf{x}_e , then this equilibrium state is called uniformly asymptotically stable in the large, or globally uniformly asymptotically stable.

The following theorem gives a sufficiency condition for an equilibrium state to be globally uniformly asymptotically stable.

Theorem 4.9 If there is a function W that is positive definite for all $\mathbf{x} \neq \mathbf{x}_e$, $W(t, \mathbf{x}_e) = 0$, as well as decrescent and radially unbounded, and if its time derivative along any solution of the system is negative definite for all \mathbf{x} , then the equilibrium state \mathbf{x}_e is globally uniformly asymptotically stable.

Proof We assume, for simplicity, that $\mathbf{x}_e = \mathbf{0}$ and follow the arguments of Willems [300, pp. 31, 32]. By Theorem 4.7, the equilibrium state \mathbf{x}_e is uniformly stable. Similarly as in Theorem 4.8, we construct the functions ψ_1 , ψ_2 , and ψ_3 . Since, by assumption, W is radially unbounded, $\psi_1(\|\mathbf{x}\|) \rightarrow \infty$ as $\|\mathbf{x}\| \rightarrow \infty$. Therefore, for any constant $d > 0$, there is a constant $b > 0$ such that $\psi_1(b) > \psi_2(d)$. See Figure 4.10 for an illustration of the procedure that can be used to find the constant b . If $\|\mathbf{x}_0\| < d$, then $\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| < b$ for all $t \geq t_0$. Indeed, because

$$\psi_1(b) > \psi_2(d) \geq W(t_0, \mathbf{x}_0) \geq W(t, \mathbf{x}(t; t_0, \mathbf{x}_0)) \geq \psi_1(\|\mathbf{x}(t; t_0, \mathbf{x}_0)\|),$$

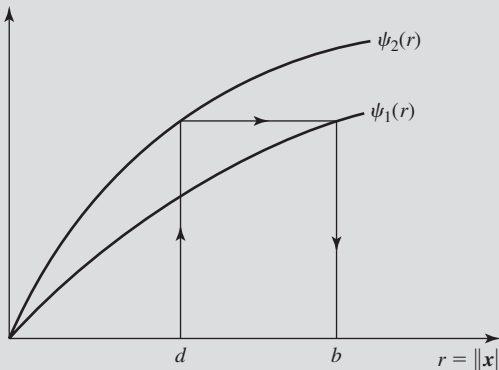


Figure 4.10 Determining a constant $b > 0$ given $d > 0$ in the proof of Theorem 4.9.

then

$$\psi_1(b) > \psi_1(\|\mathbf{x}(t; t_0, \mathbf{x}_0)\|).$$

The above implies that

$$\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| < b$$

because ψ_1 is strictly increasing. Hence all trajectories of the system are uniformly bounded. To prove their convergence to the equilibrium state, we define δ_2 , T , and b , given d and ϵ_1 , as

$$\begin{aligned}\psi_2(d) &< \psi_1(b), \\ \psi_2(\delta_2) &< \psi_1(\epsilon_1), \\ T &= \psi_1(b)/\psi_3(\delta_2).\end{aligned}$$

If $\|\mathbf{x}_0\| < d$, then for all $t > t_0 + T$, we have

$$\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| < \epsilon_1,$$

which can be shown using the same arguments as in the proof of Theorem 4.8 with b playing the role of R , and d playing the role of δ_1 . This proves that all trajectories of the system converge uniformly to the equilibrium state. Hence, the equilibrium state is globally uniformly asymptotically stable.

We now give a theorem that can be used as a test for checking the instability of nonlinear dynamical systems. This instability test is illustrated in Figure 4.11.

Theorem 4.10 Let \mathbf{x}_e be an equilibrium state of the system $\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x})$, that is, $\mathbf{f}(t, \mathbf{x}_e) = \mathbf{0}$. If there exists a function $W: \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}$ such that in the region $\mathcal{S} = \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}_e\| \leq R\}$

1. $W \in \mathcal{C}^1$,
2. $W(t, \mathbf{x}_e) = 0$,

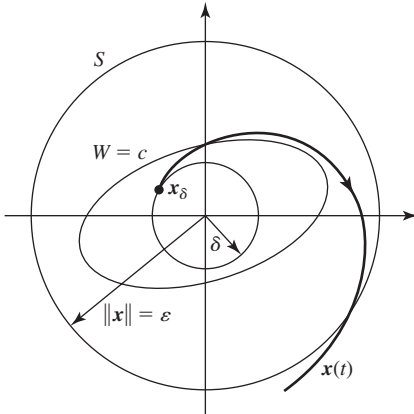


Figure 4.11 An illustration of Theorem 4.10.

3. $\dot{W}(t, \mathbf{x}) > 0$ for all $\mathbf{x} \in \mathcal{S} \setminus \{\mathbf{x}_e\}$, and
4. in every neighborhood of \mathbf{x}_e in \mathcal{S} , for any t_0 there is at least one point \mathbf{x} at which $W(t_0, \mathbf{x}) > 0$,

then the equilibrium state \mathbf{x}_e is unstable.

Proof Without loss of generality, assume that $\mathbf{x}_e = \mathbf{0}$. By assumption 4, for any t_0 and any δ , $0 < \delta < R$, there exists \mathbf{x}_δ , $0 < \|\mathbf{x}_\delta\| < \delta$, such that $W(t_0, \mathbf{x}_\delta) > 0$. The system trajectory $\mathbf{x}(t) = \mathbf{x}(t; t_0, \mathbf{x}_\delta)$ cannot approach the equilibrium state $\mathbf{x}_e = \mathbf{0}$ as $t \rightarrow \infty$ because $W(t, \mathbf{0}) = 0$ but $W(t_0, \mathbf{x}_\delta) > 0$ and $\dot{W}(t, \mathbf{x}) > 0$ for $\mathbf{x} \neq \mathbf{x}_e$. Since \dot{W} is continuous, because $W \in \mathcal{C}^1$, and the trajectory is bounded away from the equilibrium state, there is a constant $m > 0$ such that for all $t \geq t_0$,

$$\dot{W}(t, \mathbf{x}(t)) \geq m > 0,$$

as long as

$$\|\mathbf{x}(t)\| \leq R, \quad \mathbf{x} \neq \mathbf{0}.$$

Hence,

$$W(t, \mathbf{x}) = W(t_0, \mathbf{x}_\delta) + \int_{t_0}^t \dot{W}(s, \mathbf{x}(s)) ds \geq W(t_0, \mathbf{x}_\delta) + m(t - t_0). \quad (4.56)$$

The function W is continuous and therefore bounded on the compact set \mathcal{S} . However, the right-hand side of (4.56) is unbounded. Therefore the trajectory $\mathbf{x}(t; t_0, \mathbf{x}_\delta)$ cannot remain infinitely long in \mathcal{S} . Thus, for some $t_1 > t_0$ the trajectory reaches the boundary of the region $\|\mathbf{x}\| \leq R$. Hence, given any ϵ , $0 < \epsilon < R$, then for any $\delta > 0$, no matter how small, there is at least one trajectory such that $\|\mathbf{x}_\delta\| < \delta$ but $\|\mathbf{x}(t; t_0, \mathbf{x}_\delta)\| \geq \epsilon$ for some $t_1 > t_0$. Therefore, the equilibrium state is unstable.

Corollary 4.1 If there exists a function $W: \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}$ such that in the region $\mathcal{S} = \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}_e\| \leq R\}$

1. $W \in \mathcal{C}^1$,
2. $W(t, \mathbf{x}_e) = 0$,
3. $\dot{W}(t, \mathbf{x}) < 0$ for all $\mathbf{x} \in \mathcal{S} \setminus \{\mathbf{x}_e\}$, and
4. in every neighborhood of \mathbf{x}_e in \mathcal{S} , for any t_0 there is at least one point \mathbf{x} at which $W(t_0, \mathbf{x}) < 0$,

then the equilibrium state \mathbf{x}_e is unstable.

4.9 Lyapunov's Indirect Method

In this section we are concerned with the problem of investigating stability properties of an equilibrium state of a nonlinear system based on its linearization about the given equilibrium. We devise a method that allows one to determine whether the equilibrium of the nonlinear system is asymptotically stable or unstable based on the dynamics of the linearized system about this equilibrium. The method is sometimes referred to as the *Lyapunov's first method* or *Lyapunov's*

indirect method. For simplicity of notation, we assume that the equilibrium point to be tested for stability is the origin. We consider a nonlinear system model,

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}), \quad (4.57)$$

where $\mathbf{f} : \mathcal{D} \rightarrow \mathbb{R}^n$ is a continuously differentiable function from a domain $\mathcal{D} \subseteq \mathbb{R}^n$ into \mathbb{R}^n . We assume that the origin $\mathbf{x} = \mathbf{0}$ is in the interior of \mathcal{D} and is an equilibrium state of (4.57), that is, $\mathbf{f}(\mathbf{0}) = \mathbf{0}$. Let

$$\mathbf{A} = \left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{0}} = \left[\begin{array}{cccc} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \cdots & \frac{\partial f_n}{\partial x_n} \end{array} \right] \bigg|_{\mathbf{x}=\mathbf{0}} \quad (4.58)$$

be the Jacobian matrix of \mathbf{f} evaluated at $\mathbf{x} = \mathbf{0}$. Then, applying Taylor's theorem to the right-hand side of (4.57) yields

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{g}(\mathbf{x}), \quad (4.59)$$

where the function \mathbf{g} represents higher-order terms and has the property

$$\lim_{\|\mathbf{x}\| \rightarrow 0} \frac{\|\mathbf{g}(\mathbf{x})\|}{\|\mathbf{x}\|} = 0. \quad (4.60)$$

The above suggests that in a “small” neighborhood of the origin, we can approximate the nonlinear system (4.57) with its linearized system about the origin, $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$. The following theorems give conditions under which we can infer about stability properties of the origin of the nonlinear system based on stability properties of the linearized system.

Theorem 4.11 Let the origin, $\mathbf{x} = 0$, be an equilibrium state of the nonlinear system, $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x})$. Then, the origin is an asymptotically stable equilibrium of the nonlinear system if \mathbf{A} , the Jacobian matrix of \mathbf{f} evaluated at the origin, has all its eigenvalues in the open left-half plane.

Proof Our proof is based on that of Khalil [157, pp. 127–128]. To begin with, we suppose that all the eigenvalues of \mathbf{A} are in the open left-half plane. Then, by Theorem 4.1, for any positive definite symmetric matrix \mathbf{Q} , the solution \mathbf{P} to the Lyapunov equation (4.6) is also positive definite. We employ the positive definite quadratic form $V(\mathbf{x}) = \mathbf{x}^T \mathbf{P} \mathbf{x}$ as a Lyapunov function candidate for the nonlinear system $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x})$. We use (4.59) and the Lyapunov matrix equation (4.6) when evaluating the time derivative of $V(\mathbf{x})$ on the trajectories of the nonlinear system. We obtain

$$\begin{aligned} \dot{V}(\mathbf{x}) &= 2\mathbf{x}^T \mathbf{P} \mathbf{f}(\mathbf{x}) \\ &= 2\mathbf{x}^T \mathbf{P} (\mathbf{A}\mathbf{x} + \mathbf{g}(\mathbf{x})) \\ &= \mathbf{x}^T (\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A}) \mathbf{x} + 2\mathbf{x}^T \mathbf{P} \mathbf{g}(\mathbf{x}) \\ &= -\mathbf{x}^T \mathbf{Q} \mathbf{x} + 2\mathbf{x}^T \mathbf{P} \mathbf{g}(\mathbf{x}). \end{aligned} \quad (4.61)$$

The first term on the right-hand side of the above equation is negative definite, while the second term is, in general, indefinite. The function \mathbf{g} satisfies

$$\frac{\|\mathbf{g}(\mathbf{x})\|_2}{\|\mathbf{x}\|_2} \rightarrow 0 \quad \text{as } \|\mathbf{x}\|_2 \rightarrow 0.$$

Therefore, for any $\gamma > 0$ there exists $r > 0$ such that

$$\|\mathbf{g}(\mathbf{x})\|_2 < \gamma \|\mathbf{x}\|_2 \quad \text{for all } \|\mathbf{x}\|_2 < r. \quad (4.62)$$

Substituting the above into (4.61) gives

$$\begin{aligned} \dot{V}(\mathbf{x}) &< -\mathbf{x}^T \mathbf{Q} \mathbf{x} + 2\gamma \|\mathbf{P}\|_2 \|\mathbf{x}\|_2^2 \\ &\leq -(\lambda_{\min}(\mathbf{Q}) - 2\gamma \lambda_{\max}(\mathbf{P})) \|\mathbf{x}\|_2^2 \quad \text{for all } \|\mathbf{x}\|_2 < r. \end{aligned} \quad (4.63)$$

Thus, selecting

$$\gamma < \frac{\lambda_{\min}(\mathbf{Q})}{2\lambda_{\max}(\mathbf{P})}$$

ensures that $\dot{V}(\mathbf{x})$ is negative definite. By Theorem 4.8, we conclude that the origin, $\mathbf{x} = \mathbf{0}$, is an asymptotically stable equilibrium state of the nonlinear system, which completes the proof.

Theorem 4.12 Let the origin, $\mathbf{x} = \mathbf{0}$, be an equilibrium state of the nonlinear system, $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x})$. Then, the origin is unstable if \mathbf{A} has one or more eigenvalues in the right-half plane.

Proof Our proof follows that of Khalil [157, pp. 128–130]. We first consider the special case when \mathbf{A} has no eigenvalues on the imaginary axis. Suppose that the eigenvalues cluster into a group in the open right-half plane and a group in the open left-half plane. Then, we can find a nonsingular matrix \mathbf{T} such that

$$\mathbf{T} \mathbf{A} \mathbf{T}^{-1} = \begin{bmatrix} -\mathbf{A}_1 & \mathbf{O} \\ \mathbf{O} & \mathbf{A}_2 \end{bmatrix},$$

where both \mathbf{A}_1 and \mathbf{A}_2 have their eigenvalues in the open left-half plane. A possible candidate for \mathbf{T} would be the matrix that transforms \mathbf{A} into its real Jordan form. Now, let

$$\mathbf{z} = \mathbf{T} \mathbf{x} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix},$$

where the partition of \mathbf{z} is compatible with the dimensions of \mathbf{A}_1 and \mathbf{A}_2 . The change of variables, $\mathbf{z} = \mathbf{T} \mathbf{x}$, transforms the nonlinear system (4.59) into the form

$$\begin{aligned} \dot{\mathbf{z}}_1 &= -\mathbf{A}_1 \mathbf{z}_1 + \mathbf{g}_1(\mathbf{z}), \\ \dot{\mathbf{z}}_2 &= \mathbf{A}_2 \mathbf{z}_2 + \mathbf{g}_2(\mathbf{z}), \end{aligned} \quad (4.64)$$

where the functions \mathbf{g}_i , $i = 1, 2$, have the property that for any $\gamma > 0$, there exists $r > 0$ such that for $i = 1, 2$, we have

$$\|\mathbf{g}_i(\mathbf{z})\|_2 \leq \gamma \|\mathbf{z}\|_2 \quad \text{for all } \|\mathbf{z}\|_2 \leq r. \quad (4.65)$$

The origin $\mathbf{z} = \mathbf{0}$ is an equilibrium state for the nonlinear system (4.64). Because the matrix \mathbf{T} is nonsingular, stability properties of $\mathbf{z} = \mathbf{0}$ in the \mathbf{z} coordinates are the same as the stability properties of $\mathbf{x} = \mathbf{0}$ in the \mathbf{x} coordinates. We will now show that $\mathbf{z} = \mathbf{0}$ is an unstable equilibrium of the system (4.64), which implies that $\mathbf{x} = \mathbf{0}$ is an unstable equilibrium of the system (4.57). We use Theorem 4.10, which entails constructing a function V satisfying the hypotheses of that theorem. Let \mathbf{Q}_1 and \mathbf{Q}_2 be positive definite symmetric matrices of the dimensions of \mathbf{A}_1 and \mathbf{A}_2 , respectively. Because \mathbf{A}_1 and \mathbf{A}_2 have their eigenvalues in the open left-half plane, the Lyapunov matrix equations,

$$\mathbf{A}_i^T \mathbf{P}_i + \mathbf{P}_i \mathbf{A}_i = -\mathbf{Q}_i, \quad i = 1, 2,$$

have unique positive definite solutions \mathbf{P}_1 and \mathbf{P}_2 . Let

$$V(\mathbf{z}) = \mathbf{z}_1^T \mathbf{P}_1 \mathbf{z}_1 - \mathbf{z}_2^T \mathbf{P}_2 \mathbf{z}_2 = \mathbf{z}^T \begin{bmatrix} \mathbf{P}_1 & \mathbf{O} \\ \mathbf{O} & -\mathbf{P}_2 \end{bmatrix} \mathbf{z} \quad (4.66)$$

In the subspace, $\{\mathbf{z} = [\mathbf{z}_1^T \ \mathbf{z}_2^T]^T : \mathbf{z}_2 = \mathbf{0}\}$, $V(\mathbf{z}) > 0$ in every neighborhood of the origin. Let

$$\mathcal{U} = \{\mathbf{z} \in \mathbb{R}^n : \|\mathbf{z}\|_2 \leq r\}.$$

The time derivative of V evaluated on the trajectories of the system (4.64) is

$$\begin{aligned} \dot{V}(\mathbf{z}) &= 2\mathbf{z}_1^T \mathbf{P}_1 \dot{\mathbf{z}}_1 - 2\mathbf{z}_2^T \mathbf{P}_2 \dot{\mathbf{z}}_2 \\ &= 2\mathbf{z}_1^T \mathbf{P}_1 (-\mathbf{A}_1 \mathbf{z}_1 + \mathbf{g}_1(\mathbf{z})) - 2\mathbf{z}_2^T \mathbf{P}_2 (\mathbf{A}_2 \mathbf{z}_2 + \mathbf{g}_2(\mathbf{z})) \\ &= \mathbf{z}_1^T \mathbf{Q}_1 \mathbf{z}_1 + 2\mathbf{z}_1^T \mathbf{P}_1 \mathbf{g}_1(\mathbf{z}) + \mathbf{z}_2^T \mathbf{Q}_2 \mathbf{z}_2 - 2\mathbf{z}_2^T \mathbf{P}_2 \mathbf{g}_2(\mathbf{z}) \\ &\geq \lambda_{\min}(\mathbf{Q}_1) \|\mathbf{z}_1\|_2^2 + \lambda_{\min}(\mathbf{Q}_2) \|\mathbf{z}_2\|_2^2 + 2\mathbf{z}^T \begin{bmatrix} \mathbf{P}_1 \mathbf{g}_1(\mathbf{z}) \\ -\mathbf{P}_2 \mathbf{g}_2(\mathbf{z}) \end{bmatrix} \\ &\geq \lambda_{\min}(\mathbf{Q}_1) \|\mathbf{z}_1\|_2^2 + \lambda_{\min}(\mathbf{Q}_2) \|\mathbf{z}_2\|_2^2 \\ &\quad - 2\|\mathbf{z}\|_2 \sqrt{\|\mathbf{P}_1\|_2^2 \|\mathbf{g}_1(\mathbf{z})\|_2^2 + \|\mathbf{P}_2\|_2^2 \|\mathbf{g}_2(\mathbf{z})\|_2^2}. \end{aligned} \quad (4.67)$$

Let

$$\alpha = \min(\lambda_{\min}(\mathbf{Q}_1), \lambda_{\min}(\mathbf{Q}_2)) \quad \text{and} \quad \beta = \max(\|\mathbf{P}_1\|_2, \|\mathbf{P}_2\|_2).$$

With (4.65) and the above taken into account, (4.67) evaluates to

$$\dot{V}(\mathbf{z}) \geq (\alpha - 2\sqrt{2}\beta\gamma) \|\mathbf{z}\|_2^2 \quad \text{for all } \mathbf{z} \in \mathcal{U}. \quad (4.68)$$

If we select

$$\gamma < \frac{\alpha}{2\sqrt{2}\beta},$$

then $\dot{V}(\mathbf{z}) > 0$ in $\mathcal{U} \setminus \{0\}$. Thus, all hypotheses of Theorem 4.10 are satisfied and hence the origin is an unstable equilibrium of the system (4.57).

In the general case when \mathbf{A} has eigenvalues on the imaginary axis, in addition to eigenvalues in the open right-half plane, we shift the imaginary axis in order to reduce this case to the previous one. Specifically, suppose that \mathbf{A} has m eigenvalues in the open right-half plane with $\Re \lambda_i > \delta > 0$. Then, the matrix

$$\mathbf{A} - \frac{\delta}{2} \mathbf{I}$$

still has m eigenvalues in open right-half plane, but no eigenvalues on the imaginary axis. We can proceed similarly as in the previous case to conclude that the origin is an unstable equilibrium of the corresponding nonlinear system (4.57). For details of this part of the proof, we refer to Khalil [157, p. 130].

In summary, the previous two theorems provide us with a simple method for determining asymptotic stability or instability of an equilibrium state of the nonlinear system (4.57). We compute the Jacobian matrix (4.58) of the nonlinear system, evaluate it at the equilibrium state of interest, and test its eigenvalues. If $\Re \lambda_i < 0$ for all i , then the equilibrium state is asymptotically stable. If for some i , $\Re \lambda_i > 0$, then the equilibrium is unstable. Finally, if $\Re \lambda_i \leq 0$ for all i with $\Re \lambda_i = 0$ for some i , linearization fails to determine asymptotic stability of the equilibrium of the nonlinear system. In this case, the stability of the nonlinear system depends on the nature of the higher-order terms.

◆ Example 4.8

For the nonlinear system model

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} -1.1x_1 + g_2(x_2) \\ g_1(x_1) - 1.1x_2 \end{bmatrix}, \quad (4.69)$$

where

$$g_i(x_i) = \frac{2}{\pi} \arctan(0.7\pi x_i), \quad i = 1, 2,$$

we will find its equilibrium states and determine their stability properties using the linearization method. We mention here that the above system was analyzed in references 126 and 203.

To find the system equilibrium states, we solve the algebraic equations $\dot{x}_1 = 0$ and $\dot{x}_2 = 0$, that is,

$$\begin{aligned} -1.1x_1 + \frac{2}{\pi} \arctan(0.7\pi x_2) &= 0, \\ \frac{2}{\pi} \arctan(0.7\pi x_1) - 1.1x_2 &= 0. \end{aligned}$$

We obtain three equilibrium states:

$$\mathbf{x}^{(1)} = \begin{bmatrix} 0.454 \\ 0.454 \end{bmatrix}, \quad \mathbf{x}^{(2)} = \begin{bmatrix} -0.454 \\ -0.454 \end{bmatrix}, \quad \text{and} \quad \mathbf{x}^{(3)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

Let

$$\Delta \mathbf{x} = \mathbf{x} - \mathbf{x}^{(i)}, \quad i = 1, 2, 3.$$

Then, a linearized model about an equilibrium state has the form

$$\frac{d}{dt} \Delta \mathbf{x} = \mathbf{A} \Delta \mathbf{x},$$

where

$$\mathbf{A} = \left[\begin{array}{cc} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} \end{array} \right] \bigg|_{\mathbf{x}=\mathbf{x}^{(i)}} = \left[\begin{array}{cc} -1.1 & \frac{1.4}{1 + (0.7\pi x_2)^2} \\ \frac{1.4}{1 + (0.7\pi x_2)^2} & -1.1 \end{array} \right] \bigg|_{\mathbf{x}=\mathbf{x}^{(i)}}. \quad (4.70)$$

We have three linearized models:

1. About $\mathbf{x}^{(1)}$,

$$\frac{d}{dt} \Delta \mathbf{x} = \begin{bmatrix} -1.1 & 0.7073 \\ 0.7073 & -1.1 \end{bmatrix} \Delta \mathbf{x},$$

2. About $\mathbf{x}^{(2)}$,

$$\frac{d}{dt} \Delta \mathbf{x} = \begin{bmatrix} -1.1 & 0.7073 \\ 0.7073 & -1.1 \end{bmatrix} \Delta \mathbf{x},$$

3. About $\mathbf{x}^{(3)}$,

$$\frac{d}{dt} \Delta \mathbf{x} = \begin{bmatrix} -1.1 & 1.4 \\ 1.4 & -1.1 \end{bmatrix} \Delta \mathbf{x}.$$

By Theorem 4.11, the states $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$ are asymptotically stable because the corresponding linearized models have their eigenvalues in the open left-half plane: -1.8073 and -0.3927 . By Theorem 4.12, the equilibrium state $\mathbf{x}^{(3)}$ is unstable because the linearized model has an eigenvalue in the open right-half plane, at 0.3 (the other eigenvalue is at -2.5). In Figure 4.12, we show a phase portrait of the above nonlinear system, where one can see clearly the stability properties of the equilibrium states.

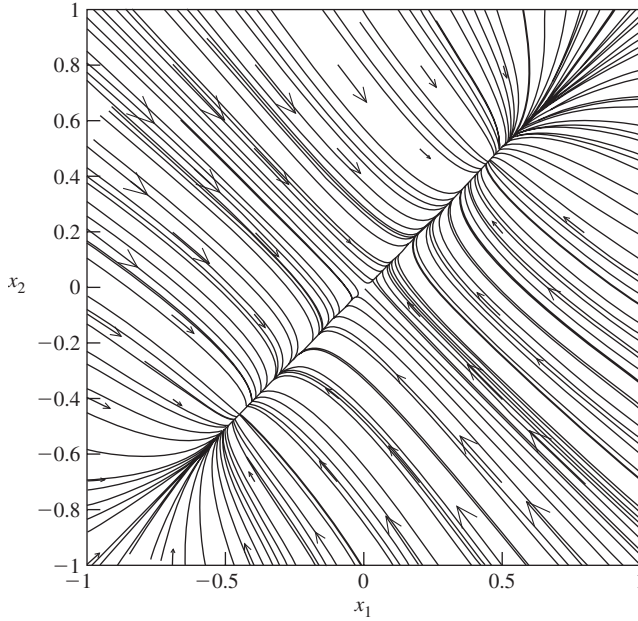


Figure 4.12 Phase portrait of the nonlinear system of Example 4.8.

4.10 Discontinuous Robust Controllers

In this section we consider the problem of constructing stabilizing state-feedback controllers for nonlinear dynamical systems whose mathematical models contain uncertainties. The design of the controllers is based on the Lyapunov theory presented in the previous sections. The class of dynamical systems considered here is modeled by

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t)) + \mathbf{G}(t, \mathbf{x}(t)) (\mathbf{u}(t) + \boldsymbol{\xi}(t, \mathbf{x}(t))), \quad \mathbf{x}(t_0) = \mathbf{x}_0, \quad (4.71)$$

where $\mathbf{x}(t) \in \mathbb{R}^n$, $\mathbf{u}(t) \in \mathbb{R}^m$, and $\boldsymbol{\xi}(t, \mathbf{x}(t)) \in \mathbb{R}^m$ represent the lumped uncertainties. We assume that the uncertain element $\boldsymbol{\xi}$ is bounded by a known continuous nonnegative function ρ , that is,

$$\|\boldsymbol{\xi}(t, \mathbf{x}(t))\| \leq \rho(t, \mathbf{x}(t)).$$

We also assume that the equilibrium state $\mathbf{x} = \mathbf{0}$ of the uncontrolled nominal system

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t)) \quad (4.72)$$

is globally uniformly asymptotically stable and that we were able to construct a Lyapunov function $W : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}$ corresponding to this equilibrium state. The Lyapunov function W is assumed to satisfy the following conditions: There exist three continuous and strictly increasing

functions $\psi_i, i = 1, 2, 3$, such that

$$\begin{aligned}\psi_i(0) &= 0, & i &= 1, 2, 3, \\ \lim_{r \rightarrow \infty} \psi_1(r) &= \infty, \\ \psi_1(\|\mathbf{x}\|) &\leq W(t, \mathbf{x}) \leq \psi_2(\|\mathbf{x}\|), \\ \psi_3(\|\mathbf{x}\|) &\leq -\dot{W}(t, \mathbf{x})\end{aligned}\tag{4.73}$$

for all t and all \mathbf{x} . We will show that if we apply the control strategy

$$\mathbf{u}(t, \mathbf{x}) = \begin{cases} -\frac{\boldsymbol{\mu}(t, \mathbf{x})}{\|\boldsymbol{\mu}(t, \mathbf{x})\|} \rho(t, \mathbf{x}) & \text{if } (t, \mathbf{x}) \notin \Omega \\ \mathbf{p} \in \mathbb{R}^m, \|\mathbf{p}\| \leq \rho(t, \mathbf{x}) & \text{if } (t, \mathbf{x}) \in \Omega, \end{cases}\tag{4.74}$$

where $\boldsymbol{\mu}(t, \mathbf{x}) = \mathbf{G}^T(t, \mathbf{x}) \nabla_{\mathbf{x}} W(t, \mathbf{x})$ and $\Omega = \{(t, \mathbf{x}) : \boldsymbol{\mu} = \mathbf{0}\}$, then the equilibrium state $\mathbf{x} = \mathbf{0}$ is a globally uniformly asymptotically stable solution of the closed-loop system. For this, consider the Lyapunov function W for the uncontrolled nominal system. The time derivative of W evaluated on the solutions of the closed-loop system is

$$\begin{aligned}\dot{W}(t, \mathbf{x}) &= \frac{\partial W}{\partial t} + \nabla_{\mathbf{x}} W^T (\mathbf{f}(t, \mathbf{x}) + \mathbf{G}(t, \mathbf{x}) (\mathbf{u} + \boldsymbol{\xi}(t, \mathbf{x}))) \\ &\leq -\psi_3(\|\mathbf{x}\|) + \nabla_{\mathbf{x}} W^T \mathbf{G}(t, \mathbf{x}) \mathbf{u} + \nabla_{\mathbf{x}} W^T \mathbf{G}(t, \mathbf{x}) \boldsymbol{\xi}(t, \mathbf{x}) \\ &\leq -\psi_3(\|\mathbf{x}\|) + \boldsymbol{\mu}^T(t, \mathbf{x}) \mathbf{u} + \|\boldsymbol{\mu}(t, \mathbf{x})\| \rho(t, \mathbf{x}).\end{aligned}\tag{4.75}$$

For $(t, \mathbf{x}) \notin \Omega$,

$$\begin{aligned}\dot{W}(t, \mathbf{x}) &\leq -\psi_3(\|\mathbf{x}\|) - \boldsymbol{\mu}^T(t, \mathbf{x}) \frac{\boldsymbol{\mu}(t, \mathbf{x})}{\|\boldsymbol{\mu}(t, \mathbf{x})\|} \rho(t, \mathbf{x}) + \|\boldsymbol{\mu}(t, \mathbf{x})\| \rho(t, \mathbf{x}) \\ &= -\psi_3(\|\mathbf{x}\|) - \frac{\|\boldsymbol{\mu}(t, \mathbf{x})\|^2}{\|\boldsymbol{\mu}(t, \mathbf{x})\|} \rho(t, \mathbf{x}) + \|\boldsymbol{\mu}(t, \mathbf{x})\| \rho(t, \mathbf{x}) \\ &= -\psi_3(\|\mathbf{x}\|) - \|\boldsymbol{\mu}(t, \mathbf{x})\| \rho(t, \mathbf{x}) + \|\boldsymbol{\mu}(t, \mathbf{x})\| \rho(t, \mathbf{x}) \\ &= -\psi_3(\|\mathbf{x}\|) \\ &< 0.\end{aligned}\tag{4.76}$$

For $(t, \mathbf{x}) \in \Omega$, we have $\boldsymbol{\mu}(t, \mathbf{x}) = \mathbf{0}$; and using (4.75); we obtain

$$\dot{W}(t, \mathbf{x}) \leq -\psi_3(\|\mathbf{x}\|) < 0.\tag{4.77}$$

Thus, W is also a Lyapunov function for the closed-loop system, which means that the uncertain system (4.71) driven by the control law (4.74) is indeed globally uniformly asymptotically stable.

◆ Example 4.9

For the uncertain system model

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \boldsymbol{\xi}_1(\mathbf{x}) + \mathbf{b}(u + \boldsymbol{\xi}_2(\mathbf{x})),$$

where

$$\mathbf{A} = \begin{bmatrix} -2 & 1 \\ 0 & -3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad \|\xi_1(\mathbf{x})\| \leq 1.75\|\mathbf{x}\|, \quad \text{and} \quad |\xi_2| \leq 10\|\mathbf{x}\| + 0.5|u|,$$

design stabilizing control strategy of the form (4.74).

The uncontrolled nominal system $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$ is asymptotically stable. Let

$$\mathbf{Q} = \begin{bmatrix} 4 & -1 \\ -1 & 6 \end{bmatrix}.$$

The eigenvalues of \mathbf{Q} are $\lambda_{\min}(\mathbf{Q}) = 3.5858$ and $\lambda_{\max}(\mathbf{Q}) = 6.4142$. Solving the Lyapunov matrix equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{Q}$ yields

$$\mathbf{P} = \mathbf{I}_2.$$

The proposed controller has the form

$$u = -\rho(\mathbf{x}) \frac{\mu(\mathbf{x})}{|\mu(\mathbf{x})|} = -\rho(\mathbf{x}) \frac{\mathbf{b}^T \mathbf{P} \mathbf{x}}{|\mathbf{b}^T \mathbf{P} \mathbf{x}|} = -\rho(\mathbf{x}) \frac{\mathbf{b}^T \mathbf{x}}{|\mathbf{b}^T \mathbf{x}|} = -\rho(\mathbf{x}) \text{sign}(\mathbf{b}^T \mathbf{x}). \quad (4.78)$$

Our goal is to determine a bound $\rho^*(\mathbf{x}) > 0$ such that for any $\rho(\mathbf{x}) \geq \rho^*(\mathbf{x})$ the uncertain closed-loop system is globally uniformly asymptotically stable. To begin, consider

$$V(\mathbf{x}) = \mathbf{x}^T \mathbf{P} \mathbf{x} = \mathbf{x}^T \mathbf{x}$$

as the Lyapunov function candidate for the closed-loop system. Evaluating the time derivative of V on the trajectories of the closed-loop system yields

$$\begin{aligned} \frac{d}{dt} V(\mathbf{x}) &= 2\mathbf{x}^T \dot{\mathbf{x}} \\ &= 2\mathbf{x}^T (\mathbf{A}\mathbf{x} + \xi_1(\mathbf{x}) + \mathbf{b}(u + \xi_2(\mathbf{x}))) \\ &= -\mathbf{x}^T \mathbf{Q} \mathbf{x} + 2\mathbf{x}^T \xi_1(\mathbf{x}) + 2\mathbf{x}^T \mathbf{b} u(\mathbf{x}) + 2\mathbf{x}^T \mathbf{b} \xi_2(\mathbf{x}). \end{aligned} \quad (4.79)$$

Substituting into the above the expression for the controller given by (4.78), taking the norms, and using bounds on the uncertain elements, we obtain

$$\begin{aligned} \frac{d}{dt} V(\mathbf{x}) &\leq -\lambda_{\min}(\mathbf{Q})\|\mathbf{x}\|^2 + 2\|\mathbf{x}\|\|\xi_1(\mathbf{x})\| - 2\rho(\mathbf{x}) \frac{|\mathbf{x}^T \mathbf{b}|^2}{|\mathbf{x}^T \mathbf{b}|} + 2|\mathbf{x}^T \mathbf{b}|\|\xi_2(\mathbf{x})\| \\ &\leq (-3.5858 + 3.5)\|\mathbf{x}\|^2 - 2\rho(\mathbf{x})|\mathbf{x}^T \mathbf{b}| + 2(10\|\mathbf{x}\| + 0.5\rho(\mathbf{x}))|\mathbf{x}^T \mathbf{b}|. \end{aligned} \quad (4.80)$$

In the above, we used the fact that

$$|u| = \rho(\mathbf{x}).$$

Collecting the terms in (4.80) yields

$$\frac{d}{dt} V(\mathbf{x}) \leq (-3.5858 + 3.5)\|\mathbf{x}\|^2 + (20\|\mathbf{x}\| - \rho(\mathbf{x}))|\mathbf{x}^T \mathbf{b}|. \quad (4.81)$$

Therefore, $\frac{dV}{dt} < 0$ for any

$$\rho(\mathbf{x}) \geq \rho^*(\mathbf{x}) = 20\|\mathbf{x}\|. \quad (4.82)$$

In summary, the uncertain closed-loop system will be globally uniformly asymptotically stable if the controller gain $\rho(\mathbf{x})$ is chosen to satisfy (4.82).

We now use the Lyapunov approach to design a model-reference tracking controller. A block diagram of the model-reference control system is shown in Figure 4.13. Model-reference control can be used for plants with nonlinearities and time-varying parameters. One of the assumptions that is made while designing a model-reference controller is that a plant model is in a controller companion form. Here, we propose a method for designing model-reference controllers for a class of multi-input systems that does not require a plant model to be in a specific form. The crux of our method is in an appropriate selection of the model-reference system. Consider a plant model of the form

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}u, \quad (4.83)$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times m}$. The model-reference signal generator is described by

$$\dot{\mathbf{x}}_d = \mathbf{A}_d\mathbf{x}_d + \mathbf{B}_d\mathbf{r}, \quad (4.84)$$

where $\mathbf{A}_d \in \mathbb{R}^{n \times n}$, $\mathbf{B}_d \in \mathbb{R}^{n \times m}$, and $\mathbf{r} \in \mathbb{R}^m$ is a command signal that is chosen so that \mathbf{x}_d represents a desired state trajectory of the plant. The matrix \mathbf{A}_d of the reference system is chosen from the set

$$\{\mathbf{A}_d : \mathbf{A}_d = \mathbf{A} - \mathbf{B}\mathbf{K}_d, \mathbf{K}_d \in \mathbb{R}^{m \times n}\}, \quad (4.85)$$

and $\mathbf{B}_d = \mathbf{B}$. Thus, the reference system is state-feedback equivalent with the plant model (4.83). We assume that the pair (\mathbf{A}, \mathbf{B}) is controllable. Therefore, we can select \mathbf{A}_d to have its eigenvalues in prescribed locations in the open left-hand plane. To proceed further, we define the tracking error \mathbf{e} as

$$\mathbf{e}(t) = \mathbf{x}(t) - \mathbf{x}_d(t). \quad (4.86)$$

Our goal is to construct a state-feedback control law so that

$$\lim_{t \rightarrow \infty} \|\mathbf{e}(t)\| = 0. \quad (4.87)$$

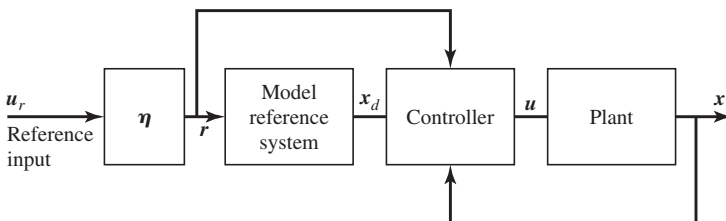


Figure 4.13 Model-reference control system block diagram with input scaling.

We now form the differential equation that describes the error dynamics. For this we evaluate the time derivative of the error vector to get

$$\begin{aligned}\dot{\mathbf{e}}(t) &= \dot{\mathbf{x}}(t) - \dot{\mathbf{x}}_d(t) \\ &= \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u} - \mathbf{A}_d\mathbf{x}_d(t) - \mathbf{B}\mathbf{r}(t).\end{aligned}\quad (4.88)$$

Recall that $\mathbf{A}_d = \mathbf{A} - \mathbf{B}\mathbf{K}_d$, hence $\mathbf{A} = \mathbf{A}_d + \mathbf{B}\mathbf{K}_d$. Substituting this expression for \mathbf{A} into the error equation (4.88), we obtain

$$\begin{aligned}\dot{\mathbf{e}}(t) &= \mathbf{A}_d(\mathbf{x}(t) - \mathbf{x}_d(t)) + \mathbf{B}(\mathbf{u} + \mathbf{K}_d\mathbf{x}(t) - \mathbf{r}(t)) \\ &= \mathbf{A}_d\mathbf{e}(t) + \mathbf{B}(\mathbf{u} + \mathbf{K}_d\mathbf{x}(t) - \mathbf{r}(t)).\end{aligned}\quad (4.89)$$

We will use a Lyapunov function approach in our further analysis that will result in a control law forcing the error $\mathbf{e}(t)$ to go asymptotically to zero. As a Lyapunov function candidate for (4.89), we take

$$V(\mathbf{e}) = \mathbf{e}^T \mathbf{P} \mathbf{e}, \quad (4.90)$$

where \mathbf{P} solves the Lyapunov matrix equation, $\mathbf{A}_d^T \mathbf{P} + \mathbf{P} \mathbf{A}_d = -\mathbf{Q}$, for some $\mathbf{Q} = \mathbf{Q}^T > 0$. Because by construction \mathbf{A}_d is asymptotically stable, the solution to the Lyapunov matrix equation is positive definite, that is, $\mathbf{P} = \mathbf{P}^T > 0$. We next evaluate the time derivative of $V(\mathbf{e}(t))$ on the trajectories of (4.89),

$$\begin{aligned}\dot{V}(\mathbf{e}) &= 2\mathbf{e}^T \mathbf{P} \dot{\mathbf{e}} \\ &= 2\mathbf{e}^T \mathbf{P}(\mathbf{A}_d\mathbf{e} + \mathbf{B}(\mathbf{u} + \mathbf{K}_d\mathbf{x} - \mathbf{r})) \\ &= \mathbf{e}^T (\mathbf{A}_d^T \mathbf{P} + \mathbf{P} \mathbf{A}_d) \mathbf{e} + 2\mathbf{e}^T \mathbf{P} \mathbf{B}(\mathbf{u} + \mathbf{K}_d\mathbf{x} - \mathbf{r}).\end{aligned}\quad (4.91)$$

Using the Lyapunov matrix equation and taking the norms of the above yields

$$\dot{V}(\mathbf{e}) \leq -\mathbf{e}^T \mathbf{Q} \mathbf{e} + 2\mathbf{e}^T \mathbf{P} \mathbf{B} \mathbf{u} + 2\|\mathbf{e}^T \mathbf{P} \mathbf{B}\| \|\mathbf{K}_d \mathbf{x}\| + 2\|\mathbf{e}^T \mathbf{P} \mathbf{B}\| \|\mathbf{r}\|. \quad (4.92)$$

If \mathbf{e} is such that $\mathbf{B}^T \mathbf{P} \mathbf{e} = \mathbf{0}$, then

$$\dot{V}(\mathbf{e}) \leq -\mathbf{e}^T \mathbf{Q} \mathbf{e} < 0. \quad (4.93)$$

We now analyze the case when $\mathbf{B}^T \mathbf{P} \mathbf{e} \neq \mathbf{0}$. Let the proposed control law be

$$\mathbf{u} = -\rho(\mathbf{x}, \mathbf{r}) \frac{\mathbf{B}^T \mathbf{P} \mathbf{e}}{\|\mathbf{B}^T \mathbf{P} \mathbf{e}\|}, \quad (4.94)$$

where $\rho(\mathbf{x}, \mathbf{r}) = \|\mathbf{K}_d \mathbf{x}\| + \|\mathbf{r}\|$. Substituting (4.94) into (4.92), we obtain

$$\begin{aligned}\dot{V}(\mathbf{e}) &\leq -\mathbf{e}^T \mathbf{Q} \mathbf{e} - 2\|\mathbf{e}^T \mathbf{P} \mathbf{B}\| (\rho(\mathbf{x}, \mathbf{r}) - \|\mathbf{K}_d \mathbf{x}\| - \|\mathbf{r}\|) \\ &= -\mathbf{e}^T \mathbf{Q} \mathbf{e} \\ &< 0.\end{aligned}$$

Since for all $\mathbf{e} \in \mathbb{R}^n$ we have $V(\mathbf{e}) > 0$ and $\dot{V}(\mathbf{e}) < 0$, the equilibrium point $\mathbf{e} = \mathbf{0}$ of the error equation (4.89) is globally uniformly asymptotically stable, which means that $\mathbf{x}(t) \rightarrow \mathbf{x}_d(t)$.

◆ Example 4.10

We now apply the above tracking controller to the vehicle model derived in Subsection 1.7.2, where $\delta_r = 0$. That is, we wish to construct a tracking controller for a vehicle with front-wheel steering. The first step is to select the poles of the model reference system. We choose the model reference system to have the following set of asymptotically stable poles:

$$\{-4 - 2j, -4 + 2j, -4.5, -5.0\}. \quad (4.95)$$

The pair (\mathbf{A}, \mathbf{b}) is controllable; hence we can find \mathbf{k}_d such that the eigenvalues of the matrix $\mathbf{A}_d = \mathbf{A} - \mathbf{b}\mathbf{k}_d$ are at the locations in (4.95). The feedback gain is

$$\mathbf{k}_d = [0.0329 \quad 2.2854 \quad 0.2429 \quad -0.1653], \quad (4.96)$$

and the resulting model reference plant matrix is

$$\mathbf{A}_d = \begin{bmatrix} -6.6638 & -107.1281 & -29.6350 & 7.7479 \\ 0 & 0 & 1 & 0 \\ -0.9205 & -65.8195 & -10.8362 & 4.7603 \\ -1 & -18.3 & 0 & 0 \end{bmatrix}. \quad (4.97)$$

The next step is to find \mathbf{P} that solves the Lyapunov equation, $\mathbf{A}_d^T \mathbf{P} + \mathbf{P} \mathbf{A}_d = -\mathbf{Q}$. We take $\mathbf{Q} = \mathbf{I}_4$. Solving the Lyapunov equation, we get

$$\mathbf{P} = \begin{bmatrix} 0.1145 & 0.7530 & -0.1809 & -0.0966 \\ 0.7530 & 34.6174 & -0.2883 & -3.3437 \\ -0.1809 & -0.2883 & 0.5144 & 0.0522 \\ -0.0966 & -3.3437 & 0.0522 & 0.6216 \end{bmatrix}. \quad (4.98)$$

Now that we have \mathbf{P} , we are ready to construct the controller. We have

$$\begin{aligned} u &= -\rho(\mathbf{x}, r) \frac{\mathbf{b}^T \mathbf{P} \mathbf{e}}{|\mathbf{b}^T \mathbf{P} \mathbf{e}|} \\ &= -(|\mathbf{k}_d \mathbf{x}| + |r|) \text{sign}(\mathbf{b}^T \mathbf{P}(\mathbf{x} - \mathbf{x}_d)) \\ &= -\left(\left\| \begin{bmatrix} 0.0329 \\ 2.2854 \\ 0.2429 \\ -0.1653 \end{bmatrix}^T \mathbf{x} \right\| + |r| \right) \text{sign} \left(\begin{bmatrix} 0.1571 \\ 26.9931 \\ 6.3324 \\ -3.0250 \end{bmatrix}^T (\mathbf{x} - \mathbf{x}_d) \right). \end{aligned} \quad (4.99)$$

Our goal is to have the reference input signal u_r to correspond to a desired vehicle path. This requires scaling the reference input. The scaling factor is

$$\eta = \left(\lim_{s \rightarrow 0} \mathbf{c}(s\mathbf{I}_4 - \mathbf{A}_d)^{-1} \mathbf{b} \right)^{-1}, \quad (4.100)$$

where $\mathbf{c} = [0 \ 0 \ 0 \ 1]$. Note that η is well-defined because the matrix \mathbf{A}_d is asymptotically stable. In our example, $\eta = -0.1653$. Thus, we have $r = \eta u_r$. The reference input is $u_r = 4h(t - 1 \text{ sec})$, where $h(t)$ is the unit step function.

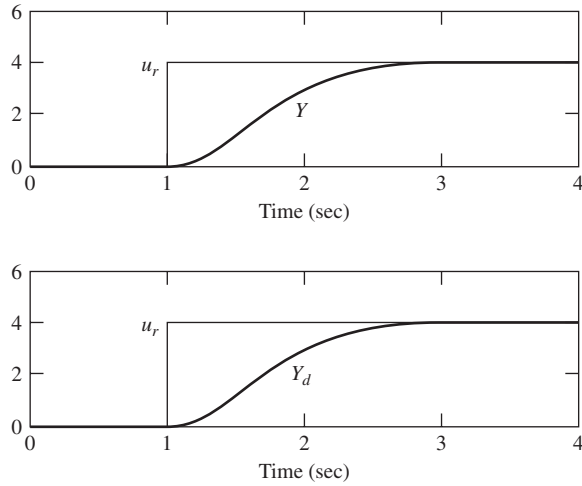


Figure 4.14 Plots of the actual output of the closed-loop system and the desired output in Example 4.10.

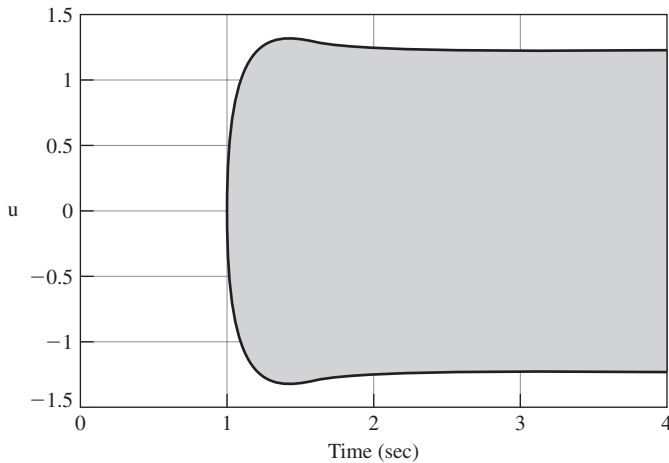


Figure 4.15 Plot of the controller action versus time in Example 4.10.

The performance of the controller, for zero initial conditions, is shown in Figures 4.14 and 4.15. In Figure 4.14, the closed-loop system output is shown in the upper subplot, while the desired output is shown in the bottom subplot. The vehicle tracks the desired path with no noticeable steady-state error. The control action is shown in Figure 4.15.

The proposed tracking controller is discontinuous in the state, which results in high activity referred to as chattering. Chattering can be reduced by using a continuous version of the above controller. The continuous controller is obtained by introducing a so-called boundary layer that smooths out the discontinuous controller. We discuss the problem of constructing a boundary layer controller in the following section.

4.11 Uniform Ultimate Boundedness

A drawback to the discontinuous controllers, discussed in the previous section, is that rapid switching tends to excite high-frequency modes of the plant, and this may degrade the performance of the controller. These problems can be alleviated by adding a so-called *boundary layer* to the discontinuous controller resulting in a continuous control strategy. However, there is a price for making a discontinuous controller a continuous one. The closed-loop system driven by a boundary layer controller may only be *uniformly ultimately bounded* rather than asymptotically stable. The notions we are concerned with in this section are described in the following three definitions.

Definition 4.18 The closed ball $\bar{B}_q(\mathbf{0}) = \{\mathbf{x} : \|\mathbf{x}\| \leq q\}$ is uniformly stable if for any $\varepsilon > q$, there is a $\delta = \delta(\varepsilon) > 0$ so that if $\|\mathbf{x}(t_0)\| \leq \delta$, then $\|\mathbf{x}(t)\| \leq \varepsilon$ for $t \geq t_0$.

For the benefit of the reader, we now rephrase the definition of uniform boundedness from Section 4.8.

Definition 4.19 The solutions of $\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x})$ are uniformly bounded if for any given $d > 0$, there exists $b = b(d) < \infty$ so that for each solution $\mathbf{x}(t) = \mathbf{x}(t; t_0, \mathbf{x}_0)$ with $\|\mathbf{x}_0\| \leq d$, we have $\|\mathbf{x}(t)\| \leq b$ for all $t \geq t_0$.

We next define the notion of uniform ultimate boundedness (UUB) of a system.

Definition 4.20 The solutions of $\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x})$ are uniformly ultimately bounded (UUB) with respect to the ball $\bar{B}_q(\mathbf{0})$ if for every $d > 0$, there exists $T(q, d)$ such that for each solution $\mathbf{x}(t) = \mathbf{x}(t; t_0, \mathbf{x}_0)$ with $\|\mathbf{x}_0\| \leq d$, we have $\mathbf{x}(t) \in \bar{B}_q(\mathbf{0})$ for $t \geq t_0 + T(q, d)$.

We now consider a dynamical system model (4.71):

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t)) + \mathbf{G}(t, \mathbf{x}(t))(\mathbf{u}(t) + \boldsymbol{\xi}(t, \mathbf{x}(t))), \quad \mathbf{x}(t_0) = \mathbf{x}_0,$$

where $\mathbf{x}(t) \in \mathbb{R}^n$, $\mathbf{u}(t) \in \mathbb{R}^m$, and $\boldsymbol{\xi}(t, \mathbf{x}(t)) \in \mathbb{R}^m$ represents the lumped uncertainties. We assume that the uncertain element $\boldsymbol{\xi}$ is bounded by a known continuous nonnegative function ρ , that is, $\|\boldsymbol{\xi}(t, \mathbf{x}(t))\| \leq \rho(t, \mathbf{x}(t))$. As in Section 4.10, we assume that the equilibrium state $\mathbf{x} = \mathbf{0}$ of the uncontrolled nominal system, $\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t))$, is globally uniformly asymptotically stable and that we were able to construct a Lyapunov function $W : \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying the conditions (4.73) for all t and all \mathbf{x} .

We showed in the previous section that if we apply the discontinuous control strategy

$$\mathbf{u}(t, \mathbf{x}) = \begin{cases} -\frac{\boldsymbol{\mu}(t, \mathbf{x})}{\|\boldsymbol{\mu}(t, \mathbf{x})\|} \rho(t, \mathbf{x}) & \text{if } \boldsymbol{\mu}(t, \mathbf{x}) \neq \mathbf{0} \\ \mathbf{p} \in \mathbb{R}^m, \quad \|\mathbf{p}\| \leq \rho(t, \mathbf{x}) & \text{if } \boldsymbol{\mu}(t, \mathbf{x}) = \mathbf{0}, \end{cases}$$

where $\boldsymbol{\mu}(t, \mathbf{x}) = \mathbf{G}^T(t, \mathbf{x}) \nabla_{\mathbf{x}} W(t, \mathbf{x})$, then the equilibrium state $\mathbf{x} = \mathbf{0}$ is a globally uniformly asymptotically stable solution of the closed-loop system. However, the above controller is

discontinuous in \mathbf{x} . Following Corless and Leitmann [56], we modify it to make it continuous by introducing a boundary layer. The proposed boundary layer controller is

$$\mathbf{u}(t, \mathbf{x}) = \begin{cases} -\frac{\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})}{\|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\|}\rho(t, \mathbf{x}) & \text{if } \|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| > \nu \\ \mathbf{p}(t, \mathbf{x}), \quad \|\mathbf{p}(t, \mathbf{x})\| \leq \rho(t, \mathbf{x}) & \text{if } \|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| \leq \nu, \end{cases} \quad (4.101)$$

where ν is the thickness of the boundary layer about the naturally defined switching surface $\{(t, \mathbf{x}) : \|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| = 0\}$, and \mathbf{p} is a continuous function such that

$$\mathbf{p} = -\frac{\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})}{\|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\|}\rho(t, \mathbf{x}) \quad \text{if } \|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| = \nu.$$

A particular example of the above type of a controller is

$$\mathbf{u}(t, \mathbf{x}) = \begin{cases} -\frac{\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})}{\|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\|}\rho(t, \mathbf{x}) & \text{if } \|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| > \nu \\ -\frac{\boldsymbol{\mu}(t, \mathbf{x})}{\nu}\rho^2(t, \mathbf{x}) & \text{if } \|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| \leq \nu. \end{cases} \quad (4.102)$$

We will now show that the closed-loop system driven by the boundary layer controller (4.101) is UUB with respect to a neighborhood of the zero state. This neighborhood can be made arbitrarily small by letting $\nu \rightarrow 0$. Before we state the results, we first compute the time derivative of a Lyapunov function candidate on the trajectories of the closed-loop system. Consider the Lyapunov function W for the uncontrolled nominal system. The time derivative of W evaluated on the solutions of the closed-loop system is

$$\begin{aligned} \dot{W}(t, \mathbf{x}) &= \frac{\partial W}{\partial t} + \nabla_{\mathbf{x}} W^T (\mathbf{f}(t, \mathbf{x}) + \mathbf{G}(t, \mathbf{x}) (\mathbf{u} + \boldsymbol{\xi}(t, \mathbf{x}))) \\ &\leq -\psi_3(\|\mathbf{x}\|) + \nabla_{\mathbf{x}} W^T \mathbf{G}(t, \mathbf{x}) \mathbf{u} + \|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\|. \end{aligned}$$

We now consider two cases. The first case is when $\|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| > \nu$. Taking into account the controller structure for this case yields

$$\begin{aligned} \dot{W}(t, \mathbf{x}) &\leq -\psi_3(\|\mathbf{x}\|) - \nabla_{\mathbf{x}} W^T \mathbf{G}(t, \mathbf{x}) \frac{\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})}{\|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\|}\rho(t, \mathbf{x}) + \|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| \\ &= -\psi_3(\|\mathbf{x}\|) - \frac{\|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\|^2}{\|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\|} + \|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| \\ &= -\psi_3(\|\mathbf{x}\|). \end{aligned}$$

If, on the other hand, $\|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| \leq \nu$, then by using (4.101) we obtain

$$\begin{aligned} \dot{W}(t, \mathbf{x}) &\leq -\psi_3(\|\mathbf{x}\|) + \nabla_{\mathbf{x}} W^T \mathbf{G}(t, \mathbf{x}) \mathbf{u} + \|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| \\ &\leq -\psi_3(\|\mathbf{x}\|) + \|\nabla_{\mathbf{x}} W^T \mathbf{G}(t, \mathbf{x})\| \|\mathbf{u}\| + \|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| \\ &\leq -\psi_3(\|\mathbf{x}\|) + 2\|\boldsymbol{\mu}(t, \mathbf{x})\rho(t, \mathbf{x})\| \\ &\leq -\psi_3(\|\mathbf{x}\|) + 2\nu. \end{aligned}$$

We conclude that for all t and all \mathbf{x} , we have

$$\dot{W}(t, \mathbf{x}) \leq -\psi_3(\|\mathbf{x}\|) + 2\nu. \quad (4.103)$$

To proceed further, we need an assumption on the function ψ_3 ; if ψ_3 is bounded and

$$\lim_{r \rightarrow \infty} \psi_3(r) = l < \infty, \quad (4.104)$$

then we assume

$$2\nu < l. \quad (4.105)$$

Proposition 4.2 The solutions of the system (4.71) driven by the boundary layer controller (4.101) are uniformly bounded.

Proof Let R be such that

$$\psi_3(R) = 2\nu.$$

Note that R is well-defined by virtue of (4.105) and the fact that ψ_3 is strictly increasing. Because ψ_3 is strictly increasing, we can write

$$R = \psi_3^{-1}(2\nu).$$

Figure 4.16 illustrates the above construction.

For a given $d > 0$, where $\|\mathbf{x}_0\| \leq d$, let

$$\tilde{d} = \max\{d, R\}.$$

We claim that

$$b = b(d) = \psi_1^{-1}(\psi_2(\tilde{d}));$$

that is, the solutions of (4.71) driven by the boundary layer controller (4.101) are uniformly bounded. In other words, if $\|\mathbf{x}_0\| \leq d$, then $\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| \leq b$ for $t \geq t_0$. First note that $\|\mathbf{x}_0\| \leq \tilde{d}$ and that $\psi_1(b) = \psi_2(\tilde{d})$. Hence, $\tilde{d} \leq b$ (see Figure 4.17 for an

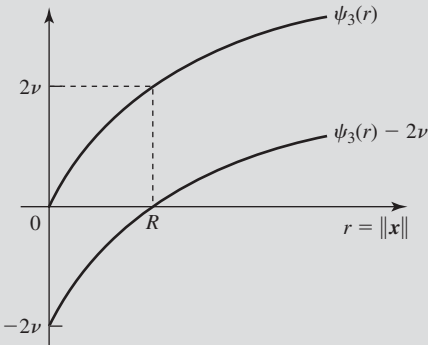


Figure 4.16 Finding R for a given ν in the proof of Proposition 4.2.

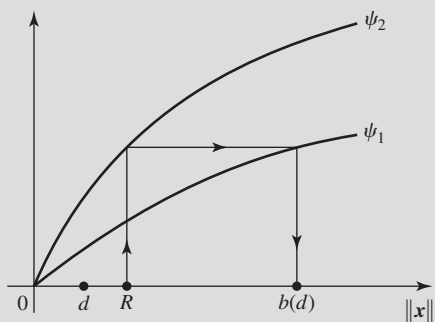


Figure 4.17 Finding b for a given d in the proof of Proposition 4.2.

illustration of this construction), and thus

$$\begin{aligned}\dot{W}(t, \mathbf{x}(t)) &\leq -\psi_3(\|\mathbf{x}(t)\|) + 2\nu \\ &= -\psi_3(\|\mathbf{x}(t)\|) + \psi_3(R).\end{aligned}\quad (4.106)$$

Suppose first that $d > R$. Then for $\|\mathbf{x}_0\| > R$ and $\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| > R$, $\dot{W} < 0$ and therefore

$$\begin{aligned}\psi_1(b) &= \psi_2(\tilde{d}) \\ &\geq W(t_0, \mathbf{x}(t_0; t_0, \mathbf{x}_0)) \\ &\geq W(t, \mathbf{x}(t; t_0, \mathbf{x}_0)) \\ &\geq \psi_1(\|\mathbf{x}(t; t_0, \mathbf{x}_0)\|),\end{aligned}$$

which implies that $\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| \leq b$ for $t \geq t_0$.

If, on the other hand, $\|\mathbf{x}_0\| \leq R$, then because $\dot{W}(t, \mathbf{x}(t)) \leq -\psi_3(\|\mathbf{x}(t)\|) + 2\nu$, the solution may reach the sphere $\|\mathbf{x}\| = R$ at some time $t_2 \geq t_0$. However, by (4.106), $\dot{W} < 0$ outside $\|\mathbf{x}\| = R$. Therefore, $\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| \leq b$. Thus the solution is uniformly bounded.

We illustrate Proposition 4.2 with Figure 4.18. The initial condition in Figure 4.18 is such that $\|\mathbf{x}_0\| < d$ and $d < R$. The trajectory is confined to stay within the ball of radius $b = \psi_1^{-1}(\psi_2(R))$.

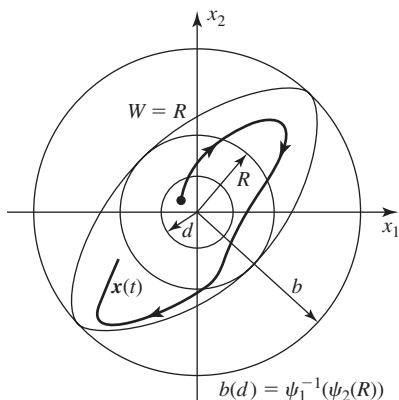


Figure 4.18 A two-dimensional illustration of an uniformly ultimately bounded (UUB) solution.

We now show, using arguments similar to those of Corless and Leitmann [56], that the solutions of the closed-loop system are UUB with respect to any ball of radius $\tilde{b} > \psi_1^{-1}(\psi_2(R))$.

Theorem 4.13 Let $\mathbf{x}(t) = \mathbf{x}(t; t_0, \mathbf{x}_0)$ be a solution of the closed-loop system for $[t_0, \infty)$ such that (4.103) is satisfied. If $\|\mathbf{x}_0\| \leq d$, then for a given

$$\tilde{b} > \psi_1^{-1}(\psi_2(R)), \quad (4.107)$$

$\|\mathbf{x}(t; t_0, \mathbf{x}_0)\| \leq \tilde{b}$ for all $t \geq t_0 + T(\tilde{b}, d)$, where

$$T(\tilde{b}, d) = \begin{cases} 0 & \text{if } d \leq \tilde{R} \\ \frac{\psi_2(d) - \psi_1(\tilde{R})}{\psi_3(\tilde{R}) - 2\nu} & \text{if } d > \tilde{R}, \end{cases}$$

and

$$\psi_2(\tilde{R}) = \psi_1(\tilde{b}). \quad (4.108)$$

Proof From assumption (4.107),

$$\psi_1(\tilde{b}) > \psi_2(R),$$

and by (4.108),

$$\psi_1(\tilde{b}) = \psi_2(\tilde{R}).$$

Hence, by (4.73), where we assumed that ψ_2 is strictly increasing,

$$\tilde{R} > R. \quad (4.109)$$

We consider two cases. The first case is when $d \leq \tilde{R}$. By uniform boundedness, for all $t \geq t_0$,

$$\|\mathbf{x}(t)\| \leq \tilde{b},$$

for any \tilde{b} such that $\psi_1(\tilde{b}) > \psi_2(\tilde{R})$. Hence in this case $T(\tilde{b}, d) = 0$.

We now consider the case when $d > \tilde{R}$. We show, by contradiction, that for some $t_1 \in [t_0, t_0 + T(\tilde{b}, d)]$,

$$\|\mathbf{x}(t_1)\| \leq \tilde{R}.$$

From this we shall conclude the result by referring to the previous case. So suppose that $\|\mathbf{x}(t)\| > \tilde{R}$ for all $t \in [t_0, t_0 + T(\tilde{b}, d)]$. Then,

$$\begin{aligned} \psi_1(\|\mathbf{x}(t_0 + T(\tilde{b}, d); t_0, \mathbf{x}_0)\|) &\leq W(t_0 + T(\tilde{b}, d), \mathbf{x}(t_0 + T(\tilde{b}, d); t_0, \mathbf{x}_0)) \\ &= W(t_0, \mathbf{x}_0) + \int_{t_0}^{t_0 + T(\tilde{b}, d)} \dot{W} dt \\ &\leq W(t_0, \mathbf{x}_0) + \int_{t_0}^{t_0 + T(\tilde{b}, d)} (-\psi_3(\tilde{R}) + 2\nu) dt \\ &\leq \psi_2(d) - \int_{t_0}^{t_0 + T(\tilde{b}, d)} (\psi_3(\tilde{R}) - 2\nu) dt. \end{aligned}$$

Using the fact that $\psi_3(\tilde{R}) - 2\nu > 0$ because by (4.109), $\tilde{R} > R$, we obtain

$$\begin{aligned}\psi_1(\|\mathbf{x}(t_0 + T(\tilde{b}, d); t_0, \mathbf{x}_0)\|) &\leq \psi_2(d) - (\psi_3(\tilde{R}) - 2\nu) T(\tilde{b}, d) \\ &= \psi_2(d) - (\psi_3(\tilde{R}) - 2\nu) \frac{\psi_2(d) - \psi_1(\tilde{R})}{\psi_3(\tilde{R}) - 2\nu} \\ &= \psi_1(\tilde{R}).\end{aligned}$$

Thus, because ψ_1 is strictly increasing, we obtain

$$\|\mathbf{x}(t_0 + T(\tilde{b}, d); t_0, \mathbf{x}_0)\| \leq \tilde{R},$$

which contradicts our assumption that for all $t \in [t_0, t_0 + T(\tilde{b}, d)]$ it should have been $\|\mathbf{x}(t)\| > \tilde{R}$. Hence, there exists $t_1 \in [t_0, t_0 + T(\tilde{b}, d)]$ such that $\|\mathbf{x}(t_1)\| \leq \tilde{R}$. Then using the same argument as in the previous case we can show that

$$\|\mathbf{x}(t)\| \leq \tilde{b}$$

for all $t \geq t_1$ and the proof is complete.

We have the following consequence of the above theorem.

Corollary 4.2 The ball

$$\{\mathbf{x} : \|\mathbf{x}\| \leq \psi_1^{-1}(\psi_2(R))\}$$

is uniformly stable.

Proof Let

$$\varepsilon > \psi_1^{-1}(\psi_2(R)),$$

and consider $\delta = \delta(\varepsilon)$ obtained from

$$\psi_1(\varepsilon) = \psi_2(\delta).$$

Obviously $\delta \geq R$. Figure 4.19 illustrates the above relations. Hence, by uniform boundedness if $\|\mathbf{x}_0\| \leq \delta$, then $\|\mathbf{x}(t)\| \leq \varepsilon$ for all $t \geq t_0$, which means that the ball is uniformly stable.

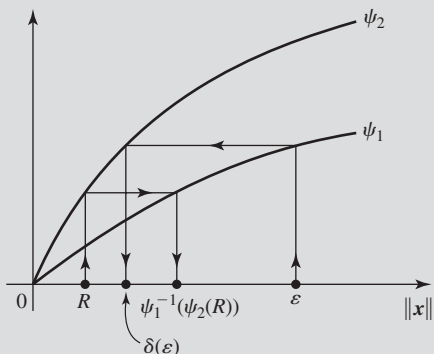


Figure 4.19 Finding δ given ε for an uniformly stable ball in the proof of Corollary 4.2.

◆ **Example 4.11**

Consider a dynamical system model

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b}(u + h(t, \mathbf{x})),$$

where

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ -1 & -1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad \text{and} \quad |h(t, \mathbf{x})| \leq 2.$$

Let \mathbf{P} be the solution of the matrix Lyapunov equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -2\mathbf{I}_2$, and let

$$u = \begin{cases} -2 \operatorname{sign}(\mathbf{b}^T \mathbf{P} \mathbf{x}) & \text{if } |\mathbf{b}^T \mathbf{P} \mathbf{x}| \geq \nu \\ -2 \frac{\mathbf{b}^T \mathbf{P} \mathbf{x}}{\nu} & \text{if } |\mathbf{b}^T \mathbf{P} \mathbf{x}| < \nu, \end{cases}$$

where $\nu = 1$. Suppose that the initial condition, $\mathbf{x}(0)$, is such that $\|\mathbf{x}(0)\|_2 \leq 10$. Find the radius, b , of a ball such that $\|\mathbf{x}(t)\|_2 \leq b$ for all $t \geq 0$.

Solving the Lyapunov equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -2\mathbf{I}_2$ yields

$$\mathbf{P} = \begin{bmatrix} 3 & 1 \\ 1 & 2 \end{bmatrix}.$$

We next evaluate $d/dt(\mathbf{x}^T \mathbf{P} \mathbf{x})$ on the trajectories of the closed-loop system. For the case when $|\mathbf{b}^T \mathbf{P} \mathbf{x}| \geq \nu$, we get

$$\frac{d}{dt}(\mathbf{x}^T \mathbf{P} \mathbf{x}) \leq -2\|\mathbf{x}\|^2.$$

For the case when $|\mathbf{b}^T \mathbf{P} \mathbf{x}| < \nu$, we obtain

$$\begin{aligned} \frac{d}{dt}(\mathbf{x}^T \mathbf{P} \mathbf{x}) &= -2\|\mathbf{x}\|^2 + 2\mathbf{x}^T \mathbf{P} \mathbf{b}(u + h) \\ &\leq -2\|\mathbf{x}\|^2 - 4 \frac{|\mathbf{x}^T \mathbf{P} \mathbf{b}|^2}{\nu} + 4|\mathbf{x}^T \mathbf{P} \mathbf{b}| \\ &\leq -2\|\mathbf{x}\|^2 + 4\nu \\ &= -2\|\mathbf{x}\|^2 + 4. \end{aligned}$$

Hence, $d/dt(\mathbf{x}^T \mathbf{P} \mathbf{x}) < 0$ outside of the ball

$$\bar{B}_R(\mathbf{0}) = \{\mathbf{x} : \|\mathbf{x}\| \leq R = \sqrt{2}\}.$$

To find the radius, b , of a ball, such that $\|\mathbf{x}(t)\|_2 \leq b$ for all $t \geq 0$ and $\|\mathbf{x}(0)\|_2 \leq 10 = d$, we first observe that

$$\lambda_{\min}(\mathbf{P})\|\mathbf{x}\|^2 = \psi_1(\|\mathbf{x}\|) \leq \mathbf{x}^T \mathbf{P} \mathbf{x} \leq \psi_2(\|\mathbf{x}\|) = \lambda_{\max}(\mathbf{P})\|\mathbf{x}\|^2.$$

Let $\tilde{d} = \max\{d, R\} = 10$. Then, the radius b is computed by solving the equation

$$\psi_1(b) = \psi_2(\tilde{d}),$$

that is,

$$\lambda_{\min}(\mathbf{P})b^2 = \lambda_{\max}(\mathbf{P})\tilde{d}^2.$$

We obtain

$$\begin{aligned} b &= \tilde{d} \sqrt{\frac{\lambda_{\max}(\mathbf{P})}{\lambda_{\min}(\mathbf{P})}} \\ &= 10 \sqrt{\frac{(5 + \sqrt{5})/2}{(5 - \sqrt{5})/2}} = 16.18. \end{aligned}$$

◆ Example 4.12

In this example, we investigate the performance of a boundary layer tracking controller on the tracking problem considered in Example 4.10. Specifically, we replace the discontinuous controller given by (4.99) with a boundary layer continuous controller obtained by replacing the sign nonlinearity with a linear saturation nonlinearity. A linear saturation nonlinearity, denoted sat_ν , can be described by

$$\text{sat}_\nu(x) = \frac{1}{2\nu} (|x + \nu| - |x - \nu|), \quad (4.110)$$

where the parameter ν controls the width of the boundary layer. In other words,

$$\text{sat}_\nu(x) = \begin{cases} 1 & \text{if } x \geq \nu, \\ \frac{1}{\nu}x & \text{if } -\nu < x < \nu, \\ -1 & \text{if } x \leq -\nu. \end{cases} \quad (4.111)$$

A plot of (4.110) for $\nu = 0.5$ is shown in Figure 4.20. A boundary layer version of the discontinuous tracking controller (4.99) is

$$\begin{aligned} u &= -\rho(\mathbf{x}, r) \text{sat}_\nu(\mathbf{b}^T \mathbf{P} \mathbf{e}) \\ &= -(|\mathbf{k}_d \mathbf{x}| + |r|) \text{sat}_\nu(\mathbf{b}^T \mathbf{P}(\mathbf{x} - \mathbf{x}_d)) \\ &= -\left(\left(\left[\begin{array}{c} 0.0329 \\ 2.2854 \\ 0.2429 \\ -0.1653 \end{array} \right]^T \mathbf{x} + |r| \right) \text{sat}_\nu \left(\left[\begin{array}{c} 0.1571 \\ 26.9931 \\ 6.3324 \\ -3.0250 \end{array} \right]^T (\mathbf{x} - \mathbf{x}_d) \right) \right). \end{aligned} \quad (4.112)$$

The plot of the actual output, Y , of the closed-loop system driven by (4.112) for $v = 0.5$ and zero initial conditions was indistinguishable from that shown in Figure 4.14 when the discontinuous controller was used. The boundary layer controller action versus time is shown in Figure 4.21. Chattering is basically eliminated in this case.

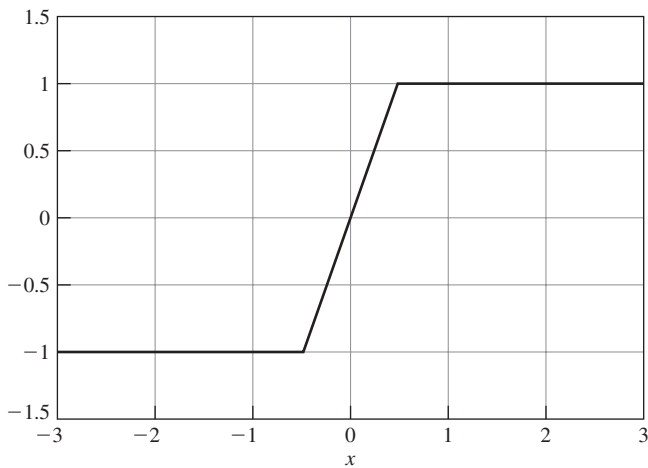


Figure 4.20 Plot of the linear saturation function described by (4.111).

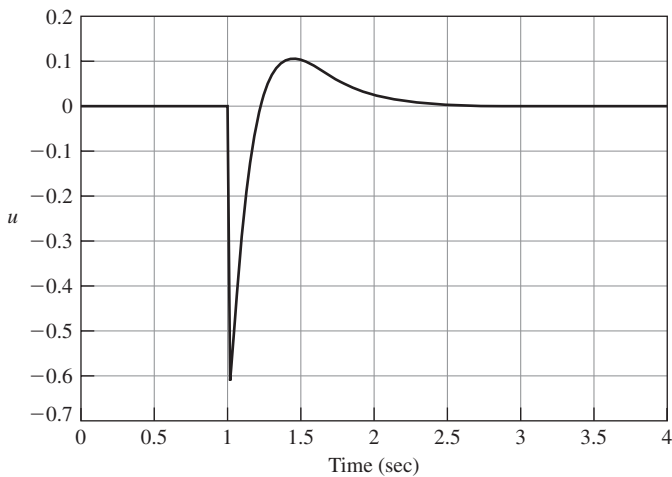


Figure 4.21 Plot of the boundary layer controller action versus time in Example 4.12.

4.12 Lyapunov-Like Analysis

In this section we present a number of results that are very useful in the analysis and design of adaptive controllers.

Lemma 4.4 (Barbalat) Given a differentiable function $f = f(t) : \mathbb{R} \rightarrow \mathbb{R}$ such that \dot{f} is uniformly continuous. If $\lim_{t \rightarrow \infty} f(t)$ exists, then

$$\lim_{t \rightarrow \infty} \dot{f}(t) = 0.$$

Proof We prove the lemma by contraposition. We assume that $\dot{f}(t)$ does not converge to 0 as $t \rightarrow \infty$. Therefore, there exists an $\varepsilon > 0$ such that for every $T > 0$ we can find $t_i \geq T$ so that $|\dot{f}(t_i)| \geq \varepsilon$. By assumption \dot{f} is uniformly continuous. Hence, given $\varepsilon/2 > 0$ there is $\delta_{\varepsilon/2} > 0$ so that for every t_i and every $\tau \in [0, \delta_{\varepsilon/2}]$, we have

$$|\dot{f}(t_i + \tau) - \dot{f}(t_i)| < \frac{\varepsilon}{2}. \quad (4.113)$$

Hence, for all $t \in [t_i, t_i + \delta_{\varepsilon/2}]$, we have

$$\begin{aligned} |\dot{f}(t)| &= |\dot{f}(t) - \dot{f}(t_i) + \dot{f}(t_i)| \\ &\geq |\dot{f}(t_i)| - |\dot{f}(t) - \dot{f}(t_i)| \\ &> \varepsilon - \frac{\varepsilon}{2} \\ &= \frac{\varepsilon}{2}. \end{aligned}$$

Suppose that $\dot{f}(t_i) > 0$. (The other case is when $\dot{f}(t_i) < 0$ and its analysis is similar to the case we analyze now.) It follows from the above that for any $t \in [t_i, t_i + \delta_{\varepsilon/2}]$,

$$\dot{f}(t) > \frac{\varepsilon}{2}.$$

Using the above, we obtain

$$\int_{t_i}^{t_i + \delta_{\varepsilon/2}} \dot{f}(t) dt > \int_{t_i}^{t_i + \delta_{\varepsilon/2}} \frac{\varepsilon}{2} dt = \frac{\varepsilon}{2} \delta_{\varepsilon/2}. \quad (4.114)$$

However,

$$\int_{t_i}^{t_i + \delta_{\varepsilon/2}} \dot{f}(t) dt = f(t_i + \delta_{\varepsilon/2}) - f(t_i). \quad (4.115)$$

Combining (4.114) and (4.115) yields

$$f(t_i + \delta_{\varepsilon/2}) - f(t_i) > \frac{\varepsilon}{2} \delta_{\varepsilon/2} \quad \text{for every } i,$$

which contradicts the assumption that $f(t)$ has a finite limit as $t \rightarrow \infty$. The proof of the lemma is complete.

A simple sufficient condition for a differentiable function to be uniformly continuous is boundedness of its derivative. Thus if a function f is twice differentiable, then its derivative \dot{f} is uniformly continuous if its second derivative \ddot{f} is bounded.

The following lemma, called by Slotine and Li [262, p. 125] the “Lyapunov-like lemma,” follows immediately from Barbalat’s lemma.

Lemma 4.5 (“Lyapunov-Like Lemma”) Given a real-valued function $W(t, \mathbf{x})$ such that

1. $W(t, \mathbf{x})$ is bounded below,
 2. $\dot{W}(t, \mathbf{x})$ is negative semidefinite, and
 3. $\dot{W}(t, \mathbf{x})$ is uniformly continuous in time,
- then

$$\dot{W}(t, \mathbf{x}) \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

◆ Example 4.13

Investigate the stability of the following dynamical system model:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} -2x_1 + x_2 u \\ -x_1 u \end{bmatrix}, \quad (4.116)$$

where the input u is bounded.

Consider the Lyapunov function candidate

$$V = \|\mathbf{x}\|^2 = x_1^2 + x_2^2.$$

Note that V is bounded below. The time derivative of V evaluated on the trajectories of the system model (4.116) is

$$\dot{V} = 2x_1 \dot{x}_1 + 2x_2 \dot{x}_2 = 2x_1(-2x_1 + x_2 u) - 2x_2 x_1 u = -4x_1^2 \leq 0.$$

Thus, the system modeled by (4.116) is uniformly stable, and hence x_1 and x_2 are bounded. We now use Lemma 4.5 to show that in addition, $x_1(t) \rightarrow 0$ as $t \rightarrow \infty$. To use this lemma, we have to show that \dot{V} is uniformly continuous in t . For this it is enough to show that the derivative of \dot{V} is bounded. The derivative of \dot{V} is

$$\ddot{V} = -8x_1 \dot{x}_1 = -8x_1(-2x_1 + x_2 u),$$

which is bounded because u is bounded by assumption and we showed that x_1 and x_2 are also bounded. Hence, \dot{V} is uniformly continuous. The function V satisfies all the conditions of Lemma 4.5, and therefore $x_1(t) \rightarrow 0$ as $t \rightarrow \infty$.

◆ Example 4.14

Consider a dynamical system model,

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b}u, \quad e = \mathbf{c}\mathbf{x}, \quad (4.117)$$

where

$$u = \boldsymbol{\phi}^T \boldsymbol{\xi}. \quad (4.118)$$

The components $\phi_i = \phi_i(t)$ of the vector-valued function $\boldsymbol{\phi} = \boldsymbol{\phi}(t) \in \mathbb{R}^K$ are adjustable parameters according to the adaptation law

$$\dot{\boldsymbol{\phi}} = -\boldsymbol{\Gamma}e\boldsymbol{\xi}, \quad (4.119)$$

where $\boldsymbol{\Gamma} = \boldsymbol{\Gamma}^T > 0$ is a design parameter matrix. The vector-valued function $\boldsymbol{\xi} = \boldsymbol{\xi}(t) \in \mathbb{R}^K$ is bounded for all t . The matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is asymptotically stable, and $\mathbf{b} \in \mathbb{R}^n$. The row vector $\mathbf{c} \in \mathbb{R}^{1 \times n}$ satisfies the equation

$$\mathbf{c} = \mathbf{b}^T \mathbf{P}, \quad (4.120)$$

where $\mathbf{P} = \mathbf{P}^T > 0$ is the solution of the matrix Lyapunov equation

$$\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -2\mathbf{Q}$$

for some $\mathbf{Q} = \mathbf{Q}^T > 0$. We will use the Lyapunov function candidate,

$$V(\mathbf{x}, \boldsymbol{\phi}) = \frac{1}{2} \mathbf{x}^T \mathbf{P} \mathbf{x} + \frac{1}{2} \boldsymbol{\phi}^T \boldsymbol{\Gamma}^{-1} \boldsymbol{\phi}, \quad (4.121)$$

and the Lyapunov-like lemma to show that \mathbf{x} and $\boldsymbol{\phi}$ are both bounded and that

$$\lim_{t \rightarrow \infty} \mathbf{x}(t) = \mathbf{0}.$$

We first evaluate the time derivative of the Lyapunov function candidate $V = V(\mathbf{x}, \boldsymbol{\phi})$ on the trajectories of the closed-loop system,

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b}(\boldsymbol{\phi}^T \boldsymbol{\xi}), \quad \dot{\boldsymbol{\phi}} = -\boldsymbol{\Gamma}e\boldsymbol{\xi}. \quad (4.122)$$

We have

$$\begin{aligned} \dot{V} &= \mathbf{x}^T \mathbf{P} \dot{\mathbf{x}} + \boldsymbol{\phi}^T \boldsymbol{\Gamma}^{-1} \dot{\boldsymbol{\phi}} \\ &= \mathbf{x}^T \mathbf{P} (\mathbf{A}\mathbf{x} + \mathbf{b}(\boldsymbol{\phi}^T \boldsymbol{\xi})) - \boldsymbol{\phi}^T e \boldsymbol{\Gamma}^{-1} \boldsymbol{\Gamma} \boldsymbol{\xi} \\ &= -\mathbf{x}^T \mathbf{Q} \mathbf{x} + (e \boldsymbol{\phi}^T - e \boldsymbol{\phi}^T) \boldsymbol{\xi} \\ &= -\mathbf{x}^T \mathbf{Q} \mathbf{x}, \end{aligned}$$

where we used the fact that $\mathbf{c} = \mathbf{b}^T \mathbf{P}$ and hence $e = \mathbf{c}\mathbf{x} = \mathbf{x}^T \mathbf{c}^T = \mathbf{x}^T \mathbf{P} \mathbf{b}$. Thus, \dot{V} is negative semidefinite in the $\begin{bmatrix} \mathbf{x} \\ \boldsymbol{\phi} \end{bmatrix}$ -space, which implies that \mathbf{x} and $\boldsymbol{\phi}$ are bounded.

To show the convergence of \mathbf{x} , we use the Lyapunov-like lemma. For this, we first show that \dot{V} is uniformly continuous. Indeed, because

$$\ddot{V} = -2\mathbf{x}^T \mathbf{Q} \dot{\mathbf{x}}$$

is bounded, since \mathbf{x} and $\dot{\mathbf{x}}$ are bounded, \dot{V} is uniformly continuous. Hence, $\dot{V} \rightarrow 0$ as $t \rightarrow \infty$. Since $\dot{V} = -\mathbf{x}^T \mathbf{Q} \mathbf{x}$ and $\mathbf{Q} = \mathbf{Q}^T > 0$, we have to have $\mathbf{x}(t) \rightarrow \mathbf{0}$ as $t \rightarrow \infty$.

We can re-state Barbalat's lemma in alternative form as follows (see, for example, Khalil [157, p. 192]).

Lemma 4.6 Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be a uniformly continuous function on $[0, \infty)$. Suppose that $\lim_{t \rightarrow \infty} \int_0^t \phi(\tau) d\tau$ exists and is finite. Then,

$$\phi(t) \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

4.13 LaSalle's Invariance Principle

In this section we generalize the Lyapunov function concept using LaSalle's invariant set. We first present this generalization, following the exposition in LaSalle [176], for nonlinear dynamical systems modeled by the difference equation

$$\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k)), \quad (4.123)$$

where $k \in \mathbb{Z}_+$, the set of nonnegative integers, and $\mathbf{x} : \mathbb{Z}_+ \rightarrow \mathbb{R}^n$; that is, $\mathbf{x}(k)$ is the state of the system model (4.123) at time k , and $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$. We assume that the vector-valued function \mathbf{f} is continuous on \mathbb{R}^n ; that is, if $\mathbf{x}^j \rightarrow \mathbf{y}$, then $\mathbf{f}(\mathbf{x}^j) \rightarrow \mathbf{f}(\mathbf{y})$. If \mathbf{x} is a vector, then $\mathbf{f}(\mathbf{x})$ is a vector and is the value of the function \mathbf{f} at \mathbf{x} .

The solution to the initial value problem

$$\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k)), \quad \mathbf{x}(0) = \mathbf{x}^0$$

is

$$\mathbf{x}(k) = \mathbf{x}^k = \mathbf{f}^k(\mathbf{x}^0),$$

where \mathbf{f}^k is the k th iterate of \mathbf{f} ; that is, $\mathbf{f}^0 = \mathbf{I}_n$, and $\mathbf{f}^k = \mathbf{f} \circ \mathbf{f}^{k-1}$.

A state $\mathbf{x}_e \in \mathbb{R}^n$ is a *fixed point* or an *equilibrium state* of the system $\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k))$ if $\mathbf{x}_e = \mathbf{f}(\mathbf{x}_e)$. Thus the solution of the system that starts at \mathbf{x}_e stays at \mathbf{x}_e for all $k \geq 0$.

A subset $\mathcal{A} \subset \mathbb{R}^n$ is called *positively invariant* if $\mathbf{f}(\mathcal{A}) \subset \mathcal{A}$; that is, if $\mathbf{x} \in \mathcal{A}$, then $\mathbf{f}(\mathbf{x}) \in \mathcal{A}$. It is clear that an equilibrium solution is invariant. The adjective "positive" is added to emphasize the fact that we will be interested in the behavior of the system for $k \geq 0$.

The solution through \mathbf{x}^0 is bounded if there is a constant M such that $\|\mathbf{f}^k(\mathbf{x}^0)\| \leq M$ for all $k \geq 0$.

Consider again a particular solution $\mathbf{f}^k(\mathbf{x}^0)$. A point $\mathbf{y} \in \mathbb{R}^n$ is a *positive limit point* or a *positive point of accumulation* of $\mathbf{f}^k(\mathbf{x}^0)$ if there is a subsequence k_i such that $\mathbf{f}^{k_i}(\mathbf{x}^0) \rightarrow \mathbf{y}$. The set of all positive limit points of $\mathbf{f}^k(\mathbf{x}^0)$ is called the *positive limit set* and denoted $\Omega(\mathbf{x}^0)$.

A sequence $\{\mathbf{x}^k\}$ is said to approach a set \mathcal{A} if

$$\text{dist}(\mathbf{x}^k, \mathcal{A}) = \inf\{\|\mathbf{x}^k - \mathbf{y}\| : \mathbf{y} \in \mathcal{A}\} \rightarrow 0,$$

where $\text{dist}(\mathbf{p}, \mathcal{S})$ denotes the distance from a point \mathbf{p} to a set \mathcal{S} .

Lemma 4.7 If \mathbf{f} is continuous and $\mathbf{f}^k(\mathbf{x}^0)$ is bounded for all $k \geq 0$, then $\Omega(\mathbf{x}^0)$ is non-empty, compact, positively invariant set. Furthermore, $\mathbf{f}^k(\mathbf{x}^0) \rightarrow \Omega(\mathbf{x}^0)$ as $k \rightarrow \infty$.

Proof Because the sequence $\{\mathbf{x}^k\}$ is bounded by the Bolzano–Weierstrass theorem, stated on page 645, there is at least one positive limit point of $\{\mathbf{x}^k\}$ and so $\Omega(\mathbf{x}^0)$ is nonempty and bounded.

We now show that the set $\Omega(\mathbf{x}^0)$ is closed; that is, if $\mathbf{y}^k \rightarrow \mathbf{y}$ with $\mathbf{y}^k \in \Omega(\mathbf{x}^0)$, then $\mathbf{y} \in \Omega(\mathbf{x}^0)$. Because $\Omega(\mathbf{x}^0)$ is nonempty and bounded, there is a sequence, say $\{\mathbf{y}^l\}$, $\mathbf{y}^l \in \Omega(\mathbf{x}^0)$, such that $\mathbf{y}^l \rightarrow \mathbf{y}$. Because each \mathbf{y}^l is in $\Omega(\mathbf{x}^0)$, for each \mathbf{y}^l there is a subsequence, $\{\mathbf{x}_l^{k_l}\}$, of $\{\mathbf{x}^k\}$ such that $\mathbf{x}_l^{k_l} \rightarrow \mathbf{y}^l$. The sequences $\{\mathbf{x}_l^{k_l}\}$ are convergent. Therefore, we can construct a sequence $\{\mathbf{x}_l\}$, a subsequence of $\{\mathbf{x}^k\}$, such that for any $\varepsilon > 0$ there exist N_1 so that

$$\|\mathbf{x}_l - \mathbf{y}^l\| < \varepsilon/2 \quad \text{for } l \geq N_1.$$

The construction of such a sequence can be accomplished as follows. For $l = 1$ select a term \mathbf{x}_1 from $\{\mathbf{x}_l^{k_l}\}$ such that

$$\|\mathbf{x}_1 - \mathbf{y}^1\| < 1.$$

Next, for $l = 2$ select a term \mathbf{x}_2 from $\{\mathbf{x}_l^{k_l}\}$ such that

$$\|\mathbf{x}_2 - \mathbf{y}^2\| < 1/2,$$

and so on. The sequence $\{\mathbf{y}^l\}$ is also convergent. Hence, for any $\varepsilon > 0$ there is N_2 such that

$$\|\mathbf{y}^l - \mathbf{y}\| < \varepsilon/2 \quad \text{for } l \geq N_2.$$

Combining the above gives

$$\|\mathbf{x}_l - \mathbf{y}\| \leq \|\mathbf{x}_l - \mathbf{y}^l\| + \|\mathbf{y}^l - \mathbf{y}\| < \varepsilon \quad \text{for } l \geq \max\{N_1, N_2\},$$

which implies that $\mathbf{y} \in \Omega(\mathbf{x}^0)$. That is, the limit of any convergent sequence with elements in $\Omega(\mathbf{x}^0)$ is also in this set, which means that the set $\Omega(\mathbf{x}^0)$ is closed. Because $\Omega(\mathbf{x}^0)$ is also bounded, it is compact.

We now show that $\Omega(\mathbf{x}^0)$ is positively invariant. For this let $\mathbf{y} \in \Omega(\mathbf{x}^0)$. Hence, there is a subsequence k_i such that $\mathbf{x}^{k_i} \rightarrow \mathbf{y}$. Because \mathbf{f} is continuous, $\mathbf{x}^{k_i+1} = \mathbf{f}(\mathbf{x}^{k_i}) \rightarrow \mathbf{f}(\mathbf{y})$ as $i \rightarrow \infty$. Therefore $\mathbf{f}(\mathbf{y}) \in \Omega(\mathbf{x}^0)$, which implies that $\Omega(\mathbf{x}^0)$ is invariant. Note that we do not have to have $\mathbf{y} = \mathbf{f}(\mathbf{y})$.

Finally, we will show that $\mathbf{f}^k(\mathbf{x}^0) \rightarrow \Omega(\mathbf{x}^0)$ as $k \rightarrow \infty$. We shall prove the statement by contradiction. First observe that $\text{dist}(\mathbf{x}^k, \Omega(\mathbf{x}^0))$ is bounded because both $\{\mathbf{x}^k\}$ and $\Omega(\mathbf{x}^0)$ are. Assume that $\text{dist}(\mathbf{x}^k, \Omega(\mathbf{x}^0))$ does not converge to zero. Then, there is a subsequence k_i such that $\mathbf{x}^{k_i} \rightarrow \mathbf{y}$ and $\rho(\mathbf{x}^{k_i}, \Omega(\mathbf{x}^0)) \rightarrow a > 0$. But because $\mathbf{y} \in \Omega(\mathbf{x}^0)$,

$$\text{dist}(\mathbf{x}^{k_i}, \Omega(\mathbf{x}^0)) \leq \|\mathbf{x}^{k_i} - \mathbf{y}\| \rightarrow 0,$$

and thus $\text{dist}(\mathbf{x}^{k_i}, \Omega(\mathbf{x}^0)) \rightarrow 0$, which is a contradiction. The proof of the lemma is now complete.

Let $V : \mathbb{R}^n \rightarrow \mathbb{R}$. Define the *Lyapunov difference* of V relative to $\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k))$, where $\mathbf{x}(k) \in \mathcal{G} \subseteq \mathbb{R}^n$, as

$$\Delta V(\mathbf{x}) = V(\mathbf{f}(\mathbf{x})) - V(\mathbf{x});$$

that is,

$$\Delta V(\mathbf{x}(k)) = V(\mathbf{x}(k+1)) - V(\mathbf{x}(k)).$$

If $\Delta V(\mathbf{x}(k)) \leq 0$, then V is nonincreasing along solutions of the system.

We say that V is a *Lyapunov function* for $\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k))$ on a set \mathcal{G} if

1. V is continuous on \mathbb{R}^n and
2. $\Delta V(\mathbf{x}(k)) \leq 0$ for all $\mathbf{x}(k) \in \mathcal{G}$.

A c -level surface of V is the set

$$V^{-1}(c) = \{\mathbf{x} : V(\mathbf{x}) = c, \mathbf{x} \in \mathbb{R}^n\}.$$

Theorem 4.14 (LaSalle's Invariance Principle) If

1. V is a Lyapunov function for $\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k))$ in a compact set \mathcal{G} and
2. the solution $\mathbf{f}^k(\mathbf{x}^0)$ is bounded and in \mathcal{G} ,

then there is a number c such that $\mathbf{f}^k(\mathbf{x}^0) \rightarrow \mathcal{M} \cap V^{-1}(c)$, where \mathcal{M} is the largest positively invariant set contained in the set $\mathcal{E} = \{\mathbf{x} \in \mathbb{R}^n : \Delta V = 0\}$.

Proof Because $\mathbf{x}^k = \mathbf{f}^k(\mathbf{x}^0)$ is bounded and in \mathcal{G} , by Lemma 4.7 the positive limit set $\Omega(\mathbf{x}^0)$ is nonempty and in \mathcal{G} . Furthermore, \mathbf{x}^k tends to $\Omega(\mathbf{x}^0)$. Now $\{V(\mathbf{x}^k)\}$ is nonincreasing and bounded below, so $V(\mathbf{x}^k) \rightarrow c$ for some c . If $\mathbf{y} \in \Omega(\mathbf{x}^0)$, then there is a subsequence k_i such that $\mathbf{x}^{k_i} \rightarrow \mathbf{y}$, and hence $V(\mathbf{x}^{k_i}) \rightarrow V(\mathbf{y})$; that is, $V(\mathbf{y}) = c$. Thus $V(\Omega(\mathbf{x}^0)) = c$, and hence $\Omega(\mathbf{x}^0) \subset V^{-1}(c)$. Because $V(\Omega(\mathbf{x}^0)) = c$ and $\Omega(\mathbf{x}^0)$ is positively invariant, for arbitrary $\mathbf{y} \in \Omega(\mathbf{x}^0)$ we have

$$\begin{aligned} \Delta V(\mathbf{y}) &= V(\mathbf{f}(\mathbf{y})) - V(\mathbf{y}) \\ &= c - c = 0; \end{aligned}$$

that is, $\Delta V(\Omega(\mathbf{x}^0)) = 0$. In summary,

$$\mathbf{x}^k \rightarrow \Omega(\mathbf{x}^0) \subset \{\mathbf{x} \in \mathbb{R}^n : \Delta V = 0\} \cap V^{-1}(c).$$

Because $\Omega(\mathbf{x}^0)$ is positively invariant, we have $\Omega(\mathbf{x}^0) \subset \mathcal{M}$, and the proof of the theorem is complete.

Using the preceding theorem we can study stability properties of difference equations.

LaSalle's invariance principle can be used to prove a result concerning the instability of equilibrium states.

Theorem 4.15 Let \mathbf{x}_e be an equilibrium state of $\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k))$, where $\{\mathbf{x}(k)\}$ is bounded. Let ΔV be positive definite with respect to \mathbf{x}_e , and let V take positive values arbitrarily close to \mathbf{x}_e . If $V(\mathbf{x}_e) = 0$, then \mathbf{x}_e is unstable.

Proof We prove the theorem by contradiction. First observe that $-V$ is a Lyapunov function because $\Delta(-V) < 0$. Because $\mathbf{x}^k = \mathbf{f}^k(\mathbf{x}^0)$ is bounded, we can apply the invariance principle to conclude that \mathbf{x}^k tends to $\{\mathbf{x} : \Delta(-V(\mathbf{x})) = 0\} = \{\mathbf{x}_e\}$, that is, $\mathbf{x}^k \rightarrow \mathbf{x}_e$. Using the above, we conclude that $V(\mathbf{x}^k) \rightarrow V(\mathbf{x}_e) = 0$. However, $\Delta V(\mathbf{x}^k) > 0$ and so

$$V(\mathbf{x}^k) > V(\mathbf{x}^{k-1}) > \cdots > V(\mathbf{x}^0) > 0.$$

Hence, for some positive ε , we have $V(\mathbf{x}^k) > \varepsilon > 0$, and thus we would have to have $V(\mathbf{x}_e) > 0$. The above contradiction completes the proof.

We now present LaSalle's invariance principle for continuous-time dynamical systems models of the form

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}), \quad (4.124)$$

where $\mathbf{f} : \mathcal{D} \rightarrow \mathbb{R}^n$ with $\mathcal{D} \subseteq \mathbb{R}^n$ being open and connected. The function \mathbf{f} is assumed to be locally Lipschitz on \mathcal{D} ; that is, each point of \mathcal{D} has a neighborhood such that \mathbf{f} satisfies the Lipschitz condition in this neighborhood with some Lipschitz constant. Each neighborhood can have a different Lipschitz constant. The above assumptions are made in order to ensure the existence and uniqueness of a solution to (4.124) for a given initial condition $\mathbf{x}(0)$.

To proceed, we need a few definitions. Suppose that $\mathbf{x}(t)$ is a solution of (4.124). Then, a point \mathbf{p} is a *positive limit point* of $\mathbf{x}(t)$ if there is a sequence $\{t_n\}$ with $t_n \rightarrow \infty$ as $n \rightarrow \infty$ such that $\mathbf{x}(t_n) \rightarrow \mathbf{p}$ as $n \rightarrow \infty$. The set of all positive limit points of $\mathbf{x}(t)$ is called the *positive limit set* of $\mathbf{x}(t)$. A set \mathcal{A} is a *positively invariant set* with respect to (4.124) if

$$\mathbf{x}(0) \in \mathcal{A} \quad \text{implies} \quad \mathbf{x}(t) \in \mathcal{A} \quad \text{for all } t \geq 0.$$

We say that $\mathbf{x}(t)$ approaches a set \mathcal{A} , as $t \rightarrow \infty$, if for each $\varepsilon > 0$ there is $T > 0$ such that

$$\text{dist}(\mathbf{x}(t), \mathcal{A}) < \varepsilon \quad \text{for all } t > T.$$

Employing arguments similar to those used in the proof of Lemma 4.7, we can prove the following lemma.

Lemma 4.8 If a solution $\mathbf{x}(t)$ to (4.124) is bounded and belongs to \mathcal{D} for $t \geq 0$, then its positive limit set Ω is non-empty, compact, and invariant. Furthermore,

$$\mathbf{x}(t) \rightarrow \Omega \quad \text{as } t \rightarrow \infty.$$

Using the above lemma, we can prove the following version of the LaSalle's invariance principle for the continuous-time case:

Theorem 4.16 If $V : \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuously differentiable function such that $\dot{V}(\mathbf{x}) \leq 0$ on the solutions of (4.124) that reside in a compact set $\mathcal{G} \subset \mathcal{D}$, then there is a number c such that $\mathbf{x}(t) \rightarrow \mathcal{M} \cap V^{-1}(c)$, where \mathcal{M} is the largest positively invariant set contained in the set $\mathcal{E} = \{\mathbf{x} \in \mathbb{R}^n : \dot{V} = 0\}$.

For further discussion of the continuous-time LaSalle's invariance principle, we refer to Khalil [157, pp. 114–115]. Applications of LaSalle's invariance principle to the stability analysis of various neural networks can be found in references 43 and 287, while in references 156, and 256 positively invariant sets were used to analyze stability of a class of nonlinear systems driven by adaptive feedback controllers.

Notes

Hahn [112] is a classical text on stability of dynamical systems. Stability of discrete dynamical systems is presented by Hahn [111] and LaSalle [176]. Extensive coverage of stability issues can also be found in Willems [300], Slotine and Li [262], Vidyasagar [289], Krstić, Kanellakopoulos, and Kokotović [171], and Khalil [157].

There is an interesting story behind the origins of the Routh test. The problem of investigating the stability of a linear dynamical system in terms of the location of the roots of its characteristic equation was posed by Maxwell [197]. On page 271 of his 1868 paper, Maxwell writes, "This condition is mathematically equivalent to the condition that all the possible roots, and all the possible parts of the impossible roots, of a certain equation shall be negative." By impossible roots, Maxwell means complex roots, and by possible parts he means real parts in today's terminology. He continues, "I have not been able completely to determine these conditions for equations of a higher degree than the third; but I hope that the subject will obtain the attention of mathematicians." A complete general solution to the problem was published by Routh in 1877. When Maxwell and Routh were classmates at Cambridge University in 1854, Routh was placed first in their final math examination ahead of Maxwell. A story says (Smith [265, p. 55]) that apparently Maxwell was so sure of coming in first on this exam that he did not even bother to rise early to hear the exam results read out in the Senate House. Maxwell sent his servant instead. When the servant came back, Maxwell asked, "Well, tell me who's second!" Obviously, he did not expect to receive the reply, "You are, sir." In 1877, Routh received the prestigious Adams Prize for solving the problem posed by Maxwell in 1868. According to an account by Truxal [283, p. 527], Routh presented his solution with the opening remark, "It has recently come to my attention that my good friend James Clerk Maxwell has had difficulty with a rather trivial problem . . ."

EXERCISES

- 4.1** (a) Consider a negative unity feedback closed-loop system whose open-loop transfer function is

$$G(s) = \frac{6s^3 + 11s^2}{s^4 + 6s + K}.$$

For what values of the parameter K the closed-loop system is asymptotically stable?

- (b) Figure 4.22 shows a closed-loop system in which the error, measured by a transducer that has a time constant of 1 sec, is applied to a proportional-plus-derivative (PD) controller. The output of the controller drives a third-order plant. Use the

Routh–Hurwitz criterion to find the region of asymptotic stability, of the closed-loop feedback system, in the (a, K) -plane.

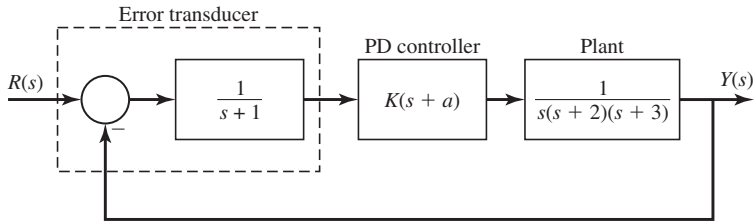


Figure 4.22 A closed-loop feedback system of Exercise 4.1.

4.2 (a) Determine if the polynomial

$$P(s) = s^7 + 2s^6 + 5s^5 + 7s^4 + 2s^3 - 5s^2 - 8s - 4$$

is Hurwitz. Use the Routh–Hurwitz criterion to find the number of the polynomial RHP zeros and determine the $j\omega$ -axis zeros.

(b) Consider a negative unity feedback system with the forward loop transfer function

$$G(s) = \frac{K}{s(s+3)(s+10)}.$$

Use the Routh–Hurwitz criterion to find the interval of K for which the closed-loop system is asymptotically stable.

(c) Consider a negative unity feedback system with the forward loop transfer function

$$G(s) = \frac{K}{s(Ts+1)(0.5s+1)}.$$

Use the Routh–Hurwitz criterion to derive the domain of asymptotic stability of the closed-loop system in the (T, K) -plane, and plot it using MATLAB.

4.3 For a dynamical system model

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ x_1 - x_2 \sin x_1 \end{bmatrix},$$

- find an equilibrium point for this system;
- linearize the system about the found equilibrium point;
- determine if the linearized system is stable, asymptotically stable, or unstable.

4.4 Consider the following model of a dynamical system,

$$\begin{aligned} \dot{x}_1 &= x_2, \\ \dot{x}_2 &= -x_1^5 - x_2^3. \end{aligned} \tag{4.125}$$

- Find the equilibrium state of the system. Then, consider the following Lyapunov

function candidate,

$$V(x_1, x_2) = \frac{1}{2}x_2^2 + \int_0^{x_1} s^5 ds.$$

Evaluate $d/dt V(x_1, x_2)$ on the trajectories of the system (4.125).

- (b) Use the results obtained in (a) to decide whether the equilibrium state is stable, asymptotically stable, or unstable.

4.5 For a dynamical system modeled by $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \boldsymbol{\zeta}(t, \mathbf{x})$, where

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ -1/2 & -3/2 \end{bmatrix} \quad \text{and} \quad \|\boldsymbol{\zeta}(t, \mathbf{x})\| \leq k\|\mathbf{x}\|,$$

find $k^* > 0$ such that for any $k < k^*$ the system is globally uniformly asymptotically stable.

4.6 (a) Find the value of

$$J_0 = 2 \int_0^\infty \|\mathbf{x}(t)\|^2 dt$$

subject to

$$\dot{\mathbf{x}} = \begin{bmatrix} 0 & 1 \\ -1 & -1 \end{bmatrix} \mathbf{x}, \quad \mathbf{x}(0) = \begin{bmatrix} 0 \\ -1 \end{bmatrix}.$$

(b) Find the value of

$$J_0 = \sum_0^\infty \|\mathbf{x}(k)\|^2$$

subject to

$$\mathbf{x}(k+1) = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/4 \end{bmatrix} \mathbf{x}(k), \quad \mathbf{x}(0) = \begin{bmatrix} 3 \\ 0 \end{bmatrix}.$$

4.7 Determine the stability or instability of the zero solution of the system

$$\begin{aligned} \dot{x}_1 &= x_2, \\ \dot{x}_2 &= x_1 - x_2 \sin x_1. \end{aligned}$$

using the second method of Lyapunov and compare the result of your analysis with the solution to Exercise 4.3.

4.8 (Based on Toptscheev and Tsypliyakov [282, p. 329]) Construct a Lyapunov function for the system modeled by

$$\begin{aligned} \dot{x}_1 &= x_2, \\ \dot{x}_2 &= -x_2 - x_1^3. \end{aligned}$$

Show that the origin is a globally stable equilibrium state of the system.

- 4.9** Investigate the stability of the following dynamical system model:

$$\begin{aligned}\dot{x}_1 &= x_2 - 2x_1x_2^2, \\ \dot{x}_2 &= -x_1 - x_2 - 2x_2^3.\end{aligned}$$

Is the null solution, $x_1(t) = x_2(t) = 0$, $t \geq 0$, locally stable, globally stable, locally asymptotically stable, globally asymptotically stable, or unstable? Justify your answer.

- 4.10** Use the second method of Lyapunov to estimate the region of asymptotic stability of the equilibrium state $[2/3 \ 2/3]^T$ of the dynamical system model

$$\begin{aligned}\dot{x}_1 &= (1 - x_1 - 0.5x_2)x_1, \\ \dot{x}_2 &= (1 - 0.5x_1 - x_2)x_2.\end{aligned}$$

- 4.11** For the nonlinear system model $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x})$, where

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} 4x_1 - 2x_1^2x_2 \\ -16x_2 + 4x_1^2x_2 \end{bmatrix}, \quad (4.126)$$

find its equilibrium points and determine their stability properties using the linearization method.

- 4.12** Consider a dynamical system modeled by

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x},$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$. The eigenvalues λ_i , $i = 1, 2, \dots, n$, of the matrix \mathbf{A} are all negative, and the matrix \mathbf{A} is diagonalizable; that is, there exists a nonsingular matrix $\mathbf{T} \in \mathbb{R}^{n \times n}$ such that

$$\mathbf{T}\mathbf{A}\mathbf{T}^{-1} = \mathbf{D} = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\}.$$

Let

$$V(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{T}^T \mathbf{T} \mathbf{x} = \mathbf{x}^T \mathbf{P} \mathbf{x}.$$

- (a) Find \dot{V} ; that is, find the derivative of $V(\mathbf{x}(t))$ with respect to time along the solutions of the system $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$.
 (b) Find the upper bound on \dot{V} such that

$$\dot{V}(\mathbf{x}) \leq c\|\mathbf{x}\|,$$

where the constant $c < 0$ depends on an eigenvalue of \mathbf{A} and an eigenvalue of \mathbf{P} .

- 4.13** Consider a system of linear equations of the form

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $m \geq n$, $\text{rank}(\mathbf{A}) = n$, and \mathbf{x} and \mathbf{b} are of appropriate dimensions. Consider also an associated iterative algorithm,

$$\mathbf{x}(k+1) = \mathbf{x}(k) + \alpha \mathbf{A}^T (\mathbf{b} - \mathbf{A}\mathbf{x}(k)), \quad \mathbf{x}(0) = \mathbf{x}_0,$$

where $\alpha > 0$. Treat the above algorithm as a discrete-time dynamical system.

- (a) Find the equilibrium state of the algorithm.
- (b) Use the conditions for asymptotic stability of discrete-time systems to find the largest range of values of α for which the algorithm converges to its equilibrium state.
- (c) Show that the algorithm, for the range of α found in (b), minimizes the objective function

$$J(\mathbf{x}) = (\mathbf{b} - \mathbf{A}\mathbf{x})^T (\mathbf{b} - \mathbf{A}\mathbf{x}).$$

4.14 Consider a system of linear equations of the form

$$\mathbf{A}\mathbf{x} = \mathbf{b},$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $m \geq n$, $\text{rank}(\mathbf{A}) = n$, and \mathbf{x} and \mathbf{b} are of appropriate dimensions. Consider also associated objective function

$$J(\mathbf{x}) = (\mathbf{b} - \mathbf{A}\mathbf{x})^T (\mathbf{b} - \mathbf{A}\mathbf{x}).$$

Write down an iterative steepest descent algorithm, in terms of \mathbf{A} and \mathbf{b} , for minimizing the above objective function. Carefully derive an expression for the step size α_k in terms of $\nabla J(\mathbf{x}^{(k)})$ and \mathbf{A} only.

4.15 (a) Given two symmetric positive definite matrices \mathbf{P} and $\tilde{\mathbf{P}}$. Show, by contradiction, that if for all \mathbf{x} ,

$$\mathbf{x}^T \mathbf{P} \mathbf{x} \geq \mathbf{x}^T \tilde{\mathbf{P}} \mathbf{x},$$

then

$$\lambda_{\max}(\mathbf{P}) \geq \lambda_{\max}(\tilde{\mathbf{P}}).$$

(b) Let $\mathbf{P} = \mathbf{P}^T > 0$ be the unique solution to

$$\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -2\mathbf{Q}, \quad (4.127)$$

where $\mathbf{Q} = \mathbf{Q}^T > 0$, and let $\hat{\mathbf{P}} = \hat{\mathbf{P}}^T > 0$ be the unique solution to

$$\mathbf{A}^T \hat{\mathbf{P}} + \hat{\mathbf{P}} \mathbf{A} = -2\hat{\mathbf{Q}}, \quad (4.128)$$

where $\hat{\mathbf{Q}} = q\mathbf{Q}$, $q > 0$. Show that

$$\frac{\lambda_{\min}(\hat{\mathbf{Q}})}{\lambda_{\max}(\hat{\mathbf{P}})} = \frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{P})}.$$

(c) Prove the following lemma of Patel and Toda [223].

Lemma 4.9 Let $\mathbf{P} = \mathbf{P}^T > 0$ be the unique solution to

$$\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -2\mathbf{Q},$$

where $\mathbf{Q} = \mathbf{Q}^T > 0$. Then, the ratio

$$\frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{P})}$$

is maximized if $\mathbf{Q} = \mathbf{I}$.

4.16 Consider the Lyapunov matrix equation,

$$\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -2\mathbf{I}_n,$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is asymptotically stable.

(a) Let σ_{\min} be the smallest singular value of \mathbf{A} ; that is,

$$\sigma_{\min} = \sqrt{\lambda_{\min}(\mathbf{A}\mathbf{A}^T)},$$

where $\lambda_{\min}(\mathbf{A}\mathbf{A}^T)$ denotes the smallest eigenvalue of $\mathbf{A}\mathbf{A}^T$. Show that $\sigma_{\min} > 0$.

(b) Let $\lambda_{\max}(\mathbf{P})$ be the largest eigenvalue of \mathbf{P} , the solution of the above Lyapunov equation. Show that

$$\lambda_{\max}(\mathbf{P}) \geq \frac{1}{\sigma_{\min}}.$$

4.17 (Based on Theorem 1 in Brockett [34]) Consider an asymptotically stable strictly proper transfer function

$$T(s) = \frac{n(s)}{d(s)}.$$

Suppose that the triple $(\mathbf{c}, \mathbf{A}, \mathbf{b})$ is a controllable and observable realization of the transfer function $T(s)$. Thus \mathbf{A} must be asymptotically stable. This and the observability assumption guarantees that we can solve the Lyapunov equation

$$\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{c}^T \mathbf{c}$$

for $\mathbf{P} = \mathbf{P}^T > 0$. Let

$$q(s) = d(s)\mathbf{b}^T \mathbf{P}(s\mathbf{I} - \mathbf{A})^{-1}\mathbf{b}.$$

Show that

$$d(s)q(-s) + q(s)d(-s) = n(s)n(-s).$$

4.18 Consider a dynamical system model

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}, \quad \mathbf{x}(0) = \mathbf{x}_0,$$

where $\mathbf{A} = \begin{bmatrix} 0 & 1 \\ -1 & -1 \end{bmatrix}$ and $\|\mathbf{x}_0\| \leq 10$. Use the Lyapunov approach to find $T \geq 0$, given $\varepsilon_1 = 20$, such that

$$\|\mathbf{x}(t; 0, \mathbf{x}_0)\| \leq \varepsilon_1 \quad \text{for all } t \geq T.$$

You can use, for example,

$$\mathbf{Q} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$$

when solving the matrix Lyapunov equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -\mathbf{Q}$.

4.19 The role of an automotive suspension is to support the vehicle body on the axles. A simple linear model of the suspension assembly associated with a single wheel, called the quarter-car model, is shown in Figure 1.27. We represent the vehicle body by

the “sprung” mass, while the tire and axle are represented by the “unsprung” mass. We model the tire as a spring with spring coefficient k_t , and we assume its damping to be negligible. The suspension consists of the spring, the shock absorber, and the variable force element. The role of the variable force element is to generate a force that compensates for the uneven road, which is the source of the system disturbances. We assume that the variable force element F_a can instantaneously provide any desired force; thus, we have a linear active suspension. The equations modeling the dynamics of the two-degree-of-freedom system depicted in Figure 4.23 are derived by applying Newton’s law to each mass of the system and by taking into account the positive direction of the z coordinates. We obtain

$$\begin{aligned} m_s \ddot{z}_s &= -k_s(z_s - z_u) - b_s(\dot{z}_s - \dot{z}_u) + F_a, \\ m_u \ddot{z}_u &= k_s(z_s - z_u) + b_s(\dot{z}_s - \dot{z}_u) - F_a - k_t(z_u - z_r). \end{aligned}$$

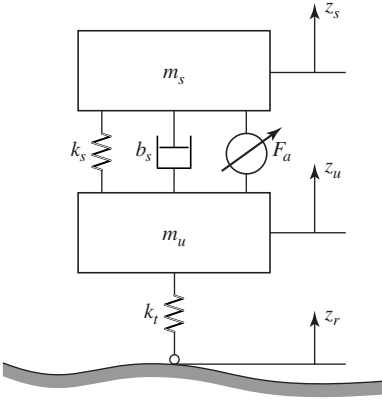


Figure 4.23 Linear quarter-car model.

Define the following state variables:

$$\begin{aligned} x_1 &= z_s - z_u = \text{suspension stroke}, \\ x_2 &= \dot{z}_s = \text{sprung mass velocity}, \\ x_3 &= z_u - z_r = \text{tire deflection}, \\ x_4 &= \dot{z}_u = \text{unsprung mass velocity}. \end{aligned}$$

Let $u = F_a$ be the system input, and let $\delta(t) = \dot{z}_r(t)$ be the disturbance input due to the road roughness.

- Represent the active suspension model in state-space format.
- Substitute the vehicle data, given in Table 4.1, into the modeling equations.
- Our goal is to isolate the sprung mass from road disturbances. Specifically, we wish to design a control law such that the sprung mass velocity $\dot{z}_s(t) = x_2$ is as close to 0 as possible in spite of road disturbances. Thus, we assume that we can measure the sprung mass velocity, which means that the output is

$$y(t) = \mathbf{c}\mathbf{x}(t) = [0 \quad 1 \quad 0 \quad 0]\mathbf{x}(t).$$

Table 4.1 Data for a Vehicle Suspension Taken from [272]

m_s	355 kg
m_u	46 kg
k_s	14384 N/m
b_s	1860 N/m/sec
k_t	300000 N/m

To investigate the effect of road disturbances on the vertical velocity of the sprung mass, set $u = 0$ and $\delta(t) = d \sin(6t)$. Generate a plot of y versus t for $d = 0.5$.

- (d) Construct a state-feedback Lyapunov-based controller of Theorem 4.4. Generate a plot of y of the closed-loop system versus t for $d = 0.5$. Discuss your strategy of the above state-feedback controller design that would minimize the influence of the disturbance $\delta(t)$ on the sprung mass velocity y .

4.20 Consider a dynamical system model

$$\dot{\mathbf{x}} = \begin{bmatrix} -2 & 0 \\ 0 & -1 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} (u + h(t, \mathbf{x})),$$

where

$$|h(t, \mathbf{x})| \leq 3.$$

Let \mathbf{P} be the solution of the matrix Lyapunov equation $\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -2\mathbf{I}_2$, and let

$$u = \begin{cases} -3 \operatorname{sign}(\mathbf{b}^T \mathbf{P} \mathbf{x}) & \text{if } |\mathbf{b}^T \mathbf{P} \mathbf{x}| \geq \nu \\ -3 \frac{\mathbf{b}^T \mathbf{P} \mathbf{x}}{\nu} & \text{if } |\mathbf{b}^T \mathbf{P} \mathbf{x}| < \nu, \end{cases}$$

where $\nu = 0.5$. Find the radius of the ball, in \mathbb{R}^2 , that is uniformly stable for the closed-loop system.



CHAPTER 5

Optimal Control

BEFORE me grew the human soul,
And after I am dead and gone,
Through grades of effort and control
The marvellous work shall still go on

—Archibald Lampman

If everything seems under control, you're just not going fast enough.

—Mario Andretti

One of the essential elements of the control problem, discussed in Section 1.1, is a means of testing the performance of any proposed control law. Whenever we use the term “best” or “optimal” to describe the effectiveness of a given control strategy, we do so with respect to some numerical index of performance called the performance index. We assume that the value of the performance index decreases as the quality of the given admissible control law increases. The admissible controller that ensures the completion of the system objective and at the same time minimizes the performance index is called an *optimal controller* for the system. In this chapter, we discuss various methods of constructing optimal controllers.

5.1 Performance Indices

Constructing a performance index—that is, choosing a means to measure the system performance—can be considered a part of the system modeling. Here, we discuss some typical choices of performance indices.

First, suppose that the objective is to control a dynamical system modeled by the equations

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t), \quad \mathbf{x}(t_0) = \mathbf{x}_0, \quad (5.1)$$

$$\mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) \quad (5.2)$$

on a fixed interval $[t_0, t_f]$ so that the components of the state vector are “small.” A suitable performance index to be minimized would be

$$J_1 = \int_{t_0}^{t_f} \mathbf{x}^T(t)\mathbf{x}(t) dt. \quad (5.3)$$

Obviously, if J_1 is small, then the state vector norm, $\|\mathbf{x}(t)\|$, is small in the sense of the above performance index.

If the objective is to control the system so that the components of the output, $\mathbf{y}(t)$, are to be small, then we could use the performance index

$$\begin{aligned} J_2 &= \int_{t_0}^{t_f} \mathbf{y}^T(t) \mathbf{y}(t) dt \\ &= \int_{t_0}^{t_f} \mathbf{x}^T(t) \mathbf{C}^T \mathbf{C} \mathbf{x}(t) dt \\ &= \int_{t_0}^{t_f} \mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) dt, \end{aligned} \quad (5.4)$$

where the weight matrix $\mathbf{Q} = \mathbf{C}^T \mathbf{C}$ is symmetric positive semidefinite.

In the case when we wish to control the system in such a manner that the components of the input, $\mathbf{u}(t)$, are “not too large,” a suitable performance index to be minimized is

$$J_3 = \int_{t_0}^{t_f} \mathbf{u}^T(t) \mathbf{u}(t) dt, \quad (5.5)$$

or

$$J_4 = \int_{t_0}^{t_f} \mathbf{u}^T(t) \mathbf{R} \mathbf{u}(t) dt, \quad (5.6)$$

where the weight matrix \mathbf{R} is symmetric positive definite. There is no loss of generality in assuming the weight matrix \mathbf{R} to be symmetric. If \mathbf{R} was not symmetric, we could represent the quadratic term $\mathbf{u}^T \mathbf{R} \mathbf{u}$ equivalently as

$$\mathbf{u}^T \mathbf{R} \mathbf{u} = \mathbf{u}^T \left(\frac{\mathbf{R} + \mathbf{R}^T}{2} \right) \mathbf{u},$$

where the matrix $\frac{1}{2}(\mathbf{R} + \mathbf{R}^T)$ is symmetric.

As pointed out by Owens [218, p. 186], we cannot simultaneously minimize the performance indices (5.3) and (5.5) because minimization of (5.3) requires large control signals, while minimization of (5.5) requires small control signals. To solve the dilemma, we could compromise between the two conflicting objectives by minimizing the performance index that is a convex combination of J_1 and J_3 ,

$$\begin{aligned} J &= \lambda J_1 + (1 - \lambda) J_3 \\ &= \int_{t_0}^{t_f} (\lambda \mathbf{x}^T(t) \mathbf{x}(t) + (1 - \lambda) \mathbf{u}^T(t) \mathbf{u}(t)) dt, \end{aligned} \quad (5.7)$$

where λ is a parameter in the range $[0, 1]$. If $\lambda = 1$, then $J = J_1$; and if $\lambda = 0$, then $J = J_3$. By trial and error, we could select λ from the interval $[0, 1]$ to compromise between the two extremes. A generalization of the performance index (5.7) is

$$J = \int_{t_0}^{t_f} (\mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) + \mathbf{u}^T(t) \mathbf{R} \mathbf{u}(t)) dt.$$

In certain applications, we may wish the final state $\mathbf{x}(t_f)$ to be as close as possible to $\mathbf{0}$. Then, a possible performance measure to be minimized is

$$\mathbf{x}^T(t_f) \mathbf{F} \mathbf{x}(t_f), \quad (5.8)$$

where \mathbf{F} is a symmetric positive definite matrix.

We can combine the performance measures (5.4), (5.6), and (5.8) when our control aim is to keep the state “small,” the control “not too large,” and the final state as near to $\mathbf{0}$ as possible. The resulting performance index is

$$J = \frac{1}{2} \mathbf{x}^T(t_f) \mathbf{F} \mathbf{x}(t_f) + \frac{1}{2} \int_{t_0}^{t_f} (\mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) + \mathbf{u}^T(t) \mathbf{R} \mathbf{u}(t)) dt, \quad (5.9)$$

where the factor $\frac{1}{2}$ is placed to simplify subsequent algebraic manipulations. Minimizing (5.9) subject to (5.1) is called the *linear quadratic regulator* problem, or the LQR problem for short.

In some cases, the controller’s goal is to force the system state to track a desired state trajectory, $\mathbf{x}_d(t)$, throughout the interval $[t_0, t_f]$ while maintaining the deviations of the actual state $\mathbf{x}(t)$ “small” from the desired trajectory with the control effort $\mathbf{u}(t)$ “not too large” and the final state $\mathbf{x}(t_f)$ being as near as possible to some desired state $\mathbf{x}_d(t_f)$. A suitable performance index to be minimized in this case might be

$$\begin{aligned} J = & \frac{1}{2} (\mathbf{x}(t_f) - \mathbf{x}_d(t_f))^T \mathbf{F} (\mathbf{x}(t_f) - \mathbf{x}_d(t_f)) \\ & + \frac{1}{2} \int_{t_0}^{t_f} ((\mathbf{x}(t) - \mathbf{x}_d(t))^T \mathbf{Q} (\mathbf{x}(t) - \mathbf{x}_d(t)) + \mathbf{u}^T(t) \mathbf{R} \mathbf{u}(t)) dt. \end{aligned} \quad (5.10)$$

5.2 A Glimpse at the Calculus of Variations

The goal of this section is to provide a brief introduction to the subject of the *calculus of variations*. The calculus of variations deals with minimization or maximization of *functionals*. A functional is a mapping or a transformation that depends on one or more functions and the values of the functional are numbers. Examples of functionals are performance indices that we introduced in the previous section. We now give two more examples of functionals.

◆ Example 5.1

The length l of the arc connecting two given points $A(x_0, y_0)$ and $B(x_1, y_1)$ in the x - y plane, as depicted in Figure 5.1, is an example of a functional. We know that

$$l = l(y(x)) = \int_{x_0}^{x_1} \sqrt{1 + (dy(x)/dx)^2} dx,$$

subject to the constraints $y(x_0) = y_0$ and $y(x_1) = y_1$.

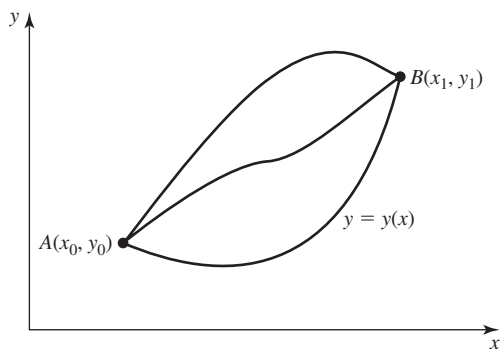


Figure 5.1 The length of arcs joining two given points in the plane is an example of a functional.

◆ Example 5.2

Another example of a functional is the area of the surface of revolution formed by a surface connecting two circular disks as illustrated in Figure 5.2. In particular, the surface of revolution can be formed by a soap film clinging to two circular disks. When a wire circle is dipped into a soap solution and withdrawn, a circular disk of soap film bounded by the circle is formed. If a second smaller circle is made to touch the disk and then moved away, the two circles will be joined by a surface of film that is a surface of revolution in the particular case when the circles are parallel and have their centers on the same axis perpendicular to their planes. If we denote the arc of intersection with the x - y plane by $y = y(x)$, then the area of the surface of revolution is

$$v = v(y(x)) = 2\pi \int_{x_0}^{x_1} y(x) \sqrt{1 + (dy(x)/dx)^2} dx,$$

subject to the constraints $y(x_0) = y_0$, $y(x_1) = y_1$. One may then want to minimize $v = v(y(x))$ to find $y = y(x)$ resulting in the surface of revolution of minimum area.

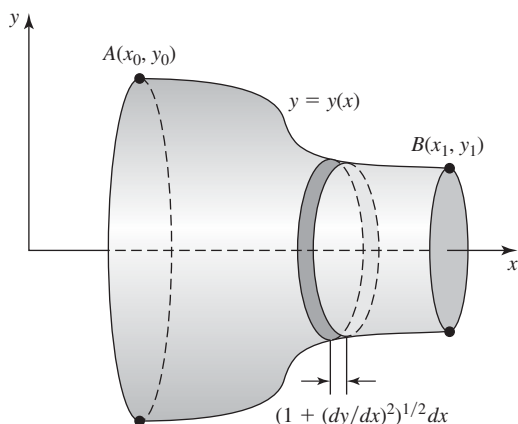


Figure 5.2 The area of surfaces of revolution connecting two circular disks is an example of a functional.

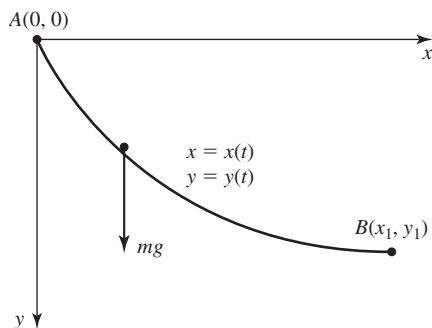


Figure 5.3 Illustration of the brachistochrone problem.

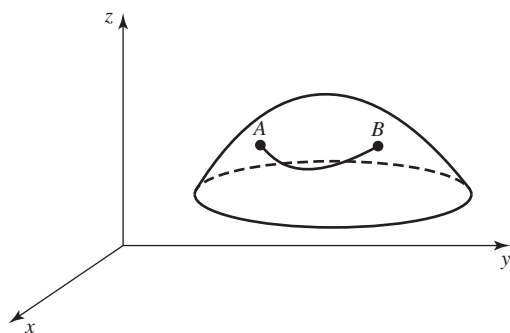


Figure 5.4 The problem of finding geodesics.

The problem of optimization of functionals dates back to Johann Bernoulli. In 1696 he published a letter in which he posed the problem, which is now known as the *brachistochrone problem*. Brachistochrone translates from Greek as the shortest time. The brachistochrone problem can be described as follows. A particle of mass m is to slide down along a frictionless track connecting two given points. The motion of the particle is to be constrained only by gravity and the shape of the curve of the track. Find the curve of the track that minimizes the time travel between the two given points. We illustrate the problem in Figure 5.3. The brachistochrone problem was solved by Johann Bernoulli himself, his brother Jacob, Isaac Newton, the Marquis de L'Hôpital, Gottfried Leibniz, and Tschirnhaus. A solution to the brachistochrone problem is presented in Example 5.23. In 1697 Johann Bernoulli posed and solved the *geodesics* problem, which can be formulated as follows: Find a line of shortest length that connects two points on a given surface; see Figure 5.4 for an illustration of the problem. The curves of the least length connecting pairs of points on a surface are called *geodesics*. A general procedure for solving the above type of problems, using the calculus of variations, was given by Leonhard Euler (who flourished between 1707 and 1783).

To acquaint ourselves with the methods of the calculus of variations, we first need to introduce concepts concerning functionals.

5.2.1 Variation and Its Properties

A *function* can be viewed as a rule, or a mapping, that assigns to each element of some set a unique element of a possibly different set. In particular, a function $x : \mathbb{R} \rightarrow \mathbb{R}$ of a real variable t

assigns to each real number a unique real number. An increment of the argument of a function of one variable t is $\Delta t = t - t_1$.

Similarly, a *functional* is a mapping that assigns to each function, from some class of functions, a unique number. Thus, we can say that a functional is a “function of a function.” Let $x : \mathbb{R} \rightarrow \mathbb{R}$ be an argument of a functional. By a *variation* δx of an argument x of a functional v we mean the difference of two functions

$$\delta x = x - x_1.$$

We assume that x can change in an arbitrary way in some class of functions.

Recall that a function $x : \mathbb{R} \rightarrow \mathbb{R}$ is continuous if a small change of its argument t corresponds to a small change of the value $x(t)$ of the function. Similarly, a functional v is said to be continuous if a “small” change of its argument x corresponds to a small change of the value of the functional.

A functional v is called linear if

$$v(ax_1 + x_2) = av(x_1) + v(x_2),$$

where a is a constant.

We will now introduce the notion of the *variation of a functional*. This notion is analogous to the notion of a function differential. To connect the two, we first consider the case of a function of one variable. Let x be a differentiable function defined on an open interval U , and let $t \in U$. The derivative of x at t is defined as

$$x'(t) = \lim_{\Delta t \rightarrow 0} \frac{x(t + \Delta t) - x(t)}{\Delta t}. \quad (5.11)$$

Let

$$\varphi(\Delta t) = \frac{x(t + \Delta t) - x(t)}{\Delta t} - x'(t). \quad (5.12)$$

The function φ is not defined at $\Delta t = 0$; however,

$$\lim_{\Delta t \rightarrow 0} \varphi(\Delta t) = 0.$$

We rewrite (5.12) as

$$x(t + \Delta t) - x(t) = x'(t)\Delta t + \varphi(\Delta t)\Delta t. \quad (5.13)$$

The above relation has meaning only when $\Delta t \neq 0$. To make it hold at $\Delta t = 0$, we define $\varphi(\Delta t)|_{\Delta t=0} = 0$. To proceed, let

$$\begin{aligned} \beta(\Delta t) &= \varphi(\Delta t) & \text{if } \Delta t > 0 \\ \beta(\Delta t) &= -\varphi(\Delta t) & \text{if } \Delta t < 0. \end{aligned} \quad (5.14)$$

Then, if x is differentiable, there exists a function β such that

$$\begin{aligned} x(t + \Delta t) - x(t) &= \Delta x = x'(t)\Delta t + \beta(\Delta t)|\Delta t| \\ &= L(t, \Delta t) + \beta(\Delta t)|\Delta t|, \end{aligned} \quad (5.15)$$

where $\lim_{\Delta t \rightarrow 0} \beta(\Delta t) = 0$, and $L(t, \Delta t) = x'(t)\Delta t$ is a linear function in Δt .

For a real-valued function $f = f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$, relation (5.15) takes the form

$$\begin{aligned} f(\mathbf{x} + \Delta \mathbf{x}) - f(\mathbf{x}) &= \Delta f = Df(\mathbf{x})\Delta \mathbf{x} + \beta(\Delta \mathbf{x})\|\Delta \mathbf{x}\| \\ &= L(\mathbf{x}, \Delta \mathbf{x}) + \beta(\Delta \mathbf{x})\|\Delta \mathbf{x}\|, \end{aligned} \quad (5.16)$$

where $\lim_{\Delta \mathbf{x} \rightarrow \mathbf{0}} \beta(\Delta \mathbf{x}) = 0$, and

$$L(\mathbf{x}, \Delta \mathbf{x}) = Df(\mathbf{x}) \Delta \mathbf{x} = \nabla f(\mathbf{x})^T \Delta \mathbf{x}$$

is a linear function in $\Delta \mathbf{x}$. Relation (5.16) holds for operators defined over Banach spaces with scalars being either real or complex numbers. The functionals we consider are operators over a Banach space of continuous functions, $\mathcal{C}([t_0, t_1])$, with the norm given by (A.99) on page 656.

If an increment $\Delta v = v(x + \delta x) - v(x)$ of a functional v can be represented as

$$\Delta v = L(x, \delta x) + \beta(\delta x) \|\delta x\|,$$

where $L(x, \delta x)$ is a linear functional with respect to δx , the term $\|\delta x\| = \max_{t \in [t_0, t_1]} |\delta x|$ denotes the maximal value of $|\delta x|$, and $\beta(\delta x) \rightarrow 0$ if $\|\delta x\| \rightarrow 0$, then the linear part L of Δv is called the *variation of the functional* and is denoted δv ; that is,

$$\delta v = L(x, \delta x).$$

◆ Example 5.3

We find the variation of the functional

$$v = \int_0^1 (2x^2(t) + x(t)) dt.$$

For this we first calculate its increment to get

$$\begin{aligned} \Delta v &= v(x + \delta x) - v(x) \\ &= \int_0^1 (2(x + \delta x)^2 + (x + \delta x)) dt - \int_0^1 (2x^2 + x) dt \\ &= \int_0^1 (2x^2 + 4x\delta x + 2(\delta x)^2 + x + \delta x - 2x^2 - x) dt \\ &= \int_0^1 (4x + 1)\delta x dt + 2 \int_0^1 (\delta x)^2 dt. \end{aligned}$$

The linear part of Δv is

$$\delta v = \int_0^1 (4x + 1)\delta x dt,$$

which is the variation of the given functional.

To proceed, note that the linear part on the right-hand side of (5.16) can be computed as

$$Df(\mathbf{x}) \Delta \mathbf{x} = L(\mathbf{x}, \Delta \mathbf{x}) = \left. \frac{d}{d\alpha} f(\mathbf{x} + \alpha \Delta \mathbf{x}) \right|_{\alpha=0} \quad (5.17)$$

As mentioned before, the above also holds for operators over Banach spaces other than Euclidean

spaces. Using this observation, we can state and prove a lemma that will allow us to calculate the variation of a given functional more efficiently than the method we used in Example 5.3.

Lemma 5.1

$$\delta v = \left. \frac{d}{d\alpha} v(x + \alpha \delta x) \right|_{\alpha=0}.$$

Proof Suppose that for a given functional v there exists its variation. This means that we can represent Δv as

$$\Delta v = v(x + \alpha \delta x) - v(x) = L(x, \alpha \delta x) + \beta(\alpha \delta x) |\alpha| \|\delta x\|.$$

Then, the derivative of $v(x + \alpha \delta x)$ with respect to α evaluated at $\alpha = 0$ is equal to

$$\begin{aligned} \lim_{\alpha \rightarrow 0} \frac{\Delta v}{\alpha} &= \lim_{\alpha \rightarrow 0} \frac{L(x, \alpha \delta x) + \beta(\alpha \delta x) |\alpha| \|\delta x\|}{\alpha} \\ &= \lim_{\alpha \rightarrow 0} \frac{L(x, \alpha \delta x)}{\alpha} + \lim_{\alpha \rightarrow 0} \frac{\beta(\alpha \delta x) |\alpha| \|\delta x\|}{\alpha} \\ &= L(x, \delta x), \end{aligned}$$

because $L(\cdot, \cdot)$ is linear with respect to the second argument, and hence

$$L(x, \alpha \delta x) = \alpha L(x, \delta x).$$

Furthermore,

$$\lim_{\alpha \rightarrow 0} \frac{\beta(\alpha \delta x) |\alpha| \|\delta x\|}{\alpha} = \lim_{\alpha \rightarrow 0} \beta(\alpha \delta x) \|\delta x\| = 0.$$

This completes the proof.

◆ Example 5.4

We will now use Lemma 5.1 to find the variation of the functional considered in Example 5.3. We have

$$\begin{aligned} \delta v &= \left. \frac{d}{d\alpha} v(x + \alpha \delta x) \right|_{\alpha=0} \\ &= \left. \frac{d}{d\alpha} \left(\int_0^1 (2(x + \alpha \delta x)^2 + (x + \alpha \delta x)) dt \right) \right|_{\alpha=0} \\ &= \left. \int_0^1 (4(x + \alpha \delta x) \delta x + \delta x) dt \right|_{\alpha=0} \\ &= \int_0^1 (4x + 1) \delta x dt, \end{aligned}$$

which is the variation of the given functional in Example 5.3.

◆ Example 5.5

We will find the variation of the functional

$$v(x_1, x_2) = v(x_1(t), x_2(t)) = \int_0^1 x_1 e^{-x_2} dt,$$

using the formula

$$\delta v = \left. \frac{d}{d\alpha} v(\mathbf{x} + \alpha \delta \mathbf{x}) \right|_{\alpha=0}.$$

We have

$$\begin{aligned} \delta v &= \left. \frac{d}{d\alpha} \left(\int_0^1 (x_1 + \alpha \delta x_1) e^{-(x_2 + \alpha \delta x_2)} dt \right) \right|_{\alpha=0} \\ &= \int_0^1 \left. \frac{d}{d\alpha} ((x_1 + \alpha \delta x_1) e^{-(x_2 + \alpha \delta x_2)}) \right|_{\alpha=0} dt \\ &= \int_0^1 (\delta x_1 e^{-(x_2 + \alpha \delta x_2)} - \delta x_2 e^{-(x_2 + \alpha \delta x_2)} (x_1 + \alpha \delta x_1)) dt \Big|_{\alpha=0} \\ &= \int_0^1 (e^{-x_2} \delta x_1 - x_1 e^{-x_2} \delta x_2) dt. \end{aligned}$$

Thus,

$$\delta v = \int_0^1 (e^{-x_2} \delta x_1 - x_1 e^{-x_2} \delta x_2) dt.$$

Definition 5.1 A functional v is maximal on some function x_0 if the value of the functional on any function x close to x_0 is less than or equal to that on x_0 , that is,

$$\Delta v = v(x) - v(x_0) \leq 0.$$

We will now give a first order necessary condition for a given functional to achieve a maximum or minimum on a curve.

Theorem 5.1 If for a functional v there exists a variation and the functional reaches a maximum or a minimum on a curve x_0 , then

$$\delta v(x_0) = 0.$$

Proof Suppose that x_0 and δx are fixed. Then

$$v(x_0 + \alpha \delta x) = \Phi(\alpha);$$

that is, $v(x_0 + \alpha \delta x)$ is a function of α that for $\alpha = 0$ achieves its maximum or minimum. This implies that $\frac{d}{d\alpha} \Phi = 0$, or, equivalently,

$$\left. \frac{d}{d\alpha} v(x_0 + \alpha \delta x) \right|_{\alpha=0} = 0.$$

Hence, by Lemma 5.1, we have $\delta v = 0$. Therefore, a functional that achieves its maximum or a minimum on a function x_0 has its variation evaluated at x_0 equal to zero.

The above results will be used in the next section to determine conditions that a curve has to satisfy to be a maximizer or a minimizer, that is, to be an extremal for a given functional.

5.2.2 Euler–Lagrange Equation

In this section, we consider the following optimization problem.

Problem 1 Suppose that we are given two points (t_0, x_0) and (t_1, x_1) in the (t, x) -plane. We wish to find a curve, or trajectory, joining the given points such that the functional

$$v(x) = \int_{t_0}^{t_1} F(t, x, \dot{x}) dt \quad (5.18)$$

along this trajectory can achieve its extremal—that is, maximal or minimal—value. The problem is illustrated in Figure 5.5. In solving Problem 1, we use the following lemma.

Lemma 5.2 Suppose that a function $\phi: \mathbb{R} \rightarrow \mathbb{R}$ is continuous on the interval $[t_0, t_1]$. Then,

$$\int_{t_0}^{t_1} \phi(t) \delta x(t) dt = 0$$

if and only if $\phi(t) = 0$ at every point of the interval $[t_0, t_1]$.

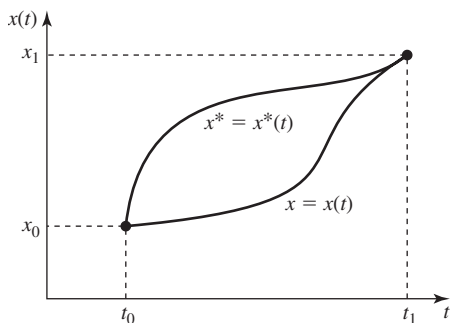


Figure 5.5 Illustration of Problem 1.

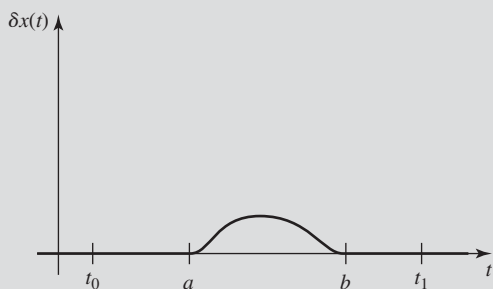


Figure 5.6 Illustration of the proof of Lemma 5.2.

Proof Sufficiency is obvious because we can satisfy $\int_{t_0}^{t_1} \phi(t) \delta x(t) dt = 0$ by choosing ϕ that vanishes at every point of the interval $[t_0, t_1]$.

We prove necessity by contraposition. We thus prove that if ϕ is nonzero on some subinterval of $[t_0, t_1]$, then $\int_{t_0}^{t_1} \phi(t) \delta x(t) dt \neq 0$. Specifically, assume, without loss of generality, that ϕ is positive on the interval $[a, b]$, where $t_0 \leq a \leq b \leq t_1$. Note that if ϕ were nonzero at only one point, then because it is continuous, it would have to be nonzero in some neighborhood of that point. Hence, we can assume that $a < b$. Consider now a variation given by

$$\delta x = \delta x(t) = \begin{cases} k^2(t-a)^2(t-b)^2 & \text{for } a \leq t \leq b, \\ 0 & \text{elsewhere.} \end{cases} \quad (5.19)$$

If we assumed ϕ to be negative on the interval $[a, b]$, then we would use δx to be negative of (5.19). We illustrate δx in Figure 5.6. Because δx is positive in $[a, b]$ and zero elsewhere in the interval $[t_0, t_1]$, the integrand of $\int_{t_0}^{t_1} \phi(t) \delta x(t) dt$ will have the same sign as ϕ in $[a, b]$ and will be zero outside this subinterval. Thus

$$\int_{t_0}^{t_1} \phi(t) \delta x(t) dt = \int_a^b \phi(t) \delta x(t) dt \neq 0,$$

and the proof is complete.

Solution to Problem 1 We assume that there is an optimal trajectory $x = x(t)$ joining the points (t_0, x_0) and (t_1, x_1) such that $\delta v(x(t)) = 0$. Let us then consider an arbitrary acceptable curve $x^* = x^*(t)$ that is close to x and from a family of curves

$$x(t, \alpha) = x(t) + \alpha(x^*(t) - x(t)) = x(t) + \alpha \delta x(t).$$

Note that for $\alpha = 0$ we get $x(t)$, and for $\alpha = 1$ we obtain $x^*(t)$. Furthermore,

$$\frac{d}{dt}(\delta x(t)) = \frac{d}{dt}(x^*(t) - x(t)) = \delta x',$$

and analogously

$$\frac{d^2}{dt^2}(\delta x(t)) = \frac{d^2}{dt^2}(x^*(t) - x(t)) = \delta x''.$$

We are now ready to investigate the behavior of the functional v given by (5.18) on the curves from the family $x(t, \alpha)$. Note that the functional v can be considered as a function of α , that is, $v(x(t, \alpha)) = \Phi(\alpha)$. It follows from the first-order necessary condition for a point to be an extremizer of a function of one variable that $\alpha = 0$ is a candidate to be an extremizer of Φ if

$$\left. \frac{d\Phi(\alpha)}{d\alpha} \right|_{\alpha=0} = 0,$$

where

$$\Phi(\alpha) = \int_{t_0}^{t_1} F(t, x(t, \alpha), x'(t, \alpha)) dt,$$

and $x'(t, \alpha) = \frac{d}{dt}x(t, \alpha)$. We now evaluate $\frac{d\Phi(\alpha)}{d\alpha}$:

$$\frac{d\Phi(\alpha)}{d\alpha} = \int_{t_0}^{t_1} \left(F_x \frac{d}{d\alpha}x(t, \alpha) + F_{x'} \frac{d}{d\alpha}x'(t, \alpha) \right) dt,$$

where $F_x = \frac{\partial}{\partial x}F(t, x(t, \alpha), x'(t, \alpha))$, and $F_{x'} = \frac{\partial}{\partial x'}F(t, x(t, \alpha), x'(t, \alpha))$. Because

$$\frac{d}{d\alpha}x(t, \alpha) = \frac{d}{d\alpha}(x(t) + \alpha\delta x(t)) = \delta x(t)$$

and

$$\frac{d}{d\alpha}x'(t, \alpha) = \frac{d}{d\alpha}(x'(t) + \alpha\delta x'(t)) = \delta x'(t),$$

we can write

$$\frac{d}{d\alpha}\Phi(\alpha) = \int_{t_0}^{t_1} (F_x(t, x(t, \alpha), x'(t, \alpha))\delta x(t) + F_{x'}(t, x(t, \alpha), x'(t, \alpha))\delta x'(t)) dt.$$

By Lemma 5.1,

$$\frac{d}{d\alpha}\Phi(0) = \delta v.$$

Therefore,

$$\delta v = \int_{t_0}^{t_1} (F_x\delta x + F_{x'}\delta x') dt,$$

where for convenience we dropped the argument t . Integrating by parts the last term, along with taking into account that $\delta x' = (\delta x)'$, yields

$$\delta v = (F_{x'}\delta x)|_{t_0}^{t_1} + \int_{t_0}^{t_1} \left(F_x - \frac{d}{dt}F_{x'} \right) \delta x dt.$$

The end points (t_0, x_0) and (t_1, x_1) are fixed. This implies that

$$\delta x|_{t=t_0} = x^*(t_0) - x(t_0) = 0,$$

and

$$\delta x|_{t=t_1} = x^*(t_1) - x(t_1) = 0.$$

Hence

$$\delta v = \int_{t_0}^{t_1} \left(F_x - \frac{d}{dt} F_{x'} \right) \delta x \, dt.$$

By Lemma 5.2,

$$\delta v = 0 \quad \text{if and only if} \quad F_x - \frac{d}{dt} F_{x'} = 0.$$

The equation $F_x - \frac{d}{dt} F_{x'} = 0$ can also be represented as

$$F_x - F_{x't} - F_{x'x}x' - F_{x'x'}x'' = 0.$$

The equation

$$\boxed{F_x - \frac{d}{dt} F_{x'} = 0} \tag{5.20}$$

is known as the Euler–Lagrange equation. The curves $x = x(t, C_1, C_2)$ that are solutions to the Euler–Lagrange equation are called *extremals*. The constants C_1 and C_2 are the integration constants. Thus, a functional can only achieve its extremum on the extremals.

◆ Example 5.6

We find the extremal for the functional

$$v = \int_0^{\pi/2} ((x')^2 - x^2) \, dt, \tag{5.21}$$

where $x(0) = 0$, and $x(\pi/2) = 1$. The Euler–Lagrange equation for the problem takes the form

$$\begin{aligned} F_x - \frac{d}{dt} F_{x'} &= \frac{\partial}{\partial x} ((x')^2 - x^2) - \frac{d}{dt} \left(\frac{\partial}{\partial x'} ((x')^2 - x^2) \right) \\ &= -2x - \frac{d}{dt} (2x') \\ &= -2x - 2x'' \\ &= 0. \end{aligned}$$

Equivalently, the Euler–Lagrange equation for the functional (5.21) can be represented as

$$x'' + x = 0.$$

The characteristic equation of the above differential equation is $s^2 + 1 = 0$.

Therefore, the general solution of the Euler–Lagrange equation has the form

$$x(t) = C_1 \cos t + C_2 \sin t.$$

Taking into account the end-point conditions gives

$$\begin{aligned} x(0) &= C_1 = 0, \\ x\left(\frac{\pi}{2}\right) &= C_2 = 1. \end{aligned}$$

Thus, we have one extremal

$$x(t) = \sin t.$$

The functional (5.21) is minimized or maximized only on the extremal $x(t) = \sin t$.

◆ Example 5.7

We now find the extremal for the functional

$$v = \int_0^1 ((x')^2 + 12xt) dt, \quad (5.22)$$

where $x(0) = 0$, and $x(1) = 1$. The Euler–Lagrange equation for the this functional is

$$\begin{aligned} F_x - \frac{d}{dt} F_{x'} &= \frac{\partial}{\partial x} ((x')^2 + 12xt) - \frac{d}{dt} \left(\frac{\partial}{\partial x'} ((x')^2 + 12xt) \right) \\ &= 12t - \frac{d}{dt} (2x') \\ &= 12t - 2x'' \\ &= 0. \end{aligned}$$

Equivalently, the Euler–Lagrange equation for the functional (5.22) can be represented as

$$x'' = 6t.$$

The general solution to the above differential equation is

$$x(t) = C_1 t + C_2 + t^3.$$

From the end-point conditions, we get

$$\begin{aligned} x(0) &= C_2 = 0, \\ x(1) &= C_1 + 1 = 1. \end{aligned}$$

Hence, $C_1 = 0$, and the extremal has the form

$$x(t) = t^3.$$

Suppose now that $A = A(t_0, x_{10}, \dots, x_{n0})$ and $B = B(t_1, x_{11}, \dots, x_{n1})$ are two given points in the $(n + 1)$ -dimensional space (t, x_1, \dots, x_n) , along with the functional

$$v = v(x_1, \dots, x_n) = \int_{t_0}^{t_1} F(t, x_1, \dots, x_n, \dot{x}_1, \dots, \dot{x}_n) dt.$$

To find a trajectory along which $\delta v = 0$, we consider only the variation in one coordinate, say x_i , while other coordinates are fixed. We perform this process for each coordinate. The functional v will depend only upon x_i , and we write

$$v(x_1, \dots, x_i, \dots, x_n) = \tilde{v}(x_i).$$

Thus, we reduced the given n -dimensional problem to n one-dimensional problems. Proceeding as before, we conclude that the functional $v = v(x_1, \dots, x_n)$ can only be extremized on the trajectories that satisfy n Euler–Lagrange equations,

$$F_{x_i} - \frac{d}{dt} F_{\dot{x}_i} = 0, \quad i = 1, \dots, n. \quad (5.23)$$

Observe that if in (5.23) we set $F = L = K - U$, where L is the Lagrangian, $K = K(\dot{q}_1, \dots, \dot{q}_n)$ is the kinetic energy of the particle, $U = U(q_1, \dots, q_n)$ is its potential energy, and $x_i = q_i$, then (5.23) becomes

$$\frac{\partial L}{\partial q_i} - \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}_i} \right) = 0, \quad i = 1, 2, \dots, n,$$

which are the Lagrange equations of motion for a particle. Thus a particle moves along a path on which the functional

$$J = \int_{t_1}^{t_2} L dt$$

is extremized, which agrees with the Hamilton principle saying that a particle moves along the path that minimizes the integral J .

In our discussion up to now, both end points were fixed. We now would like to solve a more general optimization problem, where the end points need not be fixed. For simplicity, we assume first that one of the end points is fixed, while the other is free. After we solve this problem, we can easily generalize to the case when both end points are free.

Problem 2 In the (t, x) plane the left end point A is fixed, while the right end point B is free (see Figure 5.7). It is required to find a trajectory $x = x(t)$ starting at A on which the functional $v = \int_{t_0}^{t_1} F(t, x, x') dt$ can achieve its extremal value.

Solution to Problem 2 Following Elsgolc [75, pp. 64–72], we begin by assuming that a solution to Problem 2 exists; that is, there is a trajectory $x = x(t)$ joining (t_0, x_0) and (t_1, x_1) such that $\delta v = 0$. This trajectory must satisfy the Euler–Lagrange equation. Consider an adjoining

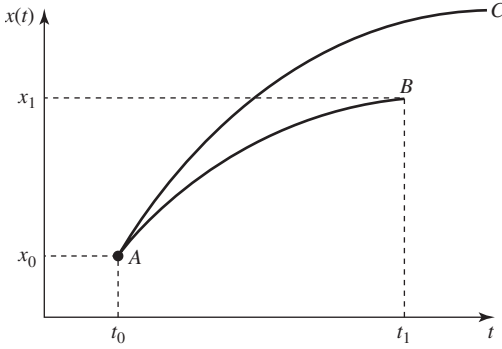


Figure 5.7 Illustration of Problem 2.

trajectory specified by $x + \delta x$ joining (t_0, x_0) and $(t_1 + \delta t_1, x_1 + \delta x_1)$. Then,

$$\begin{aligned} \Delta v &= v(x + \delta x) - v(x) \\ &= \int_{t_0}^{t_1 + \delta t_1} F(t, x + \delta x, x' + \delta x') dt - \int_{t_0}^{t_1} F(t, x, x') dt \\ &= \int_{t_1}^{t_1 + \delta t_1} F(t, x + \delta x, x' + \delta x') dt + \int_{t_0}^{t_1} (F(t, x + \delta x, x' + \delta x') - F(t, x, x')) dt. \end{aligned}$$

Our goal now is to extract from Δv its linear portion in δx , which we know is the variation of the given functional. We begin by approximating the first integral as

$$\int_{t_1}^{t_1 + \delta t_1} F(t, x + \delta x, x' + \delta x') dt \approx F(t, x, x')|_{t=t_1} \delta t_1.$$

Expanding the integrand of the second integral into a Taylor series and retaining only first-order terms yields

$$\int_{t_0}^{t_1} (F(t, x + \delta x, x' + \delta x') - F(t, x, x')) dt = \int_{t_0}^{t_1} (F_x(t, x, x') \delta x + F_{x'}(t, x, x') \delta x') dt.$$

Integrating the second term by parts and combining with the first term gives

$$\int_{t_0}^{t_1} (F_x \delta x + F_{x'} \delta x') dt = (F_{x'} \delta x) \Big|_{t_0}^{t_1} + \int_{t_0}^{t_1} \left(F_x - \frac{d}{dt} F_{x'} \right) \delta x dt,$$

where for simplicity of notation we suppressed the arguments of F_x and $F_{x'}$. We assumed that a solution to the problem exists. Hence, an extremal must satisfy the Euler–Lagrange equation $F_x - \frac{d}{dt} F_{x'} = 0$. Because the left end point is fixed, we obtain

$$\delta x|_{t=t_0} = 0.$$

Taking into account the results of the above manipulations yields

$$\Delta v \approx F(t, x, x')|_{t=t_1} \delta t_1 + (F_{x'} \delta x)|_{t=t_1}. \quad (5.24)$$

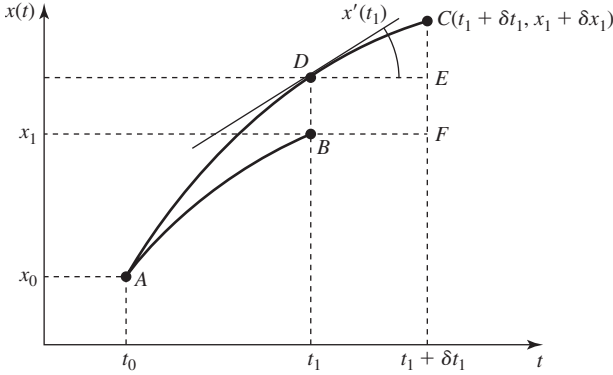


Figure 5.8 A geometric interpretation of calculating the variation of the functional for Problem 2.

The final step in obtaining δv is to evaluate $(F_{x'}\delta x)|_{t_1}$. The problem is two-dimensional, so we can illustrate it graphically as shown in Figure 5.8. Referring to Figure 5.8, we obtain

$$\begin{aligned} DB &= \delta x|_{t=t_1}, \\ FC &= \delta x_1, \\ EC &\approx x'(t_1)\delta t_1, \\ DB &= FC - EC. \end{aligned}$$

Hence,

$$\delta x|_{t=t_1} \approx \delta x_1 - x'(t_1)\delta t_1.$$

Substituting the above into (5.24) gives

$$\begin{aligned} \delta v &= F(t, x, x')|_{t=t_1}\delta t_1 + F_{x'}|_{t=t_1}(\delta x_1 - x'(t_1)\delta t_1) \\ &= (F - x'F_{x'})|_{t=t_1}\delta t_1 + F_{x'}|_{t=t_1}\delta x_1. \end{aligned}$$

If the variations δt_1 and δx_1 are arbitrary, we obtain $\delta v = 0$ if in addition to the Euler–Lagrange equation the following conditions are also satisfied:

$$(F - x'F_{x'})|_{t=t_1} = 0 \quad \text{and} \quad F_{x'}|_{t=t_1} = 0.$$

The above end-point relations are called the *transversality conditions*. In conclusion, the optimal trajectory must satisfy both the Euler–Lagrange equation and the transversality conditions.

◆ Example 5.8

Consider the functional of Example 5.7, where now the left end point is fixed as before—that is, $x(0) = 0$ —while the right end point $x(1)$ is free. Here, we have $t_1 = 1$. Because $\delta t_1 = 0$ and δx_1 is arbitrary, the transversality conditions for the right end point are reduced to one condition,

$$F_{x'}|_{t=t_1} = 0.$$

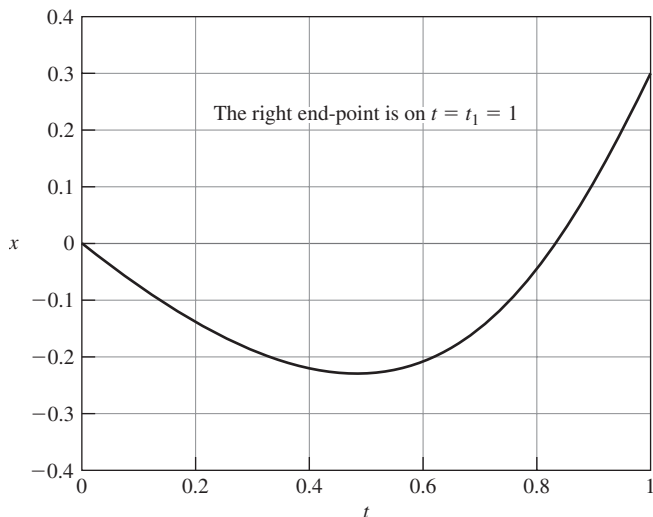


Figure 5.9 A plot of the extremal, $x(t) = -3t + t^3$, of Example 5.8.

Taking into account the solution to the Euler–Lagrange equation obtained in Example 5.7 and then using the fact that the left end point is fixed, we obtain

$$\begin{aligned}
 0 &= F_{x'}|_{t=t_1} \\
 &= 2x'(t_1) \\
 &= 2(C_1 + 3t^2)|_{t=1} \\
 &= 2(C_1 + 3).
 \end{aligned}$$

Hence, $C_1 = -3$ and the extremal for this case is

$$x(t) = -3t + t^3.$$

In Figure 5.9, we show a plot of the above extremal.

◆ Example 5.9

We consider again the same functional as in Example 5.7, where, this time, the left end point is fixed, as before, while at the right end point t_1 is free and $x(t_1) = 1$. Hence, in this case, $\delta x_1 = 0$ while δt_1 is arbitrary. The transversality conditions reduce to one condition,

$$(F - x'F_{x'})|_{t=t_1} = 0.$$

We perform simple calculations to obtain

$$F - x'F_{x'} = 12xt - (x')^2,$$

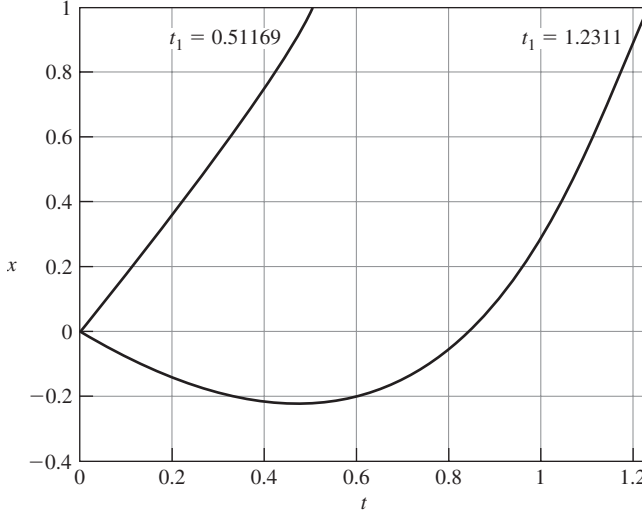


Figure 5.10 Plots of the extremals of Example 5.9.

where

$$x' = (C_1 t + t^3)' = C_1 + 3t^2.$$

Note that because $x(t_1) = C_1 t_1 + t_1^3 = 1$, we have

$$C_1 = \frac{1 - t_1^3}{t_1}.$$

Taking into account the above and the fact that $x(t_1) = 1$, we perform simple algebraic manipulations to arrive at the following equation for t_1 :

$$(F - x'F_{x'})|_{t=t_1} = (12x(t_1)t_1 - (x'(t_1))^2) = -4t_1^6 + 8t_1^3 - 1 = 0.$$

Solving the above equation, we obtain the following two solutions:

$$t_1 = \left(1 \pm \sqrt{\frac{3}{4}}\right)^{1/3}.$$

Plots of the two extremals for this example are shown in Figure 5.10.

If the position of the end point B is restricted to a curve $x = g_1(t)$, then

$$\delta x_1 \approx g_1'(t_1)\delta t_1,$$

and hence

$$\delta v = (F + (g_1' - x')F_{x'})\delta t_1 = 0$$

if

$$(F + (g'_1 - x')F_{x'})|_{t=t_1} = 0. \quad (5.25)$$

The above is the transversality condition for the case when the end point B is restricted to the curve $x = g_1(t)$.

In the next section we discuss a specific problem of minimizing a functional subject to constraints in the form of differential equations.

5.3 Linear Quadratic Regulator

5.3.1 Algebraic Riccati Equation (ARE)

Consider a linear system modeled by

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}, \quad \mathbf{x}(0) = \mathbf{x}_0$$

and the associated performance index

$$J = \int_0^\infty (\mathbf{x}^T \mathbf{Q}\mathbf{x} + \mathbf{u}^T \mathbf{R}\mathbf{u}) dt,$$

where $\mathbf{Q} = \mathbf{Q}^T \geq 0$ and $\mathbf{R} = \mathbf{R}^T > 0$. Our goal is to construct a stabilizing linear state-feedback controller of the form $\mathbf{u} = -\mathbf{K}\mathbf{x}$ that minimizes the performance index J . We denote such a linear control law by \mathbf{u}^* .

First, we assume that a linear state-feedback optimal controller exists such that the optimal closed-loop system

$$\dot{\mathbf{x}} = (\mathbf{A} - \mathbf{B}\mathbf{K})\mathbf{x}$$

is asymptotically stable. The above assumption implies that there is a Lyapunov function $V = \mathbf{x}^T \mathbf{P}\mathbf{x}$ for the closed-loop system; that is, for some positive definite matrix \mathbf{P} the time derivative dV/dt evaluated on the trajectories of the closed-loop system is negative definite.

Theorem 5.2 If the state-feedback controller $\mathbf{u}^* = -\mathbf{K}\mathbf{x}$ is such that

$$\min_{\mathbf{u}} \left(\frac{dV}{dt} + \mathbf{x}^T \mathbf{Q}\mathbf{x} + \mathbf{u}^T \mathbf{R}\mathbf{u} \right) = 0, \quad (5.26)$$

for some $V = \mathbf{x}^T \mathbf{P}\mathbf{x}$, then the controller is optimal.

Proof We can rewrite the condition of the theorem as

$$\left. \frac{dV}{dt} \right|_{\mathbf{u}=\mathbf{u}^*} + \mathbf{x}^T \mathbf{Q}\mathbf{x} + \mathbf{u}^{*T} \mathbf{R}\mathbf{u}^* = 0.$$

Hence,

$$\left. \frac{dV}{dt} \right|_{\mathbf{u}=\mathbf{u}^*} = -\mathbf{x}^T \mathbf{Q}\mathbf{x} - \mathbf{u}^{*T} \mathbf{R}\mathbf{u}^*.$$

Integrating both sides of the resulting equation with respect to time from 0 to ∞ , we obtain

$$V(\mathbf{x}(\infty)) - V(\mathbf{x}(0)) = - \int_0^\infty (\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^{*T} \mathbf{R} \mathbf{u}^*) dt.$$

Because by assumption the closed-loop system is asymptotically stable, $\mathbf{x}(\infty) = \mathbf{0}$, and therefore

$$V(\mathbf{x}(0)) = \mathbf{x}_0^T \mathbf{P} \mathbf{x}_0 = \int_0^\infty (\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^{*T} \mathbf{R} \mathbf{u}^*) dt.$$

Thus, we showed that if a linear state-feedback controller satisfies the assumption of the theorem, then the value of the performance index for such a controller is

$$J(\mathbf{u}^*) = \mathbf{x}_0^T \mathbf{P} \mathbf{x}_0.$$

To show that such a controller is indeed optimal, we use a proof by contradiction. We assume that (5.26) holds and that \mathbf{u}^* is not optimal. Suppose that $\tilde{\mathbf{u}}$ yields a smaller value of J , that is,

$$J(\tilde{\mathbf{u}}) < J(\mathbf{u}^*).$$

It follows from (5.26) that

$$\left. \frac{dV}{dt} \right|_{\mathbf{u}=\tilde{\mathbf{u}}} + \mathbf{x}^T \mathbf{Q} \mathbf{x} + \tilde{\mathbf{u}}^T \mathbf{R} \tilde{\mathbf{u}} \geq 0,$$

that is,

$$\left. \frac{dV}{dt} \right|_{\mathbf{u}=\tilde{\mathbf{u}}} \geq -\mathbf{x}^T \mathbf{Q} \mathbf{x} - \tilde{\mathbf{u}}^T \mathbf{R} \tilde{\mathbf{u}}.$$

Integrating the above with respect to time from 0 to ∞ yields

$$V(\mathbf{x}(0)) \leq \int_0^\infty (\mathbf{x}^T \mathbf{Q} \mathbf{x} + \tilde{\mathbf{u}}^T \mathbf{R} \tilde{\mathbf{u}}) dt$$

implying that

$$J(\mathbf{u}^*) \leq J(\tilde{\mathbf{u}}),$$

which is a contradiction, and the proof is complete.

It follows from the above theorem that the synthesis of the optimal control law involves finding an appropriate Lyapunov function, or equivalently, the matrix \mathbf{P} . The appropriate \mathbf{P} is found by minimizing

$$f(\mathbf{u}) = \frac{dV}{dt} + \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u}, \quad (5.27)$$

We first apply to (5.27) the necessary condition for unconstrained minimization,

$$\left. \frac{\partial}{\partial \mathbf{u}} \left(\frac{dV}{dt} + \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u} \right) \right|_{\mathbf{u}=\mathbf{u}^*} = \mathbf{0}^T.$$

Differentiating the above yields

$$\begin{aligned}
 \frac{\partial}{\partial \mathbf{u}} \left(\frac{dV}{dt} + \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u} \right) &= \frac{\partial}{\partial \mathbf{u}} (2\mathbf{x}^T \mathbf{P} \dot{\mathbf{x}} + \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u}) \\
 &= \frac{\partial}{\partial \mathbf{u}} (2\mathbf{x}^T \mathbf{P} \mathbf{A} \mathbf{x} + 2\mathbf{x}^T \mathbf{P} \mathbf{B} \mathbf{u} + \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u}) \\
 &= 2\mathbf{x}^T \mathbf{P} \mathbf{B} + 2\mathbf{u}^T \mathbf{R} \\
 &= \mathbf{0}^T.
 \end{aligned}$$

Hence, the candidate for an optimal control law has the form

$$\mathbf{u}^* = -\mathbf{R}^{-1} \mathbf{B}^T \mathbf{P} \mathbf{x} = -\mathbf{K} \mathbf{x}, \quad (5.28)$$

where $\mathbf{K} = \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P}$. Note that

$$\begin{aligned}
 \frac{\partial^2}{\partial \mathbf{u}^2} \left(\frac{dV}{dt} + \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u} \right) &= \frac{\partial^2}{\partial \mathbf{u}^2} (2\mathbf{x}^T \mathbf{P} \mathbf{A} \mathbf{x} + 2\mathbf{x}^T \mathbf{P} \mathbf{B} \mathbf{u} + \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u}) \\
 &= \frac{\partial}{\partial \mathbf{u}} (2\mathbf{x}^T \mathbf{P} \mathbf{B} + 2\mathbf{u}^T \mathbf{R}) \\
 &= 2\mathbf{R} \\
 &> 0.
 \end{aligned}$$

Thus the second-order sufficiency condition for \mathbf{u}^* to minimize (5.27) is satisfied.

We now turn our attention to finding an appropriate matrix \mathbf{P} . The optimal closed-loop system has the form

$$\dot{\mathbf{x}} = (\mathbf{A} - \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P}) \mathbf{x}, \quad \mathbf{x}(0) = \mathbf{x}_0.$$

Our optimal controller satisfies the equation

$$\left. \frac{dV}{dt} \right|_{\mathbf{u}=\mathbf{u}^*} + \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^{*T} \mathbf{R} \mathbf{u}^* = 0,$$

that is,

$$2\mathbf{x}^T \mathbf{P} \mathbf{A} \mathbf{x} + 2\mathbf{x}^T \mathbf{P} \mathbf{B} \mathbf{u}^* + \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^{*T} \mathbf{R} \mathbf{u}^* = 0.$$

We substitute the expression for \mathbf{u}^* , given by (5.28), into the above equation and represent it as

$$\mathbf{x}^T (\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A}) \mathbf{x} - 2\mathbf{x}^T \mathbf{P} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P} \mathbf{x} + \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{x}^T \mathbf{P} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P} \mathbf{x} = 0.$$

Factoring out \mathbf{x} yields

$$\mathbf{x}^T (\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} + \mathbf{Q} - \mathbf{P} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P}) \mathbf{x} = 0.$$

The above equation should hold for any \mathbf{x} . For this to be true we have to have

$$\boxed{\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} + \mathbf{Q} - \mathbf{P} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P} = \mathbf{0}.} \quad (5.29)$$

The above equation is referred to as the *algebraic Riccati equation* (ARE). In conclusion, the synthesis of the optimal linear state feedback controller minimizing the performance index

$$J = \int_0^\infty (\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u}) dt$$

subject to

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}, \quad \mathbf{x}(0) = \mathbf{x}_0$$

requires solving the ARE given by (5.29).

◆ Example 5.10

Consider the following model of a dynamical system:

$$\dot{x} = 2u_1 + 2u_2, \quad x(0) = 3,$$

along with the associated performance index

$$J = \int_0^\infty (x^2 + ru_1^2 + ru_2^2) dt,$$

where $r > 0$ is a parameter.

1. We first find the solution to the ARE corresponding to the linear state-feedback optimal controller. We have

$$\mathbf{A} = 0, \quad \mathbf{B} = [2 \quad 2], \quad \mathbf{Q} = 1, \quad \mathbf{R} = r \mathbf{I}_2.$$

The ARE for this problem is

$$\mathbf{O} = \mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} + \mathbf{Q} - \mathbf{P} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P} = 1 - \frac{8}{r} p^2,$$

whose solution is

$$p = \sqrt{\frac{r}{8}}.$$

2. We now write the equation of the closed-loop system driven by the optimal controller. The optimal controller has the form

$$\mathbf{u} = -\mathbf{R}^{-1} \mathbf{B}^T \mathbf{P} \mathbf{x} = -\frac{1}{\sqrt{2r}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} x.$$

Hence, the closed-loop optimal system is described by

$$\dot{x} = [2 \quad 2] \mathbf{u} = -\frac{4}{\sqrt{2r}} x.$$

3. Finally, we find the value of J for the optimal closed-loop system. We have

$$J = \mathbf{x}^T(0) \mathbf{P} \mathbf{x}(0) = \frac{9}{2} \sqrt{\frac{r}{2}}.$$

We can directly verify that $J = \frac{9}{2}\sqrt{\frac{r}{2}}$ as follows:

$$\begin{aligned}\min J &= \min \int_0^\infty (x^2 + ru_1^2 + ru_2^2) dt \\ &= \int_0^\infty 9 \exp\left(-\frac{8t}{\sqrt{2r}}\right) (1 + 1) dt = 9 \frac{\sqrt{2r}}{4} = \frac{9}{2}\sqrt{\frac{r}{2}}.\end{aligned}$$

A major difficulty when solving the ARE is that it is a nonlinear equation. In the next section we present an efficient method for solving the ARE.

5.3.2 Solving the ARE Using the Eigenvector Method

We now present a method for solving the ARE referred to as the MacFarlane and Potter method or the eigenvector method. In our discussion of the method, we follow the development of Power and Simpson [237]. We begin by representing the ARE in the form

$$\begin{bmatrix} P & -I_n \end{bmatrix} \begin{bmatrix} A & -BR^{-1}B^T \\ -Q & -A^T \end{bmatrix} \begin{bmatrix} I_n \\ P \end{bmatrix} = O. \quad (5.30)$$

The $2n \times 2n$ matrix in the middle is called the *Hamiltonian matrix*. We use the symbol H to denote the Hamiltonian matrix, that is,

$$H = \begin{bmatrix} A & -BR^{-1}B^T \\ -Q & -A^T \end{bmatrix}.$$

Thus, the ARE can be represented as

$$\begin{bmatrix} P & -I_n \end{bmatrix} H \begin{bmatrix} I_n \\ P \end{bmatrix} = O.$$

We premultiply the above equation by X^{-1} and then postmultiply it by X , where X is a nonsingular $n \times n$ matrix,

$$\begin{bmatrix} X^{-1}P & -X^{-1} \end{bmatrix} H \begin{bmatrix} X \\ PX \end{bmatrix} = O. \quad (5.31)$$

Observe that if we can find matrices X and PX such that

$$H \begin{bmatrix} X \\ PX \end{bmatrix} = \begin{bmatrix} X \\ PX \end{bmatrix} \Lambda,$$

then equation (5.31) becomes

$$\begin{bmatrix} X^{-1}P & -X^{-1} \end{bmatrix} \begin{bmatrix} X \\ PX \end{bmatrix} \Lambda = O.$$

We have thus reduced the problem of solving the ARE to that of constructing appropriate matrices X and PX . To proceed further, let \mathbf{v}_i be an eigenvector of \mathbf{H} and let s_i be the corresponding eigenvalue; then

$$\mathbf{H}\mathbf{v}_i = s_i\mathbf{v}_i.$$

In our further discussion, we assume that \mathbf{H} has at least n distinct real eigenvalues among its $2n$ eigenvalues. (The results obtained can be generalized for the case when the eigenvalues of \mathbf{H} are complex or nondistinct.) Thus, we can write

$$\mathbf{H}[\mathbf{v}_1 \quad \mathbf{v}_2 \quad \cdots \quad \mathbf{v}_n] = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \cdots \quad \mathbf{v}_n] \begin{bmatrix} s_1 & 0 & \cdots & 0 \\ 0 & s_2 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & s_n \end{bmatrix}.$$

Let

$$\begin{bmatrix} X \\ PX \end{bmatrix} = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \cdots \quad \mathbf{v}_n]$$

and

$$\Lambda = \begin{bmatrix} s_1 & 0 & \cdots & 0 \\ 0 & s_2 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & s_n \end{bmatrix}.$$

The above choice of X and PX constitutes a possible solution to the equation

$$[X^{-1}P \quad -X^{-1}] \begin{bmatrix} X \\ PX \end{bmatrix} \Lambda = \mathbf{O}.$$

To construct P , we partition the $2n \times n$ eigenvector matrix $[\mathbf{v}_1 \quad \mathbf{v}_2 \quad \cdots \quad \mathbf{v}_n]$ into two $n \times n$ submatrices as follows:

$$[\mathbf{v}_1 \quad \mathbf{v}_2 \quad \cdots \quad \mathbf{v}_n] = \begin{bmatrix} \mathbf{W} \\ \mathbf{Z} \end{bmatrix}.$$

Thus,

$$\begin{bmatrix} X \\ PX \end{bmatrix} = \begin{bmatrix} \mathbf{W} \\ \mathbf{Z} \end{bmatrix}.$$

Taking $X = \mathbf{W}$ and $PX = \mathbf{Z}$ and assuming that \mathbf{W} is invertible, we obtain

$$\boxed{P = \mathbf{Z}\mathbf{W}^{-1}}. \quad (5.32)$$

We now have to decide what set of n eigenvalues should be chosen from those of \mathbf{H} to construct the particular P . In the case when all $2n$ eigenvalues of \mathbf{H} are distinct, the number of different

matrices \mathbf{P} generated by the above-described method is given by

$$\frac{(2n)!}{(n!)^2}.$$

Let $\mathbf{Q} = \mathbf{C}^T \mathbf{C}$ be a full rank factorization of \mathbf{Q} , as defined on page 636. It follows from Theorem 3 of Kučera [172, p. 346] that the Hamiltonian matrix \mathbf{H} has n eigenvalues in the open left-half complex plane and n in the open right-half plane if and only if the pair (\mathbf{A}, \mathbf{B}) is stabilizable and the pair (\mathbf{A}, \mathbf{C}) is detectable. The matrix \mathbf{P} that we seek corresponds to the asymptotically stable eigenvalues of \mathbf{H} . The eigenvalues of the Hamiltonian matrix come in pairs $\pm s_i$ (see Exercise 5.15). The characteristic polynomial of \mathbf{H} contains only even powers of s . With such constructed matrix \mathbf{P} , we have the following result:

Theorem 5.3 The poles of the closed-loop system

$$\dot{\mathbf{x}}(t) = (\mathbf{A} - \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P})\mathbf{x}(t)$$

are those eigenvalues of \mathbf{H} having negative real parts.

Proof Because $[\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_n] = \begin{bmatrix} \mathbf{W} \\ \mathbf{Z} \end{bmatrix}$, we can write

$$\begin{bmatrix} \mathbf{A} & -\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T \\ -\mathbf{Q} & -\mathbf{A}^T \end{bmatrix} \begin{bmatrix} \mathbf{W} \\ \mathbf{Z} \end{bmatrix} = \begin{bmatrix} \mathbf{W} \\ \mathbf{Z} \end{bmatrix} \mathbf{\Lambda}.$$

Performing appropriate multiplication of $n \times n$ block matrices yields

$$\mathbf{A}\mathbf{W} - \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{Z} = \mathbf{W}\mathbf{\Lambda},$$

or

$$\mathbf{A} - \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{Z}\mathbf{W}^{-1} = \mathbf{A} - \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P} = \mathbf{W}\mathbf{\Lambda}\mathbf{W}^{-1},$$

because $\mathbf{P} = \mathbf{Z}\mathbf{W}^{-1}$. Thus, the matrix $\mathbf{A} - \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}$ is similar to the matrix $\mathbf{\Lambda}$ whose eigenvalues are the asymptotically stable eigenvalues of \mathbf{H} . The proof is complete.

◆ Example 5.11

Consider the following model of a dynamical system:

$$\dot{x} = 2x + u,$$

along with the associated performance index

$$J = \int_0^\infty (x^2 + ru^2) dt.$$

Find the value of r such that the optimal closed-loop system has its pole at -3 .

We form the associated Hamiltonian matrix

$$\mathbf{H} = \begin{bmatrix} \mathbf{A} & -\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T \\ -\mathbf{Q} & -\mathbf{A}^T \end{bmatrix} = \begin{bmatrix} 2 & -\frac{1}{r} \\ -1 & -2 \end{bmatrix}.$$

The characteristic equation of \mathbf{H} is

$$\det(s\mathbf{I}_2 - \mathbf{H}) = s^2 - 4 - \frac{1}{r} = 0.$$

Hence,

$$r = \frac{1}{5}$$

results in the optimal closed-loop system having its pole located at -3 .

The eigenvector method for solving the ARE can be generalized to the case when the eigenvalues of the Hamiltonian matrix \mathbf{H} are not necessarily distinct or real.

5.3.3 Optimal Systems with Prescribed Poles

We can always find a state feedback $u = -\mathbf{K}\mathbf{x}$ such that the closed-loop system $\dot{\mathbf{x}} = (\mathbf{A} - \mathbf{B}\mathbf{K})\mathbf{x}$ will have its poles in prescribed locations (symmetric with respect to the real axis) if and only if the pair (\mathbf{A}, \mathbf{B}) is reachable. In the multi-input systems the gain matrix \mathbf{K} that shifts the poles to the prescribed locations is not unique. We would like to use the remaining degrees of freedom to construct a gain matrix \mathbf{K} that shifts the system poles to the desired locations and at the same time minimizes a performance index. Specifically, we are given a system model,

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}u(t),$$

the desired set of closed-loop poles,

$$\{s_1, s_2, \dots, s_n\},$$

and an $m \times m$ positive definite matrix $\mathbf{R} = \mathbf{R}^T$. We wish to construct a state-feedback control law $\mathbf{u}(t) = -\mathbf{K}\mathbf{x}(t)$ and an $n \times n$ positive semidefinite matrix $\mathbf{Q} = \mathbf{Q}^T$ such that the closed-loop system has its poles in the desired locations, and at the same time the performance index

$$J = \int_0^\infty (\mathbf{x}^T(t)\mathbf{Q}\mathbf{x}(t) + \mathbf{u}^T(t)\mathbf{R}\mathbf{u}(t)) dt$$

is minimized. In other words, we are interested in designing an optimal control system with prescribed poles. We illustrate the problem with a very simple example.

◆ Example 5.12

For a dynamical system modeled by

$$\begin{aligned}\dot{x} &= 2x + u \\ &= \mathbf{A}x + \mathbf{B}u,\end{aligned}$$

and the performance index

$$\begin{aligned} J &= \int_0^\infty (qx^2 + u^2) dt \\ &= \int_0^\infty (\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{u}^T \mathbf{R} \mathbf{u}) dt, \end{aligned}$$

find q so that the eigenvalue of

$$\mathbf{A} - \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P}$$

is equal to -3 .

The corresponding Hamiltonian matrix is

$$\mathbf{H} = \begin{bmatrix} 2 & -1 \\ -q & -2 \end{bmatrix}.$$

The characteristic polynomial of \mathbf{H} is

$$\det[s\mathbf{I}_2 - \mathbf{H}] = \det \begin{bmatrix} s-2 & 1 \\ q & s+2 \end{bmatrix} = s^2 - 4 - q.$$

We want

$$\det[s\mathbf{I}_2 - \mathbf{H}] = (s+3)(s-3) = s^2 - 9.$$

Hence, the desired value of q is $q = 5$.

We now discuss a method for designing optimal control systems with prescribed poles. In constructing the gain matrix \mathbf{K} and the weight matrix \mathbf{Q} , we repeatedly use a state variable transformation. To proceed, we analyze the system model, the performance index, and the Hamiltonian matrix in a new basis. Let

$$\mathbf{x} = \mathbf{T} \tilde{\mathbf{x}},$$

where $\det \mathbf{T} \neq 0$. Then, the system model in the new coordinates has the form

$$\begin{aligned} \dot{\tilde{\mathbf{x}}}(t) &= \mathbf{T}^{-1} \mathbf{A} \mathbf{T} \tilde{\mathbf{x}}(t) + \mathbf{T}^{-1} \mathbf{B} \mathbf{u}(t) \\ &= \tilde{\mathbf{A}} \tilde{\mathbf{x}}(t) + \tilde{\mathbf{B}} \mathbf{u}(t). \end{aligned}$$

The performance index J , after state variable transformation, takes the form

$$\tilde{J} = \int_0^\infty (\tilde{\mathbf{x}}^T(t) \mathbf{T}^T \mathbf{Q} \mathbf{T} \tilde{\mathbf{x}}(t) + \mathbf{u}^T(t) \mathbf{R} \mathbf{u}(t)) dt.$$

The Hamiltonian matrix is transformed to

$$\begin{aligned} \tilde{\mathbf{H}} &= \begin{bmatrix} \mathbf{T}^{-1} \mathbf{A} \mathbf{T} & -\mathbf{T}^{-1} \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{T}^{-T} \\ -\mathbf{T}^T \mathbf{Q} \mathbf{T} & -\mathbf{T}^T \mathbf{A}^T \mathbf{T}^{-T} \end{bmatrix} \\ &= \begin{bmatrix} \tilde{\mathbf{A}} & -\tilde{\mathbf{B}} \mathbf{R}^{-1} \tilde{\mathbf{B}}^T \\ -\tilde{\mathbf{Q}} & -\tilde{\mathbf{A}}^T \end{bmatrix}, \end{aligned}$$

where $T^{-T} = (T^{-1})^T = (T^T)^{-1}$. Observe that

$$\tilde{H} = \begin{bmatrix} T^{-1} & \mathbf{0} \\ \mathbf{0} & T^T \end{bmatrix} H \begin{bmatrix} T & \mathbf{0} \\ \mathbf{0} & T^{-T} \end{bmatrix};$$

that is, \tilde{H} and H are similarity matrices because

$$\begin{bmatrix} T^{-1} & \mathbf{0} \\ \mathbf{0} & T^T \end{bmatrix} \begin{bmatrix} T & \mathbf{0} \\ \mathbf{0} & T^{-T} \end{bmatrix} = I_{2n}.$$

Hence, the eigenvalues of H and \tilde{H} are identical. Our next step is to compute the characteristic equation of \tilde{H} as a function of \tilde{Q} . Computing $\det(sI_{2n} - \tilde{H})$ is facilitated after we perform some manipulations on the characteristic matrix $sI_{2n} - \tilde{H}$ that do not change its determinant. Specifically, we premultiply the characteristic matrix $sI_{2n} - \tilde{H}$ by a $2n \times 2n$ matrix whose determinant equals 1. Then, the determinant of the product is the same as $\det(sI_{2n} - \tilde{H})$, but it is much easier to evaluate. We have

$$\begin{aligned} & \begin{bmatrix} I_n & \mathbf{0} \\ -\tilde{Q}(sI_n - \tilde{A})^{-1} & I_n \end{bmatrix} \begin{bmatrix} sI_n - \tilde{A} & \tilde{B}R^{-1}\tilde{B}^T \\ \tilde{Q} & sI_n + \tilde{A}^T \end{bmatrix} \\ &= \begin{bmatrix} sI_n - \tilde{A} & \tilde{B}R^{-1}\tilde{B}^T \\ \mathbf{0} & sI_n + \tilde{A}^T - \tilde{Q}(sI_n - \tilde{A})^{-1}\tilde{B}R^{-1}\tilde{B}^T \end{bmatrix}. \end{aligned}$$

Let

$$F = \tilde{B}R^{-1}\tilde{B}^T,$$

then

$$\begin{aligned} \det(sI_{2n} - \tilde{H}) &= \det \begin{bmatrix} sI_n - \tilde{A} & F \\ \mathbf{0} & sI_n + \tilde{A}^T - \tilde{Q}(sI_n - \tilde{A})^{-1}F \end{bmatrix} \\ &= \det(sI_n - \tilde{A}) \det(sI_n + \tilde{A}^T - \tilde{Q}(sI_n - \tilde{A})^{-1}F). \end{aligned}$$

Suppose now that the matrix A has n linearly independent eigenvectors; hence it is diagonalizable. Let

$$T = [t_1 \quad t_2 \quad \cdots \quad t_n]$$

be a matrix composed of the eigenvectors of A . Therefore,

$$\tilde{A} = T^{-1}AT = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix}.$$

Note that $\tilde{\mathbf{A}}^T = \tilde{\mathbf{A}}$. Assume that

$$\tilde{\mathbf{Q}} = \tilde{\mathbf{Q}}_j = \begin{bmatrix} 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & \tilde{q}_{jj} & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix};$$

that is, $\tilde{\mathbf{Q}}_j$ has only one nonzero element \tilde{q}_{jj} . Under the above assumptions the matrix $(s\mathbf{I}_n + \tilde{\mathbf{A}}^T - \tilde{\mathbf{Q}}_j(s\mathbf{I}_n - \tilde{\mathbf{A}})^{-1}\mathbf{F})$ takes the form

$$\begin{aligned} s\mathbf{I}_n + \tilde{\mathbf{A}}^T - \tilde{\mathbf{Q}}_j(s\mathbf{I}_n - \tilde{\mathbf{A}})^{-1}\mathbf{F} &= \begin{bmatrix} s + \lambda_1 & 0 & \cdots & 0 \\ 0 & s + \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & s + \lambda_n \end{bmatrix} \\ &\quad - \begin{bmatrix} 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & \tilde{q}_{jj} & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{s - \lambda_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{s - \lambda_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{s - \lambda_n} \end{bmatrix} \begin{bmatrix} f_{11} & f_{12} & \cdots & f_{1n} \\ f_{21} & f_{22} & \cdots & f_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ f_{n1} & f_{n2} & \cdots & f_{nn} \end{bmatrix} \\ &= \begin{bmatrix} s + \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & s + \lambda_2 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & s + \lambda_j & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & s + \lambda_n \end{bmatrix} - \begin{bmatrix} 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ \frac{\tilde{q}_{jj}f_{j1}}{s - \lambda_j} & \frac{\tilde{q}_{jj}f_{j2}}{s - \lambda_j} & \frac{\tilde{q}_{jj}f_{jj}}{s - \lambda_j} & \cdots & \frac{\tilde{q}_{jj}f_{jn}}{s - \lambda_j} \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix} \\ &= \begin{bmatrix} s + \lambda_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & s + \lambda_2 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ -\frac{\tilde{q}_{jj}f_{j1}}{s - \lambda_j} & -\frac{\tilde{q}_{jj}f_{j2}}{s - \lambda_j} & \cdots & \cdots & s + \lambda_j & -\frac{\tilde{q}_{jj}f_{jj}}{s - \lambda_j} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & s + \lambda_n \end{bmatrix}. \end{aligned}$$

Hence,

$$\det(s\mathbf{I}_n + \tilde{\mathbf{A}}^T - \tilde{\mathbf{Q}}_j(s\mathbf{I}_n - \tilde{\mathbf{A}})^{-1}\mathbf{F}) = \left(\prod_{\substack{i=1 \\ i \neq j}}^n (s + \lambda_i) \right) \left(s + \lambda_j - \frac{\tilde{q}_{jj}f_{jj}}{s - \lambda_j} \right).$$

We can now evaluate the characteristic polynomial of the Hamiltonian matrix,

$$\begin{aligned}
 \det(s\mathbf{I}_{2n} - \mathbf{H}) &= \det(s\mathbf{I}_{2n} - \tilde{\mathbf{H}}) \\
 &= \left(\prod_{i=1}^n (s - \lambda_i) \right) \left(\prod_{\substack{i=1 \\ i \neq j}}^n (s + \lambda_i) \right) \left(s + \lambda_j - \frac{\tilde{q}_{jj} f_{jj}}{s - \lambda_j} \right) \\
 &= \prod_{\substack{i=1 \\ i \neq j}}^n (s + \lambda_i)(s - \lambda_i) ((s + \lambda_j)(s - \lambda_j) - \tilde{q}_{jj} f_{jj}).
 \end{aligned}$$

From the form of the characteristic polynomial of the Hamiltonian matrix, we conclude that selecting an appropriate value of \tilde{q}_{jj} allows us to shift the j th pole of the closed-loop system into a prespecified location while leaving the other poles in their original positions. To see this we represent the term that involves \tilde{q}_{jj} as

$$\begin{aligned}
 (s + \lambda_j)(s - \lambda_j) - \tilde{q}_{jj} f_{jj} &= s^2 - \lambda_j^2 - \tilde{q}_{jj} f_{jj} \\
 &= s^2 - (\lambda_j^2 + \tilde{q}_{jj} f_{jj}) \\
 &= (s - s_j)(s + s_j),
 \end{aligned}$$

where

$$s_j = \pm \sqrt{\lambda_j^2 + \tilde{q}_{jj} f_{jj}}.$$

From the above expression for s_j , we can calculate \tilde{q}_{jj} that results in the j th pole being shifted to $-\sqrt{\lambda_j^2 + \tilde{q}_{jj} f_{jj}}$. We obtain

$$\tilde{q}_{jj} = \frac{s_j^2 - \lambda_j^2}{f_{jj}}. \quad (5.33)$$

Note that while selecting the desired location for s_j , we have to make sure that $\tilde{q}_{jj} > 0$ to ensure positive semidefiniteness of the weight matrix $\tilde{\mathbf{Q}}_j$. Having $\tilde{\mathbf{Q}}_j$, we compute $\mathbf{Q}_j = \mathbf{T}^{-T} \tilde{\mathbf{Q}}_j \mathbf{T}^{-1}$ and solve the Riccati equation

$$\mathbf{A}^T \mathbf{P}_j + \mathbf{P}_j \mathbf{A} + \mathbf{Q}_j - \mathbf{P}_j \mathbf{B} \mathbf{R}^{-1} \mathbf{B}^T \mathbf{P}_j = \mathbf{O} \quad (5.34)$$

for \mathbf{P}_j . We compute $\mathbf{u} = -\mathbf{R}^{-1} \mathbf{B}^T \mathbf{P}_j \mathbf{x} + \mathbf{u}_j = -\mathbf{K}_j \mathbf{x} + \mathbf{u}_j$ that shifts only one pole of the closed-loop system

$$\begin{aligned}
 \dot{\mathbf{x}} &= (\mathbf{A} - \mathbf{B} \mathbf{K}_j) \mathbf{x} + \mathbf{B} \mathbf{u}_j \\
 &= \mathbf{A}_j \mathbf{x} + \mathbf{B} \mathbf{u}_j
 \end{aligned}$$

to its specified position. We then start with the new system

$$\dot{\mathbf{x}} = \mathbf{A}_j \mathbf{x} + \mathbf{B} \mathbf{u}_j$$

and shift the next pole to a specified location, and then we combine the obtained feedback gains. The above method is based on the result, established by Solheim [267], that different feedback

matrices K_j and the corresponding weight matrices Q_j can be combined to give the final optimal gain K and the corresponding weight matrix Q . Indeed, suppose we apply the above method to the system $\dot{x} = (A - BK_j)x + Bu_j$. After finding an appropriate weight matrix Q_{j+1} , we solve the Riccati equation

$$(A - BK_j)^T P_{j+1} + P_{j+1}(A - BK_j) + Q_{j+1} - P_{j+1}BR^{-1}B^T P_{j+1} = 0 \quad (5.35)$$

for P_{j+1} to compute the gain K_{j+1} . Adding the two Riccati equations (5.34) and (5.35), and performing simple manipulations yields

$$A^T(P_j + P_{j+1}) + (P_j + P_{j+1})A + Q_j + Q_{j+1} - (P_j + P_{j+1})BR^{-1}B^T(P_j + P_{j+1}) = 0.$$

Let $P = P_j + P_{j+1}$ and $Q = Q_j + Q_{j+1}$, then

$$\begin{aligned} K &= K_j + K_{j+1} \\ &= R^{-1}B^T P_j + R^{-1}B^T P_{j+1}. \end{aligned}$$

In a similar manner we can show that the method works for more than just two iterations.

◆ Example 5.13

Consider a simple model of one-link robot manipulator shown in Figure 5.11. The motion of the robot arm is controlled by a DC motor via a gear. The DC motor is armature-controlled and its schematic is shown in Figure 5.12. We assume that the motor moment of inertia is negligible compared with that of the robot arm. We model the arm as a point mass m attached to the end of a massless rod of length l . Hence the arm inertia is $I_a = ml^2$. We assume that the gear train has no backlash, and all connecting shafts are rigid. As can be seen from Figure 5.11, counterclockwise rotation of the arm is defined as positive, and clockwise rotation

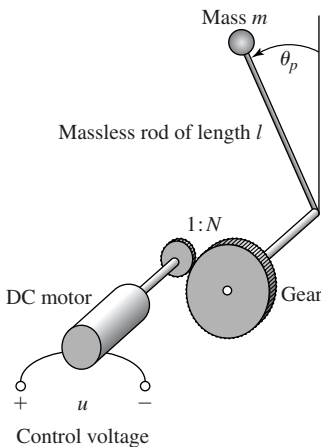


Figure 5.11 One-link manipulator controlled by a DC motor via a gear.

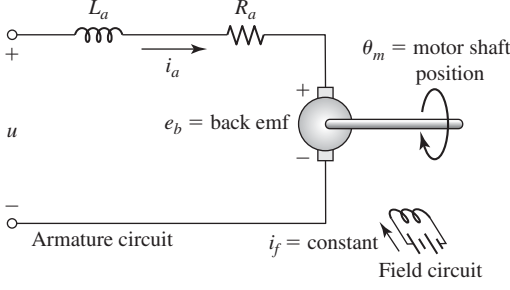


Figure 5.12 Schematic of an armature-controlled DC motor of Example 5.13.

of the arm is defined as negative, whereas counterclockwise rotation of the motor shaft is defined as negative, and clockwise rotation of the shaft is defined as positive. The torque delivered by the motor is

$$T_m = K_m i_a,$$

where K_m is the motor-torque constant, and i_a is the armature current. Let N denote the gear ratio. Then we have

$$\frac{\theta_p}{\theta_m} = \frac{\text{radius of motor gear}}{\text{radius of arm gear}} = \frac{\text{number of teeth of motor gear}}{\text{number of teeth of arm gear}} = \frac{1}{N}.$$

This is because the gears are in contact and thus

$$\theta_p \times \text{radius of arm gear} = \theta_m \times \text{radius of motor gear},$$

and the radii of the gears are proportional to their number of teeth. The work done by the gears must be equal. Let T_p denote the torque applied to the robot arm. Then,

$$T_p \theta_p = T_m \theta_m.$$

Thus, the torque applied to the pendulum is

$$T_p = N T_m = N K_m i_a.$$

We use Newton's second law to write the equation modeling the arm dynamics,

$$I_a \frac{d^2 \theta_p}{dt^2} = mgl \sin \theta_p + T_p. \quad (5.36)$$

Substituting into (5.36) the expressions for I_a and T_p and then rearranging yields

$$ml^2 \frac{d^2 \theta_p}{dt^2} = mgl \sin \theta_p + N K_m i_a, \quad (5.37)$$

where $g = 9.8 \text{ m/sec}^2$ is the gravitational constant. Applying Kirchhoff's voltage law to the armature circuit yields

$$L_a \frac{di_a}{dt} + R_a i_a + K_b N \frac{d\theta_p}{dt} = u,$$

where K_b is the back emf constant. We assume that $L_a \approx 0$. Hence,

$$u = R_a i_a + K_b N \frac{d\theta_p}{dt}, \quad (5.38)$$

We next compute i_a from (5.38) and substitute the result into (5.37) to get

$$ml^2 \frac{d^2\theta_p}{dt^2} = mgl \sin \theta_p + NK_m \left(\frac{u}{R_a} - \frac{K_b N \frac{d\theta_p}{dt}}{R_a} \right). \quad (5.39)$$

We can now construct a state-space model of the one-link robot. For this we choose the following state and output variables:

$$x_1 = \theta_p, \quad x_2 = \frac{d\theta_p}{dt} = \omega_p, \quad \text{and} \quad y = x_1.$$

Then, using (5.39), we obtain the following simple state-space model of the robot manipulator:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ \frac{g}{l} \sin x_1 - \frac{K_b K_m N^2}{ml^2 R_a} x_2 + \frac{NK_m}{ml^2 R_a} u \end{bmatrix},$$

$$y = x_1.$$

Reasonable parameters for the robot are: $l = 1 \text{ m}$, $m = 1 \text{ kg}$, $N = 10$, $K_m = 0.1 \text{ Nm/A}$, $K_b = 0.1 \text{ Vsec/rad}$, $R_a = 1 \Omega$. With the above parameters the robot model takes the form:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ 9.8 \sin x_1 - x_2 + u \end{bmatrix},$$

$$y = x_1.$$

Time histories of state trajectories of the uncontrolled nonlinear system, $u = 0$, are shown in Figure 5.13 for the initial conditions $x_1(0) = 1$ and $x_2(0) = 0$. A phase portrait of the uncontrolled nonlinear system is shown in Figure 5.14. The linearized model about $\mathbf{x} = \mathbf{0}$, $u = 0$ has the form

$$\frac{d}{dt} \Delta \mathbf{x} = \begin{bmatrix} 0 & 1 \\ 9.8 & -1 \end{bmatrix} \Delta \mathbf{x} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \Delta u, \quad (5.40)$$

$$\Delta y = [1 \quad 0] \Delta \mathbf{x}.$$

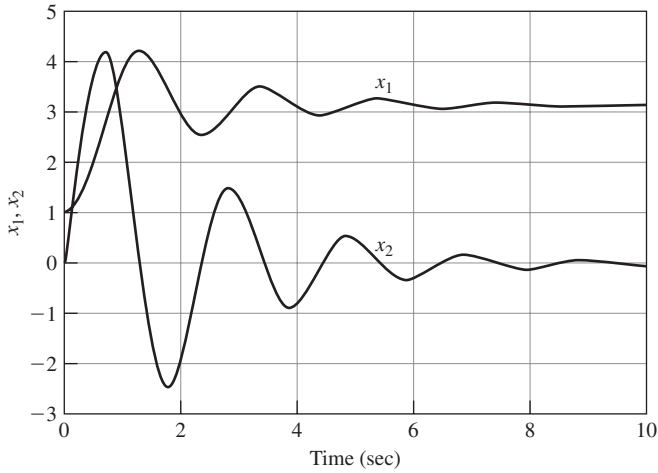


Figure 5.13 Plots of $y = x_1$ and x_2 versus time of the uncontrolled nonlinear system of Example 5.13.

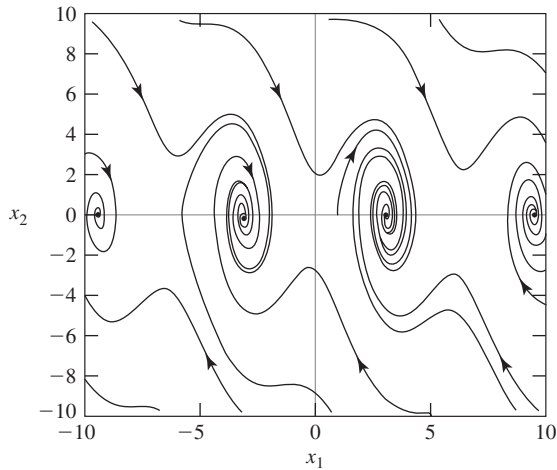


Figure 5.14 A phase portrait of the uncontrolled nonlinear system of Example 5.13.

Plots of $y = x_1$ and x_2 versus time are shown in Figure 5.15. A phase portrait of the linearized system is shown in Figure 5.16. Let

$$J = \int_0^{\infty} (y^2 + u^2) dt.$$

We shall find the linear state-feedback control law $u = -\mathbf{k}\mathbf{x}$ that minimizes J subject to the equations given by (5.40). We have

$$\mathbf{Q} = \mathbf{c}^T \mathbf{c} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{R} = [1].$$

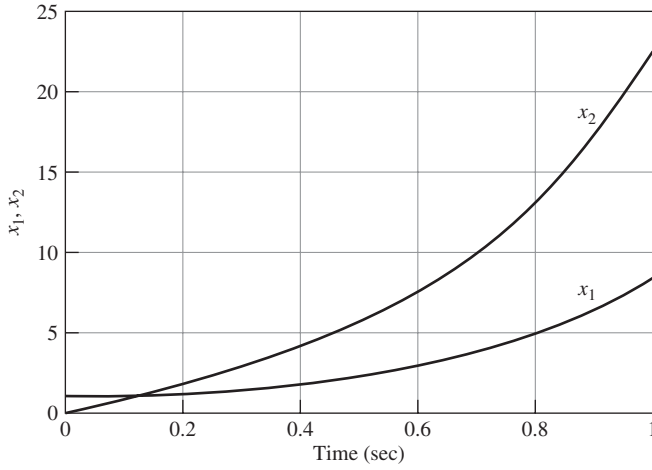


Figure 5.15 Plots of x_1 and x_2 versus time of the linearized system of Example 5.13.

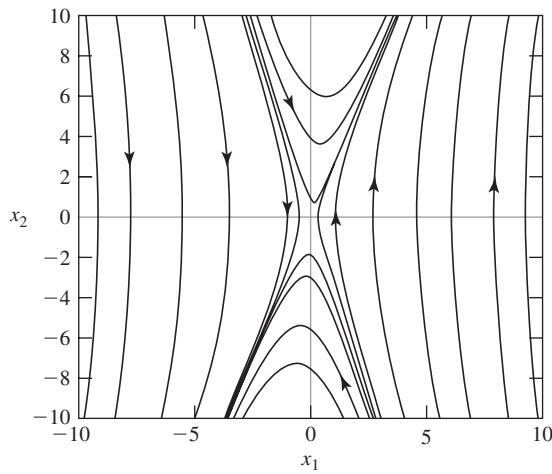


Figure 5.16 A phase portrait of the linearized system of Example 5.13.

Solving the associated Riccati equation yields

$$\mathbf{P} = \begin{bmatrix} 72.3371 & 19.6509 \\ 19.6509 & 5.3484 \end{bmatrix}. \quad (5.41)$$

Hence,

$$\mathbf{k} = [19.6509 \quad 5.3484] \quad (5.42)$$

and

$$\mathbf{A}_c = \begin{bmatrix} 0 & 1 \\ -9.8509 & -6.3484 \end{bmatrix}.$$

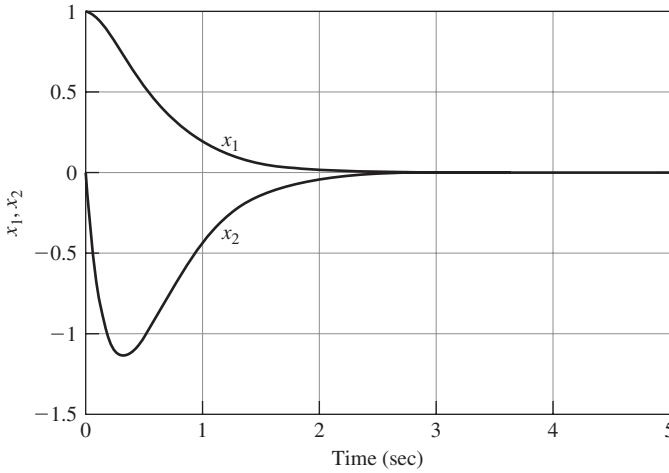


Figure 5.17 Plots of x_1 and x_2 versus time of the linear closed-loop system of Example 5.13.

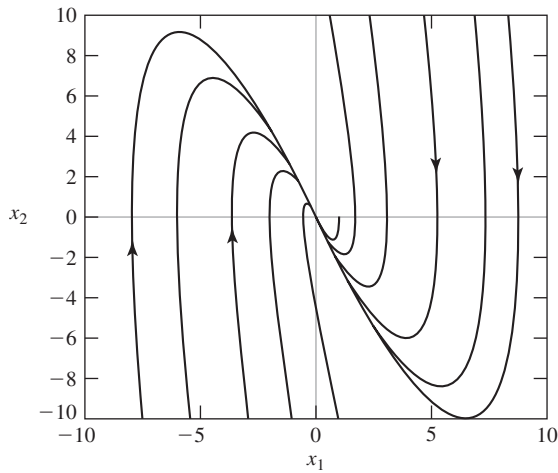


Figure 5.18 A phase portrait of the linear closed-loop system of Example 5.13.

Plots of x_1 and x_2 versus time of the closed-loop linear system

$$\dot{\mathbf{x}} = (\mathbf{A} - \mathbf{bk})\mathbf{x} = \mathbf{A}_c\mathbf{x}$$

$$y = x_1,$$

when the initial conditions are $x_1(0) = 1$ and $x_2(0) = 0$, are shown in Figure 5.17. A phase portrait of the closed-loop linear system is shown in Figure 5.18. We now apply the above optimal controller to the nonlinear model. Plots of x_1 and x_2 versus time of the closed-loop nonlinear systems are shown in Figure 5.19. A phase portrait of the nonlinear closed-loop system is shown in Figure 5.20. The closed-loop poles of the linearized system—that is, the eigenvalues of \mathbf{A}_c —are $\lambda_1 = -2.7003$ and $\lambda_2 = -3.6481$.

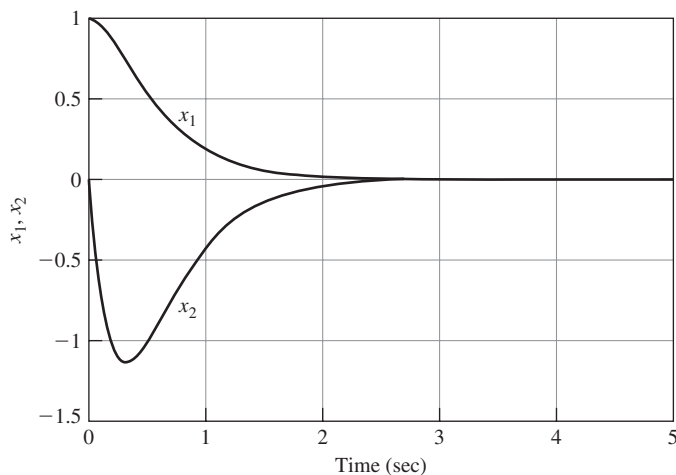


Figure 5.19 Plots of x_1 and x_2 of the nonlinear closed-loop system of Example 5.13.

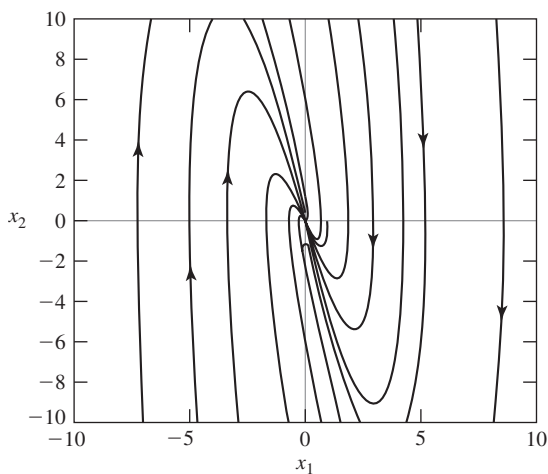


Figure 5.20 A phase portrait of the nonlinear closed-loop system of Example 5.13.

◆ Example 5.14

We use the data from the previous example to find the weight matrix \mathbf{Q}_2 and the corresponding \mathbf{k}_2 so that λ_2 is shifted to $s_2 = -5$; that is, the new poles of the closed-loop system

$$\dot{\mathbf{x}} = (\mathbf{A}_c - \mathbf{b}\mathbf{k}_2)\mathbf{x}$$

are $s_1 = \lambda_1$ and $s_2 = -5$. For the control law $u = -\mathbf{k}\mathbf{x} + r$, the linear closed-loop system is

$$\begin{aligned}\dot{\mathbf{x}} &= \begin{bmatrix} 0 & 1 \\ -9.8509 & -6.3484 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u \\ &= \mathbf{A}_c \mathbf{x} + \mathbf{b}r,\end{aligned}$$

where for simplicity we used \mathbf{x} rather than $\Delta\mathbf{x}$, and u instead of Δu . The eigenvalues of \mathbf{A}_c are $\lambda_1 = -2.7003$ and $\lambda_2 = -3.6481$. We have $\lambda_1 \neq \lambda_2$ and thus the matrix \mathbf{A}_c is diagonalizable. Note that different eigenvalues are only a sufficient condition for diagonalizability. We construct the similarity transformation using the eigenvectors of \mathbf{A}_c . The eigenvectors of \mathbf{A}_c are

$$\mathbf{v}_1 = \begin{bmatrix} 0.3473 \\ -0.9378 \end{bmatrix} \quad \text{and} \quad \mathbf{v}_2 = \begin{bmatrix} -0.2644 \\ -0.9644 \end{bmatrix}.$$

Let $\mathbf{x} = \mathbf{T}\mathbf{z}$, where $\mathbf{T} = [\mathbf{v}_1 \ \mathbf{v}_2]$. Then, \mathbf{A}_c in the new coordinates has the form

$$\tilde{\mathbf{A}}_c = \mathbf{T}^{-1} \mathbf{A}_c \mathbf{T} = \begin{bmatrix} -2.7003 & 0 \\ 0 & -3.6481 \end{bmatrix}$$

and

$$\tilde{\mathbf{b}} = \mathbf{T}^{-1} \mathbf{b} = \begin{bmatrix} 3.0383 \\ 3.9912 \end{bmatrix}.$$

Next, we compute

$$\mathbf{T}^{-1} \mathbf{b} \mathbf{R}^{-1} \mathbf{b}^T \mathbf{T}^{-T} = \begin{bmatrix} 9.2312 & 12.1264 \\ 12.1264 & 15.9296 \end{bmatrix}.$$

Hence,

$$\tilde{q}_{22} = \frac{s_2^2 - \lambda_2^2}{15.9296} = 0.7339.$$

Therefore,

$$\tilde{\mathbf{Q}}_2 = \begin{bmatrix} 0 & 0 \\ 0 & 0.7339 \end{bmatrix} = \mathbf{T}^T \mathbf{Q}_2 \mathbf{T}$$

and

$$\mathbf{Q}_2 = \mathbf{T}^{-T} \tilde{\mathbf{Q}}_2 \mathbf{T}^{-1} = \begin{bmatrix} 85.2450 & 31.5686 \\ 31.5686 & 11.6907 \end{bmatrix}.$$

Solving the ARE with \mathbf{A}_c , \mathbf{b} , \mathbf{Q}_2 , and \mathbf{R} yields

$$\mathbf{P}_2 = \begin{bmatrix} 9.8572 & 3.6504 \\ 3.6504 & 1.3518 \end{bmatrix}$$

and

$$\mathbf{k}_2 = [3.6504 \quad 1.3518].$$

As a check, if we take

$$\mathbf{Q}_{new} = \mathbf{Q} + \mathbf{Q}_2 = \begin{bmatrix} 86.2450 & 31.5686 \\ 31.5686 & 11.6907 \end{bmatrix},$$

where $\mathbf{Q} = \mathbf{c}^T \mathbf{c}$ as in Example 5.13, and $\mathbf{R}_{new} = \mathbf{R}$, then solving the associated ARE using \mathbf{A} , rather than \mathbf{A}_C , we obtain

$$\mathbf{P}_{new} = \begin{bmatrix} 82.1943 & 23.3013 \\ 23.3013 & 6.7002 \end{bmatrix} = \mathbf{P} + \mathbf{P}_2,$$

where \mathbf{P} is given by (5.41), and

$$\mathbf{k}_{new} = [23.3013 \quad 6.7002] = \mathbf{k} + \mathbf{k}_2,$$

where \mathbf{k} is the optimal gain given by (5.42).

◆ Example 5.15

For the linearized model given by (5.40) and $\mathbf{R} = [2]$, we will now determine \mathbf{Q} and the corresponding \mathbf{k} that result in the closed-loop system $\dot{\mathbf{x}} = (\mathbf{A} - \mathbf{b}\mathbf{k})\mathbf{x}$ with poles located at $s_1 = -4$ and $s_2 = -5$, where

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 9.8 & -1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

We shift poles into the desired locations one by one. We first diagonalize \mathbf{A} using similarity transformation constructed from eigenvectors of \mathbf{A} . The eigenvalues of \mathbf{A} are $\lambda_1 = 2.6702$ and $\lambda_2 = -3.6702$. We denote the corresponding eigenvectors by \mathbf{v}_1 and \mathbf{v}_2 , respectively. We obtain

$$\mathbf{T}^{-1} \mathbf{A} \mathbf{T} = \begin{bmatrix} 2.6702 & 0 \\ 0 & -3.6702 \end{bmatrix},$$

where

$$\mathbf{T} = [\mathbf{v}_1 \quad \mathbf{v}_2] = \begin{bmatrix} 0.3507 & -0.2629 \\ 0.9365 & 0.9648 \end{bmatrix}.$$

We will now construct \mathbf{k}_1 so that λ_1 is shifted to $s_1 = -4$, while λ_2 remains unchanged; that is, the eigenvalues of $\mathbf{A} - \mathbf{b}\mathbf{k}_1$ become s_1 and λ_2 . For this, we first compute

$$\mathbf{F} = \mathbf{T}^{-1} \mathbf{b} \mathbf{R}^{-1} \mathbf{T}^T \mathbf{b}^T = \begin{bmatrix} 0.1011 & 0.1349 \\ 0.1349 & 0.1800 \end{bmatrix}.$$

Hence, $f_{11} = 0.1011$,

$$\tilde{q}_{11} = \frac{4^2 - (2.6702)^2}{0.1011} = 87.7352,$$

and

$$\tilde{Q}_1 = \begin{bmatrix} 87.7352 & 0 \\ 0 & 0 \end{bmatrix} = T^T Q_1 T.$$

Hence,

$$Q_1 = T^{-T} \tilde{Q}_1 T^{-1} = \begin{bmatrix} 239.0024 & 65.1202 \\ 65.1202 & 17.7431 \end{bmatrix}.$$

Solving the associated ARE, using $R = [2]$, yields

$$P_1 = \begin{bmatrix} 179.7014 & 48.9626 \\ 48.9626 & 13.3407 \end{bmatrix}.$$

Therefore,

$$k_1 = [24.4813 \quad 6.6703]$$

and

$$A_1 = A - bk_1 = \begin{bmatrix} 0 & 1 \\ -14.6813 & -7.6703 \end{bmatrix}.$$

We will now construct k_2 that shifts λ_2 to s_2 leaving s_1 unchanged. We first diagonalize A_1 using its eigenvectors. Note that A_1 and A share the common eigenvector v_2 . The other eigenvector of A_1 corresponding to its eigenvalue s_1 is $u_1 = [-0.2425 \quad 0.9701]^T$. We have

$$T^{-1} A_1 T = \begin{bmatrix} s_1 & 0 \\ 0 & \lambda_2 \end{bmatrix},$$

where

$$T = [u_1 \quad v_2] = \begin{bmatrix} -0.2425 & 0.2629 \\ 0.9701 & -0.9648 \end{bmatrix}.$$

We next compute

$$T^{-1} b R^{-1} b^T T^{-T} = \begin{bmatrix} 78.0607 & 72.0157 \\ 72.0157 & 66.4389 \end{bmatrix}.$$

Thus, $\tilde{f}_{22} = 66.4389$,

$$\tilde{q}_{22} = \frac{5^2 - (-3.6702)^2}{66.4389} = 0.1735,$$

and

$$\tilde{Q}_2 = \begin{bmatrix} 0 & 0 \\ 0 & 0.1735 \end{bmatrix} = T^T Q_2 T.$$

From the above, we obtain

$$Q_2 = T^{-T} \tilde{Q}_2 T^{-1} = \begin{bmatrix} 368.9003 & 92.2211 \\ 92.2211 & 23.0543 \end{bmatrix}.$$

Solving the associated ARE, we get

$$\mathbf{P}_2 = \begin{bmatrix} 42.5495 & 10.6369 \\ 10.6369 & 2.6591 \end{bmatrix}$$

and

$$\mathbf{k}_2 = [5.3185 \quad 1.3296].$$

The final weight matrix \mathbf{Q} that results in shifting poles into the desired locations is

$$\mathbf{Q} = \mathbf{Q}_1 + \mathbf{Q}_2 = \begin{bmatrix} 607.9027 & 157.3413 \\ 157.3413 & 40.7974 \end{bmatrix}.$$

The corresponding \mathbf{P} is

$$\mathbf{P} = \mathbf{P}_1 + \mathbf{P}_2 = \begin{bmatrix} 222.2509 & 59.5996 \\ 59.5996 & 15.9998 \end{bmatrix}.$$

The corresponding optimal gain vector \mathbf{k} is

$$\mathbf{k} = \mathbf{k}_1 + \mathbf{k}_2 = [29.7998 \quad 7.9999].$$

5.3.4 Optimal Saturating Controllers

Asymptotic stability of an equilibrium state of a nonlinear system is a local property. This means that the system's trajectories starting in a neighborhood of the equilibrium state will approach this equilibrium asymptotically. Often, we are interested in determining how far from the equilibrium state the system's trajectories can start and still converge to the equilibrium asymptotically. A region around an equilibrium state within which all trajectories approach the equilibrium is called a *region of attraction*, or *region of asymptotic stability* (RAS). Finding the exact RAS analytically can be a challenging task. A tool often used to estimate the exact RAS is a Lyapunov function. Specifically, if $V = V(\mathbf{x})$ is a Lyapunov function of the equilibrium state of interest—that is, V is positive definite with respect to the equilibrium state in a domain \mathcal{S} and the time derivative of V evaluated on the trajectories of the system is negative definite in \mathcal{S} —then the set

$$\Omega_c = \{\mathbf{x} : \mathbf{x} \in \mathcal{S}, V(\mathbf{x}) \leq c\}$$

is an estimate of the exact RAS. In other words, a region of attraction is a set of states \mathbf{x} for which $\dot{V}(\mathbf{x}) < 0$ and $V(\mathbf{x}) \leq c$ for some positive constant c . The set Ω_c may be a conservative estimate, often a crude underestimate, of the true RAS. In order to obtain the least conservative estimate of the domain of attraction, corresponding to the given Lyapunov function, we could try to maximize the size of Ω_c by finding the largest V contour that can just fit on the equilibrium side of the surface described by $\dot{V}(\mathbf{x}) = 0$. This can be accomplished by solving the optimization problem for c such that

$$c = \min_{\dot{V}(\mathbf{x})=0} V(\mathbf{x}). \quad (5.43)$$

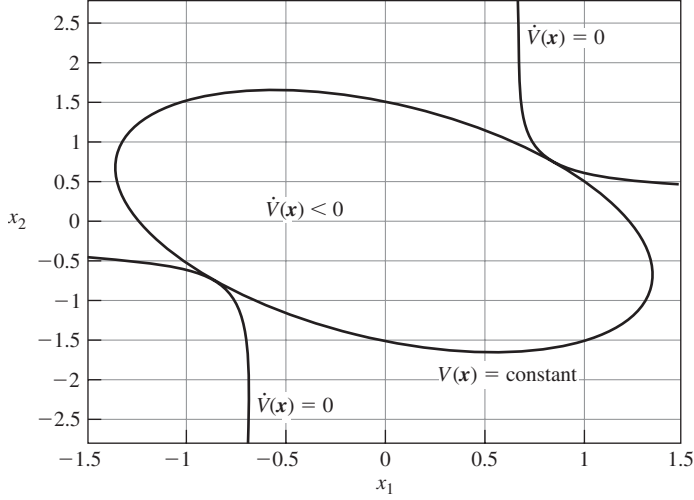


Figure 5.21 Estimating the region of asymptotic stability.

Solving the above optimization problem corresponds to maximizing the size of the region $\{\mathbf{x} : V(\mathbf{x}) \leq c\}$ in which $\dot{V}(\mathbf{x}) < 0$ by making the boundary of $\{\mathbf{x} : V(\mathbf{x}) \leq c\}$ to be tangent to the boundary of $\{\mathbf{x} : \dot{V}(\mathbf{x}) < 0\}$. This is illustrated in Figure 5.21.

We next estimate the RAS for single-input linear dynamical systems with saturating linear optimal control. Our development follows that of Friedland [91, Section 5.1]. We will use such controllers in Chapter 11 to stabilize chaotic systems. We consider a single-input linear system model,

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b}U_{\max}\text{sat}(u/U_{\max}), \quad (5.44)$$

where $\mathbf{x} \in \mathbb{R}^n$ and

$$\text{sat}(z) = \begin{cases} 1 & \text{if } z \geq 1 \\ z & \text{if } |z| < 1 \\ -1 & \text{if } z \leq -1. \end{cases}$$

In the linear region, $|u| \leq U_{\max}$, (5.44) becomes $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b}u$. We will consider using the control law of the form

$$u = U_{\max}\text{sat}(-\mathbf{k}\mathbf{x}/U_{\max}). \quad (5.45)$$

When $|\mathbf{k}\mathbf{x}| \leq U_{\max}$, the control law given by (5.45) becomes a linear state-feedback controller, $u = -\mathbf{k}\mathbf{x}$. For simplicity of notation, we can assume that $U_{\max} = 1$. Indeed, we can always scale the input vector \mathbf{b} so that $U_{\max} = 1$. The gain vector \mathbf{k} is obtained by solving the associated LQR problem with the performance index

$$J = \int_0^\infty (\mathbf{x}^T \mathbf{Q}\mathbf{x} + ru^2) dt,$$

where we assume that the pair (\mathbf{A}, \mathbf{b}) is reachable, $\mathbf{Q} = \mathbf{Q}^T > 0$, and $r > 0$. Solving the algebraic

Riccati equation,

$$A^T P + P A + Q - P b b^T P / r = O, \quad (5.46)$$

we use its solution, $P = P^T > 0$, to construct the gain vector,

$$k = b^T P / r. \quad (5.47)$$

We then proceed with evaluating of $\frac{d}{dt} V = \frac{d}{dt} x^T P x$ on the trajectories of the closed-loop system,

$$\begin{aligned} \dot{V} &= 2x^T P \dot{x} \\ &= x^T (A^T P + P A)x - 2x^T P b \text{sat}(kx). \end{aligned} \quad (5.48)$$

By (5.47),

$$b^T P = rk. \quad (5.49)$$

Using the above and (5.46) gives

$$A^T P + P A = rk^T k - Q. \quad (5.50)$$

Substituting (5.50) and (5.49) into (5.48) yields

$$\dot{V} = -x^T Q x - 2rkx \text{sat}(kx) + r(kx)^2 \quad (5.51)$$

For x such that $|kx| \leq 1$, the above becomes

$$\dot{V} = -x^T Q x - r(kx)^2,$$

which is negative definite. Thus, the origin is (locally) asymptotically stable for the closed-loop system.

An estimate of the RAS is a set of states bounded by the surface $V = \text{constant}$ for which $\dot{V} < 0$. To find as large an estimate as possible, we will attempt to find the surface $\dot{V} = 0$ and find $c = \min V(x)$ on that surface. Then, the region bounded by $V(x) = c$ is an estimate of the RAS. Solving the above problem amounts to solving the optimization problem,

$$\begin{aligned} &\min_x x^T P x \\ &\text{subject to } -x^T Q x - 2rkx \text{sat}(kx) + r(kx)^2 = 0. \end{aligned} \quad (5.52)$$

Solving the above problem is not a simple task. Instead, we will solve an easier problem resulting in a more conservative estimate of the RAS. We first rewrite (5.52) as

$$y^2 - 2y \text{sat}(y) - x^T Q x / r = 0, \quad (5.53)$$

where $y = kx$. We will investigate the case when $|y| \geq 1$ because for $|y| < 1$, we have $\dot{V} < 0$ and we are interested in the solution to the above when $\dot{V} = 0$. When $y \geq 1$, (5.53) becomes

$$y^2 - 2y - x^T Q x / r = 0. \quad (5.54)$$

because for $y \geq 1$, $\text{sat}(y) = 1$. The solution to (5.54) for $y \geq 1$ is

$$y = 1 + \sqrt{1 + x^T Q x / r} > 2, \quad (5.55)$$

while for $y \leq -1$ the solution to (5.53) has the form

$$y = -1 - \sqrt{1 + \mathbf{x}^T \mathbf{Q} \mathbf{x} / r} < -2. \quad (5.56)$$

The above solutions satisfy the inequality

$$|y| \geq 2.$$

Thus, we have $\dot{V} < 0$ for $|y| = |\mathbf{k}\mathbf{x}| \leq 2$ and hence the RAS includes the region

$$\{\mathbf{x} : V(\mathbf{x}) \leq \tilde{c}\},$$

where \tilde{c} is obtained by solving the optimization problem

$$\tilde{c} = \max_{|\mathbf{k}\mathbf{x}| \leq 2} \mathbf{x}^T \mathbf{P} \mathbf{x}$$

We will use the Lagrange multipliers method to solve the above constrained optimization problem. We will solve the equivalent problem,

$$\begin{aligned} & \max \mathbf{x}^T \mathbf{P} \mathbf{x} \\ & \text{subject to } (\mathbf{k}\mathbf{x})^2 = 4. \end{aligned}$$

The Lagrangian for the above problem is

$$l(\mathbf{x}, \lambda) = \mathbf{x}^T \mathbf{P} \mathbf{x} + \lambda(4 - (\mathbf{k}\mathbf{x})^2).$$

Applying to the Lagrangian l the first-order necessary conditions for unconstrained minimization gives

$$D_{\mathbf{x}} l(\mathbf{x}, \lambda) = 2\mathbf{x}^T \mathbf{P} - 2\lambda(\mathbf{k}\mathbf{x})\mathbf{k} = \mathbf{0}^T, \quad (5.57)$$

$$D_{\lambda} l(\mathbf{x}, \lambda) = 4 - (\mathbf{k}\mathbf{x})^2 = 0. \quad (5.58)$$

From (5.57), we obtain

$$\mathbf{x}^T \mathbf{P} \mathbf{x} = \lambda(\mathbf{k}\mathbf{x})^2.$$

Because $(\mathbf{k}\mathbf{x})^2 = 4$,

$$\lambda = \frac{1}{4} \mathbf{x}^T \mathbf{P} \mathbf{x},$$

and therefore

$$\tilde{c} = \mathbf{x}^T \mathbf{P} \mathbf{x} = 4\lambda.$$

Again, using (5.57), we obtain

$$\mathbf{x}^T \mathbf{P} = \lambda(\mathbf{k}\mathbf{x})\mathbf{k}.$$

Postmultiplying both sides of the above by $\mathbf{P}^{-1}\mathbf{k}^T$ gives

$$\mathbf{x}^T \mathbf{k}^T = \mathbf{k}\mathbf{x} = \lambda(\mathbf{k}\mathbf{x})\mathbf{k}\mathbf{P}^{-1}\mathbf{k}^T. \quad (5.59)$$

Note that $\mathbf{k}\mathbf{x} \neq 0$ because $(\mathbf{k}\mathbf{x})^2 = 4$. Hence, from (5.59),

$$\lambda = \frac{1}{\mathbf{k}\mathbf{P}^{-1}\mathbf{k}^T}.$$

Substituting (5.47) into the above gives

$$\lambda = \frac{1}{\mathbf{k} \mathbf{P}^{-1} \mathbf{k}^T} = \frac{r^2}{\mathbf{b}^T \mathbf{P} \mathbf{b}}.$$

Hence,

$$\tilde{c} = 4\lambda = \frac{4r^2}{\mathbf{b}^T \mathbf{P} \mathbf{b}}$$

5.3.5 Linear Quadratic Regulator for Discrete Systems on an Infinite Time Interval

Consider a discrete linear plant model,

$$\mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k), \quad k = 0, 1, 2, \dots, \quad (5.60)$$

with a specified initial condition $\mathbf{x}(0) = \mathbf{x}_0$, where $\mathbf{x}(k) \in \mathbb{R}^n$ and $\mathbf{u}(k) \in \mathbb{R}^m$. We assume that the pair (\mathbf{A}, \mathbf{B}) is reachable. Our objective is to construct a stabilizing linear state-feedback controller,

$$\mathbf{u}(k) = -\mathbf{K}\mathbf{x}(k),$$

that minimizes the quadratic performance index

$$J(\mathbf{u}) = \sum_{k=0}^{\infty} \{ \mathbf{x}^T(k) \mathbf{Q} \mathbf{x}(k) + \mathbf{u}^T(k) \mathbf{R} \mathbf{u}(k) \}, \quad (5.61)$$

where $\mathbf{Q} = \mathbf{Q}^T \geq 0$, and $\mathbf{R} = \mathbf{R}^T > 0$. We assume that the components of the control vector are unconstrained. An optimal control law will be denoted \mathbf{u}^* .

Because we seek the controller among the ones that stabilize the system (5.60), the closed-loop system,

$$\mathbf{x}(k+1) = (\mathbf{A} - \mathbf{B}\mathbf{K})\mathbf{x}(k), \quad (5.62)$$

is asymptotically stable; that is, the eigenvalues of the matrix $\mathbf{A} - \mathbf{B}\mathbf{K}$ are in the open unit disk. This implies that there exists a Lyapunov function, $V(\mathbf{x}(k)) = \mathbf{x}^T(k) \mathbf{P} \mathbf{x}(k)$, for the system (5.62). Therefore the first forward difference, $\Delta V(\mathbf{x}(k)) = V(\mathbf{x}(k+1)) - V(\mathbf{x}(k))$, evaluated on the trajectories of (5.62) is negative definite. We now state and prove a sufficient condition for $\mathbf{u}(k) = -\mathbf{K}\mathbf{x}(k)$ to be optimal.

Theorem 5.4 If the state feedback controller, $\mathbf{u}(k) = -\mathbf{K}\mathbf{x}(k)$, is such that

$$\min_{\mathbf{u}} (\Delta V(\mathbf{x}(k)) + \mathbf{x}^T(k) \mathbf{Q} \mathbf{x}(k) + \mathbf{u}^T(k) \mathbf{R} \mathbf{u}(k)) = 0, \quad (5.63)$$

then it is optimal.

Proof We can represent (5.63) as

$$\Delta V(\mathbf{x}(k))|_{\mathbf{u}=\mathbf{u}^*} + \mathbf{x}^T(k) \mathbf{Q}\mathbf{x}(k) + \mathbf{u}^{*T}(k) \mathbf{R}\mathbf{u}^*(k) = 0.$$

Hence,

$$\Delta V(\mathbf{x}(k))|_{\mathbf{u}=\mathbf{u}^*} = (V(\mathbf{x}(k+1)) - V(\mathbf{x}(k))|_{\mathbf{u}=\mathbf{u}^*}) = -\mathbf{x}^T(k) \mathbf{Q}\mathbf{x}(k) - \mathbf{u}^{*T}(k) \mathbf{R}\mathbf{u}^*(k).$$

We sum both sides from $k = 0$ to $k = \infty$ to get

$$V(\mathbf{x}(\infty)) - V(\mathbf{x}(0)) = -\sum_{k=0}^{\infty} (\mathbf{x}^T(k) \mathbf{Q}\mathbf{x}(k) + \mathbf{u}^{*T}(k) \mathbf{R}\mathbf{u}^*(k)).$$

Because by assumption the closed-loop system is asymptotically stable, we have $\mathbf{x}(\infty) = \mathbf{0}$. Hence,

$$V(\mathbf{x}(0)) = \mathbf{x}_0^T \mathbf{P} \mathbf{x}_0 = \sum_{k=0}^{\infty} (\mathbf{x}^T(k) \mathbf{Q}\mathbf{x}(k) + \mathbf{u}^{*T}(k) \mathbf{R}\mathbf{u}^*(k)).$$

Thus, if a linear stabilizing controller satisfies the hypothesis of the theorem, then the value of the performance index (5.61) resulting from applying this controller is

$$J = \mathbf{x}_0^T \mathbf{P} \mathbf{x}_0.$$

To show that such a controller is an optimal one, we will use a proof by contradiction. We assume that the hypothesis of the theorem holds but the controller \mathbf{u}^* that satisfies this hypothesis is not optimal; that is, there is a controller $\tilde{\mathbf{u}}$ such that

$$J(\tilde{\mathbf{u}}) < J(\mathbf{u}^*). \quad (5.64)$$

The hypothesis of the theorem implies that

$$\Delta V(\mathbf{x}(k))|_{\mathbf{u}=\tilde{\mathbf{u}}} + \mathbf{x}^T(k) \mathbf{Q}\mathbf{x}(k) + \tilde{\mathbf{u}}^T(k) \mathbf{R}\tilde{\mathbf{u}}(k) \geq 0. \quad (5.65)$$

We represent (5.65) as

$$\Delta V(\mathbf{x}(k))|_{\mathbf{u}=\tilde{\mathbf{u}}} \geq -\mathbf{x}^T(k) \mathbf{Q}\mathbf{x}(k) - \tilde{\mathbf{u}}^T(k) \mathbf{R}\tilde{\mathbf{u}}(k).$$

Summing the above from $k = 0$ to $k = \infty$ yields

$$V(\mathbf{x}(0)) \leq \sum_{k=0}^{\infty} (\mathbf{x}^T(k) \mathbf{Q}\mathbf{x}(k) + \tilde{\mathbf{u}}^T(k) \mathbf{R}\tilde{\mathbf{u}}(k));$$

that is,

$$J(\tilde{\mathbf{u}}) \geq J(\mathbf{u}^*),$$

which contradicts (5.64). The proof is complete.

Finding an optimal controller involves finding an appropriate quadratic Lyapunov function $V(\mathbf{x}) = \mathbf{x}^T \mathbf{P} \mathbf{x}$, which then is used to construct the optimal controller. We first find \mathbf{u}^* that minimizes the function

$$f = f(\mathbf{u}(k)) = \Delta V(\mathbf{x}(k)) + \mathbf{x}^T(k) \mathbf{Q}\mathbf{x}(k) + \mathbf{u}^T(k) \mathbf{R}\mathbf{u}(k). \quad (5.66)$$

To proceed, we perform preliminary manipulations on the function f . In particular, we use the

definition of ΔV and substitute into (5.66) the system equation (5.60) to obtain

$$\begin{aligned} f(\mathbf{u}(k)) &= \mathbf{x}^T(k+1)\mathbf{P}\mathbf{x}(k+1) - \mathbf{x}^T(k)\mathbf{P}\mathbf{x}(k) + \mathbf{x}^T(k)\mathbf{Q}\mathbf{x}(k) + \mathbf{u}^T(k)\mathbf{R}\mathbf{u}(k) \\ &= (\mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k))^T\mathbf{P}(\mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k)) - \mathbf{x}^T(k)\mathbf{P}\mathbf{x}(k) + \mathbf{x}^T(k)\mathbf{Q}\mathbf{x}(k) \\ &\quad + \mathbf{u}^T(k)\mathbf{R}\mathbf{u}(k). \end{aligned} \quad (5.67)$$

Applying to the above function the first-order necessary condition for a relative minimum gives

$$\begin{aligned} \frac{\partial f(\mathbf{u}(k))}{\partial \mathbf{u}(k)} &= 2(\mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k))^T\mathbf{P}\mathbf{B} + 2\mathbf{u}^T(k)\mathbf{R} \\ &= 2\mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{B} + 2\mathbf{u}^T(k)(\mathbf{R} + \mathbf{B}^T\mathbf{P}\mathbf{B}) \\ &= \mathbf{0}^T. \end{aligned}$$

The matrix $\mathbf{R} + \mathbf{B}^T\mathbf{P}\mathbf{B}$ is positive definite because \mathbf{R} is positive definite, and therefore the matrix $\mathbf{R} + \mathbf{B}^T\mathbf{P}\mathbf{B}$ is invertible. Hence,

$$\mathbf{u}^*(k) = -(\mathbf{R} + \mathbf{B}^T\mathbf{P}\mathbf{B})^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) = -\mathbf{K}\mathbf{x}(k), \quad (5.68)$$

where

$$\mathbf{K} = (\mathbf{R} + \mathbf{B}^T\mathbf{P}\mathbf{B})^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}.$$

Let

$$\mathbf{S} = \mathbf{R} + \mathbf{B}^T\mathbf{P}\mathbf{B}.$$

Then, we represent (5.68) as

$$\mathbf{u}^*(k) = -\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}\mathbf{x}(k). \quad (5.69)$$

We next check if \mathbf{u}^* satisfies the second-order sufficient condition for a relative minimizer of the function (5.66). We have

$$\begin{aligned} &\frac{\partial^2}{\partial \mathbf{u}^2(k)}(\mathbf{x}^T(k+1)\mathbf{P}\mathbf{x}(k+1) - \mathbf{x}^T(k)\mathbf{P}\mathbf{x}(k) + \mathbf{x}^T(k)\mathbf{Q}\mathbf{x}(k) + \mathbf{u}^T(k)\mathbf{R}\mathbf{u}(k)) \\ &= \frac{\partial}{\partial \mathbf{u}(k)}(2\mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{B} + 2\mathbf{u}^T(k)(\mathbf{R} + \mathbf{B}^T\mathbf{P}\mathbf{B})) \\ &= 2(\mathbf{R} + \mathbf{B}^T\mathbf{P}\mathbf{B}) \\ &> 0, \end{aligned} \quad (5.70)$$

that is, \mathbf{u}^* satisfies the second-order sufficient condition for a relative minimizer of the function (5.66).

The optimal controller (5.69) can be constructed if we have found an appropriate positive definite matrix \mathbf{P} . Our next task is to devise a method that would allow us to compute the desired matrix \mathbf{P} . For this, we first find the equation describing the closed-loop system driven by the optimal controller given by (5.69):

$$\mathbf{x}(k+1) = (\mathbf{A} - \mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A})\mathbf{x}(k). \quad (5.71)$$

Our controller satisfies the hypothesis of Theorem 5.4, that is,

$$\mathbf{x}^T(k+1)\mathbf{P}\mathbf{x}(k+1) - \mathbf{x}^T(k)\mathbf{P}\mathbf{x}(k) + \mathbf{x}^T(k)\mathbf{Q}\mathbf{x}(k) + \mathbf{u}^{*T}(k)\mathbf{R}\mathbf{u}^*(k) = 0.$$

Substituting into (5.71) the expression for \mathbf{u}^* given by (5.69), along with performing simple manipulations, gives

$$\begin{aligned}
& \mathbf{x}^T(k)(\mathbf{A} - \mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A})^T\mathbf{P}(\mathbf{A} - \mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A})\mathbf{x}(k) - \mathbf{x}^T(k)\mathbf{P}\mathbf{x}(k) \\
& \quad + \mathbf{x}^T(k)\mathbf{Q}\mathbf{x}(k) + \mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}\mathbf{R}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) \\
& = \mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) - \mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) \\
& \quad - \mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) + \mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) \\
& \quad - \mathbf{x}^T(k)\mathbf{P}\mathbf{x}(k) + \mathbf{x}^T(k)\mathbf{Q}\mathbf{x}(k) + \mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}\mathbf{R}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) \\
& = \mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) - \mathbf{x}^T(k)\mathbf{P}\mathbf{x}(k) + \mathbf{x}^T(k)\mathbf{Q}\mathbf{x}(k) \\
& \quad - 2\mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) \\
& \quad + \mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}(\mathbf{R} + \mathbf{B}^T\mathbf{P}\mathbf{B})\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) \\
& = \mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) - \mathbf{x}^T(k)\mathbf{P}\mathbf{x}(k) + \mathbf{x}^T(k)\mathbf{Q}\mathbf{x}(k) \\
& \quad - 2\mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) \\
& \quad + \mathbf{x}^T(k)\mathbf{A}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A}\mathbf{x}(k) \\
& = \mathbf{x}^T(k)(\mathbf{A}^T\mathbf{P}\mathbf{A} - \mathbf{P} + \mathbf{Q} - \mathbf{A}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A})\mathbf{x}(k) \\
& = 0.
\end{aligned} \tag{5.72}$$

The above equation should hold for any \mathbf{x} . For this to be true we have to have

$$\boxed{\mathbf{A}^T\mathbf{P}\mathbf{A} - \mathbf{P} + \mathbf{Q} - \mathbf{A}^T\mathbf{P}\mathbf{B}\mathbf{S}^{-1}\mathbf{B}^T\mathbf{P}\mathbf{A} = \mathbf{O}} \tag{5.73}$$

The above equation is called the *discrete algebraic Riccati equation*.

5.4 Dynamic Programming

The objective of this section is to introduce the reader to the method of dynamic programming invented by Richard E. Bellman (1920–1984) in 1957. Dynamic programming can be used to synthesize optimal controllers for nonlinear, time-varying dynamical systems. We first present dynamic programming for discrete-time systems and then discuss the continuous version of dynamic programming.

5.4.1 Discrete-Time Systems

Dynamic programming is based on the principle of optimality (PO), which was formulated by Bellman as follows:

An optimal policy has the property that whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision.

In other words, an optimal trajectory has the property that at an intermediate point, no matter how it was reached, the rest of the trajectory must coincide with an optimal trajectory as computed from this intermediate point as the initial point. We will justify the PO by contradiction. Suppose that we are given an optimal trajectory starting from initial state $\mathbf{x}_0 = \mathbf{x}(t_0)$ at time t_0 and

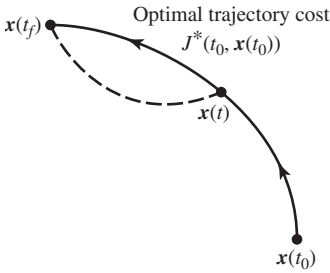


Figure 5.22 Illustration of the Principle of Optimality.

terminating at the final state $\mathbf{x}(t_f)$ at time t_f . Let $\mathbf{x}(t)$ be an intermediate state on the optimal trajectory and let $J^*(t, \mathbf{x}(t))$ be the optimal cost of traveling from $\mathbf{x}(t)$ to the final state $\mathbf{x}(t_f)$. Thus, $J^*(t_0, \mathbf{x}(t_0))$ is the optimal cost of traveling from $\mathbf{x}_0 = \mathbf{x}(t_0)$ to $\mathbf{x}(t_f)$. It then follows that the portion of the optimal trajectory from $\mathbf{x}(t)$ to $\mathbf{x}(t_f)$ is also optimal. For if some other path, indicated by the dashed curve in Figure 5.22, connecting $\mathbf{x}(t)$ and $\mathbf{x}(t_f)$ resulted in a smaller cost than the solid path, then this would provide a less expensive route from \mathbf{x}_0 to $\mathbf{x}(t_f)$, thus contradicting the optimality of the original trajectory.

We note that the PO seems to be already known to Johann Bernoulli, who stated in 1706 that “. . . all curves, which should give a maximum, preserve also in all their sections the laws of the same maximum.” Here, the reader should interpret maximum as optimum.

Bellman’s principle of optimality serves to “limit the number of potentially optimal control strategies that must be investigated” [182, p. 315]. The optimal control strategy is determined, using the PO, by working backward from the final stage.

We illustrate this “backwardness” with an example of a sequential decision-making process. We define a decision as the choice of alternative paths leaving a given node.

◆ Example 5.16

We use the PO to find the minimum cost path that starts at the node **a** and ends at any one of the nodes **k**, **l**, **m** in the routing network depicted in Figure 5.23. The travel costs are shown beside each path segment. The line segments can only be traversed from left to right. Applying the PO to this routing problem, we find the optimal path by performing a backward pass through the network. Let J_q^* be the minimum cost from a generic node **q** to any one of the three possible final nodes. Then, we have

$$J_g^* = 16, \quad J_h^* = 8, \quad J_i^* = 7, \quad J_j^* = 6.$$

The next stage yields

$$J_d^* = \min\{3 + J_g^*, 11 + J_h^*\} = \min\{19, 19\} = 19.$$

Thus, in this case, we have two optimal paths from node **d**. We can take either **dgk** or **dhk** paths. We then compute optimal paths from nodes **e** and **f**. We have

$$J_e^* = \min\{6 + J_h^*, 10 + J_i^*\} = \min\{6 + 8, 10 + 7\} = 14$$

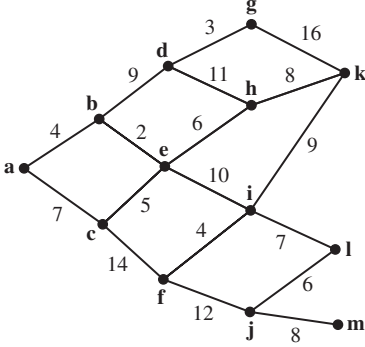


Figure 5.23 A routing network.

and

$$J_f^* = \min\{4 + J_i^*, 12 + J_j^*\} = \min\{4 + 7, 12 + 6\} = 11.$$

Thus, the optimal path emanating from node **e** is **ehk**, and the path from node **f** is **fil**. The next stage yields

$$J_b^* = \min\{9 + J_d^*, 2 + J_e^*\} = \min\{9 + 19, 2 + 14\} = 16,$$

and

$$J_c^* = \min\{5 + J_e^*, 14 + J_f^*\} = \min\{5 + 14, 14 + 11\} = 19.$$

Therefore, the optimal path starting from **b** is **behk**, and starting from **c** the optimal path is **cehk**. Finally, the initial, and hence final, stage yields

$$J_a^* = \min\{4 + J_b^*, 7 + J_c^*\} = \min\{4 + 16, 7 + 19\} = 20.$$

Hence, the optimal, minimum total cost route from **a** to any one of the three possible termination nodes **k**, **l**, or **m** is **abehk**. The total cost of this path is 20.

We now consider a class of discrete-time dynamical systems modeled by the difference equation

$$\mathbf{x}(k+1) = \mathbf{f}(k, \mathbf{x}(k), \mathbf{u}(k)), \quad (5.74)$$

with a given initial condition $\mathbf{x}(k_0) = \mathbf{x}_0$, and the associated performance index

$$J = J_0 = \Phi(N, \mathbf{x}(N)) + \sum_{k=0}^{N-1} F(k, \mathbf{x}(k), \mathbf{u}(k)).$$

Let $J_k^*(\mathbf{x}(k))$ denote the minimum cost of transferring the system from $\mathbf{x}(k)$ to the terminal point $\mathbf{x}(N)$. Because $\mathbf{x}(N)$ is the terminal point, we have $J_N^*(\mathbf{x}(N)) = \Phi(N, \mathbf{x}(N))$. Using the PO, we obtain the cost of transfer from $\mathbf{x}(N-1)$ to the terminal point:

$$\begin{aligned} J_{N-1}^*(\mathbf{x}(N-1)) &= \min_{\mathbf{u}(N-1)} \{F(N-1, \mathbf{x}(N-1), \mathbf{u}(N-1)) + J_N^*(\mathbf{x}(N))\} \\ &= \min_{\mathbf{u}(N-1)} \{F(N-1, \mathbf{x}(N-1), \mathbf{u}(N-1)) + \Phi(N, \mathbf{x}(N))\}. \end{aligned}$$

When applying the PO to solving optimal control problems for discrete-time dynamical systems, we usually perform two stage-by-stage passes through the time stages. We begin with a backward pass. For this, we use (5.74) to eliminate $\mathbf{x}(N)$. Next, we carry out the minimization to find $\mathbf{u}^*(N-1)$. We then repeat the process. A typical step of the backward pass is described as

$$\begin{aligned} J_k^*(\mathbf{x}(k)) &= \min_{\mathbf{u}(k)} \{F(k, \mathbf{x}(k), \mathbf{u}(k)) + J_{k+1}^*(\mathbf{x}(k+1))\} \\ &= \min_{\mathbf{u}(k)} \{F(k, \mathbf{x}(k), \mathbf{u}(k)) + J_{k+1}^*(f(k, \mathbf{x}(k), \mathbf{u}(k)))\}. \end{aligned}$$

The backward pass is completed when the initial time k_0 is reached. Now, because $\mathbf{x}_0 = \mathbf{x}(k_0)$ is known, we find $\mathbf{u}_0 = \mathbf{u}(k_0)$ in terms of this state. We can then proceed with the *forward pass*. We use \mathbf{x}_0 , \mathbf{u}_0 , and (5.74) to compute $\mathbf{x}(k_0+1)$. The state $\mathbf{x}(k_0+1)$ is then used to compute $\mathbf{u}(k_0+1)$ from $\mathbf{x}(k_0+2) = f(k_0+1, \mathbf{x}(k_0+1), \mathbf{u}(k_0+1))$, and so on.

◆ Example 5.17

Consider a scalar discrete-time dynamical system modeled by

$$x(k+1) = 2x(k) - 3u(k), \quad x(0) = 4,$$

and the performance index

$$J = J_0 = (x(2) - 10)^2 + \frac{1}{2} \sum_{k=0}^1 (x^2(k) + u^2(k)).$$

Note that the final state, in this example, is free. We use the PO to find optimal $u(0)$ and $u(1)$. There are no constraints on $u(k)$. We begin with the backward pass. We have

$$J^*(x(2)) = (x(2) - 10)^2.$$

Then,

$$\begin{aligned} J^*(x(1)) &= \min_{u(1)} \left\{ \frac{1}{2}(x^2(1) + u^2(1)) + J^*(x(2)) \right\} \\ &= \min_{u(1)} \left\{ \frac{1}{2}(x^2(1) + u^2(1)) + (2x(1) - 3u(1) - 10)^2 \right\}. \end{aligned}$$

Because there are no constraints on $u(1)$, we find optimal $u(1)$ as a function of $x(1)$ by applying the first-order necessary condition for unconstrained optimization to get

$$\frac{\partial}{\partial u(1)} \left\{ \frac{1}{2}(x^2(1) + u^2(1)) + (2x(1) - 3u(1) - 10)^2 \right\} = 0.$$

Hence,

$$u(1) = \frac{1}{19}(12x(1) - 60),$$

and therefore

$$J^*(x(1)) = \frac{1}{2} \left(x^2(1) + \left(\frac{12x(1) - 60}{19} \right)^2 \right) + \left(2x(1) - 3 \frac{12x(1) - 60}{19} - 10 \right)^2.$$

Next, we compute $u(0)$. For this observe that

$$J^*(x(0)) = \min_{u(0)} \left\{ \frac{1}{2} (x^2(0) + u^2(0)) + J^*(x(1)) \right\}.$$

Taking into account that

$$x(1) = 2x(0) - 3u(0), \quad x(0) = 4,$$

and the fact that there are no constraints on $u(0)$, we again use the first-order necessary condition for unconstrained optimization to find optimal $u(0)$. We compute

$$\frac{\partial}{\partial u(0)} \left\{ \frac{1}{2} (x^2(0) + u^2(0)) + J^*(x(1)) \right\} = 0$$

to get

$$\begin{aligned} \frac{\partial}{\partial u(0)} \left\{ 8 + \frac{1}{2} u^2(0) + \frac{1}{2} \left(\frac{12(8 - 3u(0)) - 60}{19} \right)^2 \right\} \\ + \frac{\partial}{\partial u(0)} \left\{ 2(8 - 3u(0)) - 3 \frac{12(8 - 3u(0)) - 60}{19} - 10 \right\}^2 = 0. \end{aligned}$$

Performing some manipulations yields $13.7895u(0) = 27.7895$. Hence, $u(0) = 2.0153$. This completes the backward pass. We are now ready for the forward pass. We use the system difference equation, $x(0)$, and $u(0)$ to compute $x(1) = 1.9541$. Using the above, we find

$$u(1) = \frac{12x(1) - 60}{19} = -1.9237.$$

Therefore,

$$x(2) = 2x(1) - 3u(1) = 9.6794,$$

which completes the forward pass.

5.4.2 Discrete Linear Quadratic Regulator Problem

In Example 5.17 above, we solved a simple discrete linear quadratic regulator problem where we found the open-loop optimal control strategy. We will now apply the PO to find the optimal control of a linear state-feedback form for the discrete linear quadratic regulator problem. This

problem is stated as follows. Given a discrete linear plant model,

$$\mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k), \quad 0 \leq k \leq N-1,$$

with a specified initial condition $\mathbf{x}(0) = \mathbf{x}_0$. We wish to calculate the optimal control sequence

$$\{\mathbf{u}^*(0), \mathbf{u}^*(1), \dots, \mathbf{u}^*(N-1)\},$$

that minimizes the quadratic performance index

$$J = J_0 = \frac{1}{2}\mathbf{x}^T(N)\mathbf{H}_N\mathbf{x}(N) + \frac{1}{2}\sum_{k=0}^{N-1}\{\mathbf{x}^T(k)\mathbf{Q}\mathbf{x}(k) + \mathbf{u}^T(k)\mathbf{R}\mathbf{u}(k)\},$$

where $\mathbf{H}_N = \mathbf{H}_N^T \geq 0$, $\mathbf{Q} = \mathbf{Q}^T \geq 0$, and $\mathbf{R} = \mathbf{R}^T > 0$. We assume that the components of the control vector are unconstrained. To solve this problem we use the PO. To begin, we note that

$$J^*(\mathbf{x}(N)) = \frac{1}{2}\mathbf{x}^T(N)\mathbf{H}_N\mathbf{x}(N).$$

The above is the penalty for being in a state $\mathbf{x}(N)$ at a time N .

We now decrement k to $N-1$ to get

$$J_{N-1}(\mathbf{x}(N-1)) = \frac{1}{2}\mathbf{x}^T(N)\mathbf{H}_N\mathbf{x}(N) + \frac{1}{2}\mathbf{x}^T(N-1)\mathbf{Q}\mathbf{x}(N-1) + \frac{1}{2}\mathbf{u}^T(N-1)\mathbf{R}\mathbf{u}(N-1).$$

We use the state equation to eliminate $\mathbf{x}(N)$ from J_{N-1} to obtain

$$\begin{aligned} J_{N-1} &= \frac{1}{2}(\mathbf{A}\mathbf{x}(N-1) + \mathbf{B}\mathbf{u}(N-1))^T \mathbf{H}_N (\mathbf{A}\mathbf{x}(N-1) + \mathbf{B}\mathbf{u}(N-1)) \\ &\quad + \frac{1}{2}\mathbf{x}^T(N-1)\mathbf{Q}\mathbf{x}(N-1) + \frac{1}{2}\mathbf{u}^T(N-1)\mathbf{R}\mathbf{u}(N-1). \end{aligned}$$

Because there are no constraints on control, we can apply the first-order necessary condition from static optimization to find $\mathbf{u}^*(N-1)$ as a function of $\mathbf{x}(N-1)$,

$$\begin{aligned} \frac{\partial J_{N-1}}{\partial \mathbf{u}(N-1)} &= (\mathbf{A}\mathbf{x}(N-1) + \mathbf{B}\mathbf{u}(N-1))^T \mathbf{H}_N \mathbf{B} + \mathbf{u}^T(N-1)\mathbf{R} \\ &= \mathbf{0}^T. \end{aligned}$$

Solving for the control satisfying the first-order necessary condition, we obtain

$$\mathbf{u}^*(N-1) = -(\mathbf{R} + \mathbf{B}^T \mathbf{H}_N \mathbf{B})^{-1} \mathbf{B}^T \mathbf{H}_N \mathbf{A} \mathbf{x}_{N-1}.$$

Let

$$\mathbf{K}_{N-1} = (\mathbf{R} + \mathbf{B}^T \mathbf{H}_N \mathbf{B})^{-1} \mathbf{B}^T \mathbf{H}_N \mathbf{A},$$

then

$$\mathbf{u}^*(N-1) = -\mathbf{K}_{N-1}\mathbf{x}(N-1).$$

Computing the optimal cost of transferring the system from $N-1$ to N yields

$$\begin{aligned} J_{N-1}^* &= \frac{1}{2}\mathbf{x}^T(N-1)(\mathbf{A} - \mathbf{B}\mathbf{K}_{N-1})^T \mathbf{H}_N (\mathbf{A} - \mathbf{B}\mathbf{K}_{N-1})\mathbf{x}(N-1) \\ &\quad + \frac{1}{2}\mathbf{x}^T(N-1)(\mathbf{K}_{N-1}^T \mathbf{R} \mathbf{K}_{N-1} + \mathbf{Q})\mathbf{x}(N-1). \end{aligned}$$

Let

$$\mathbf{H}_{N-1} = (\mathbf{A} - \mathbf{B}\mathbf{K}_{N-1})^T \mathbf{H}_N (\mathbf{A} - \mathbf{B}\mathbf{K}_{N-1}) + \mathbf{K}_{N-1}^T \mathbf{R} \mathbf{K}_{N-1} + \mathbf{Q},$$

then

$$J_{N-1}^* = \frac{1}{2} \mathbf{x}^T (N-1) \mathbf{H}_{N-1} \mathbf{x} (N-1).$$

Decrementing k to $N-2$ yields

$$J_{N-2} = \frac{1}{2} \mathbf{x}^T (N-1) \mathbf{H}_{N-1} \mathbf{x} (N-1) + \frac{1}{2} \mathbf{x}^T (N-2) \mathbf{Q} \mathbf{x} (N-2) + \frac{1}{2} \mathbf{u}^T (N-2) \mathbf{R} \mathbf{u} (N-2).$$

Note that J_{N-2} has the same form as J_{N-1} . Thus, we obtain analogous optimal feedback gain where N is being replaced with $N-1$. Continuing in this fashion, we get the following results for each $k = N-1, N-2, \dots, 0$:

$$\begin{aligned} \mathbf{K}_k &= (\mathbf{R} + \mathbf{B}^T \mathbf{H}_{k+1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{H}_{k+1} \mathbf{A}, \\ \mathbf{u}^*(k) &= -\mathbf{K}_k \mathbf{x}(k), \\ \mathbf{H}_k &= (\mathbf{A} - \mathbf{B}\mathbf{K}_k)^T \mathbf{H}_{k+1} (\mathbf{A} - \mathbf{B}\mathbf{K}_k) + \mathbf{K}_k^T \mathbf{R} \mathbf{K}_k + \mathbf{Q}, \end{aligned}$$

and

$$J_k^* = \frac{1}{2} \mathbf{x}^T (k) \mathbf{H}_k \mathbf{x}(k).$$

The above control scheme can be implemented by computing the sequence of gain matrices $\{\mathbf{K}_k\}$ off-line and stored. Then, we can implement the controller $\mathbf{u}^*(k) = -\mathbf{K}_k \mathbf{x}(k)$. Note that the time-varying nature of the feedback controller may make an implementation of the control strategy quite complicated, especially for large n and N . However, as we go to the limit as $N \rightarrow \infty$, we can simplify the solution to the problem. Specifically, let us consider a (reachable) system model

$$\mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k), \quad k \geq 0, \quad \mathbf{x}(0) = \mathbf{x}_0$$

and the performance index

$$J = \frac{1}{2} \sum_{k=0}^{\infty} \{ \mathbf{x}^T(k) \mathbf{Q} \mathbf{x}(k) + \mathbf{u}^T(k) \mathbf{R} \mathbf{u}(k) \}$$

obtained from the previously considered performance index by letting $N \rightarrow \infty$ and setting $\mathbf{H}_N = \mathbf{0}$. Then, the optimal controller takes the form

$$\mathbf{u}^*(k) = -\mathbf{K}_{\infty} \mathbf{x}(k), \quad k \geq 0,$$

where

$$\mathbf{K}_{\infty} = (\mathbf{R} + \mathbf{B}^T \mathbf{H}_{\infty} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{H}_{\infty} \mathbf{A} \quad (5.75)$$

and \mathbf{H}_{∞} is defined as the limit obtained by continuing the recursion

$$\mathbf{H}_k = (\mathbf{A} - \mathbf{B}\mathbf{K}_k)^T \mathbf{H}_{k+1} (\mathbf{A} - \mathbf{B}\mathbf{K}_k) + \mathbf{K}_k^T \mathbf{R} \mathbf{K}_k + \mathbf{Q} \quad (5.76)$$

back to $k = -\infty$, that is,

$$\lim_{k \rightarrow -\infty} \mathbf{H}_k = \mathbf{H}_\infty.$$

Alternatively, letting $k \rightarrow -\infty$ in (5.76), we obtain

$$\mathbf{H}_\infty = (\mathbf{A} - \mathbf{B}\mathbf{K}_\infty)^T \mathbf{H}_\infty (\mathbf{A} - \mathbf{B}\mathbf{K}_\infty) + \mathbf{K}_\infty^T \mathbf{R} \mathbf{K}_\infty + \mathbf{Q}, \quad (5.77)$$

where \mathbf{K}_∞ is given by (5.75). Substituting (5.75) into (5.77) yields

$$\begin{aligned} \mathbf{H}_\infty = & (\mathbf{A} - \mathbf{B}(\mathbf{R} + \mathbf{B}^T \mathbf{H}_\infty \mathbf{B})^{-1} \mathbf{B}^T \mathbf{H}_\infty \mathbf{A})^T \mathbf{H}_\infty (\mathbf{A} - \mathbf{B}(\mathbf{R} + \mathbf{B}^T \mathbf{H}_\infty \mathbf{B})^{-1} \mathbf{B}^T \mathbf{H}_\infty \mathbf{A}) \\ & + ((\mathbf{R} + \mathbf{B}^T \mathbf{H}_\infty \mathbf{B})^{-1} \mathbf{B}^T \mathbf{H}_\infty \mathbf{A})^T \mathbf{R} ((\mathbf{R} + \mathbf{B}^T \mathbf{H}_\infty \mathbf{B})^{-1} \mathbf{B}^T \mathbf{H}_\infty \mathbf{A}) + \mathbf{Q}. \end{aligned} \quad (5.78)$$

Thus, we can compute \mathbf{H}_∞ by solving the above algebraic equation. We note that the above formidable looking equation is just the discrete algebraic Riccati equation. Indeed, setting $\mathbf{P} = \mathbf{H}_\infty$, using the notation

$$\mathbf{S} = \mathbf{R} + \mathbf{B}^T \mathbf{H}_\infty \mathbf{B} = \mathbf{R} + \mathbf{B}^T \mathbf{P} \mathbf{B},$$

and performing simple manipulations, we obtain (5.73).

We next discuss a continuous form of dynamic programming. We begin our discussion by analyzing the minimum time regulator problem. Then, we focus our attention on a more general optimal control problem.

5.4.3 Continuous Minimum Time Regulator Problem

In this subsection we consider a control system modeled by

$$\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x}, \mathbf{u})$$

with a given initial state $\mathbf{x}(t_0) = \mathbf{x}_0$. We restrict the controller $\mathbf{u}(t)$ to be restrained by the compact set $U \subset \mathbb{R}^m$. For example, we may restrict the components u_i , $i = 1, 2, \dots, m$, of the control vector \mathbf{u} to satisfy the amplitude constraints

$$|u_i| \leq k_i,$$

where k_i are fixed real numbers. A piecewise continuous control vector \mathbf{u} satisfying the above constraints is called admissible. The set of admissible controllers is denoted U . Our goal is to determine an admissible control law \mathbf{u}^* that steers the system's trajectory from a given initial state $\mathbf{x}(t_0) = \mathbf{x}_0$ to a specified final state \mathbf{x}_f in the shortest possible time. This is equivalent to minimizing the performance index

$$J = \int_{t_0}^{t_f} dt = t_f - t_0.$$

We call such a controller *time optimal*. We assume that for any $\mathbf{x} \neq \mathbf{x}_f$ there exists admissible, time optimal control transferring the system from \mathbf{x} to \mathbf{x}_f in minimum time. We denote by $J^*(t, \mathbf{x}) = J^*(t, \mathbf{x}(t))$ the time of the optimal transfer from $\mathbf{x}(t)$ to the specified \mathbf{x}_f . Observe that $J^*(t, \mathbf{x}) = J^*(t, x_1, x_2, \dots, x_n)$, that is,

$$J^*: \mathbb{R}^{n+1} \rightarrow \mathbb{R},$$

which means that J^* is a real-valued function of $n + 1$ variables. We assume that the function J^* is continuous and has continuous partial derivatives everywhere except at the point \mathbf{x}_f .

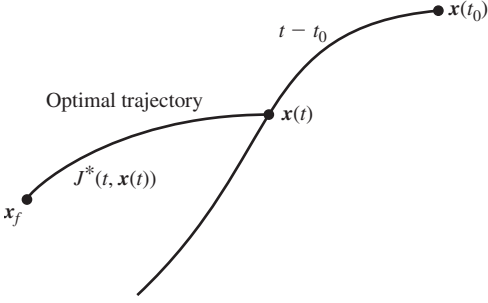


Figure 5.24 Illustration of the derivation of the condition that a controller must satisfy when transferring a dynamical system from an initial state to the final state in an optimal manner.

Now, let $\mathbf{x}_0 = \mathbf{x}(t_0) \neq \mathbf{x}_f$ be an arbitrary initial state in \mathbb{R}^n . Assume that we applied some admissible $\mathbf{u}_0 \in U$ on the time interval $[t_0, t]$. As a result of this control action, the system is transferred from $\mathbf{x}(t_0)$ to $\mathbf{x}(t)$. The time spent on the motion from $\mathbf{x}(t_0)$ to $\mathbf{x}(t)$ is $t - t_0$. We then apply an admissible optimal control that transfers the system in minimum time $J^*(t, \mathbf{x}(t))$ from $\mathbf{x}(t)$ to \mathbf{x}_f . The above considerations are illustrated in Figure 5.24. Denoting the optimal time of transferring the system from $\mathbf{x}(t_0)$ to \mathbf{x}_f by $J^*(t_0, \mathbf{x}(t_0))$, we obtain

$$(t - t_0) + J^*(t, \mathbf{x}(t)) \geq J^*(t_0, \mathbf{x}(t_0)).$$

Rearranging the terms of the above equation and dividing by $(t - t_0)$ yields

$$1 + \frac{J^*(t, \mathbf{x}(t)) - J^*(t_0, \mathbf{x}(t_0))}{t - t_0} \geq 0.$$

Letting $t \rightarrow t_0$ and taking into account that $\dot{x}_i = f_i(t, \mathbf{x}, \mathbf{u})$, $i = 1, 2, \dots, n$, we obtain

$$\begin{aligned} 1 + \frac{\partial J^*(t, \mathbf{x})}{\partial t} + \sum_{i=1}^n \frac{\partial J^*(t, \mathbf{x})}{\partial x_i} \dot{x}_i &= 1 + \frac{\partial J^*(t, \mathbf{x})}{\partial t} + \sum_{i=1}^n \frac{\partial J^*(t, \mathbf{x})}{\partial x_i} f_i \\ &= 1 + \frac{\partial J^*}{\partial t} + \frac{\partial J^*}{\partial \mathbf{x}} \mathbf{f} \\ &\geq 0. \end{aligned}$$

We conclude that for any $\mathbf{x}_0 \neq \mathbf{x}_f$ and for any $\mathbf{u} \in U$,

$$1 + \frac{\partial J^*}{\partial t} + \frac{\partial J^*}{\partial \mathbf{x}} \mathbf{f} \geq 0.$$

We next define the Hamiltonian function for the time optimal control problem to be

$$H(t, \mathbf{x}, \mathbf{u}, \mathbf{p}) = 1 + \mathbf{p}^T \mathbf{f}(t, \mathbf{x}, \mathbf{u}),$$

where

$$\mathbf{p} = \left[\frac{\partial J^*}{\partial x_1} \quad \frac{\partial J^*}{\partial x_2} \quad \dots \quad \frac{\partial J^*}{\partial x_n} \right]^T.$$

Thus, we can say that for any $\mathbf{x}_0 \neq \mathbf{x}_f$ and for any $\mathbf{u} \in U$,

$$\frac{\partial J^*(t, \mathbf{x})}{\partial t} + H(t, \mathbf{x}, \mathbf{u}) \geq 0.$$

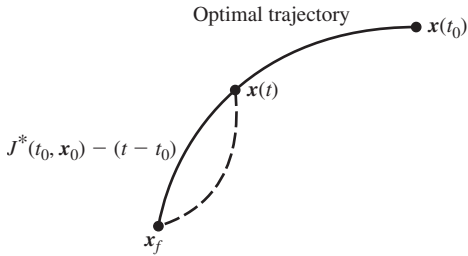


Figure 5.25 Illustration of the derivation of the optimality condition of Theorem 5.5.

We will now show that

$$\frac{\partial J^*(t, \mathbf{x})}{\partial t} + H(t, \mathbf{x}, \mathbf{u}) = 0$$

for any optimal transfer, that is,

$$\frac{\partial J^*(t, \mathbf{x})}{\partial t} + \min_{\mathbf{u} \in U} H(t, \mathbf{x}, \mathbf{u}) = 0 \quad \text{for any } \mathbf{x} \neq \mathbf{x}_f.$$

In the following, we use the PO. Suppose that $\mathbf{u}(t) \in U$ is an optimal controller that transfers the system from the initial state $\mathbf{x}(t_0) \neq \mathbf{x}_f$ to the given final state \mathbf{x}_f in a time t_f . Thus, we have $J^*(t_0, \mathbf{x}(t_0)) = t_f - t_0$. The transfer from $\mathbf{x}(t)$ to \mathbf{x}_f requires $(t_f - t)$ time units. See Figure 5.25 for an illustration of the arguments. Hence,

$$t_f - t = J^*(t_0, \mathbf{x}(t_0)) - (t - t_0),$$

and it is not possible to go from $\mathbf{x}(t)$ to \mathbf{x}_f in a time less than $J^*(t_0, \mathbf{x}(t_0)) - (t - t_0)$. For if such a more rapid motion were to exist, indicated by the dashed curve in Figure 5.25, then by moving from \mathbf{x}_0 to $\mathbf{x}(t)$ during the time $t - t_0$ and then moving from $\mathbf{x}(t)$ to \mathbf{x}_f during a time less than $J^*(t_0, \mathbf{x}(t_0)) - (t - t_0)$, we would accomplish the transfer from \mathbf{x}_0 to \mathbf{x}_f during a time less than $J^*(t_0, \mathbf{x}(t_0))$, which is impossible. Therefore, we can write

$$\begin{aligned} t_f - t &= J^*(t, \mathbf{x}(t)) \\ &= J^*(t_0, \mathbf{x}(t_0)) - (t - t_0). \end{aligned}$$

Rearranging terms, dividing by $t - t_0$, and letting $t \rightarrow t_0$ yields

$$\frac{\partial J^*(t, \mathbf{x})}{\partial t} + \min_{\mathbf{u} \in U} H(t, \mathbf{x}, \mathbf{u}) = 0 \quad \text{for any } \mathbf{x} \neq \mathbf{x}_f.$$

We summarize our considerations in the following theorem.

Theorem 5.5 Let $\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x}, \mathbf{u})$ be a control system model, let U be a set of admissible controllers, and let \mathbf{x}_f be a prescribed final state. If we suppose that

1. for any $\mathbf{x} \neq \mathbf{x}_f$ there exists an admissible optimal control transferring \mathbf{x} to \mathbf{x}_f , and
2. the function $J^*(t, \mathbf{x}(t))$, which is the time of optimal transfer of the system from \mathbf{x} , at time t , to \mathbf{x}_f , is continuously differentiable except at the point \mathbf{x}_f ,

then

$$\begin{aligned} \frac{\partial J^*(t, \mathbf{x})}{\partial t} + H(t, \mathbf{x}, \mathbf{u}) &\geq 0 && \text{for any } \mathbf{u} \in U \text{ and } \mathbf{x} \neq \mathbf{x}_f, \\ \frac{\partial J^*(t, \mathbf{x})}{\partial t} + \min_{\mathbf{u} \in U} H(t, \mathbf{x}, \mathbf{u}) &= 0 && \text{for any optimal transfer where } \mathbf{x} \neq \mathbf{x}_f. \end{aligned}$$

The above theorem is a special case of a more general result known as the continuous form of the dynamic programming. We now use the PO to derive a continuous form of dynamic programming for a class of problems with more general performance indices.

5.4.4 The Hamilton–Jacobi–Bellman Equation

We consider a control process model of the form

$$\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x}, \mathbf{u}), \quad \mathbf{x}(t_0) = \mathbf{x}_0, \quad (5.79)$$

and we consider the performance index to minimize

$$J(t_0, \mathbf{x}(t_0), \mathbf{u}) = \Phi(t_f, \mathbf{x}(t_f)) + \int_{t_0}^{t_f} F(\tau, \mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau. \quad (5.80)$$

Define

$$J(t, \mathbf{x}(t), \mathbf{u}(\tau)) = \Phi(t_f, \mathbf{x}(t_f)) + \int_t^{t_f} F(\tau, \mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau$$

with $t \leq \tau \leq t_f$, and let

$$J^*(t, \mathbf{x}(t)) = \min_{\mathbf{u}} J(t, \mathbf{x}(t), \mathbf{u}(\tau)).$$

We subdivide the interval $[t, t_f]$ as

$$[t, t_f] = [t, t + \Delta t] \cup [t + \Delta t, t_f].$$

Then, we can write

$$J^*(t, \mathbf{x}(t), \mathbf{u}(\tau)) = \min_{\mathbf{u}} \left\{ \int_t^{t+\Delta t} F d\tau + \int_{t+\Delta t}^{t_f} F d\tau + \Phi(t_f, \mathbf{x}(t_f)) \right\}.$$

However, by the PO the trajectory on the interval $[t + \Delta t, t_f]$ must be optimal. Hence,

$$J^*(t, \mathbf{x}(t)) = \min_{\mathbf{u}} \left\{ \int_t^{t+\Delta t} F d\tau + J^*(t + \Delta t, \mathbf{x}(t + \Delta t)) \right\}.$$

We next expand $J^*(t + \Delta t, \mathbf{x}(t + \Delta t))$ in a Taylor series about the point $(t, \mathbf{x}(t))$ to get

$$J^*(t, \mathbf{x}(t)) = \min_{\mathbf{u}} \left\{ \int_t^{t+\Delta t} F d\tau + J^* + \frac{\partial J^*}{\partial t} \Delta t + \frac{\partial J^*}{\partial \mathbf{x}} (\mathbf{x}(t + \Delta t) - \mathbf{x}(t)) + \text{H.O.T.} \right\},$$

where H.O.T. stands for higher-order terms. Because $J^*(t, \mathbf{x}(t))$ is independent of \mathbf{u} , we cancel the terms $J^*(t, \mathbf{x}(t))$ out and use that fact that $\mathbf{x}(t + \Delta t) - \mathbf{x}(t) \approx \dot{\mathbf{x}} \Delta t$ to obtain

$$0 = \min_{\mathbf{u}} \left\{ \int_t^{t+\Delta t} F d\tau + \frac{\partial J^*}{\partial t} \Delta t + \frac{\partial J^*}{\partial \mathbf{x}} \dot{\mathbf{x}} \Delta t + \text{H.O.T.} \right\}.$$

By assumption, Δt is small; hence we can write

$$0 = \min_{\mathbf{u}} \left\{ F \Delta t + \frac{\partial J^*}{\partial t} \Delta t + \frac{\partial J^*}{\partial \mathbf{x}} \mathbf{f} \Delta t + \text{H.O.T.} \right\}.$$

Dividing the above equation by Δt and letting $\Delta t \rightarrow 0$ yields

$$0 = \frac{\partial J^*}{\partial t} + \min_{\mathbf{u}} \left\{ F + \frac{\partial J^*}{\partial \mathbf{x}} \mathbf{f} \right\}. \quad (5.81)$$

Let $H = F + \frac{\partial J^*}{\partial \mathbf{x}} \mathbf{f}$ denote the Hamiltonian function. Then, we write the partial differential equation (5.81) as

$$0 = \frac{\partial J^*}{\partial t} + \min_{\mathbf{u}} H \quad (5.82)$$

subject to the boundary condition

$$J^*(t_f, \mathbf{x}(t_f)) = \Phi(t_f, \mathbf{x}(t_f)).$$

The partial differential equation (5.82) for the optimal cost $J^*(t, \mathbf{x}(t))$ is called the *Hamilton–Jacobi–Bellman* (HJB) equation. It provides the solution to the optimal control problem for general nonlinear dynamical systems. However, analytical solution to the HJB equation is difficult to obtain in most cases. We now illustrate, with three examples, how the HJB equation can be used to construct optimal controllers.

◆ Example 5.18

We consider a general class of dynamical systems modeled by the scalar differential equation

$$\dot{x} = ax + bu, \quad x(t_0) = x_0,$$

with the associated cost index

$$J(t_0) = \frac{1}{2} f x^2(t_f) + \frac{1}{2} \int_{t_0}^{t_f} (q x^2(t) + r u^2(t)) dt,$$

where the initial time t_0 equals zero, the final time $t_f < \infty$ is fixed, and the final state $x(t_f)$ is free. Our goal is construct the control u^* that minimizes $J(0)$.

We first form the Hamiltonian function for the problem:

$$H\left(t, x, u, \frac{\partial J^*}{\partial x}\right) = \frac{1}{2} q x^2 + \frac{1}{2} r u^2 + \frac{\partial J^*}{\partial x} (ax + bu).$$

Because there are no constraints on the control u , we can determine the form of the optimal control by applying the first-order necessary condition for static optimization. We compute

$$\frac{\partial H}{\partial u} = ru + b \frac{\partial J^*}{\partial x} = 0.$$

The optimal control, therefore, has the form

$$u^* = -\frac{1}{r} b \frac{\partial J^*}{\partial x}.$$

This control indeed minimizes the Hamiltonian function H because

$$\frac{\partial^2 H}{\partial u^2} = r > 0;$$

that is, the second-order sufficient condition for static minimization is satisfied. We thus reduced the problem of constructing the optimal control to that of finding $\partial J^* / \partial x$. Substituting u^* into the HJB equation yields

$$\begin{aligned} \frac{\partial J^*}{\partial t} + \frac{1}{2} q x^2 + \frac{b^2}{2r} \left(\frac{\partial J^*}{\partial x} \right)^2 + a x \frac{\partial J^*}{\partial x} - \frac{b^2}{r} \left(\frac{\partial J^*}{\partial x} \right)^2 \\ = \frac{\partial J^*}{\partial t} + \frac{1}{2} q x^2 - \frac{b^2}{2r} \left(\frac{\partial J^*}{\partial x} \right)^2 + a x \frac{\partial J^*}{\partial x} \\ = 0. \end{aligned}$$

Let us now assume that J^* is a quadratic function of the state for all $t \leq t_f$, that is,

$$J^* = \frac{1}{2} p(t) x^2,$$

where $p(t)$ is a scalar function of time to be determined. Partial derivatives of J^* are

$$\frac{\partial J^*}{\partial t} = \frac{1}{2} \dot{p}(t) x^2 \quad \text{and} \quad \frac{\partial J^*}{\partial x} = p(t) x.$$

Substituting the above into the HJB equation and performing some manipulations gives

$$\left(\frac{1}{2} \dot{p}(t) + \frac{1}{2} q + a p(t) - \frac{b^2 p^2(t)}{2r} \right) x^2 = 0.$$

The above holds for all $x = x(t)$ if and only if

$$\frac{1}{2} \dot{p}(t) + \frac{1}{2} q + a p(t) - \frac{b^2 p^2(t)}{2r} = 0,$$

or, equivalently,

$$\dot{p}(t) + q + 2ap(t) - \frac{b^2 p^2(t)}{r} = 0 \quad (5.83)$$

subject to the boundary condition $p(t_f) = f$ that we obtain by comparing $J^* = \frac{1}{2} p(t) x^2$ and $J(t_f) = \frac{1}{2} f x^2(t_f)$. To solve the above equation we assume that

$$p(t) = \rho \frac{\dot{w}(t)}{w(t)},$$

where the parameter ρ is to be determined. Note that

$$\dot{p}(t) = \rho \frac{\ddot{w}(t) w(t) - \dot{w}^2(t)}{w^2(t)}.$$

Substituting the above into (5.83), we obtain

$$\begin{aligned} \dot{p} + q + 2ap - \frac{b^2 p^2}{r} &= \rho \frac{\ddot{w} w - \dot{w}^2}{w^2} + q \frac{w^2}{w^2} + 2a\rho \frac{\dot{w}}{w} - \frac{b^2 \rho^2 \dot{w}^2 / r}{w^2} \\ &= \frac{\rho \ddot{w} w - \rho \dot{w}^2 + q w^2 + 2a\rho \dot{w} w - b^2 \rho^2 \dot{w}^2 / r}{w^2} \\ &= 0. \end{aligned}$$

We now choose ρ so that the terms nonlinear in \dot{w} cancel out. For this it is enough that

$$-\rho \dot{w}^2 - b^2 \rho^2 \dot{w}^2 / r = 0.$$

Hence, if we set

$$\rho = -\frac{r}{b^2},$$

then the equation to be solved takes the form

$$\frac{\left(-\frac{r}{b^2} \ddot{w} - \frac{2ar}{b^2} \dot{w} + qw \right) w}{w^2} = 0.$$

The above will be satisfied if

$$-\frac{r}{b^2} \ddot{w} - \frac{2ar}{b^2} \dot{w} + qw = 0,$$

or, equivalently,

$$\ddot{w} + 2a\dot{w} - \frac{qb^2}{r} w = 0,$$

which is just a second-order linear differential equation. Its characteristic equation has two real roots:

$$s_{1,2} = -a \pm \sqrt{a^2 + \frac{qb^2}{r}}.$$

Hence

$$w(t) = C_1 e^{s_1 t} + C_2 e^{s_2 t},$$

and

$$p(t) = \rho \frac{\dot{w}(t)}{w(t)} = \rho \frac{s_1 C_1 e^{s_1 t} + s_2 C_2 e^{s_2 t}}{C_1 e^{s_1 t} + C_2 e^{s_2 t}}.$$

Using the boundary value $p(t_f) = f$, we express C_1 in terms of C_2 to obtain

$$C_1 = \frac{\rho s_2 C_2 e^{s_1 t_f} - f C_2 e^{s_2 t_f}}{f e^{s_1 t_f} - \rho s_1 e^{s_1 t_f}}.$$

We substitute the above expression for C_1 into $p(t)$ and cancel out C_2 . Then, we use the obtained result to construct the optimal state-feedback controller:

$$u^* = -\frac{1}{r} b \frac{\partial J^*}{\partial x} = -\frac{1}{r} b p(t) x$$

◆ Example 5.19

Consider the schematic of an armature controlled DC motor, given in Figure 5.26, where the system parameters are $R_a = 2 \Omega$, $L_a \approx 0$ H, $K_b = 2$ V/rad/sec, and $K_i = 2$ Nm/A, and the equivalent moment of inertia referred to the motor shaft is $I_{eq} = 1 \text{ kg} \cdot \text{m}^2$. The friction is assumed to be negligible.

1. Construct the first-order differential equation modeling the DC motor.
2. Use the Hamilton–Jacobi–Bellman equation to find the optimal state-feedback

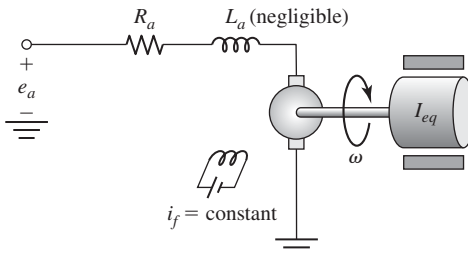


Figure 5.26 Schematic of an armature-controlled DC motor system of Example 5.19.

controller, $e_a = k(t)\omega$, which minimizes the performance index

$$J = \frac{1}{2}\omega^2(10) + \int_0^{10} R_a i_a^2(t) dt,$$

where i_a is the armature current. The final state is free. There are no constraints on e_a . Assume $J^* = \frac{1}{2}p(t)\omega^2$.

Applying Kirchhoff's voltage law to the armature circuit yields

$$R_a i_a + K_b \omega = e_a.$$

The torque developed by the motor, T_m , is

$$I_{eq}\dot{\omega} = T_m = K_i i_a.$$

Substituting i_a from the first equation into the second and dividing both sides by I_{eq} yields

$$\dot{\omega} = -\frac{K_i K_b}{R_a I_{eq}}\omega + \frac{K_i}{I_{eq} R_a}e_a.$$

Substituting into the above the parameter values, we obtain

$$\dot{\omega} = -2\omega + e_a. \quad (5.84)$$

We now represent the performance index, J , in terms of ω and e_a . Applying Ohm's law to the armature circuit, we get

$$e_a - K_b \omega = R_a i_a.$$

Hence,

$$J = \frac{1}{2}\omega^2(10) + \int_0^{10} \frac{(e_a - K_b \omega)^2}{R_a} dt.$$

The Hamiltonian function is

$$\begin{aligned} H &= \frac{(e_a - K_b \omega)^2}{R_a} + \frac{\partial J^*}{\partial \omega} \left(-\frac{K_i K_b}{R_a I_{eq}}\omega + \frac{K_i}{I_{eq} R_a}e_a \right) \\ &= \frac{(e_a - 2\omega)^2}{2} + \frac{\partial J^*}{\partial \omega}(-2\omega + e_a). \end{aligned}$$

Because there are no constraints on e_a , we can find the optimal control by solving the equation

$$0 = \frac{\partial H}{\partial e_a} = \frac{2(e_a - K_b \omega)}{R_a} + \frac{K_i}{I_{eq} R_a} \frac{\partial J^*}{\partial \omega}.$$

Hence,

$$e_a = 2\omega - \frac{\partial J^*}{\partial \omega}. \quad (5.85)$$

Substituting (5.85) into the HJB equation yields

$$0 = \frac{\partial J^*}{\partial t} + \frac{1}{2} \left(-\frac{\partial J^*}{\partial \omega} \right)^2 + \frac{\partial J^*}{\partial \omega} \left(-\frac{\partial J^*}{\partial \omega} \right) = \frac{\partial J^*}{\partial t} - \frac{1}{2} \left(\frac{\partial J^*}{\partial \omega} \right)^2.$$

By assumption $J^* = \frac{1}{2} p(t) \omega^2$. Hence,

$$\frac{\partial J^*}{\partial t} = \frac{1}{2} \dot{p} \omega^2 \quad \text{and} \quad \frac{\partial J^*}{\partial \omega} = p \omega.$$

Substituting the above into the HJB equation, we get

$$\frac{1}{2} (\dot{p} - p^2) \omega^2 = 0.$$

We next solve the nonlinear differential equation

$$\dot{p} - p^2 = 0. \tag{5.86}$$

We assume that $p = \rho \frac{\dot{w}}{w}$. Therefore,

$$\dot{p} = \rho \frac{\ddot{w}w - \dot{w}^2}{w^2}.$$

We substitute the expressions for p and \dot{p} into (5.86) to get

$$\frac{\rho \ddot{w}w - \rho \dot{w}^2 - \rho^2 \dot{w}^2}{w^2} = 0.$$

Selecting $\rho = -1$ eliminates the nonlinear terms in the numerator, and we are left with the equation

$$\ddot{w} = 0,$$

whose solution is

$$w = w(t) = C_1 t + C_2,$$

where C_1 and C_2 are integration constants. Hence,

$$p = -\frac{\dot{w}}{w} = -\frac{C_1}{C_1 t + C_2}.$$

Because

$$J^*(10) = \frac{1}{2} \omega^2(10) = \frac{1}{2} p(10) \omega^2(10),$$

we conclude that $p(10) = 1$. We use this information to eliminate one of the integration constants. We obtain

$$p(10) = 1 = -\frac{C_1}{10C_1 + C_2}.$$

Hence,

$$C_2 = -11C_1,$$

and

$$p = p(t) = -\frac{C_1}{C_1 t + C_2} = -\frac{C_1}{C_1 t - 11C_1} = -\frac{1}{t - 11}.$$

Therefore,

$$e_a^* = 2\omega - \frac{\partial J^*}{\partial \omega} = 2\omega - p\omega = \left(2 + \frac{1}{t - 11}\right)\omega$$

In the above examples, the final time t_f was fixed and $t_f < \infty$. In the following example, we consider the linear quadratic regulator design for $t_f = \infty$ using the HJB equation.

◆ Example 5.20

For the dynamical system model

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u$$

and the performance index

$$J = \frac{1}{2} \int_0^\infty (x_1^2 + x_2^2 + u^2) dt,$$

we will use the HJB equation to find the optimal state-feedback controller. We first form the Hamiltonian function

$$H = \frac{1}{2}(x_1^2 + x_2^2 + u^2) + \frac{\partial J^*}{\partial x_1} x_2 + \frac{\partial J^*}{\partial x_2} (-x_2 + u).$$

We assume that J^* has the form

$$J^* = \frac{1}{2}(\alpha x_1^2 + 2\beta x_1 x_2 + \gamma x_2^2).$$

Hence,

$$\frac{\partial J^*}{\partial t} = 0,$$

and the HJB equation takes the form

$$\min_u H = 0,$$

with the boundary condition

$$J^*(\mathbf{x}(t_f)) = J^*(\mathbf{x}(\infty)) = 0.$$

Because there are no constraints on u , we apply the first-order necessary condition for unconstrained optimization to find u minimizing H :

$$\frac{\partial H}{\partial u} = 0.$$

We denote by u^* the solution to the above equation. Differentiating, we obtain

$$u^* + \frac{\partial J^*}{\partial x_2} = 0.$$

Hence,

$$u^* = -\frac{\partial J^*}{\partial x_2}.$$

We substitute u^* into the HJB equation to get

$$\frac{1}{2} \left(x_1^2 + x_2^2 + \left(\frac{\partial J^*}{\partial x_2} \right)^2 \right) + \frac{\partial J^*}{\partial x_1} x_2 - \frac{\partial J^*}{\partial x_2} x_2 - \left(\frac{\partial J^*}{\partial x_2} \right)^2 = 0.$$

Note that

$$\frac{\partial J^*}{\partial x_1} = \alpha x_1 + \beta x_2 \quad \text{and} \quad \frac{\partial J^*}{\partial x_2} = \beta x_1 + \gamma x_2.$$

Substituting the above into the HJB equation yields

$$\frac{1}{2} (x_1^2 + x_2^2 + (\beta x_1 + \gamma x_2)^2) + (\alpha x_1 + \beta x_2) x_2 - (\beta x_1 + \gamma x_2) x_2 - (\beta x_1 + \gamma x_2)^2 = 0.$$

We rearrange the above equation to obtain

$$\frac{1}{2} x_1^2 (1 - \beta^2) + x_1 x_2 (\alpha - \beta - \beta \gamma) + \frac{1}{2} x_2^2 (1 - 2\gamma - \gamma^2 + 2\beta) = 0.$$

For this equation to hold, we have to have

$$\begin{aligned} 1 - \beta^2 &= 0, \\ \alpha - \beta - \beta \gamma &= 0, \\ 1 - 2\gamma - \gamma^2 + 2\beta &= 0. \end{aligned}$$

There are three solutions to the above system of equations:

1. $\alpha = 2, \beta = \gamma = 1.$
2. $\alpha = -2, \beta = 1, \gamma = -3.$
3. $\alpha = 0, \beta = \gamma = -1.$

Thus, we have three candidate solutions to the problem:

$$\begin{aligned} u^1 &= -x_1 - x_2, \\ u^2 &= -x_1 + 3x_2, \\ u^3 &= x_1 + x_2. \end{aligned}$$

However, only u^1 stabilizes the plant. Hence, the optimal control law for the problem is

$$u^* = -x_1 - x_2.$$

Note that only for this controller the boundary condition, $J^*(\mathbf{x}(t_f)) = J^*(\mathbf{x}(\infty)) = 0$, is satisfied.

5.5 Pontryagin's Minimum Principle

5.5.1 Optimal Control with Constraints on Inputs

In this section we provide an informal derivation of the minimum principle of Pontryagin, which is a generalization of the Euler–Lagrange equations that also includes problems with constraints on the control inputs. Only a special case of the minimum principle will be stated, but this special case covers a large class of control problems. We consider the problem of minimizing the performance index given by

$$J = \Phi(t_f, \mathbf{x}(t_f)) + \int_{t_0}^{t_f} F(\mathbf{x}(t), \mathbf{u}(t)) dt \quad (5.87)$$

subject to

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, \mathbf{u}), \quad \mathbf{x}(t_0) = \mathbf{x}_0, \quad \text{and} \quad \mathbf{x}(t_f) = \mathbf{x}_f. \quad (5.88)$$

We consider two cases: fixed final state and free final state.

It follows from our discussion of the HJB equation that the optimal controller, \mathbf{u}^* , minimizes the Hamiltonian function, that is,

$$\mathbf{u}^* = \arg \min_{\mathbf{u} \in U} H(\mathbf{x}, \mathbf{u}, \mathbf{p}),$$

where $H(\mathbf{x}, \mathbf{u}, \mathbf{p}) = H = F + \mathbf{p}^T \mathbf{f}$, and

$$\mathbf{p} = [p_1 \quad p_2 \quad \cdots \quad p_n]^T = \left[\frac{\partial J^*}{\partial x_1} \quad \frac{\partial J^*}{\partial x_2} \quad \cdots \quad \frac{\partial J^*}{\partial x_n} \right]^T.$$

To proceed, we need one more assumption in addition to the assumptions we made while discussing the HJB equation. We assume that $J^*(\mathbf{x})$ has continuous second partial derivatives $\partial^2 J^* / \partial x_i \partial x_j$; that is, the functions p_i have continuous first derivatives $\partial p_i / \partial x_j$. This implies that

$$\frac{\partial \mathbf{p}}{\partial \mathbf{x}} = \left(\frac{\partial \mathbf{p}}{\partial \mathbf{x}} \right)^T. \quad (5.89)$$

We also assume that the components of \mathbf{f} have continuous first derivatives $\partial f_i / \partial x_j$. We then

compute

$$\begin{aligned}\frac{\partial H}{\partial \mathbf{x}} &= \frac{\partial}{\partial \mathbf{x}}(F + \mathbf{p}^T \mathbf{f}) \\ &= \frac{\partial F}{\partial \mathbf{x}} + \mathbf{f}^T \frac{\partial \mathbf{p}}{\partial \mathbf{x}} + \mathbf{p}^T \frac{\partial \mathbf{f}}{\partial \mathbf{x}},\end{aligned}\quad (5.90)$$

where $\partial F / \partial \mathbf{x} = \nabla_{\mathbf{x}}^T F$. From our previous discussion it follows that the Hamiltonian function attains its minimum on the optimal trajectory. Taking this into account, we conclude that for the optimal process we obtain

$$\frac{\partial H}{\partial \mathbf{x}} = \mathbf{0}^T \quad \text{for } (\mathbf{x}, \mathbf{u}) = (\mathbf{x}^*, \mathbf{u}^*).$$

Hence, for the optimal process $(\mathbf{x}^*, \mathbf{u}^*)$ we have

$$\frac{\partial F}{\partial \mathbf{x}} + \mathbf{f}^T \frac{\partial \mathbf{p}}{\partial \mathbf{x}} + \mathbf{p}^T \frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \mathbf{0}^T. \quad (5.91)$$

We now differentiate \mathbf{p} with respect to time to get

$$\dot{\mathbf{p}} = \frac{\partial \mathbf{p}}{\partial \mathbf{x}} \dot{\mathbf{x}} = \frac{\partial \mathbf{p}}{\partial \mathbf{x}} \mathbf{f}. \quad (5.92)$$

Combining (5.91) and (5.92) and taking into account (5.89), we obtain

$$\dot{\mathbf{p}} = - \left(\frac{\partial F}{\partial \mathbf{x}} \right)^T - \left(\frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right)^T \mathbf{p}. \quad (5.93)$$

Thus, \mathbf{p} is the solution to the differential equation (5.93), and hence $\mathbf{p} = \mathbf{p}(t)$. Therefore $\partial \mathbf{p} / \partial \mathbf{x} = \mathbf{0}$, and by (5.90) we have

$$\frac{\partial H}{\partial \mathbf{x}} = \frac{\partial F}{\partial \mathbf{x}} + \mathbf{p}^T \frac{\partial \mathbf{f}}{\partial \mathbf{x}}. \quad (5.94)$$

Comparing (5.93) and (5.94) gives

$$\dot{\mathbf{p}} = - \left(\frac{\partial H}{\partial \mathbf{x}} \right)^T \quad (5.95)$$

The above equation is called in the literature the *adjoint* or *costate* equation.

We summarize our development in the following theorem.

Theorem 5.6 Necessary conditions for $\mathbf{u} \in U$ to minimize (5.87) subject to (5.88) are

$$\dot{\mathbf{p}} = - \left(\frac{\partial H}{\partial \mathbf{x}} \right)^T,$$

where $H = H(\mathbf{x}, \mathbf{u}, \mathbf{p}) = F(\mathbf{x}, \mathbf{u}) + \mathbf{p}^T \mathbf{f}(\mathbf{x}, \mathbf{u})$, and

$$H(\mathbf{x}^*, \mathbf{u}^*, \mathbf{p}^*) = \min_{\mathbf{u} \in U} H(\mathbf{x}, \mathbf{u}, \mathbf{p}).$$

If the final state, $\mathbf{x}(t_f)$, is free, then in addition to the above conditions it is required that the following end-point condition is satisfied:

$$\mathbf{p}(t_f) = \nabla_{\mathbf{x}} \Phi|_{t=t_f}.$$

We now illustrate the above theorem with the following well-known example.

◆ Example 5.21

We consider a controlled object modeled by

$$\begin{aligned}\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} &= \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u \\ &= \mathbf{A}\mathbf{x} + \mathbf{b}u\end{aligned}$$

with the performance index

$$J = \int_0^{t_f} dt.$$

The control is required to satisfy

$$|u(t)| \leq 1$$

for all $t \in [0, t_f]$. This constraint means that the control must have magnitude no greater than 1. Our objective is to find admissible control that minimizes J and transfers the system from a given initial state \mathbf{x}_0 to the origin.

We begin by finding the Hamiltonian function for the problem:

$$\begin{aligned}H &= 1 + \mathbf{p}^T (\mathbf{A}\mathbf{x} + \mathbf{b}u) \\ &= 1 + [p_1 \quad p_2] \left(\begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u \right) \\ &= 1 + p_1 x_2 + p_2 u.\end{aligned}$$

The costate equations are

$$\begin{bmatrix} \dot{p}_1 \\ \dot{p}_2 \end{bmatrix} = - \begin{bmatrix} \frac{\partial H}{\partial x_1} \\ \frac{\partial H}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 0 \\ -p_1 \end{bmatrix}.$$

Solving the costate equations yields

$$p_1 = d_1 \quad \text{and} \quad p_2 = -d_1 t + d_2,$$

where d_1 and d_2 are integration constants. We now find an admissible control minimizing the Hamiltonian:

$$\begin{aligned}\arg_u \min H &= \arg_u \min (1 + p_1 x_2 + p_2 u) \\ &= \arg_u \min (p_2 u) \\ &= \begin{cases} u(t) = 1 & \text{if } p_2 < 0, \\ u(t) = -1 & \text{if } p_2 > 0. \end{cases}\end{aligned}$$

Hence,

$$u^*(t) = -\text{sign}(p_2^*) = -\text{sign}(-d_1 t + d_2),$$

where

$$\text{sign}(z) = \begin{cases} 1 & \text{if } z > 0, \\ -1 & \text{if } z < 0. \end{cases}$$

Thus, the optimal control law is piecewise constant taking the values 1 or -1 . This control law has at most two intervals of constancy because the argument of the sign function is a linear function, $-d_1 t + d_2$, that changes its sign at most once. This type of control is called a *bang-bang control* because it switches back and forth between its extreme values. System trajectories for $u = 1$ and $u = -1$ are families of parabolas, and their equations were derived in Example 2.2. A phase portrait for $u = 1$ is shown in Figure 2.4, while that for $u = -1$ is shown in Figure 2.5. Note that only one parabola from each family passes through the specified terminal point $\mathbf{x} = [0 \ 0]^T$ in the state plane. Segments of the two parabolas through the origin form the switching curve,

$$x_1 = -\frac{1}{2}x_2^2 \operatorname{sign}(x_2).$$

This means that if an initial state is above the switching curve, then $u = -1$ is used until the switching curve is reached. Then, $u = 1$ is used to reach the origin. For an initial state below the switching curve, the control $u = 1$ is used first to reach the switching curve, and then the control is switched to $u = -1$. The above control action can be described as

$$u = \begin{cases} 1 & \text{if } v < 0, \\ -\operatorname{sign}(x_2) & \text{if } v = 0, \\ -1 & \text{if } v > 0, \end{cases}$$

where $v = v(x_1, x_2) = x_1 + \frac{1}{2}x_2^2 \operatorname{sign}(x_2) = 0$ is the equation describing the switching curve. We can use the above equation to synthesize a closed-loop system such that starting at an arbitrary initial state in the state plane, the trajectory will always be moving in an optimal fashion toward the origin. Once the origin is reached, the trajectory will stay there. A phase portrait of the closed-loop time optimal system is given in Figure 5.27.

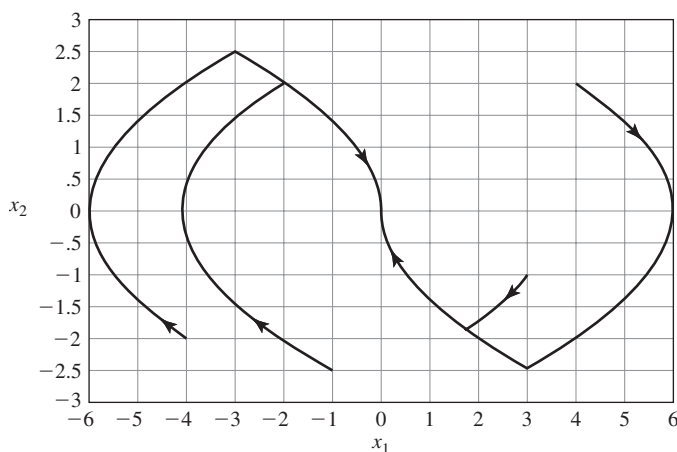


Figure 5.27 A phase-plane portrait of the time optimal closed-loop system of Example 5.21.

◆ Example 5.22

Consider the schematic of an armature-controlled DC motor, given in Figure 5.26. Use the minimum principle to find the armature voltage $e_a(t)$ such that the motor angular velocity ω changes from $\omega(0) = 0$ rad/sec at $t = 0$ to $\omega(1) = 10$ rad/sec while minimizing the energy dissipated in the armature resistor, R_a . There are no constraints on e_a . Note that the energy dissipated in the armature resistor is

$$J = \int_0^1 R_a i_a^2(t) dt,$$

where i_a is the armature current.

The modeling equation of the DC motor was derived in Example 5.19. It is given by (5.84) on page 288. We now represent the expression for the energy dissipated in the armature resistor in terms of ω and e_a . Applying Ohm's law to the armature circuit, we get

$$e_a - K_b \omega = R_a i_a.$$

Hence,

$$J = \int_0^1 R_a i_a^2(t) dt = \int_0^1 \frac{(e_a - K_b \omega)^2}{R_a} dt.$$

The Hamiltonian function is

$$H = \frac{(e_a - K_b \omega)^2}{R_a} + p \left(-\frac{K_i K_b}{R_a I_{eq}} \omega + \frac{K_i}{I_{eq} R_a} e_a \right).$$

Because there are no constraints on e_a , we can find the optimal control by solving the equation

$$0 = \frac{\partial H}{\partial e_a} = \frac{2(e_a - K_b \omega)}{R_a} + \frac{K_i}{I_{eq} R_a} p.$$

Hence,

$$e_a = 2\omega - p. \quad (5.96)$$

The costate equation is

$$\dot{p} = -\frac{\partial H}{\partial \omega} = \frac{2K_b(e_a - K_b \omega)}{R_a} + \frac{K_i K_b}{R_a I_{eq}} p = 2(e_a - 2\omega) + 2p. \quad (5.97)$$

We substitute (5.96) into (5.97) to obtain

$$\dot{p} = 0.$$

Hence,

$$p(t) = \text{constant} = c.$$

Therefore,

$$e_a = 2\omega - c. \quad (5.98)$$

Substituting (5.98) into (5.84), we get

$$\dot{\omega} = -c.$$

Solving the above yields

$$\omega(t) = -ct + \omega(0) = -ct,$$

because $\omega(0) = 0$. We next use the boundary condition $\omega(1) = 10$ to find $c = -10$. Hence, the optimal armature voltage is

$$e_a^*(t) = 20t + 10$$

◆ Example 5.23

In this example, we will solve Johann Bernoulli's brachistochrone problem, introduced in Section 5.2, using the minimum principle. We first recall the problem as restated by Johann Bernoulli in 1697 (quoted from Dunham [70, p. 200]):

"Among the infinitely many curves which join the two given points . . . choose the one such that, if the curve is replaced by a thin tube or groove, and a small sphere placed in it and released, then this will pass from one point to the other in the shortest time."

It is difficult not to mention the reward that one obtains for solving the problem (quoted from Dunham [70, p. 200]):

"Let who can seize quickly the prize which we have promised to the solver. Admittedly this prize is neither of gold nor silver, for these appeal only to base and venal souls. . . . Rather, since virtue itself is its own most desirable reward and a fame is a powerful incentive, we offer the prize fitting for the man of noble blood, compounded of honor, praise, and approbation . . ." We add that Johann Bernoulli wrote the above after he already knew the solution to the problem!

Model and the Performance Index

The field is conservative, hence the sums of the kinetic and potential energy is constant,

$$\frac{1}{2}mv^2(t) - mgy(t) = \frac{1}{2}mv^2(t_0) - mgy(t_0) = 0,$$

where we assumed that the particle is at rest at t_0 . Using the above, we can find the particle's speed at any time $t \geq t_0$:

$$v(t) = \sqrt{2gy(t)}.$$

Then, the brachistochrone problem can be formulated as an optimal control problem of the form

$$\begin{aligned} &\text{minimize} \quad J(t_0) = \int_{t_0}^T dt \\ &\text{subject to} \quad \begin{cases} \dot{x} = v \cos \theta, \\ \dot{y} = v \sin \theta, \end{cases} \end{aligned}$$

where θ is the angle that the particle velocity vector forms with the positive x axis, while \dot{x} and \dot{y} are the horizontal and vertical components of the particle velocity vector. The path angle θ is considered as the control input to be determined.

Particle Trajectory

To proceed, we form the Hamiltonian function

$$H = 1 + p_1 v \cos \theta + p_2 v \sin \theta, \quad (5.99)$$

where p_1 and p_2 are the costate variables. The costate equations are

$$\begin{bmatrix} \dot{p}_1 \\ \dot{p}_2 \end{bmatrix} = - \begin{bmatrix} \frac{\partial H}{\partial x} \\ \frac{\partial H}{\partial y} \end{bmatrix} = \begin{bmatrix} 0 \\ -\frac{g}{v}(p_1 \cos \theta + p_2 \sin \theta) \end{bmatrix}.$$

It follows from the above that $p_1 = K = \text{constant}$. Because there are no constraints on the control input θ , the necessary condition for the control optimality is

$$\frac{\partial H}{\partial \theta} = 0,$$

where

$$\frac{\partial H}{\partial \theta} = -p_1 v \sin \theta + p_2 v \cos \theta.$$

Hence,

$$p_2 = p_1 \tan \theta = K \tan \theta.$$

We next find an expression for $\dot{\theta}$. First note that

$$\dot{p}_2 = \frac{\partial p_2}{\partial \theta} \dot{\theta} = \frac{\partial (K \tan \theta)}{\partial \theta} \dot{\theta} = \frac{K}{\cos^2 \theta} \dot{\theta}. \quad (5.100)$$

On the other hand,

$$\begin{aligned} \dot{p}_2 &= -\frac{g}{v}(p_1 \cos \theta + p_2 \sin \theta) \\ &= -\frac{g}{v}(K \cos \theta + K \tan \theta \sin \theta) \\ &= -\frac{Kg}{v} \left(\cos \theta + \frac{\sin^2 \theta}{\cos \theta} \right) \\ &= -\frac{Kg}{v} \left(\frac{\cos^2 \theta + \sin^2 \theta}{\cos \theta} \right) \\ &= -\frac{Kg}{v \cos \theta}. \end{aligned} \quad (5.101)$$

Combining (5.100) and (5.101) yields

$$\dot{\theta} = -\frac{g}{v} \cos \theta. \quad (5.102)$$

We now use the boundary conditions to express x and y as functions of time and θ . First, using the above we can write

$$\dot{y} = \frac{dy}{d\theta} \dot{\theta} = \frac{dy}{d\theta} \left(-\frac{g}{v} \cos \theta \right).$$

On the other hand,

$$\dot{y} = v \sin \theta.$$

Equating the right-hand sides of the above two equations yields

$$\frac{dy}{v^2} = \frac{dy}{2gy} = -\frac{1}{g} \tan \theta d\theta.$$

Integrating both sides gives

$$\frac{1}{2g} \int \frac{dy}{y} = -\frac{1}{g} \int \tan \theta d\theta,$$

that is,

$$\frac{1}{2g} \ln |y| = \frac{1}{g} \ln |\cos \theta|.$$

Performing manipulations, we obtain

$$y = C \cos^2 \theta,$$

where C is an integration constant. Evaluating the above at $t = T$ gives

$$y(T) = y_1 = C \cos^2 \theta(T).$$

Hence,

$$C = \frac{y_1}{\cos^2 \theta(T)}$$

and

$$\boxed{y(t) = \frac{y_1}{\cos^2 \theta(T)} \cos^2 \theta(t)} \quad (5.103)$$

In a similar way, we can obtain an expression for x as a function of time and θ . Indeed, on the one hand we have

$$\dot{x} = \frac{dx}{d\theta} \dot{\theta} = -\frac{dx}{d\theta} \frac{g}{v} \cos \theta,$$

while on the other hand we obtain

$$\dot{x} = v \cos \theta.$$

Comparing the right-hand sides of the above two equations gives

$$\frac{dx}{d\theta} = -\frac{v^2}{g} = -2y = -2\frac{y_1}{\cos^2 \theta(T)} \cos^2 \theta.$$

We represent the above as

$$dx = -2\frac{y_1}{\cos^2 \theta(T)} \cos^2 \theta d\theta.$$

Integrating both sides, using the identity, $\cos^2 \theta = \frac{1}{2}(1 + \cos 2\theta)$, yields

$$x = -2\frac{y_1}{\cos^2 \theta(T)} \left(\frac{\theta}{2} + \frac{\sin 2\theta}{4} \right) + \tilde{C},$$

where \tilde{C} is the integration constant. Evaluating the above at $t = T$ gives

$$x_1 = -2\frac{y_1}{\cos^2 \theta(T)} \left(\frac{\theta(T)}{2} + \frac{\sin 2\theta(T)}{4} \right) + \tilde{C}.$$

Hence,

$$\tilde{C} = x_1 + \frac{y_1}{\cos^2 \theta(T)} \left(\theta(T) + \frac{\sin 2\theta(T)}{2} \right).$$

Using the above, we obtain

$$x(t) = x_1 + \frac{y_1}{2 \cos^2 \theta(T)} (2(\theta(T) - \theta(t)) + \sin 2\theta(T) - \sin 2\theta(t)) \quad (5.104)$$

We point out that the equation for $x(t)$, $y(t)$ and $\dot{\theta}$ specify an implicit feedback control law in the sense that the control θ can be calculated given the state $(x(t), y(t))$.

We now use the boundary conditions to find the description of the time-optimal closed loop system. We can assume, without loss of generality, that $x(t_0) = y(t_0) = x(0) = y(0) = 0$. Then, using the expression for $y(t)$ we note that $y(t_0) = 0 = \frac{y_1}{\cos^2 \theta(T)} \cos^2 \theta(t_0)$. Thus, the angle of departure is $\theta(t_0) = 90^\circ$. With the above initial path angle, y , and hence v , is increased at the outset as quickly as possible. Using these data, we next evaluate x at $t = t_0$ to form an equation that can be used to find the angle of arrival $\theta(T)$. The equation is

$$2x_1 \cos^2 \theta(T) + y_1(2\theta(T) - \pi + \sin 2\theta(T)) = 0.$$

The above equation can be solved easily using MATLAB's Symbolic Toolbox or the MATLAB's `fsolve` function. Once $\theta(T)$ is found, we use the expression for y to find $\cos \theta$, where

$$\cos \theta(t) = \sqrt{\frac{y(t)}{y_1}} \cos \theta(T).$$

Hence,

$$\sin \theta(t) = \sqrt{1 - \cos^2 \theta(t)}.$$

We substitute the above expression into $\dot{x} = v \cos \theta$ and $\dot{y} = v \sin \theta$ to obtain the state-space model of the time optimal motion of the particle moving under the force of gravity from a given point to another specified point. The equations are

$$\begin{aligned}\dot{x}(t) &= \sqrt{2gy(t)} \sqrt{\frac{y(t)}{y_1}} \cos \theta(T), \\ \dot{y}(t) &= \sqrt{2gy(t)} \sqrt{1 - \frac{y(t)}{y_1} \cos^2 \theta(T)}\end{aligned}$$

subject to the boundary conditions $[x(t_0) \ y(t_0)] = [x_0 \ y_0]$ and $[x(T) \ y(T)] = [x_1 \ y_1]$.

The Time of Quickest Descent

We now derive an expression that can be used to compute the time of the particle travel. We first establish an auxiliary result that we will use in our derivation. Specifically, we will show that $\dot{\theta}(t)$ is constant on the interval $[t_0, T]$. Differentiating both sides of equation (5.102) gives

$$\begin{aligned}\frac{d\dot{\theta}}{dt} &= \frac{d}{dt} \left(-\frac{g}{\sqrt{2gy}} \cos \theta \right) \\ &= -\frac{g}{\sqrt{2g}} \left(\frac{d}{dt} \left(\frac{1}{\sqrt{y}} \cos \theta \right) \right) \\ &= -\frac{g}{\sqrt{2g}} \left(-\frac{1}{2} \frac{\dot{y} \cos \theta}{y^{3/2}} - \frac{\sin \theta}{\sqrt{y}} \dot{\theta} \right) \\ &= \frac{g}{\sqrt{2g}} \left(\frac{1}{2} \frac{\dot{y} \cos \theta}{y^{3/2}} + \frac{\sin \theta}{\sqrt{y}} \dot{\theta} \right).\end{aligned}$$

Substituting into the above the expressions for $\dot{\theta}$ and \dot{y} gives

$$\begin{aligned}\frac{d\dot{\theta}}{dt} &= \frac{g}{\sqrt{2g}} \left(\frac{1}{2} \frac{\sqrt{2gy} \sin \theta \cos \theta}{y^{3/2}} - g \frac{\sin \theta \cos \theta}{\sqrt{y} \sqrt{2gy}} \right) \\ &= \frac{g}{\sqrt{2g}} \left(\frac{\sqrt{g}}{y\sqrt{2}} \sin \theta \cos \theta - \frac{\sqrt{g}}{y\sqrt{2}} \sin \theta \cos \theta \right) \\ &= 0.\end{aligned}$$

Because $\dot{\theta}$ is constant on the interval $[t_0, T]$, we use (5.102) to write

$$\frac{d\theta}{dt} = -\frac{g}{\sqrt{2gy(T)}} \cos \theta(T).$$

Separating the variables yields

$$d\theta = -\frac{g}{\sqrt{2gy(T)}} \cos \theta(T) dt.$$

Integrating both sides of the above, we obtain

$$\int_{\theta(t)}^{\theta(T)} d\theta = -\frac{g}{\sqrt{2gy(T)}} \cos \theta(T) \int_{t_0}^T dt.$$

Hence,

$$\theta(T) - \theta(t) = -\frac{g}{\sqrt{2gy(T)}} \cos \theta(T) (T - t_0).$$

To obtain an expression for the particle travel time, we set in the above equation $t = t_0 = 0$ and perform simple manipulations to obtain

$$T = \frac{\pi/2 - \theta(T)}{\cos \theta(T)} \sqrt{2y(T)/g}$$

It's an Upside-Down Cycloid!

The path of quickest descent is an upside-down cycloid. A cycloid is a curve described by a point on the circumference of a circle that rolls without slipping (see Figure 5.28). The equations describing a cycloid are

$$\begin{aligned} x(\phi) &= R(\phi - \sin \phi), \\ y(\phi) &= R(1 - \cos \phi). \end{aligned}$$

We will now show that equations (5.104) and (5.103) describe a cycloid. To proceed, we define

$$R = \frac{y_1}{2 \cos^2 \theta(T)} \quad \text{and} \quad \phi(t) = \pi - 2\theta(t). \quad (5.105)$$

Using the above notation and taking into account that $\sin(\pi - \phi) = \sin \phi$, we rewrite (5.104) as

$$x(t) - x_1 + R(\phi(T) - \sin \phi(T)) = R(\phi(t) - \sin \phi(t)).$$

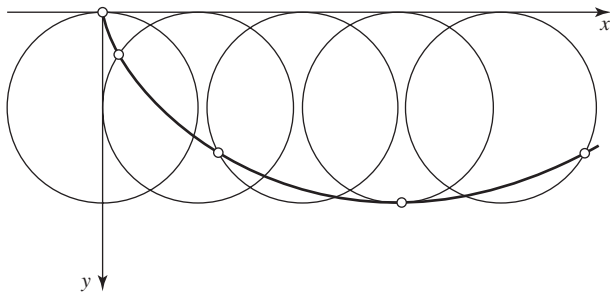


Figure 5.28 Constructing an upside-down cycloid—the solution to the brachistochrone problem.

Using the notation given in (5.105) and the fact that $\cos 2\phi = 2 \cos^2 \phi - 1$, we represent (5.103) as

$$\begin{aligned} y(t) &= \frac{y_1}{\cos^2 \theta(T)} \cos^2 \theta(t) \\ &= \frac{y_1}{\cos^2 \theta(T)} \frac{\cos 2\theta(t) + 1}{2} \\ &= R(1 + \cos(\pi - \phi(t))) \\ &= R(1 - \cos \phi(t)), \end{aligned}$$

because $\cos(\pi - \phi(t)) = -\cos \phi$.

In the following, we discuss a general minimum time regulator problem for a class of multi-input dynamical systems modeled by

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}.$$

The goal is to transfer the system from some initial state to the origin in minimum time subject to

$$|u_i| \leq k_i, \quad k_i > 0, \quad i = 1, 2, \dots, m.$$

Thus, the performance index for the problem is $J = \int_{t_0}^{t_f} dt$. We begin our analysis by forming the associated Hamiltonian function

$$\begin{aligned} H &= 1 + \mathbf{p}^T (\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}) \\ &= 1 + \mathbf{p}^T \mathbf{A}\mathbf{x} + \sum_{i=1}^m (\mathbf{p}^T \mathbf{b}_i) u_i, \end{aligned}$$

where \mathbf{b}_i is the i th column of \mathbf{B} . Let

$$s_i(t) = \mathbf{p}^{*T}(t) \mathbf{b}_i.$$

Then, applying Pontryagin's minimum principle, we obtain the following form of the optimal control:

$$u_i^*(t) = -k_i \text{sign}(s_i(t)), \quad i = 1, 2, \dots, m.$$

Thus, the optimal control for the i th control channel has the bang-bang form, $u_i^* = \pm k_i$. Note that $s_i(t) = \mathbf{p}^{*T}(t) \mathbf{b}_i = 0$ is the switching curve for the i th control channel. We compute the costate vector from the costate equation,

$$\begin{aligned} \dot{\mathbf{p}}^{*T} &= -\frac{\partial H}{\partial \mathbf{x}} \\ &= -\frac{\partial}{\partial \mathbf{x}} (\mathbf{p}^{*T} \mathbf{A}\mathbf{x}) \\ &= -\mathbf{p}^{*T} \mathbf{A}, \end{aligned}$$

or

$$\dot{\mathbf{p}}^* = -\mathbf{A}^T \mathbf{p}^*.$$

The solution to the costate equation is

$$\mathbf{p}^*(t) = e^{-\mathbf{A}^T t} \mathbf{p}^*(0).$$

Thus, the i th switching curve can be represented as

$$s_i(t) = \mathbf{p}^{*T}(0) e^{-\mathbf{A}^T t} \mathbf{b}_i = 0.$$

Observe that if $s_i(t) = 0$ on some time interval, then $u_i^*(t)$ is indefinite on this interval. We now show that under some mild assumptions there is no time interval on which $s_i(t) = 0$; that is, there is no time interval on which $u_i^*(t)$ is indefinite. First, we assume that $\mathbf{b}_i \neq \mathbf{0}$. Note that in our problem here, $\min_{\mathbf{u} \in U} H(\mathbf{x}, \mathbf{u}) = 0$ because $J^*(t, \mathbf{x}) = J^*(\mathbf{x})$ and therefore $\partial J^*/\partial t = 0$. Hence, $1 + \mathbf{p}^{*T}(t)(\mathbf{A}\mathbf{x}^* + \mathbf{B}\mathbf{u}^*) = 0$ for all t for the optimal process. Thus, the optimal costate vector $\mathbf{p}^*(t)$ cannot be zero for any value of t . Finally, if $s_i = \mathbf{p}^{*T} \mathbf{b}_i$ is zero on some time interval—that is, s_i is constant on a time interval—then this implies that on the above time interval, we obtain

$$\dot{s}_i = \dot{\mathbf{p}}^{*T} \mathbf{b}_i = -\mathbf{p}^{*T} \mathbf{A} \mathbf{b}_i = 0.$$

Similarly, higher derivatives of $s_i(t) = \mathbf{p}^{*T} \mathbf{b}_i$ will also be zero on the time interval on which $s_i(t) = \mathbf{p}^{*T} \mathbf{b}_i$ is zero, that is, $s_i^{(j)}(t) = (-1)^j \mathbf{p}^{*T} \mathbf{A}^j \mathbf{b}_i = 0$, $j = 1, 2, \dots$. We represent the above relations as

$$\mathbf{p}^{*T} [\mathbf{b}_i \quad \mathbf{A} \mathbf{b}_i \quad \dots \quad \mathbf{A}^{n-1} \mathbf{b}_i \quad \dots] = \mathbf{0}^T.$$

If the system is controllable by the i th input acting alone, then

$$\text{rank}[\mathbf{b}_i \quad \mathbf{A} \mathbf{b}_i \quad \dots \quad \mathbf{A}^{n-1} \mathbf{b}_i \quad \dots] = n,$$

and we would have to have $\mathbf{p}^* = \mathbf{p}^*(t) = \mathbf{0}$. However, we already ruled out this case and thus $s_i(t) = \mathbf{p}^{*T} \mathbf{b}_i$ cannot be zero on any time interval. In conclusion, if the system is controllable by each input acting alone, then there is no time interval on which optimal control vector is indefinite.

5.5.2 A Two-Point Boundary-Value Problem

We consider the problem of minimizing the quadratic performance index

$$J = \frac{1}{2} \mathbf{x}^T(t_f) \mathbf{F} \mathbf{x}(t_f) + \frac{1}{2} \int_{t_0}^{t_f} (\mathbf{x}^T(t) \mathbf{Q} \mathbf{x}(t) + \mathbf{u}^T(t) \mathbf{R} \mathbf{u}(t)) dt \quad (5.106)$$

subject to

$$\dot{\mathbf{x}}(t) = \mathbf{A} \mathbf{x}(t) + \mathbf{B} \mathbf{u}(t), \quad \mathbf{x}(t_0) = \mathbf{x}_0, \quad \mathbf{x}(t_f) \text{ is free.} \quad (5.107)$$

We assume that there are no constraints on the control input \mathbf{u} . The weight matrices \mathbf{F} and \mathbf{Q} are real symmetric positive semidefinite, and the matrix \mathbf{R} is real symmetric positive definite. The Hamiltonian function for the above linear quadratic control problem is

$$H = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u} + \mathbf{p}^T (\mathbf{A} \mathbf{x} + \mathbf{B} \mathbf{u}), \quad (5.108)$$

where $\mathbf{p} \in \mathbb{R}^n$ is the costate vector. It follows from Theorem 5.6 that the optimal controller must minimize the Hamiltonian function. Because the control vector is unconstrained, the necessary

condition for optimality of the control \mathbf{u} is

$$\frac{\partial H}{\partial \mathbf{u}} = \mathbf{0}^T. \quad (5.109)$$

Evaluating (5.109) yields

$$\frac{\partial H}{\partial \mathbf{u}} = \mathbf{u}^T \mathbf{R} + \mathbf{p}^T \mathbf{B} = \mathbf{0}^T. \quad (5.110)$$

Hence, the optimal controller has the form

$$\mathbf{u} = -\mathbf{R}^{-1} \mathbf{B}^T \mathbf{p}. \quad (5.111)$$

Substituting the above into (5.107), we obtain

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} - \mathbf{B}\mathbf{R}^{-1} \mathbf{B}^T \mathbf{p}, \quad \mathbf{x}(t_0) = \mathbf{x}_0. \quad (5.112)$$

We next find the costate equation,

$$\begin{aligned} \dot{\mathbf{p}} &= - \left(\frac{\partial H}{\partial \mathbf{x}} \right)^T \\ &= -\mathbf{Q}\mathbf{x} - \mathbf{A}^T \mathbf{p}. \end{aligned} \quad (5.113)$$

Combining (5.112) and the above yields $2n$ linear differential equations:

$$\begin{bmatrix} \dot{\mathbf{x}} \\ \dot{\mathbf{p}} \end{bmatrix} = \begin{bmatrix} \mathbf{A} & -\mathbf{B}\mathbf{R}^{-1} \mathbf{B}^T \\ -\mathbf{Q} & -\mathbf{A}^T \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{p} \end{bmatrix}. \quad (5.114)$$

By Theorem 5.6, the state vector at time t_f must satisfy

$$\mathbf{p}(t_f) = \frac{1}{2}(\nabla \mathbf{x}^T \mathbf{F} \mathbf{x})|_{t=t_f} = \mathbf{F} \mathbf{x}(t_f). \quad (5.115)$$

In summary, we reduced the linear quadratic control problem to solving $2n$ linear differential equations with mixed boundary conditions, which is an example of a *two-point boundary value problem*, or TPBVP for short. We shall now solve the above TPBVP. If we had the initial conditions $\mathbf{x}(t_0)$ and $\mathbf{p}(t_0)$, then the solution to (5.114) would have the form

$$\begin{bmatrix} \mathbf{x}(t_f) \\ \mathbf{p}(t_f) \end{bmatrix} = e^{\mathbf{H}(t_f-t_0)} \begin{bmatrix} \mathbf{x}(t_0) \\ \mathbf{p}(t_0) \end{bmatrix}, \quad (5.116)$$

where

$$\mathbf{H} = \begin{bmatrix} \mathbf{A} & -\mathbf{B}\mathbf{R}^{-1} \mathbf{B}^T \\ -\mathbf{Q} & -\mathbf{A}^T \end{bmatrix}.$$

However, in this particular TPBVP, $\mathbf{p}(t_0)$ is unknown and instead we are given $\mathbf{p}(t_f)$. To proceed, let

$$\begin{bmatrix} \mathbf{x}(t_f) \\ \mathbf{p}(t_f) \end{bmatrix} = e^{\mathbf{H}(t_f-t)} \begin{bmatrix} \mathbf{x}(t) \\ \mathbf{p}(t) \end{bmatrix}, \quad (5.117)$$

where

$$e^{H(t_f-t)} = \begin{bmatrix} \Phi_{11}(t_f, t) & \Phi_{12}(t_f, t) \\ \Phi_{21}(t_f, t) & \Phi_{22}(t_f, t) \end{bmatrix}.$$

Each of the blocks $\Phi_{ij}(t_f, t)$, $i, j = 1, 2$, is $n \times n$. Thus we have

$$\begin{bmatrix} \mathbf{x}(t_f) \\ \mathbf{p}(t_f) \end{bmatrix} = \begin{bmatrix} \Phi_{11}(t_f, t) & \Phi_{12}(t_f, t) \\ \Phi_{21}(t_f, t) & \Phi_{22}(t_f, t) \end{bmatrix} \begin{bmatrix} \mathbf{x}(t) \\ \mathbf{p}(t) \end{bmatrix}. \quad (5.118)$$

The above matrix equation can be represented as

$$\mathbf{x}(t_f) = \Phi_{11}(t_f, t)\mathbf{x}(t) + \Phi_{12}(t_f, t)\mathbf{p}(t), \quad (5.119)$$

$$\mathbf{p}(t_f) = \Phi_{21}(t_f, t)\mathbf{x}(t) + \Phi_{22}(t_f, t)\mathbf{p}(t). \quad (5.120)$$

Using the boundary condition (5.115) and (5.119) gives

$$\mathbf{p}(t_f) = \mathbf{F}\mathbf{x}(t_f) = \mathbf{F}\Phi_{11}(t_f, t)\mathbf{x}(t) + \mathbf{F}\Phi_{12}(t_f, t)\mathbf{p}(t). \quad (5.121)$$

Subtracting (5.120) from (5.121), we obtain

$$\mathbf{0} = (\mathbf{F}\Phi_{11}(t_f, t) - \Phi_{21}(t_f, t))\mathbf{x}(t) + (\mathbf{F}\Phi_{12}(t_f, t) - \Phi_{22}(t_f, t))\mathbf{p}(t). \quad (5.122)$$

Using the above, we compute

$$\begin{aligned} \mathbf{p}(t) &= (\Phi_{22}(t_f, t) - \mathbf{F}\Phi_{12}(t_f, t))^{-1}(\mathbf{F}\Phi_{11}(t_f, t) - \Phi_{21}(t_f, t))\mathbf{x}(t) \\ &= \mathbf{P}(t)\mathbf{x}(t), \end{aligned} \quad (5.123)$$

where

$$\mathbf{P}(t) = (\Phi_{22}(t_f, t) - \mathbf{F}\Phi_{12}(t_f, t))^{-1}(\mathbf{F}\Phi_{11}(t_f, t) - \Phi_{21}(t_f, t)). \quad (5.124)$$

Substituting (5.123) into (5.111) yields the optimal state-feedback controller

$$\mathbf{u}(t) = -\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}(t)\mathbf{x}(t). \quad (5.125)$$

We will now show that $\mathbf{P}(t)$ satisfies a matrix differential equation. Indeed, differentiating (5.123) gives

$$\dot{\mathbf{P}}\mathbf{x} + \mathbf{P}\dot{\mathbf{x}} - \dot{\mathbf{p}} = \mathbf{0}. \quad (5.126)$$

Substituting into (5.112), (5.113), and (5.123) yields

$$(\dot{\mathbf{P}} + \mathbf{P}\mathbf{A} - \mathbf{P}\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P} + \mathbf{Q} + \mathbf{A}^T\mathbf{P})\mathbf{x} = \mathbf{0}. \quad (5.127)$$

The preceding must hold throughout $t_0 \leq t \leq t_f$. Hence, \mathbf{P} must satisfy

$$\boxed{\dot{\mathbf{P}} = \mathbf{P}\mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P} - \mathbf{A}^T\mathbf{P} - \mathbf{P}\mathbf{A} - \mathbf{Q}} \quad (5.128)$$

subject to the boundary condition

$$\mathbf{P}(t_f) = \mathbf{F}. \quad (5.129)$$

Equation (5.128) is called the matrix *Riccati differential equation*. Note that because \mathbf{F} is

symmetric, so is $P = P(t)$. Observe that $P = P(t)$ given by (5.124) is a solution to (5.128) subject to the boundary condition (5.129).

The Riccati differential equation (5.128) associated with a class of one-dimensional plants was solved in Example 5.18.

Notes

Sussmann and Willems [273] made the following statement, in 1997, regarding the origins of optimal control: “Optimal control was born in 1697—300 years ago—in Groningen, a university town in the north of The Netherlands, when Johann Bernoulli, professor of mathematics at the local university from 1695 to 1705, published his solution of the *brachystochrone problem*.”

For a serious student of optimal control, the seminal text by Lee and Markus [177] is recommended. Less advanced than Lee and Markus are Kirk [160], Bryson and Ho [36], Owens [218], and Pierre [234]. The classical text by Kwakernaak and Sivan [173] is also worth perusing because of its clarity and its treatment of the stochastic aspects of control problems. Informative, and at the same time entertaining, accounts of the origins of the calculus of variations are given by Dunham [70, Chapter 8] and by Sussmann and Willems [273]. Elsgolc [75] provides a lucid introduction, illustrated with many examples, to the calculus of variations. Applications of the calculus of variations in physics and engineering are provided by Weinstock [297]. Interrelations between the calculus of variations and optimal control are discussed by Hestenes [121]. Frequency-domain aspects of optimal control and robust optimal control design are discussed by Burl [39]. Informal discussion of the Pontryagin minimum principle in Section 5.5 was adapted from Boltyanskii [30]. For more numerical examples related to optimal control, see Lewis and Syrmos [182], and Bryson [37]. Exciting applications of optimal control to a number of control problems of technology and science can be found in Lee and Markus [177].

EXERCISES

5.1 Find the variations of the functionals:

$$(a) \ v(x_1, x_2) = \int_0^{\pi/2} (x_1^4 + 2x_1^2 x_2^2 + x_2^2) dt,$$

$$(b) \ v(x_1, x_2) = \int_0^1 (e^{x_1 - x_2} - e^{x_1 + x_2}) dt.$$

5.2 Find the equations of the curves that are extremals for the functional

$$v(x) = \int_0^{t_f} \left(\frac{1}{4} \dot{x}^2 - t\dot{x} - x + \frac{1}{2} x\dot{x} \right) dt$$

for the boundary conditions

$$(a) \ x(0) = 1, t_f = 2, \text{ and } x(2) = 10;$$

$$(b) \ x(0) = 1, t_f = 10, \text{ and } x(10) \text{ is free.}$$

5.3 Find the equation of the curve that is an extremal for the functional

$$v(x) = \int_0^{t_f} \left(\frac{1}{2} \dot{x}^2 - t\dot{x} + 2x^2 \right) dt$$

subject to the boundary conditions

$$x(0) = 3/4, \quad t_f = 1, \quad \text{and} \quad x(1) \text{ is free.}$$

5.4 Given a dynamical system model

$$\begin{aligned}\dot{\mathbf{x}} &= \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u, \\ y &= [2 \quad 0] \mathbf{x},\end{aligned}$$

and the associated performance index

$$J = \int_0^\infty (y^2 + u^2) dt.$$

Assuming that $u = -\mathbf{k}\mathbf{x}$ is such that J is minimized, find the poles of the closed-loop system.

5.5 Given a dynamical system model

$$\begin{aligned}\dot{\mathbf{x}} &= \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u, \\ y &= [2 \quad 0] \mathbf{x}\end{aligned}$$

and the performance index

$$J = \int_0^\infty (y^2 + u^2) dt.$$

Use the Kronecker product to represent the associated algebraic Riccati equation as the nonlinear vector algebraic equation. Solve it, and construct the optimal state-feedback control law $u = -\mathbf{k}\mathbf{x}$. Assume that

$$\mathbf{x}(0) = \begin{bmatrix} 1 \\ 2 \end{bmatrix}.$$

Find the optimal value of J .

5.6 Find a state-feedback control law u so as to minimize

$$J = \int_0^1 (x^2 + u^2) dt$$

subject to

$$\dot{x} = -x + u, \quad x(0) = 2.$$

Assume $J^* = p(t)x^2$.

5.7 Use the principle of optimality to find the set of numbers x_1, x_2, \dots, x_n satisfying $x_1 + x_2 + \dots + x_n = c$, where c is a positive constant such that the product $x_1 x_2 \dots x_n$ is maximized (see [20]).

5.8 Use dynamic programming to find the maximum cost path between node **a** and node **t** in the routing network depicted in Figure 5.29, where the travel costs are shown

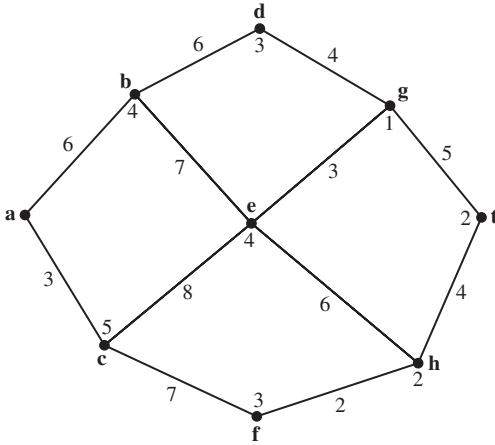


Figure 5.29 A routing network of Exercise 5.8.

beside each segment and toll charges are shown by each node. The network can only be traversed from left to right.

- 5.9** Use dynamic programming to find $u(0)$ and $u(1)$ that minimize

$$J = (x(2) - 2)^2 + \sum_{k=0}^1 u^2(k)$$

subject to

$$x(k+1) = x(k) + u(k), \quad x(0) = 1.$$

- 5.10** Consider the following optimal control problem:

$$\text{minimize } J_2(u) = \frac{1}{2} \sum_{k=0}^1 u(k)^2$$

$$\text{subject to } x(k+1) = x(k) + 2u(k), \quad x(0) = 3, \quad x(2) = 0.$$

Convert the above optimal control problem into a quadratic optimization problem with an equality constraint using the composite input vector

$$z = \begin{bmatrix} u(0) \\ u(1) \end{bmatrix}.$$

Then, employ Lagrange's multiplier theorem to find the optimal sequence $\{u^*(0), u^*(1)\}$.

- 5.11** Consider the following optimal control problem:

$$\text{minimize } J_2(u) = \frac{1}{2} x(2)^2 + \frac{1}{2} \sum_{k=0}^1 (x(k)^2 + u(k)^2)$$

$$\text{subject to } x(k+1) = x(k) + 2u(k), \quad x(0) = 3.$$

Convert the above optimal control problem into a quadratic optimization problem with an equality constraint using the composite input-state vector

$$z = \begin{bmatrix} u(0) \\ x(1) \\ u(1) \\ x(2) \end{bmatrix}.$$

Then, employ Lagrange's multiplier theorem to find the optimal sequence $\{u^*(0), u^*(1)\}$.

5.12 Find $u = u(t)$ so as to minimize

$$J = \int_0^1 u^2 dt$$

subject to

$$\dot{x} = -2x + u, \quad x(0) = 10, \quad x(1) = 0.$$

5.13 For the circuit shown in Figure 5.30, find the input voltage, $u(t)$, that changes the current in the inductor from $i = 10$ A at $t = 0$ to $i = 0$ at $t = 1$ sec while minimizing the performance index,

$$J = \int_0^1 u^2(t) dt.$$

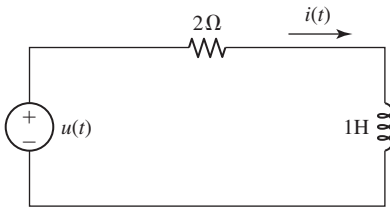


Figure 5.30 A circuit for Exercise 5.13.

5.14 Use the minimum principle to find the input voltage $u(t)$ that charges the capacitor, in Figure 5.31, from $x(0) = 2$ V at $t = 0$ to $x(10) = 12$ V while minimizing the energy dissipated in the resistor. There are no constraints on $u(t)$. Note that, by Kirchhoff's voltage law,

$$Ri(t) + x(t) = u(t),$$

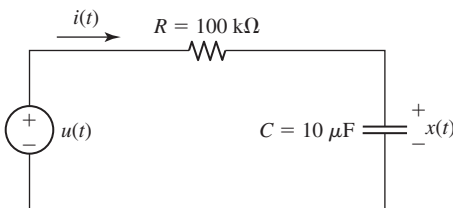


Figure 5.31 A capacitor charging circuit of Exercise 5.14.

where $i(t) = C \frac{dx(t)}{dt}$, and the energy dissipated is

$$J = \int_0^{10} R i^2(t) dt.$$

5.15 Show that if μ is an eigenvalue of the Hamiltonian matrix \mathbf{H} , then so is $-\mu$.

5.16 Consider the Hamiltonian function defined by

$$H = F + \mathbf{p}^T \mathbf{f}.$$

Show that if the functions \mathbf{f} and F do not explicitly depend upon t , then

$$H|_{u=u^*} = \text{constant}, \quad t_0 \leq t \leq t_f.$$

5.17 Given the Hamiltonian matrix corresponding to a dynamical system and its associated performance index. The Hamiltonian matrix is

$$\mathbf{H} = \begin{bmatrix} 2 & -5 \\ -1 & -2 \end{bmatrix}.$$

The MATLAB command `[V,D] = eig(H)` resulted in the following data:

$$\mathbf{V} = \begin{bmatrix} 0.9806 & 0.7071 \\ -0.1961 & 0.7071 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} 3 & 0 \\ 0 & -3 \end{bmatrix}.$$

- (a) Write the equation of the closed loop system driven by the optimal controller $u = -kx$.
- (b) Find the solution to the algebraic Riccati equation corresponding to the optimal controller.

5.18 (a) Suppose that s_i is an eigenvalue of the matrix

$$\mathbf{A} - \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P},$$

and $\mathbf{v}_i \neq \mathbf{0}$ is the corresponding eigenvector, where \mathbf{P} is the symmetric positive definite solution of the algebraic Riccati equation. Show that

$$\begin{bmatrix} \mathbf{v}_i \\ \mathbf{P}\mathbf{v}_i \end{bmatrix}$$

is an eigenvector of the associated Hamiltonian matrix corresponding to the eigenvalue s_i .

- (b) Show that if $\mathbf{Q} = \mathbf{Q}^T > 0$, then the real parts of the eigenvalues of the matrix

$$\mathbf{A} - \mathbf{B}\mathbf{R}^{-1}\mathbf{B}^T\mathbf{P}$$

are all negative, where $\mathbf{P} = \mathbf{P}^T > 0$ is the solution of the associated algebraic Riccati equation.

5.19 Given the following model of a dynamical system:

$$\dot{x}_1 = x_2 + 5,$$

$$\dot{x}_2 = u,$$

where

$$|u| \leq 1.$$

The performance index to be minimized is

$$J = \int_0^{t_f} dt.$$

Find the state-feedback control law $u = u(x_1, x_2)$ that minimizes J and drives the system from a given initial condition $\mathbf{x}(0) = [x_1(0), x_2(0)]^T$ to the final state $\mathbf{x}(t_f) = \mathbf{0}$. Proceed as indicated below.

- Derive the equations of the optimal trajectories.
- Derive the equation of the switching curve.
- Write the expression for the optimal state-feedback controller.

5.20 (Based on Rugh and Murphy [246]) Given a dynamical system model

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b}u,$$

where

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & & & \cdots & & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ -a_0 & -a_1 & -a_2 & \cdots & a_{n-2} & a_{n-1} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix},$$

and the associated performance index

$$J = \frac{1}{2} \int_0^\infty (\mathbf{x}^T \mathbf{Q} \mathbf{x} + u^2) dt, \quad \mathbf{Q} = \mathbf{Q}^T \geq 0.$$

We label the above problem as (\mathbf{A}, \mathbf{b}) . Consider now the dynamical system model

$$\dot{\tilde{\mathbf{x}}} = \tilde{\mathbf{A}}\tilde{\mathbf{x}} + \tilde{\mathbf{b}}\tilde{u},$$

where $\mathbf{b} = \tilde{\mathbf{b}}$,

$$\tilde{\mathbf{A}} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & & & \cdots & & \vdots \\ 0 & 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & 0 & \cdots & 0 & 0 \end{bmatrix},$$

and the associated performance index has the form

$$\tilde{J} = \frac{1}{2} \int_0^\infty (\tilde{\mathbf{x}}^T \tilde{\mathbf{Q}} \tilde{\mathbf{x}} + \tilde{u}^2) dt,$$

where $\tilde{\mathbf{Q}} = \mathbf{Q} + \mathbf{D}\tilde{\mathbf{A}} + \tilde{\mathbf{A}}^T \mathbf{D} + \mathbf{D}\mathbf{b}\mathbf{b}^T \mathbf{D}$, and

$$\mathbf{D} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & -a_0 \\ 0 & 0 & 1 & \cdots & 0 & -a_1 \\ \vdots & & & \cdots & & \vdots \\ 0 & 0 & 0 & \cdots & 0 & -a_{n-2} \\ -a_0 & -a_1 & -a_2 & \cdots & -a_{n-2} & -a_{n-1} \end{bmatrix}.$$

We refer to the second problem as (\tilde{A}, \tilde{b}) .

- (a) Show that the symmetric matrix P is a solution to the (A, b) problem if, and only if, the symmetric matrix $\tilde{P} = P - D$ is a solution to the (\tilde{A}, \tilde{b}) problem.
- (b) Show that $A - bb^T P = \tilde{A} - \tilde{b}\tilde{b}^T \tilde{P}$.

Hint Note that $\tilde{A} = A - bb^T D$ and $\tilde{b} = b$. Write down the algebraic Riccati equation for the (\tilde{A}, \tilde{b}) problem and perform appropriate substitutions.

- 5.21** Consider the following variation of the brachistochrone problem: A particle of mass m is to slide down a frictionless track connecting two given points that are not located on the vertical line. The particle's initial speed is v_0 . The motion of the particle is constrained only by gravity and the shape of the track. Find the curve of the track that minimizes the time travel between the two given points. Solve this brachistochrone problem using the minimum principle of Pontryagin. Illustrate your derivations with a numerical example and simulations using MATLAB.
- 5.22** (a) Consider the ground vehicle model described in Subsection 1.7.2. Let the input signal be $u = \delta_f$ (the front wheel steer angle) and $\delta_r = 0$. Write a MATLAB script that animates the lateral motion Y versus the longitudinal motion X of the vehicle, the steer angle δ_f , and the heading angle ψ . Use the steer angle shown in Figure 5.32.
- (b) Use the linear model and the data from Table 1.1 to design a linear quadratic regulator (LQR). In this part, the control has the form

$$\delta_f = -kx + Y_{ref},$$

where Y_{ref} is the command input shown in Figure 5.33. That is, the closed-loop

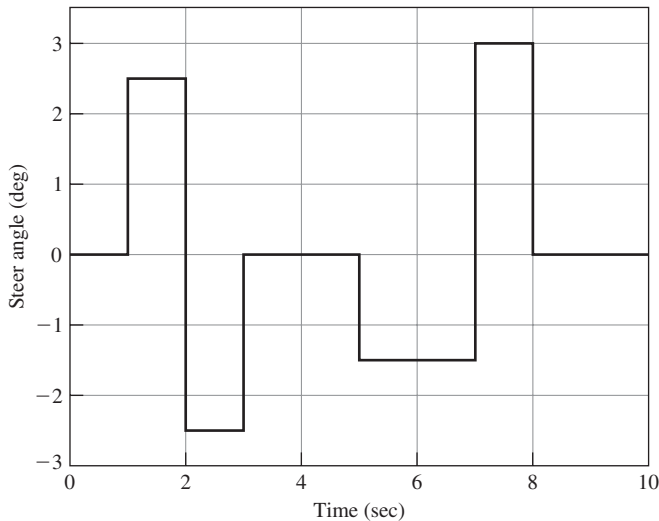


Figure 5.32 Plot of the steer angle δ_f versus time for Exercise 5.22.

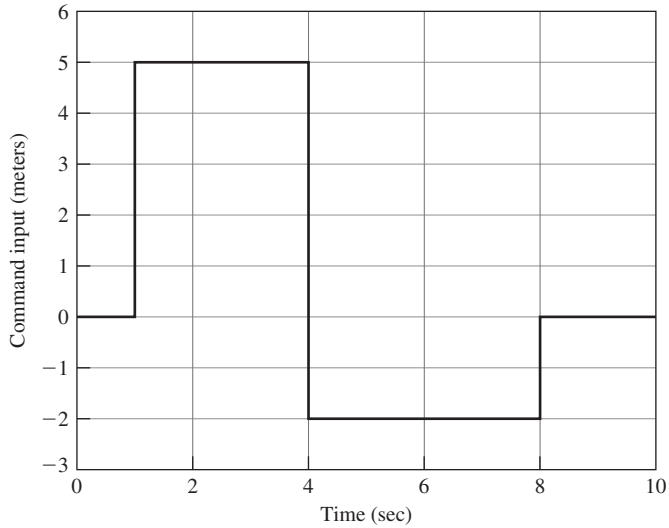


Figure 5.33 Plot of the command input Y_{ref} versus time for Exercise 5.22.

system here will have the form

$$\dot{\mathbf{x}} = (\mathbf{A} - \mathbf{b}\mathbf{k})\mathbf{x} + \mathbf{b}Y_{ref}.$$

You as a designer will select the weights \mathbf{Q} and r . See how the linear vehicle tracks the command input. Observe the sign difference in the command input and the heading angle. Your LQR gain should be calculated for each chosen v_x . In other words, if in your interactive simulation a specific value of v_x is chosen, then your program should calculate the corresponding gain for this particular value of v_x . Implement the LQR state feedback employing the vehicle model that uses the nonlinear cornering characteristic shown in Figure 1.15 and test the performance of the closed-loop system utilizing the command input shown in Figure 5.33. You may need to impose constraints on δ_f —for example, $|\delta_f| \leq 20^\circ$. (Do not forget about converting degrees into radians.)



CHAPTER 6

Sliding Modes

As philosopher of science Karl Popper has emphasized, a good theory is characterized by the fact that it makes a number of predictions that could in principle be disproved or falsified by observation. Each time new experiments are observed to agree with the predictions the theory survives, and our confidence in it is increased; but if ever a new observation is found to disagree, we have to abandon or modify the theory. At least that is what is supposed to happen, but you can always question the competence of the person who carried out the observation.

—Stephen Hawking, *A Brief History of Time* [117]

6.1 Simple Variable Structure Systems

The purpose of this section is to informally introduce the reader to variable structure sliding mode control systems. A formal analysis is presented in the following sections.

A *variable structure system* is a dynamical system whose structure changes in accordance with the current value of its state. A variable structure system can be viewed as a system composed of independent structures together with a switching logic between each of the structures. With appropriate switching logic, a variable structure system can exploit the desirable properties of each of the structures the system is composed of. Even more, a variable structure system may have a property that is not a property of any of its structures. We illustrate the above ideas with two numerical examples.

◆ Example 6.1

This example was adapted from Utkin [284]. We consider a dynamical system model

$$\begin{aligned}\dot{x}_1 &= x_2, \\ \dot{x}_2 &= -ux_1,\end{aligned}\tag{6.1}$$

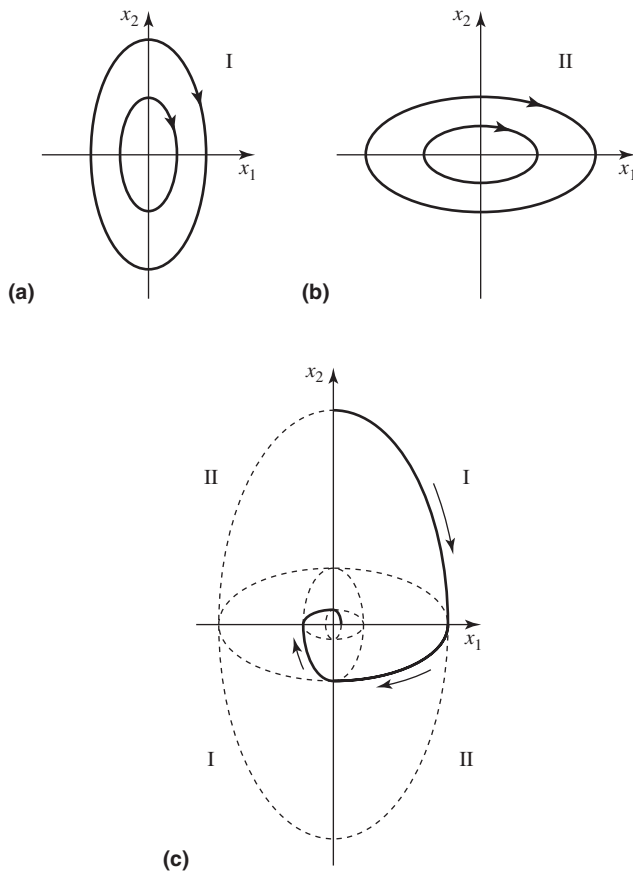


Figure 6.1 (a) and (b) Phase-plane portraits of structures making up a variable structure system. (c) Phase-plane portrait of the variable structure system itself, where the structure of the system is changed any time the system's trajectory crosses either of the coordinate axis. (Adapted from the survey paper by Utkin [284].)

having two structures corresponding to $u = 1/a$ and $u = a$, where a is a positive constant greater than 1; in our simulations we take $a = 5$. The phase-plane portraits of the structures are shown in Figures 6.1(a) and 6.1(b). The phase-plane portraits of the individual structures are families of ellipses. Neither of the structures is asymptotically stable—they are only stable. However, by choosing a suitable switching logic between the structures, we can make the resulting variable structure system asymptotically stable. Indeed, suppose the structure of the system is changed any time the system's trajectory crosses either of the coordinate axes of the state plane, that is,

$$a = \begin{cases} 1/5 & \text{if } x_1 x_2 < 0, \\ 5 & \text{if } x_1 x_2 > 0. \end{cases} \quad (6.2)$$

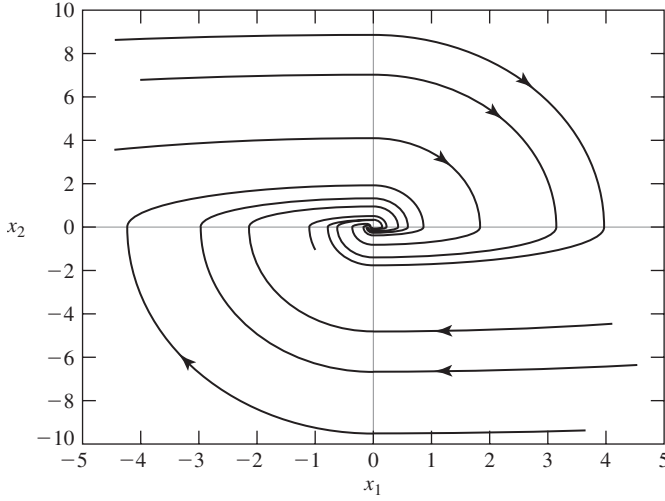


Figure 6.2 Phase-plane portrait of the variable structure system described by (6.1) and (6.2).

A phase plane portrait of the system (6.1) with the switching logic specified by (6.2) is shown in Figure 6.2. We add that the trajectories spiral clockwise towards the origin.

◆ Example 6.2

Consider the double integrator

$$\ddot{x} = u,$$

whose state-space model is

$$\begin{aligned}\dot{x}_1 &= x_2, \\ \dot{x}_2 &= u.\end{aligned}\tag{6.3}$$

Let $u = \pm 2$. Typical trajectories of the above system are shown in Figure 6.3. We thus have two independent structures, neither of them being stable. Let us now use the following switching logic between each of the structures:

$$u = -2 \operatorname{sign}(x_1 + x_2).\tag{6.4}$$

With the switching logic in (6.4), the closed-loop system changes its structure anytime the state trajectory crosses the line

$$x_1 + x_2 = 0.$$

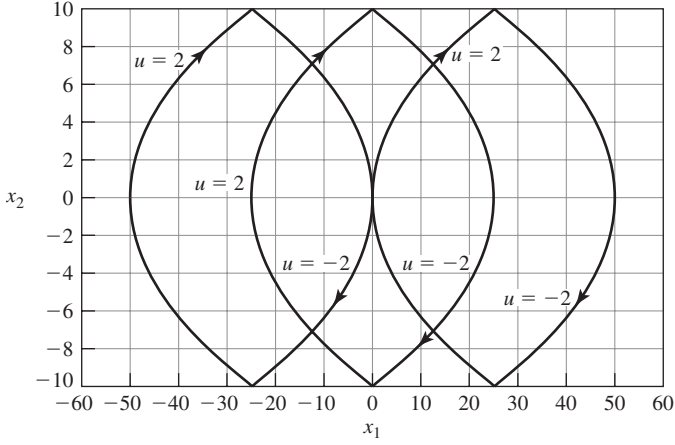


Figure 6.3 Typical trajectories of the double integrator system.

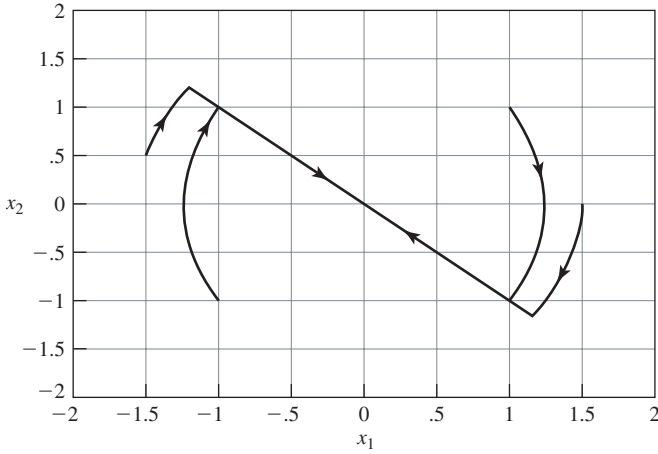


Figure 6.4 Typical trajectories of the closed-system (6.3), (6.4) for initial conditions “close” to the origin of the state-plane. The applied control law is given by (6.4).

We introduce the notation

$$\sigma(\mathbf{x}) = \sigma(x_1, x_2) = x_1 + x_2.$$

Then, we can describe the switching logic action as

$$u = \begin{cases} -2 & \text{if } \sigma(\mathbf{x}) > 0, \\ 2 & \text{if } \sigma(\mathbf{x}) < 0. \end{cases} \quad (6.5)$$

Typical trajectories of the closed-loop system described by (6.3) and (6.4) for initial conditions located “close” to the origin of the state-plane are shown in Figure 6.4. Close to the origin, once the trajectories intercept the switching line, they are

confined to it thereafter. If we increase the integration step size, then we can observe oscillations, or *chattering*, of the trajectory about the switching surface as it is explained in Figure 6.5. A plot of two trajectories with an increased integration time is shown in Figure 6.6.

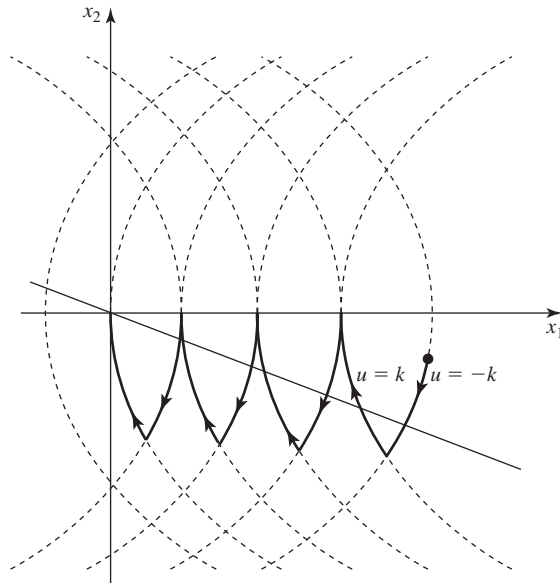


Figure 6.5 Chattering in a variable structure sliding mode system—a magnified view.

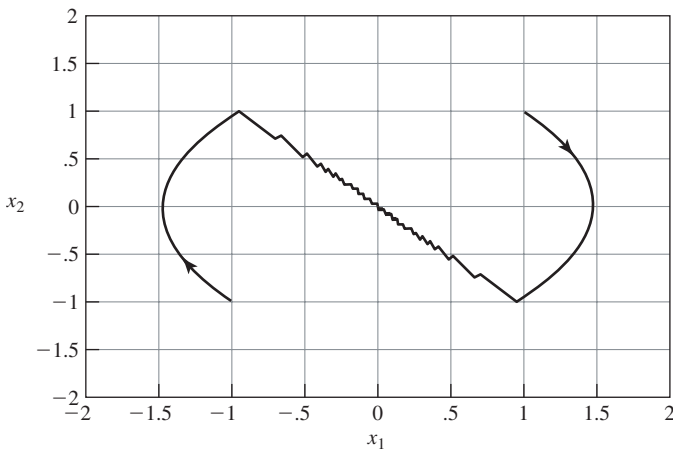


Figure 6.6 An illustration of chattering in a variable structure sliding mode system consisting of the double integrator and the control law, $u = -\text{sign}(x_1 + x_2)$. Two trajectories starting from two different initial conditions are plotted.

6.2 Sliding Mode Definition

The motion of the system while confined to the switching line or a surface is referred to as *sliding*. A sliding mode will exist if in the vicinity of the switching surface the state velocity vectors are directed toward the surface. In such a case, the switching surface attracts trajectories when they are in its vicinity; and once a trajectory intersects the switching surface, it will stay on it thereafter. A surface $\sigma(\mathbf{x}) = 0$ is attractive if

1. any trajectory starting on the surface remains there, and
2. any trajectory starting outside the surface tends to it at least asymptotically.

Thus, for a sliding motion to occur we need

$$\lim_{\sigma \rightarrow 0^+} \dot{\sigma} < 0 \quad \text{and} \quad \lim_{\sigma \rightarrow 0^-} \dot{\sigma} > 0.$$

The above conditions are illustrated in Figure 6.7. They ensure that the motion of the state trajectory \mathbf{x} on either side of the switching surface $\sigma(\mathbf{x}) = 0$ is toward the switching surface. The two conditions may be combined to give

$$\sigma \dot{\sigma} < 0 \tag{6.6}$$

in the neighborhood of the switching surface. Note, however, that the closed-loop system is now modeled by differential equations with discontinuous right-hand sides. The modeling differential equations are discontinuous on the switching surface $\sigma(\mathbf{x}) = 0$. This is the reason we often refer to the switching surface as the *discontinuity surface*. A rigorous analysis of such systems cannot be carried out using methods of classical theory of differential equations. A solution to the differential equation $\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x})$ with a continuous right-hand side is a function $\mathbf{x}(t)$ that has a derivative and satisfies the differential equation everywhere on a given interval. In the case when the right-hand side of $\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x})$ is discontinuous, we have to use more general methods such as one, for example, developed by Filippov [84]. We illustrate these points with a numerical example taken from Filippov [84, p. 1].

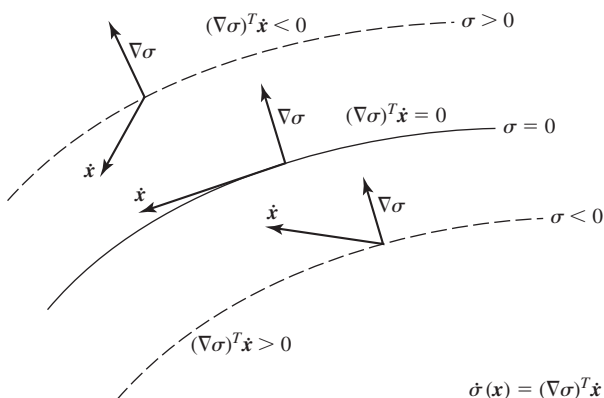


Figure 6.7 An illustration of the sufficient condition of the attractivity to the switching surface.

◆ Example 6.3

Consider the following differential equation with the discontinuous right-hand side:

$$\dot{x} = 1 - 2 \operatorname{sign}(x). \quad (6.7)$$

For $x < 0$, we have $\dot{x} = 3$, while for $x > 0$, $\dot{x} = -1$. In Figure 6.8, plots of x versus time for different initial conditions are given. As t increases, the solution x of (6.7) reaches the line $x = 0$ and stays there. The function $x(t)$, however, does not satisfy (6.7) in the usual sense because if $\dot{x}(t) = 0$ and $x(t) = 0$, the right-hand side of (6.7) becomes

$$1 - 2 \operatorname{sign}(0) = 1 \neq 0,$$

where we defined $\operatorname{sign}(0) = 0$.

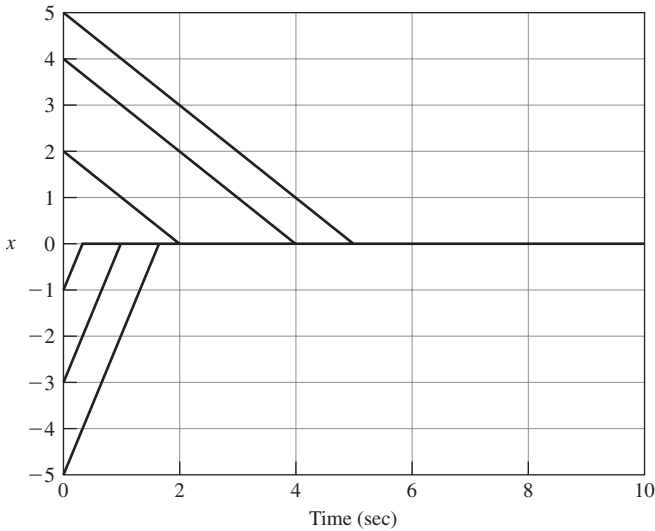


Figure 6.8 Plots of x versus time for different initial conditions for the differential equation (6.7).

We now briefly describe Filippov's approach. Filippov [84, pp. 2–3] recommends that the generalized solution to an ordinary differential equation with a discontinuous right-hand side should meet the following criteria:

1. For a differential equation with a continuous right-hand side, the generalized solution should be equivalent to the usual one.
2. For the equation $\dot{x} = f(t)$, the generalized solution should be of the form

$$x(t) = \int f(t) dt + \text{constant}.$$

3. For any initial condition $x(t_0) = x_0$, the solution must exist at least for $t > t_0$ and should be continuable.
4. The generalized solution should be applicable to a broad class of equations that model dynamical systems.
5. The limit of an uniformly convergent sequence of solutions should be a solution.
6. After applying to a solution a common change of variables, the result of the transformation should be a solution.

We next describe the proposed generalized solution in the sense of Filippov. We begin with an example. Suppose that we are given a single-input dynamical system model, $\dot{x} = f(x, u)$ along with a switching surface $\sigma(x) = 0$. Let $f^- = f(x, u^-)$ and $f^+ = f(x, u^+)$, where $u = u^-$ if $\sigma < 0$ and $u = u^+$ if $\sigma > 0$. We next join the velocity vectors f^- and f^+ . The resulting velocity vector f^0 is obtained by drawing the line tangent to the switching surface that intersects the line segment joining the velocity vectors f^- and f^+ . Thus,

$$f^0 = \alpha f^- + (1 - \alpha) f^+,$$

where $\alpha \in [0, 1]$. The above construction is illustrated in Figure 6.9. Observe that f^0 belongs to the smallest convex set containing f^- and f^+ . Recall that a set S is said to be convex if for every $x, y \in S$ and every real number α , $0 < \alpha < 1$, we have $\alpha x + (1 - \alpha)y \in S$.

Following Itkis [140], we give a simple interpretation of the above construction. The number α may be interpreted as the time the system's trajectory remains below the switching surface, while $(1 - \alpha)$ may be interpreted as the time the trajectory remains above the switching surface during the trajectory oscillation period about the switching surface. Specifically, if the trajectory performs one oscillation about $\sigma = 0$ in a "small" time interval $\Delta\tau$, then it will traverse a distance $\|f^-\| \alpha \Delta\tau$ in time $\alpha \Delta\tau$ with velocity f^- . During the remaining time of oscillation the trajectory will traverse a distance $\|f^+\| (1 - \alpha) \Delta\tau$ in time $(1 - \alpha) \Delta\tau$ with velocity f^+ . The average displacement of the trajectory during the time $\Delta\tau$ will then be $\|\alpha f^- + (1 - \alpha) f^+\| \Delta\tau$ with the average velocity f^0 .

Suppose now that there are m discontinuity surfaces of the form

$$\sigma_i = \{x : \sigma_i(x) = 0\}, \quad i = 1, 2, \dots, m.$$

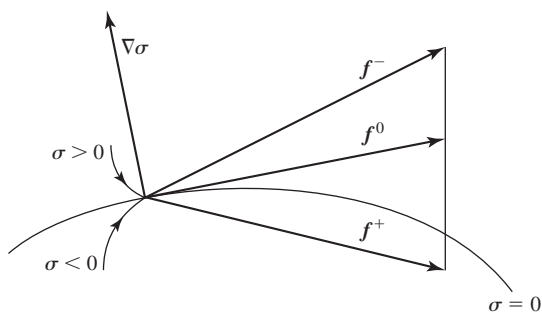


Figure 6.9 Constructing a solution in the sense of Filippov in \mathbb{R}^2 .

Let

$$H = \bigcup_{i=1}^m \sigma_i.$$

Then, we take a solution of the differential equation $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x})$ with the discontinuous right-hand side to be a solution of the differential inclusion

$$\dot{\mathbf{x}}(t) \in G(\mathbf{x}(t)),$$

where for each \mathbf{x} the set $G(\mathbf{x})$ is the smallest closed convex set containing the cluster values, limits, of the function $\mathbf{f}(\mathbf{y})$ as $\mathbf{y} \rightarrow \mathbf{x}$, $\mathbf{y} \notin H$; that is, a solution of the dynamical system $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x})$ is an absolutely continuous function $\mathbf{x}(t)$, in the sense of Definition A.2, defined on an interval, or segment, L for which $\dot{\mathbf{x}}(t) \in G(\mathbf{x}(t))$ almost everywhere on L (Filippov [84]). Observe that if the function \mathbf{f} is continuous at a point \mathbf{x} —that is, $\mathbf{x} \notin H$ —then $G(\mathbf{x})$ consists of a single point, namely $\mathbf{f}(\mathbf{x})$, and the solution satisfies $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x})$ in the usual sense. However, if $\mathbf{x} \in H$, then \mathbf{x} lies on one or more of the surfaces $\sigma_1, \sigma_2, \dots, \sigma_m$. Suppose that \mathbf{x} lies on the surfaces $\sigma_{i_1}, \sigma_{i_2}, \dots, \sigma_{i_k}$. It then follows that for a sufficiently small $\delta > 0$ the surfaces $\sigma_{i_1}, \sigma_{i_2}, \dots, \sigma_{i_k}$ partition the δ -neighborhood of \mathbf{x} into $\ell \leq 2^k$ regions $\mathcal{R}_1, \mathcal{R}_2, \dots, \mathcal{R}_\ell$, in each of which the function \mathbf{f} is continuous. Let \mathbf{f}_j denote the function \mathbf{f} restricted to the region \mathcal{R}_j . Then, $G(\mathbf{x})$ is the smallest closed convex set containing the set of points

$$\left\{ \lim_{\substack{\tilde{\mathbf{x}} \rightarrow \mathbf{x} \\ \tilde{\mathbf{x}} \in \mathcal{R}_j}} \mathbf{f}_j(\tilde{\mathbf{x}}) : j = 1, 2, \dots, \ell \right\}. \quad (6.8)$$

The velocity vector, $\dot{\mathbf{x}} = \mathbf{f}^0$, lies then on the intersection of the tangent plane to the surfaces of discontinuity and the convex set described by (6.8). This is illustrated in Figures 6.9 and 6.10.

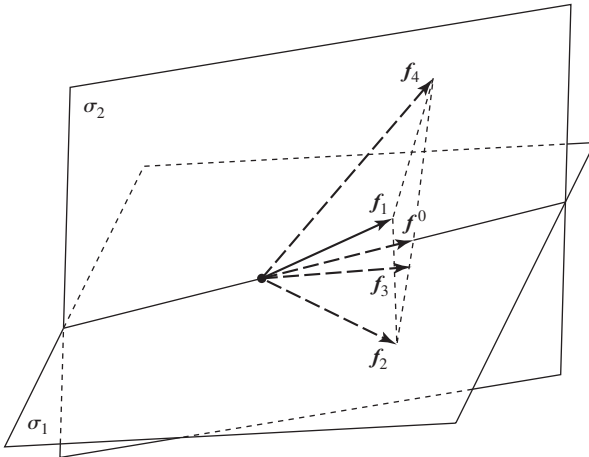


Figure 6.10 Constructing a solution in the sense of Filippov in \mathbb{R}^3 .

6.3 A Simple Sliding Mode Controller

In this section we discuss the problem of designing sliding mode controllers for single-input plants. The design of a variable structure sliding mode controller consists of two phases:

- Sliding (switching) surface design so as to achieve the desired system behavior, like the stability to the origin, when restricted to the surface.
- Selecting feedback gains of the controller so that the closed-system is stable to the sliding surface.

There are a number of benefits when the system is in sliding. In particular:

- While in sliding, the system is not affected by matched uncertainties.
- While in sliding, the system behavior is governed by a reduced set of differential equations. This is what we refer to as the *order reduction*. It is very useful in designing variable structure sliding mode controllers.

We now discuss a sliding surface design phase for linear single-input plants. Without loss of generality, we assume that the plant model is already in its controller companion form,

$$\dot{\mathbf{x}} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & \cdots & 1 \\ -a_0 & -a_1 & -a_2 & \cdots & -a_{n-1} \end{bmatrix} \mathbf{x} + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} u.$$

Let the sliding surface be

$$\begin{aligned} \sigma(\mathbf{x}) &= \mathbf{s}\mathbf{x} \\ &= [s_1 \quad s_2 \quad \cdots \quad s_{n-1} \quad 1] \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\ &= s_1 x_1 + s_2 x_2 + \cdots + s_{n-1} x_{n-1} + x_n \\ &= 0. \end{aligned}$$

Note that, without loss of generality, we assumed that $s_n = 1$. If this was not the case, we could divide both sides of $\mathbf{s}\mathbf{x} = 0$ by $s_n \neq 0$ to ensure that $s_n = 1$. Using the above equation of sliding surface, we can write

$$x_n = -s_1 x_1 - s_2 x_2 - \cdots - s_{n-1} x_{n-1}.$$

We now combine the equations of the plant model and that of the sliding surface. First note that

$$\dot{x}_{n-1} = x_n.$$

Hence, we can write

$$\begin{aligned}\dot{x}_1 &= x_2, \\ \dot{x}_2 &= x_3, \\ &\vdots \\ \dot{x}_{n-1} &= -s_1x_1 - s_2x_2 - \cdots - s_{n-1}x_{n-1},\end{aligned}\tag{6.9}$$

or, equivalently in matrix form,

$$\dot{\mathbf{z}} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & \cdots & 1 \\ -s_1 & -s_2 & -s_3 & \cdots & -s_{n-1} \end{bmatrix} \mathbf{z},\tag{6.10}$$

where $\mathbf{z} = [x_1 \ x_2 \ \cdots \ x_{n-1}]^T$. The above equation describes the system dynamics in sliding. Observe the order reduction of the system dynamics in sliding. We now have $n - m = n - 1$ first-order differential equations. The poles of the reduced-order system are the roots of the equation

$$\lambda^{n-1} + s_{n-1}\lambda^{n-2} + \cdots + s_2\lambda + s_1 = 0.$$

The response of the system in sliding is completely specified by an appropriate choice of the components s_1, s_2, \dots, s_{n-1} of the switching surface.

Consider again Example 6.2. The sliding surface for this example can be represented as

$$\sigma(\mathbf{x}) = [s_1 \quad 1]\mathbf{x} = 0,$$

or, equivalently,

$$s_1x_1 + x_2 = 0.$$

The dynamics of the system in sliding are described by

$$\dot{z} = -s_1z,$$

which is just a first-order system. Note that to ensure the reduced-system stability, we select $s_1 > 0$.

After designing a sliding surface, we construct a feedback controller. In this section we discuss a controller design for linear single-input plants. The controller objective is to drive the plant state to the sliding surface and maintain it on the surface for all subsequent time. We use a generalized Lyapunov approach in constructing the controller. Specifically, we use a distance, $V = 0.5\sigma^2$, from the sliding surface σ as a Lyapunov function candidate. Then, we select the controllers' gains so that the time derivative of the chosen Lyapunov function candidate evaluated on the solution of the controlled system is negative definite with respect to the switching surface, thus ensuring the motion of the state trajectory to the surface as it is illustrated in Figure 6.7. Our goal is to find $u = u(\mathbf{x})$ so that

$$\frac{d}{dt} \left(\frac{1}{2}\sigma^2 \right) = \sigma \dot{\sigma} < 0.$$

We consider the controller structure of the form

$$u = k_1^\pm x_1 + k_2^\pm x_2 + \cdots + k_n^\pm x_n, \quad (6.11)$$

where k_i^\pm are the feedback gains to be determined so that the condition $\sigma \dot{\sigma} < 0$ is satisfied. To determine the feedback gains, we substitute u , given by (6.11), into the expression $\sigma \dot{\sigma}$ and utilize (6.9) to get

$$\begin{aligned} \sigma \dot{\sigma} &= \sigma [s_1 \quad s_2 \quad \cdots \quad 1] \dot{\mathbf{x}} \\ &= \sigma (s_1 \dot{x}_1 + \cdots + s_{n-1} \dot{x}_{n-1} + \dot{x}_n) \\ &= \sigma (s_1 x_2 + \cdots + s_{n-1} x_n + \dot{x}_n) \\ &= \sigma (s_1 x_2 + \cdots + s_{n-1} x_n - a_0 x_1 - \cdots - a_{n-1} x_n + u) \\ &= \sigma (s_1 x_2 + \cdots + s_{n-1} x_n - a_0 x_1 - \cdots - a_{n-1} x_n + k_1^\pm x_1 + \cdots + k_n^\pm x_n). \end{aligned}$$

Rearranging the terms yields

$$\sigma \dot{\sigma} = \sigma x_1 (-a_0 + k_1^\pm) + \sigma x_2 (-a_1 + s_1 + k_2^\pm) + \cdots + \sigma x_n (-a_{n-1} + s_{n-1} + k_n^\pm).$$

A possible way to ensure that $\sigma \dot{\sigma} < 0$ is to select the gains to force each term in $\sigma \dot{\sigma}$ to be negative, that is,

$$\begin{aligned} \sigma x_1 (-a_0 + k_1^\pm) &< 0, \\ \sigma x_2 (-a_1 + s_1 + k_2^\pm) &< 0, \\ &\vdots \\ \sigma x_n (-a_{n-1} + s_{n-1} + k_n^\pm) &< 0. \end{aligned}$$

A possible choice of gains that leads to satisfaction of the attractivity to the surface condition, $\sigma \dot{\sigma} < 0$, is

$$\begin{aligned} k_1^+ &< -|a_0| & \text{if } \sigma x_1 > 0, \\ k_1^- &> |a_0| & \text{if } \sigma x_1 < 0, \\ &\vdots \\ k_n^+ &< -(|a_{n-1}| + |s_{n-1}|) & \text{if } \sigma x_n > 0, \\ k_n^- &> |a_{n-1}| + |s_{n-1}| & \text{if } \sigma x_n < 0. \end{aligned} \quad (6.12)$$

We can also use the following controller gains:

$$\begin{aligned} k_1^\pm &= -(|a_0| + \varepsilon) \text{sign}(x_1 \sigma), \\ &\vdots \\ k_n^\pm &= -(|a_{n-1}| + |s_{n-1}| + \varepsilon) \text{sign}(x_n \sigma), \end{aligned} \quad (6.13)$$

where $\varepsilon > 0$ is a design parameter. Note that the controllers' gains given by (6.13) satisfy conditions (6.12). The larger the value of ε the faster the trajectory converges to the sliding surface.

6.4 Sliding in Multi-Input Systems

In this section we describe a sliding mode in linear multi-input plants. When the number of the plant inputs is m , we can use sliding surface of the form

$$\{\mathbf{x} : \mathbf{S}\mathbf{x} = \mathbf{0}, \mathbf{S} \in \mathbb{R}^{m \times n}\}.$$

We also use $\mathbf{S}\mathbf{x} = \mathbf{0}$ to denote the above sliding surface. In many applications, nonlinear sliding surfaces arise. Here, we use linear sliding surfaces only, although one can use nonlinear sliding surfaces as well. Let $\boldsymbol{\sigma}(\mathbf{x}) = \mathbf{S}\mathbf{x}$. Then,

$$\boldsymbol{\sigma}(\mathbf{x}) = \begin{bmatrix} \sigma_1(\mathbf{x}) \\ \sigma_2(\mathbf{x}) \\ \vdots \\ \sigma_m(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} s_1\mathbf{x} \\ s_2\mathbf{x} \\ \vdots \\ s_m\mathbf{x} \end{bmatrix},$$

where $s_i \in \mathbb{R}^{1 \times n}$. Recall, that the system is in sliding mode if $\boldsymbol{\sigma}(\mathbf{x}(t)) = \mathbf{0}$ for $t \geq t_0$, where \mathbf{x} is the state trajectory and t_0 is a specified time. It follows that in a sliding mode the velocity vector $\dot{\mathbf{x}}(t)$ is tangent to the switching surface. Equivalently,

$$(\nabla \sigma_i)^T \dot{\mathbf{x}}(t) = s_i \dot{\mathbf{x}}(t) = 0 \quad \text{for } i = 1, 2, \dots, m,$$

that is,

$$\dot{\boldsymbol{\sigma}}(\mathbf{x}(t)) = \mathbf{S}\dot{\mathbf{x}}(t) = \mathbf{0}.$$

We can thus characterize the system in sliding mode as

$$\begin{aligned} \boldsymbol{\sigma}(\mathbf{x}(t)) &= \mathbf{0}, \\ \dot{\boldsymbol{\sigma}}(\mathbf{x}(t)) &= \mathbf{0}. \end{aligned} \tag{6.14}$$

Consider now the plant model $\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$. Combining $\mathbf{S}\dot{\mathbf{x}}(t) = \mathbf{0}$ and the modeling equation of the plant, we get

$$\mathbf{S}\mathbf{A}\mathbf{x} + \mathbf{S}\mathbf{B}\mathbf{u} = \mathbf{0}.$$

If the matrix $\mathbf{S}\mathbf{B}$ is nonsingular, then we can solve the above algebraic equation for \mathbf{u} . The result is called the *equivalent control*, and it is denoted using the symbol \mathbf{u}_{eq} . We have

$$\mathbf{u}_{eq} = -(\mathbf{S}\mathbf{B})^{-1}\mathbf{S}\mathbf{A}\mathbf{x}.$$

Substituting \mathbf{u}_{eq} into the plant model equation yields

$$\dot{\mathbf{x}} = (\mathbf{I}_n - \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\mathbf{S})\mathbf{A}\mathbf{x}. \tag{6.15}$$

The above system is sometimes referred to as the *equivalent system*. Thus, we can describe the behavior of the system in sliding alternatively as

$$\begin{aligned} \dot{\mathbf{x}} &= (\mathbf{I}_n - \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\mathbf{S})\mathbf{A}\mathbf{x}, \\ \mathbf{S}\mathbf{x} &= \mathbf{0}. \end{aligned} \tag{6.16}$$

Note that in sliding, the system is governed by a reduced set of differential equations that are obtained by combining the above two sets of equations.

6.5 Sliding Mode and System Zeros

In this section we show that the poles of the reduced-order system in sliding are the so-called *system zeros* of the dynamical system described by

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}, \\ \boldsymbol{\sigma} &= \mathbf{S}\mathbf{x},\end{aligned}\tag{6.17}$$

where $\mathbf{S} \in \mathbb{R}^{m \times n}$ and we assume that $\det \mathbf{S}\mathbf{B} \neq 0$. The system zeros are defined using the *system matrix*. The system matrix, $\mathbf{P}(s)$, corresponding to the plant (6.17) has the form

$$\mathbf{P}(s) = \begin{bmatrix} s\mathbf{I}_n - \mathbf{A} & -\mathbf{B} \\ \mathbf{S} & \mathbf{O} \end{bmatrix}.\tag{6.18}$$

The system matrix $\mathbf{P}(s)$ corresponding to the plant (6.17) is a square matrix of order $n + m$. The zeros of the determinant of the system matrix (6.18) are the system zeros of (6.17). The system zeros are invariant under the following set of transformations:

1. Nonsingular state variable transformation
2. Nonsingular transformation of the inputs
3. Nonsingular transformation of the outputs
4. State feedback to the inputs
5. Output feedback to the rates of change of the states

We now evaluate the determinant of the system matrix (6.18). To facilitate this task, we premultiply the system matrix by the matrix

$$\begin{bmatrix} \mathbf{I}_n & \mathbf{O} \\ -\mathbf{S}(s\mathbf{I}_n - \mathbf{A})^{-1} & \mathbf{I}_m \end{bmatrix}.$$

The determinant of the above matrix is unity and so the determinant of the product is the same as the determinant of the system matrix,

$$\begin{aligned}\det \mathbf{P}(s) &= \det \left(\begin{bmatrix} \mathbf{I}_n & \mathbf{O} \\ -\mathbf{S}(s\mathbf{I}_n - \mathbf{A})^{-1} & \mathbf{I}_m \end{bmatrix} \begin{bmatrix} s\mathbf{I}_n - \mathbf{A} & -\mathbf{B} \\ \mathbf{S} & \mathbf{O} \end{bmatrix} \right) \\ &= \det \begin{bmatrix} s\mathbf{I}_n - \mathbf{A} & -\mathbf{B} \\ \mathbf{O} & \mathbf{S}(s\mathbf{I}_n - \mathbf{A})^{-1}\mathbf{B} \end{bmatrix} \\ &= \det(s\mathbf{I}_n - \mathbf{A}) \det(\mathbf{S}(s\mathbf{I}_n - \mathbf{A})^{-1}\mathbf{B}).\end{aligned}\tag{6.19}$$

We next evaluate the characteristic polynomial of the equivalent system (6.15). While performing manipulations, we use the result of Exercise A.3 to obtain

$$\begin{aligned}\det(s\mathbf{I}_n - \mathbf{A} + \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\mathbf{S}\mathbf{A}) \\ &= \det(s\mathbf{I}_n - \mathbf{A}) \det(\mathbf{I}_n + (s\mathbf{I}_n - \mathbf{A})^{-1}\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\mathbf{S}\mathbf{A}) \\ &= \det(s\mathbf{I}_n - \mathbf{A}) \det(\mathbf{I}_m + \mathbf{S}\mathbf{A}(s\mathbf{I}_n - \mathbf{A})^{-1}\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}).\end{aligned}\tag{6.20}$$

In our further manipulations, we use the identity

$$\mathbf{A}(\text{Adj}(s\mathbf{I}_n - \mathbf{A})) = s \text{Adj}(s\mathbf{I}_n - \mathbf{A}) - \det(s\mathbf{I}_n - \mathbf{A})\mathbf{I}_n.\tag{6.21}$$

The above follows from the fact that

$$(s\mathbf{I}_n - \mathbf{A})(s\mathbf{I}_n - \mathbf{A})^{-1} = \mathbf{I}_n,$$

which is equivalent to

$$(s\mathbf{I}_n - \mathbf{A})\text{Adj}(s\mathbf{I}_n - \mathbf{A}) = \det(s\mathbf{I}_n - \mathbf{A})\mathbf{I}_n. \quad (6.22)$$

We use (6.21) to get

$$\begin{aligned} & \det(\mathbf{I}_m + \mathbf{S}\mathbf{A}(s\mathbf{I}_n - \mathbf{A})^{-1}\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}) \\ &= \det\left(\frac{\det(s\mathbf{I}_n - \mathbf{A})\mathbf{I}_m + \mathbf{S}(s\text{Adj}(s\mathbf{I}_n - \mathbf{A}) - \det(s\mathbf{I}_n - \mathbf{A})\mathbf{I}_n)\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}}{\det(s\mathbf{I}_n - \mathbf{A})}\right) \\ &= \det\left(\frac{s\mathbf{S}\text{Adj}(s\mathbf{I}_n - \mathbf{A})\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}}{\det(s\mathbf{I}_n - \mathbf{A})}\right) \\ &= s^m \det(\mathbf{S}(s\mathbf{I}_n - \mathbf{A})^{-1}\mathbf{B}) \det(\mathbf{S}\mathbf{B})^{-1}. \end{aligned}$$

Using (6.19) and the above expression, we represent (6.20) as

$$\begin{aligned} & \det(s\mathbf{I}_n - \mathbf{A} + \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\mathbf{S}\mathbf{A}) \\ &= s^m \det(\mathbf{S}\mathbf{B})^{-1} \det(s\mathbf{I}_n - \mathbf{A}) \det(\mathbf{S}(s\mathbf{I}_n - \mathbf{A})^{-1}\mathbf{B}), \end{aligned}$$

and hence

$$\boxed{\det(s\mathbf{I}_n - \mathbf{A} + \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\mathbf{S}\mathbf{A}) = \det(\mathbf{S}\mathbf{B})^{-1} s^m \det \mathbf{P}(s)} \quad (6.23)$$

We conclude that the dynamics—that is, the poles—of the system modeled by (6.17) in sliding along $\sigma = \mathbf{S}\mathbf{x}$ are determined by its system zeros. This fact leads to the observation that the switching surface design can be viewed as choosing an output matrix \mathbf{S} so that the system (6.17) has a desired set of system zero locations, which in turn govern the dynamics of the system in sliding along $\mathbf{S}\mathbf{x} = \mathbf{0}$.

6.6 Nonideal Sliding Mode

We consider a nonlinear dynamical system modeled by

$$\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x}) + \mathbf{B}(t, \mathbf{x})\mathbf{u}, \quad (6.24)$$

where $\mathbf{x}(t) \in \mathbb{R}^n$, $\mathbf{f}(t, \mathbf{x}(t)) \in \mathbb{R}^n$, $\mathbf{B}(t, \mathbf{x}(t)) \in \mathbb{R}^{n \times m}$, and $\mathbf{u} \in \mathbb{R}^m$. The functions \mathbf{f} and \mathbf{B} are continuous functions of all their arguments. The switching surface is $\{\mathbf{x} : \sigma(\mathbf{x}) = \mathbf{0}\}$, where $\sigma(\mathbf{x}(t)) \in \mathbb{R}^m$. Let

$$\mathbf{S} = \frac{\partial \sigma}{\partial \mathbf{x}},$$

where we assume that $\mathbf{S}\mathbf{B}$ is invertible. Ideally, once the system enters sliding mode, an appropriately constructed controller is switched with an infinite frequency to keep the system on the surface $\{\sigma = \mathbf{0}\}$. Thus, in the ideal case, once the state trajectory hits the switching surface, say at t_0 , then for $t \geq t_0$ we have $\sigma = \mathbf{0}$ and $\dot{\sigma}(\mathbf{x}) = \mathbf{0}$. Solving the algebraic equation

$$\dot{\sigma}(\mathbf{x}) = \mathbf{S}\dot{\mathbf{x}} = \mathbf{S}(\mathbf{f}(t, \mathbf{x}) + \mathbf{B}(t, \mathbf{x})\mathbf{u}) = \mathbf{0}$$

for \mathbf{u} yields

$$\mathbf{u} = -(\mathbf{SB})^{-1}\mathbf{Sf}.$$

Recall that the above expression for \mathbf{u} is called the equivalent control and is denoted \mathbf{u}_{eq} . Substituting \mathbf{u}_{eq} into to the modeling equation (6.24), we obtain

$$\dot{\mathbf{x}} = \mathbf{f} - \mathbf{B}(\mathbf{SB})^{-1}\mathbf{Sf} = (\mathbf{I}_n - \mathbf{B}(\mathbf{SB})^{-1}\mathbf{S})\mathbf{f}.$$

The above equation together with the condition $\boldsymbol{\sigma} = \mathbf{0}$ describes the system in ideal sliding.

We now study a nonideal sliding mode. In nonideal sliding mode we no longer have $\dot{\boldsymbol{\sigma}}(\mathbf{x}) = \mathbf{0}$. Hence, solving the equation

$$\dot{\boldsymbol{\sigma}}(\mathbf{x}) = \mathbf{S}(\mathbf{f}(t, \mathbf{x}) + \mathbf{B}(t, \mathbf{x})\mathbf{u})$$

for \mathbf{u} , which we now refer to as $\tilde{\mathbf{u}}$, yields

$$\tilde{\mathbf{u}} = -(\mathbf{SB})^{-1}\mathbf{Sf} + (\mathbf{SB})^{-1}\dot{\boldsymbol{\sigma}}.$$

Substituting the above into the plant dynamics gives

$$\dot{\tilde{\mathbf{x}}} = \mathbf{f} - \mathbf{B}(\mathbf{SB})^{-1}\mathbf{Sf} + \mathbf{B}(\mathbf{SB})^{-1}\dot{\boldsymbol{\sigma}}, \quad (6.25)$$

where we used the notation $\tilde{\mathbf{x}}$ to indicate the state of the system in nonideal sliding mode. The above describes the plant modeled by (6.24) in nonideal sliding mode. The term $\mathbf{B}(\mathbf{SB})^{-1}\dot{\boldsymbol{\sigma}}$ on the right-hand side of (6.25) can be viewed as representing nonidealities. In real sliding the state trajectory $\tilde{\mathbf{x}}$ will reside in some vicinity of the manifold $\{\tilde{\mathbf{x}} : \boldsymbol{\sigma}(\tilde{\mathbf{x}}) = \mathbf{0}\}$ so that

$$\|\boldsymbol{\sigma}(\tilde{\mathbf{x}})\| \leq \Delta,$$

where $\Delta > 0$. We now show, following Utkin [285, pp. 45–48], that if $\Delta \rightarrow 0$, then the real sliding mode behavior approaches that of the ideal sliding mode behavior.

Theorem 6.1 (Utkin, 1974) Suppose the following:

1. A solution to $\dot{\tilde{\mathbf{x}}} = \mathbf{f}(t, \tilde{\mathbf{x}}) + \mathbf{B}(t, \tilde{\mathbf{x}})\tilde{\mathbf{u}}$, on the interval $[0, T]$, satisfies the condition

$$\|\boldsymbol{\sigma}(\tilde{\mathbf{x}})\| \leq \Delta,$$

where $\tilde{\mathbf{u}}$ may incorporate nonidealities corresponding to $\dot{\boldsymbol{\sigma}} \neq \mathbf{0}$.

2. There is a Lipschitz constant L for the right-hand side of

$$\dot{\tilde{\mathbf{x}}} = \mathbf{f} - \mathbf{B}(\mathbf{SB})^{-1}\mathbf{Sf}.$$

3. Partial derivatives of the function

$$\mathbf{B}(\mathbf{SB})^{-1}$$

exist and are bounded in any bounded domain of \mathbb{R}^n .

4. For the function

$$\mathbf{f}(t, \tilde{\mathbf{x}}) + \mathbf{B}(t, \tilde{\mathbf{x}})\tilde{\mathbf{u}},$$

denoted $\mathbf{f} + \mathbf{B}\tilde{\mathbf{u}}$, there exist nonnegative constants M_1 and M_2 such that

$$\|\mathbf{f} + \mathbf{B}\tilde{\mathbf{u}}\| \leq M_1 + M_2\|\tilde{\mathbf{x}}\|.$$

Then, for any pair of state trajectories that are solutions to the ideal and nonideal sliding models for which the initial conditions satisfy

$$\|\mathbf{x}(0) - \tilde{\mathbf{x}}(0)\| \leq P\Delta, \quad P \geq 0,$$

there exists a nonnegative H such that

$$\|\mathbf{x}(t) - \tilde{\mathbf{x}}(t)\| \leq H\Delta$$

for all $t \in [0, T]$.

Proof The trajectory of the system in an ideal sliding mode satisfies

$$\mathbf{x}(t) = \mathbf{x}(0) + \int_0^t (\mathbf{f} - \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\mathbf{S}\mathbf{f})(\mathbf{x}(\tau))d\tau, \quad (6.26)$$

while the trajectory of the system in nonideal sliding satisfies

$$\tilde{\mathbf{x}}(t) = \tilde{\mathbf{x}}(0) + \int_0^t (\mathbf{f} - \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\mathbf{S}\mathbf{f})(\tilde{\mathbf{x}}(\tau))d\tau + \int_0^t \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\dot{\sigma}(\tilde{\mathbf{x}}(\tau))d\tau.$$

Integrating the last term of the above by parts gives

$$\begin{aligned} \tilde{\mathbf{x}}(t) = & \tilde{\mathbf{x}}(0) + \int_0^t (\mathbf{f} - \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\mathbf{S}\mathbf{f})(\tilde{\mathbf{x}}(\tau))d\tau \\ & + \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\sigma(\tilde{\mathbf{x}}(\tau))\Big|_0^t - \int_0^t \left(\frac{d}{d\tau}\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\right)\sigma(\tilde{\mathbf{x}}(\tau))d\tau. \end{aligned}$$

Subtracting the above from (6.26), taking the norms of both sides, and applying the triangle inequality yields

$$\begin{aligned} \|\mathbf{x}(t) - \tilde{\mathbf{x}}(t)\| \leq & \|\mathbf{x}(0) - \tilde{\mathbf{x}}(0)\| + \|\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\sigma(\tilde{\mathbf{x}}(\tau))\Big|_0^t\| \\ & + \int_0^t \left\| \left(\frac{d}{d\tau}\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\right)\sigma(\tilde{\mathbf{x}}(\tau)) \right\| d\tau \\ & + \int_0^t \|(\mathbf{f} - \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\mathbf{S}\mathbf{f})(\mathbf{x}(\tau)) - (\mathbf{f} - \mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\mathbf{S}\mathbf{f})(\tilde{\mathbf{x}}(\tau))\| d\tau. \end{aligned}$$

Taking into account the hypothesis $\|\mathbf{x}(0) - \tilde{\mathbf{x}}(0)\| \leq P\Delta$ and assumption 2, we obtain

$$\begin{aligned} \|\mathbf{x}(t) - \tilde{\mathbf{x}}(t)\| \leq & P\Delta + \|\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\sigma(\tilde{\mathbf{x}}(\tau))\Big|_0^t\| \\ & + \int_0^t \left\| \left(\frac{d}{d\tau}\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}\right)\sigma(\tilde{\mathbf{x}}(\tau)) \right\| d\tau + \int_0^t L\|\mathbf{x}(\tau) - \tilde{\mathbf{x}}(\tau)\| d\tau. \end{aligned}$$

We next use assumption 1 to get

$$\begin{aligned}\|\mathbf{x}(t) - \tilde{\mathbf{x}}(t)\| &\leq \left(P + \|\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}(\tilde{\mathbf{x}}(t))\| + \|\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}(\tilde{\mathbf{x}}(0))\| \right. \\ &\quad \left. + \int_0^t \left\| \left(\frac{d}{d\tau} \mathbf{B}(\mathbf{S}\mathbf{B})^{-1} \right) (\tilde{\mathbf{x}}(\tau)) \right\| d\tau \right) \Delta \\ &\quad + \int_0^t L \|\mathbf{x}(\tau) - \tilde{\mathbf{x}}(\tau)\| d\tau.\end{aligned}\tag{6.27}$$

The trajectory of the system in nonideal sliding is bounded on the interval $[0, T]$. Indeed, using assumption 4 to the system $\dot{\tilde{\mathbf{x}}} = \mathbf{f}(t, \tilde{\mathbf{x}}) + \mathbf{B}(t, \tilde{\mathbf{x}})\tilde{\mathbf{u}}$, we get

$$\begin{aligned}\|\tilde{\mathbf{x}}(t)\| &\leq \|\tilde{\mathbf{x}}(0)\| + \int_0^t \|\mathbf{f}(\tau, \tilde{\mathbf{x}}) + \mathbf{B}(\tau, \tilde{\mathbf{x}})\tilde{\mathbf{u}}\| d\tau \\ &\leq \|\tilde{\mathbf{x}}(0)\| + \int_0^t (M_1 + M_2 \|\tilde{\mathbf{x}}\|) d\tau \\ &\leq \|\tilde{\mathbf{x}}(0)\| + M_1 T + \int_0^t M_2 \|\tilde{\mathbf{x}}\| d\tau.\end{aligned}$$

Applying now the Bellman–Gronwall lemma (from 663), to the above integral inequality, we obtain

$$\begin{aligned}\|\tilde{\mathbf{x}}(t)\| &\leq (\|\tilde{\mathbf{x}}(0)\| + M_1 T) e^{M_2 t} \\ &\leq (\|\tilde{\mathbf{x}}(0)\| + M_1 T) e^{M_2 T}, \quad t \in [0, T].\end{aligned}$$

The above means that $\tilde{\mathbf{x}}$ is bounded on the interval $[0, T]$ and hence the function $\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}$ is bounded on the same interval. By assumption 3, partial derivatives of the function $\mathbf{B}(\mathbf{S}\mathbf{B})^{-1}$ exist and are bounded. Applying the above to (6.27), we conclude that for some $S \geq 0$ that depends on the initial conditions, time T , and constant P , we have

$$\|\mathbf{x}(t) - \tilde{\mathbf{x}}(t)\| \leq S\Delta + L \int_0^t \|\mathbf{x}(\tau) - \tilde{\mathbf{x}}(\tau)\| d\tau.$$

Application of the Bellman–Gronwall lemma to the above inequality yields

$$\|\mathbf{x}(t) - \tilde{\mathbf{x}}(t)\| \leq S\Delta e^{L T}, \quad t \in [0, T].$$

Let $H = Se^{LT}$. Then, we can write

$$\|\mathbf{x}(t) - \tilde{\mathbf{x}}(t)\| \leq H\Delta,$$

which completes the proof of the theorem.

6.7 Sliding Surface Design

We now prove a sufficient condition for the existence of a sliding surface $\sigma(\mathbf{x}) = \mathbf{S}\mathbf{x} = \mathbf{0}$ such that the plant model $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$ restricted to it has arbitrarily prescribed poles. A general method for designing a sliding surface is contained implicitly in the proof of the following theorem that can be found in Hui and Žak [132]. In our discussion, we use the following notation. For a full-column-rank matrix \mathbf{X} , its left inverse will be denoted \mathbf{X}^g . Thus, $\mathbf{X}^g\mathbf{X} = \mathbf{I}$, where \mathbf{I} is the identity matrix whose size is equal to the number of columns of \mathbf{X} . The imaginary part of a complex number z is denoted $\Im(z)$. The set of eigenvalues of a square matrix \mathbf{A} , also called its spectrum, is denoted $\text{eig}(\mathbf{A})$.

Theorem 6.2 Consider the plant model

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}, \quad (6.28)$$

where $\mathbf{x}(t) \in \mathbb{R}^n$, $\mathbf{u}(t) \in \mathbb{R}^m$, $\text{rank}(\mathbf{B}) = m$, and the pair (\mathbf{A}, \mathbf{B}) is controllable. Let

$$\{\lambda_1, \lambda_2, \dots, \lambda_{n-m}\} \in \mathbb{C},$$

where \mathbb{C} denotes the set of complex numbers, be an arbitrary set of $n - m$ complex numbers such that any λ_i with $\Im(\lambda_i) \neq 0$ appears in this set in conjugate pair. Then, there exist full-rank matrices $\mathbf{W} \in \mathbb{R}^{n \times (n-m)}$ and $\mathbf{W}^g \in \mathbb{R}^{(n-m) \times n}$ so that $\mathbf{W}^g\mathbf{W} = \mathbf{I}_{n-m}$, $\mathbf{W}^g\mathbf{B} = \mathbf{O}$, and the spectrum of the matrix $\mathbf{W}^g\mathbf{A}\mathbf{W}$ is

$$\text{eig}(\mathbf{W}^g\mathbf{A}\mathbf{W}) = \{\lambda_1, \lambda_2, \dots, \lambda_{n-m}\}.$$

Furthermore, there exists a matrix $\mathbf{S} \in \mathbb{R}^{m \times n}$ such that $\mathbf{S}\mathbf{B} = \mathbf{I}_m$ and the system (6.28) restricted to the surface $\mathbf{S}\mathbf{x} = \mathbf{0}$ has the poles $\lambda_1, \lambda_2, \dots, \lambda_{n-m}$.

Proof We first prove the existence of the desired matrices \mathbf{W} and \mathbf{W}^g . We can find $\mathbf{T} \in \mathbb{R}^{n \times (n-m)}$ so that $[\mathbf{T} \ \mathbf{B}]$ is invertible because \mathbf{B} has full rank. Let

$$[\mathbf{T} \ \mathbf{B}]^{-1} = \begin{bmatrix} \mathbf{T}^g \\ \mathbf{B}^g \end{bmatrix}.$$

We next introduce a new coordinate \mathbf{z} such that $\mathbf{x} = [\mathbf{T} \ \mathbf{B}]\mathbf{z}$. Then, in the new coordinates the system $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$ has the form

$$\dot{\mathbf{z}} = \begin{bmatrix} \mathbf{T}^g\mathbf{A}\mathbf{T} & \mathbf{T}^g\mathbf{A}\mathbf{B} \\ \mathbf{B}^g\mathbf{A}\mathbf{T} & \mathbf{B}^g\mathbf{A}\mathbf{B} \end{bmatrix} \mathbf{z} + \begin{bmatrix} \mathbf{O} \\ \mathbf{I}_m \end{bmatrix} \mathbf{u}.$$

Because the pair (\mathbf{A}, \mathbf{B}) is controllable, we can use the PBH rank test, given on page 106, to conclude that

$$\text{rank}[\lambda \mathbf{I}_{n-m} - \mathbf{T}^g\mathbf{A}\mathbf{T} \quad \mathbf{T}^g\mathbf{A}\mathbf{B}] = n - m \quad \text{for all } \lambda \in \mathbb{C},$$

which means that the pair $(T^g AT, T^g AB)$ is controllable. The set $\{\lambda_1, \lambda_2, \dots, \lambda_{n-m}\}$ is symmetric with respect to the real axis; therefore we can find a matrix $F \in \mathbb{R}^{m \times (n-m)}$ such that

$$\text{eig}(T^g AT - T^g ABF) = \{\lambda_1, \lambda_2, \dots, \lambda_{n-m}\}.$$

Let $W^g = T^g$ and let $W = T - BF$. Then, because by construction $T^g T = I_{n-m}$ and $T^g B = O$, we have $W^g W = I_{n-m}$ and $W^g B = O$. Note that

$$W^g AW = T^g AT - T^g ABF.$$

Hence,

$$\text{eig}(W^g AW) = \text{eig}(T^g AT - T^g ABF) = \{\lambda_1, \lambda_2, \dots, \lambda_{n-m}\}.$$

We next prove the existence of the sliding surface with the desired properties. We first note that the matrix $[W \ B]$ is invertible because $W^g W = I_{n-m}$, $W^g B = O$, and W, B are full rank. Hence, $\text{Range}(W) \cap \text{Range}(B) = \{0\}$ and we write

$$[W \ B]^{-1} = \begin{bmatrix} W^g \\ S \end{bmatrix}$$

for some $S \in \mathbb{R}^{m \times n}$. Note that $SB = I_m$ and $SW = O$. We now introduce new coordinates

$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} W^g \\ S \end{bmatrix} x,$$

where $z_2 = Sx$. In the new coordinates, the equivalent system $\dot{x} = (I_n - B(SB)^{-1}S)Ax$ takes the form

$$\begin{aligned} \dot{z} &= \begin{bmatrix} W^g \\ S \end{bmatrix} (I_n - B(SB)^{-1}S) A [W \ B] z \\ &= \begin{bmatrix} W^g AW & W^g AB \\ O & O \end{bmatrix} z. \end{aligned}$$

The restriction to $z_2 = Sx = 0$ is

$$\dot{z}_1 = W^g AW z_1,$$

and hence the system restricted to the sliding surface $Sx = 0$ has the same eigenvalues as the matrix $W^g AW$. The proof of the theorem is now complete.

We use the results of the proof of the above theorem to propose the following sliding surface design algorithm.

SWITCHING SURFACE DESIGN ALGORITHM

Given a controllable pair (A, B) with $\text{rank } B = m$, along with the set of desired poles $\{\lambda_i, i = 1, 2, \dots, n - m\}$ of the system restricted to a sliding surface, do the following:

STEP 1 Find a matrix $T \in \mathbb{R}^{(n-m) \times n}$ such that $[T \ B]$ is invertible.

STEP 2 Compute

$$[T \ B]^{-1} = \begin{bmatrix} T^g \\ B^g \end{bmatrix}.$$

STEP 3 Find a matrix $F \in \mathbb{R}^{m \times (n-m)}$ such that

$$\text{eig}(T^g A T - T^g A B F) = \{\lambda_1, \lambda_2, \dots, \lambda_{n-m}\}.$$

STEP 4 Take $W = T - B F$.

STEP 5 Compute

$$[W \ B]^{-1} = \begin{bmatrix} W^g \\ S \end{bmatrix}.$$

Form sliding surface $Sx = 0$ using the last m rows of $[W \ B]^{-1}$.

We illustrate the above algorithm with a simple numerical example.

◆ Example 6.4

Given a dynamical system model of the form (6.28), where

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

We design a sliding surface

$$\sigma = Sx = 0$$

so that the system restricted to it will behave like a linear system described by

$$\dot{z} = -z.$$

Following the steps of the of the design algorithm, we first find a matrix T such that the augmented matrix $[T \ B]$ is invertible. For example, we can select

$$T = [1 \ 0 \ 0]^T.$$

Note that

$$[T \ B] = I_3.$$

Hence, the inverse is

$$[\mathbf{T} \quad \mathbf{B}]^{-1} = \begin{bmatrix} \mathbf{T}^g \\ \mathbf{B}^g \end{bmatrix} = \mathbf{I}_3.$$

Next, we find $\mathbf{F} \in \mathbb{R}^{m \times (n-m)}$ such that

$$\mathbf{T}^g \mathbf{A} \mathbf{T} - \mathbf{T}^g \mathbf{A} \mathbf{B} \mathbf{F} = -\mathbf{I}.$$

We have

$$\mathbf{F} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

We then compute $\mathbf{W} = \mathbf{T} - \mathbf{B} \mathbf{F} = [1 \quad -1 \quad 0]^T$ and find

$$[\mathbf{W} \quad \mathbf{B}]^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{W}^g \\ \mathbf{S} \end{bmatrix}.$$

Hence,

$$\mathbf{S} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

6.8 State Estimation of Uncertain Systems

6.8.1 Discontinuous Estimators

So far in our discussion of variable structure sliding mode control system design we assumed that the entire state vector of the plant to be controlled is available to us through measurement. If the state vector cannot be measured, then we can try to determine its suitable estimate that can be used by the controller instead of the true state vector. If we use a state vector estimate for the inaccessible state, then we split a control design problem into three phases. The first phase is the design of a controller assuming the availability of the entire state vector of the plant. In the second phase, we construct a state estimator based on the plant input and output. Finally, in the third phase, we combine the controller and the estimator into the combined controller–estimator compensator, where the controller uses the plant state estimate instead of the actual state.

We discuss the problem of constructing a state estimator for a class of dynamical systems with matched uncertainties modeled by

$$\begin{aligned} \dot{\mathbf{x}} &= \mathbf{A} \mathbf{x} + \mathbf{B}(\mathbf{u} + \boldsymbol{\xi}), \\ \mathbf{y} &= \mathbf{C} \mathbf{x}, \end{aligned}$$

where $\mathbf{C} \in \mathbb{R}^{p \times n}$, and $p \geq m$, which means that the number of input channels does not exceed the number of output channels. The vector function $\boldsymbol{\xi}$ models the lumped uncertainties or

nonlinearities in the plant. We assume that there is a nonnegative real-valued function $\rho = \rho(y)$ such that

$$\|\xi\| \leq \rho.$$

We assume that the pair (A, C) is observable. Our goal is to construct the plant state estimator. Let \hat{x} be an estimate of x . Let e denote the estimation error, that is,

$$e(t) = \hat{x}(t) - x(t).$$

The observability of (A, C) implies the existence of a matrix $L \in \mathbb{R}^{n \times p}$ such that $A - LC$ has prescribed (symmetric with respect to the real axis) eigenvalues in the open left-half plane. Because $A - LC$ is asymptotically stable, by the Lyapunov theorem, for any $Q = Q^T > 0$, there is a unique $P = P^T > 0$ such that

$$(A - LC)^T P + P(A - LC) = -Q.$$

We choose Q , if possible, so that for some $F \in \mathbb{R}^{m \times p}$ we have

$$FC = B^T P.$$

We need this technical condition for the construction of the state estimator. This condition has an interesting system theoretic interpretation. Its satisfaction means that the transfer function matrix

$$G(s) = FC(sI - A + LC)^{-1}B$$

is positive real; that is, the elements of $G(s)$ are analytic in the open right-half plane and $G(s) + G^*(s) \geq 0$ in the open right-half plane. Here the superscript asterisk denotes the conjugate transpose.

To proceed, define the vector function

$$E(e, \eta) = \begin{cases} \frac{FCe}{\|FCe\|} \eta & \text{for } FCe \neq 0, \\ v \in \mathbb{R}^m, \|v\| \leq \eta & \text{for } FCe = 0, \end{cases}$$

where $\eta \geq 0$ is a design parameter. Note that if we have a single-input single-output plant, then we can write

$$E(e, \eta) = \eta \operatorname{sign}(FCe).$$

We now introduce a state estimator described in Walcott and Žak [292]. It has the form

$$\dot{\hat{x}} = (A - LC)\hat{x} - BE(e, \eta) + Ly + Bu.$$

If $\eta \geq \rho$, the function \hat{x} is an asymptotic estimate of the state x , that is,

$$\lim_{t \rightarrow \infty} e(t) = \lim_{t \rightarrow \infty} (\hat{x}(t) - x(t)) = 0.$$

To prove the above statement using Lyapunov's type of arguments, construct the differential equation satisfied by the estimation error e :

$$\dot{e} = \dot{\hat{x}} - \dot{x} = (A - LC)e - BE(e, \eta) - B\xi. \quad (6.29)$$

Then, show that

$$\frac{d}{dt}(\mathbf{e}^T \mathbf{P} \mathbf{e}) = -\mathbf{e}^T \mathbf{Q} \mathbf{e} < 0,$$

which implies

$$\lim_{t \rightarrow \infty} \mathbf{e}(t) = \mathbf{0}.$$

It follows from the above that the estimation error is insensitive to matched uncertainty.

6.8.2 Boundary Layer Estimators

We can view the error differential equation (6.29) as a model of a closed-loop system, with matched uncertainties, controlled by $-\mathbf{E}(\mathbf{e}, \eta)$. This “controller” is discontinuous and varies its structure with respect to the switching surface $\{\mathbf{e} : \mathbf{F} \mathbf{C} \mathbf{e} = \mathbf{0}\}$. In our analysis of (6.29), we assumed that $\mathbf{E}(\mathbf{e}, \eta)$ can switch with infinitely high frequency. In practice, the switching frequency is finite, albeit very high. High-frequency switching is not always desirable because it may degrade the system performance. In mechanical systems, it may lead to excessive wear and tear of the system’s components. To alleviate the problem, we smooth out the discontinuous part of the estimator by introducing a boundary layer around the surface $\{\mathbf{e} : \mathbf{F} \mathbf{C} \mathbf{e} = \mathbf{0}\}$. This results in replacing $\mathbf{E}(\mathbf{e}, \eta)$ with its continuous approximation $\tilde{\mathbf{E}}(\mathbf{e}, \eta)$. An example of such a continuous approximation is

$$\tilde{\mathbf{E}}(\mathbf{e}, \eta) = \begin{cases} \frac{\mathbf{F} \mathbf{C} \mathbf{e}}{\|\mathbf{F} \mathbf{C} \mathbf{e}\|} \eta & \text{for } \|\mathbf{F} \mathbf{C} \mathbf{e}\| \geq \nu, \\ \frac{\mathbf{F} \mathbf{C} \mathbf{e}}{\nu} \eta & \text{for } \|\mathbf{F} \mathbf{C} \mathbf{e}\| < \nu, \end{cases}$$

where ν is the width of the boundary layer. However, the error equation, or more precisely $\mathbf{e} = \mathbf{0}$, is not anymore asymptotically stable when a boundary layer estimator is used. To describe the behavior of the error differential equation resulting from the application of a boundary layer estimator, we recall the following definitions.

Definition 6.1 The closed ball $\bar{B}_q(\mathbf{0}) = \{\mathbf{e} : \|\mathbf{e}\| < q\}$ is uniformly stable if for any $\varepsilon > q$, there is $\delta(\varepsilon) > 0$ so that if $\|\mathbf{e}(t_0)\| \leq \delta(\varepsilon)$, then $\|\mathbf{e}(t)\| \leq \varepsilon$ for $t \geq t_0$.

Definition 6.2 The equilibrium $\mathbf{e} = \mathbf{0}$ is uniformly bounded if for any given $r > 0$, there exists $d(r) < \infty$ so that for each solution $\mathbf{e}(t) = \mathbf{e}(t; t_0, \mathbf{e}_0)$ with $\|\mathbf{e}_0\| \leq r$, we have $\|\mathbf{e}(t)\| \leq d(r)$ for all $t \geq t_0$.

We next adapt the definition of uniform ultimate boundedness, abbreviated UUB, to the estimation error, \mathbf{e} .

Definition 6.3 The error \mathbf{e} is uniformly ultimately bounded (UUB) with respect to a closed ball \bar{B} if for every $r > 0$, there is $T(r) \geq 0$ such that for each solution $\mathbf{e}(t)$ with $\|\mathbf{e}(t_0)\| \leq r$, we have $\mathbf{e}(t) \in \bar{B}$ for $t \geq t_0 + T(r)$.

We can show that the error of the boundary layer estimator is UUB. Specifically, the estimation error of the boundary layer estimator ends up in a ball whose radius depends on ν . The ball radius tends to zero with ν tending to zero. To see this consider, as before, the Lyapunov function

candidate $V(\mathbf{e}) = \mathbf{e}^T \mathbf{P} \mathbf{e}$. For $\|\mathbf{F} \mathbf{C} \mathbf{e}\| \geq \nu$, we have

$$\begin{aligned} \frac{d}{dt} V(\mathbf{e}) &= 2\mathbf{e}^T \mathbf{P} \dot{\mathbf{e}} \\ &= 2\mathbf{e}^T \mathbf{P} ((\mathbf{A} - \mathbf{L} \mathbf{C})\mathbf{e} - \mathbf{B} \tilde{\mathbf{E}}(\mathbf{e}, \eta) - \mathbf{B} \xi) \\ &= -\mathbf{e}^T \mathbf{Q} \mathbf{e} - 2\mathbf{e}^T \mathbf{P} \mathbf{B} \frac{\mathbf{F} \mathbf{C} \mathbf{e}}{\|\mathbf{F} \mathbf{C} \mathbf{e}\|} \eta - 2\mathbf{e}^T \mathbf{P} \mathbf{B} \xi. \end{aligned} \quad (6.30)$$

Because by assumption, $\mathbf{F} \mathbf{C} = \mathbf{B}^T \mathbf{P}$, therefore

$$\mathbf{e}^T \mathbf{P} \mathbf{B} \mathbf{F} \mathbf{C} \mathbf{e} = \|\mathbf{F} \mathbf{C} \mathbf{e}\|^2 = \|\mathbf{e}^T \mathbf{P} \mathbf{B}\|^2.$$

Substituting the above into (6.30) and taking into account that $\|\mathbf{e}^T \mathbf{P} \mathbf{B} \xi\| \leq \|\mathbf{e}^T \mathbf{P} \mathbf{B}\| \rho$, we obtain

$$\frac{d}{dt} V(\mathbf{e}) \leq -\mathbf{e}^T \mathbf{Q} \mathbf{e} - 2\|\mathbf{F} \mathbf{C} \mathbf{e}\| \eta + 2\|\mathbf{e}^T \mathbf{P} \mathbf{B}\| \eta = -\mathbf{e}^T \mathbf{Q} \mathbf{e}.$$

For $\|\mathbf{F} \mathbf{C} \mathbf{e}\| < \nu$, we have

$$\begin{aligned} \frac{d}{dt} V(\mathbf{e}) &= 2\mathbf{e}^T \mathbf{P} \dot{\mathbf{e}} \\ &= -\mathbf{e}^T \mathbf{Q} \mathbf{e} - 2\mathbf{e}^T \mathbf{P} \mathbf{B} \frac{\mathbf{F} \mathbf{C} \mathbf{e}}{\nu} \eta - 2\mathbf{e}^T \mathbf{P} \mathbf{B} \xi \\ &\leq -\mathbf{e}^T \mathbf{Q} \mathbf{e} - 2\frac{\|\mathbf{F} \mathbf{C} \mathbf{e}\|^2}{\nu} \eta + 2\|\mathbf{e}^T \mathbf{P} \mathbf{B}\| \rho \\ &\leq -\mathbf{e}^T \mathbf{Q} \mathbf{e} + 2\nu\rho \\ &\leq -\lambda_{\min}(\mathbf{Q})\|\mathbf{e}\|^2 + 2\nu\rho. \end{aligned} \quad (6.31)$$

In summary,

$$\begin{aligned} \frac{d}{dt} V(\mathbf{e}) &= -\mathbf{e}^T \mathbf{Q} \mathbf{e} - 2\mathbf{e}^T \mathbf{P} \mathbf{B} \tilde{\mathbf{E}}(\mathbf{e}, \eta) - 2\mathbf{e}^T \mathbf{P} \mathbf{B} \xi \\ &\leq \begin{cases} -\mathbf{e}^T \mathbf{Q} \mathbf{e} & \text{for } \|\mathbf{F} \mathbf{C} \mathbf{e}\| \geq \nu \\ -\lambda_{\min}(\mathbf{Q})\|\mathbf{e}\|^2 + 2\nu\rho & \text{for } \|\mathbf{F} \mathbf{C} \mathbf{e}\| < \nu. \end{cases} \end{aligned}$$

Hence, we can conclude that $\frac{d}{dt} V(\mathbf{e}) < 0$ for all $\mathbf{e} \in \mathbb{R}^n$ such that

$$\|\mathbf{e}\| > \sqrt{\frac{2\nu\rho}{\lambda_{\min}(\mathbf{Q})}} \stackrel{\text{def}}{=} R.$$

We illustrate the regions used in the discussion of the boundary layer state estimator in Figure 6.11. Referring to Figure 6.11 and Definition 6.2, we have for $r \leq R$,

$$d(r) = \sqrt{\frac{\lambda_{\max}(\mathbf{P})}{\lambda_{\min}(\mathbf{P})}} R.$$

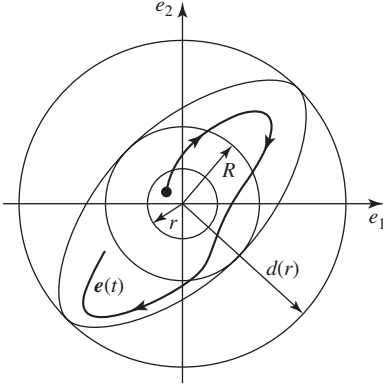


Figure 6.11 A two-dimensional illustration of an UUB error.

If, on the other hand, $r > R$, then

$$d(r) = \sqrt{\frac{\lambda_{\max}(\mathbf{P})}{\lambda_{\min}(\mathbf{P})}} r.$$

We conclude that the estimation error is UUB with respect to any ball with a radius greater than

$$\sqrt{\frac{\lambda_{\max}(\mathbf{P})}{\lambda_{\min}(\mathbf{P})}} R.$$

6.9 Sliding Modes in Solving Optimization Problems

In this section we use sliding modes to analyze dynamics of an analog circuit, proposed by Rodríguez-Vázquez et al. [242], that can solve convex optimization problems with linear constraints. This class of networks for solving optimization problems is particularly attractive because it can be realized using switched-capacitor technology. The network we analyze is modeled by a set of nonlinear ordinary differential equations with discontinuous right-hand sides. In carrying out the analysis, we use concepts from the theory of differential equations with discontinuous right-hand sides and the Lyapunov stability theory. We show that the equilibrium points of the differential equations modeling the network coincide with the optimizers of the underlying optimization problem. We then prove that irrespective of the initial state of the network, its state converges to the solution set of the optimization problem. Our analysis is based on the results of the paper by Glazos, Hui, and Žak [98].

6.9.1 Optimization Problem Statement

We consider the constrained optimization problem

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b}, \end{aligned} \tag{6.32}$$

where $\mathbf{x} \in \mathbb{R}^n$, $f: \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, and $\mathbf{b} \in \mathbb{R}^m$. Let

$$\Omega = \{\mathbf{x} : \mathbf{A}\mathbf{x} \leq \mathbf{b}\}. \tag{6.33}$$

denote the feasible region for the above optimization problem. The feasible set Ω is the intersection of m closed half-spaces and is therefore a convex polytope. Let

$$h_i = \{\mathbf{x} : \mathbf{a}_i \mathbf{x} = b_i\}, \quad i = 1, 2, \dots, m,$$

denote the hyperplanes associated with the above m closed half-spaces, where \mathbf{a}_i is the i th row of \mathbf{A} and b_i is the i th component of the vector \mathbf{b} . Also, we define

$$H = \bigcup_{i=1}^m h_i.$$

We let Γ denote the collection of all local minimizers of the optimization problem (6.32). Lastly, we introduce the index sets

$$I(\mathbf{x}) = \{i : \mathbf{a}_i \mathbf{x} = b_i\}$$

and

$$J(\mathbf{x}) = \{i : \mathbf{a}_i \mathbf{x} > b_i\}.$$

Recall that a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ defined on a convex set Ω is said to be convex if, for every $\mathbf{x}, \mathbf{y} \in \Omega$ and every $0 \leq \alpha \leq 1$, there holds

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}).$$

Let $f \in \mathcal{C}^1$. Then, f is convex over a convex set Ω if and only if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f^T(\mathbf{x})(\mathbf{y} - \mathbf{x})$$

for all $\mathbf{x}, \mathbf{y} \in \Omega$.

We make the following assumptions:

- A1.** The objective function f is convex over \mathbb{R}^n and has continuous first partial derivatives on \mathbb{R}^n ; that is, $f \in \mathcal{C}^1$.
- A2.** The feasible set Ω is a nonempty, bounded, and closed polyhedron, which means that Ω is a compact set.
- A3.** Each point in Ω is a regular point of the constraints; that is, for any $\mathbf{x} \in \Omega$ the vectors $\mathbf{a}_i, i \in I(\mathbf{x})$, are linearly independent.

It follows from Assumption A1 that the optimization problem (6.32) is a convex programming problem. Indeed, we are asked to minimize a convex objective function f over a polyhedron Ω .

Problem (6.32) is well-posed in the sense that the feasible set Ω is nonempty and the optimal value of f is not $-\infty$ because by assumptions A1 and A2 the set $f(\Omega)$ is compact.

An elementary overview of results concerning convex optimization problems can be found, for example, in Chong and Zak [50, Chapter 21]. We state two results related to problem (6.32) that are relevant in our analysis.

Lemma 6.1 Any local minimizer of the optimization problem (6.32) is a global minimizer.

Furthermore, the Karush–Kuhn–Tucker condition for optimality is both necessary and sufficient for problem (6.32). Specifically, we have the following lemma:

Lemma 6.2 Suppose that $\mathbf{x}^* \in \Omega$. Then, $\mathbf{x}^* \in \Gamma$ if and only if there exist nonnegative constants $\mu_1, \mu_2, \dots, \mu_m$ such that

$$\nabla f(\mathbf{x}^*) + \sum_{i \in I(\mathbf{x}^*)} \mu_i \mathbf{a}_i^T = \mathbf{0}.$$

The network applied to the optimization problem (6.32) has the form

$$\begin{aligned} \tau \dot{\mathbf{x}}(t) &= \mathbf{g}(\mathbf{x}(t)) \\ &= \begin{cases} -\mu \sum_{i \in J(\mathbf{x}(t))} \mathbf{a}_i^T & \text{if } \mathbf{x}(t) \notin \Omega \\ -\nabla f(\mathbf{x}(t)) & \text{if } \mathbf{x}(t) \in \Omega, \end{cases} \\ \mathbf{x}(0) &= \mathbf{x}_0, \end{aligned} \quad (6.34)$$

where τ and μ are positive design constants. The function $\mathbf{g}(\mathbf{x})$ is discontinuous on the surfaces h_1, h_2, \dots, h_m and is continuous elsewhere. A phenomenon commonly occurring in systems such as (6.34) is the so-called sliding mode that we discussed in Section 6.2. An analysis of the behavior of the above dynamical system using a sliding mode approach is presented in the following sections.

We remark that we can view (6.34) as a dynamical gradient system for the unconstrained problem

$$\text{minimize } \kappa_0 f(\mathbf{x}) + \mu \sum_{i=1}^m \kappa_i (\mathbf{a}_i \mathbf{x} - b_i),$$

where

$$\kappa_i = \begin{cases} 0 & \text{if } \mathbf{a}_i \mathbf{x} \leq b_i, \\ 1 & \text{if } \mathbf{a}_i \mathbf{x} > b_i \end{cases}$$

and

$$\kappa_0 = \begin{cases} 1 & \text{if } \kappa_i = 0 \text{ for all } i = 1, 2, \dots, m, \\ 0 & \text{otherwise;} \end{cases}$$

that is, $\kappa_0 = \prod_{i=1}^m (1 - \kappa_i)$. Using the preceding notation, we can represent system (6.34) as

$$\tau \dot{\mathbf{x}}(t) = \mathbf{g}(\mathbf{x}(t)) = -\kappa_0 \nabla f(\mathbf{x}(t)) - \mu \sum_{i=1}^m \kappa_i \mathbf{a}_i^T, \quad \mathbf{x}(0) = \mathbf{x}_0. \quad (6.35)$$

6.9.2 Penalty Function Method

An approach that can be used to construct an analog optimizer to solve a constrained optimization problem is to first convert, or approximate, the constrained optimization problem into an associated unconstrained optimization problem and then design an analog network that solves the unconstrained problem. Such a network is typically an implementation of the dynamical gradient system for minimizing the objective function of the associated unconstrained problem.

One method of approximating a constrained optimization problem into an associated unconstrained optimization problem is the penalty function method. For more details on the subject of the penalty method see, for example, Luenberger [190]. The idea behind the penalty function method is to approximate the constrained optimization problem

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) \\ &\text{subject to} && \mathbf{x} \in \Omega, \end{aligned} \tag{6.36}$$

where $\Omega \subset \mathbb{R}^n$ is the feasible set, with an associated unconstrained problem of the form

$$\text{minimize} \quad f(\mathbf{x}) + \zeta p(\mathbf{x}), \tag{6.37}$$

where the objective function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous on \mathbb{R}^n and ζ is a positive constant. The function $p: \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous on \mathbb{R}^n and satisfies the following:

1. $p(\mathbf{x}) \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$.
2. $p(\mathbf{x}) = 0$ if and only if $\mathbf{x} \in \Omega$.

If for a finite value of the penalty parameter ζ the solution of the unconstrained problem (6.37) coincides with the solution of the constrained problem (6.36), then the penalty function p is said to be exact. Bertsekas [26] has shown that except for trivial cases, an exact penalty function cannot be everywhere differentiable—see Exercise 6.9 for an illustration of this statement. Quite often we can find an exact penalty function that results in the function $f(\mathbf{x}) + \zeta p(\mathbf{x})$ being piecewise smooth. However, the dynamical gradient system for such a function will have a discontinuous right-hand side. As such, a rigorous analysis of its behavior cannot be carried out using only methods derived from the classical theory of differential equations. Rather, the analysis of the system must be carried out using other methods—for example, those reported by Filippov [84] and by Aubin and Cellina [14]. Such an approach was undertaken in 1967 by Karpinskaya [150] and was later undertaken by Korovin and Utkin [164], Utkin [286], Flåm [87], Flåm and Zowe [86], Žak et al. [314], and Chong, Hui, and Žak [49], among others. A differential geometric approach to the analysis of dynamical systems that solve optimization problems is reported by Faybusovich [80] and by Helmke and Moore [119].

6.9.3 Dynamical Gradient Circuit Analysis

We will now show that every trajectory of system (6.34), irrespective of the initial state, converges to the solution set of the optimization problem (6.32). We do this by showing that every trajectory of system (6.34) has the following properties:

1. Irrespective of the initial state \mathbf{x}_0 , the resulting trajectory reaches the feasible set Ω in finite time and is thereafter confined to Ω . This portion of the trajectory is hereafter referred to as the *reaching phase*.
2. While confined to the feasible set Ω , the trajectory converges asymptotically to the solution set Γ . This portion of the trajectory is hereafter referred to as the *convergence phase*.

Reaching Phase We begin by introducing the function $V: \mathbb{R}^n \rightarrow \mathbb{R}$ of the form

$$V(\mathbf{x}) = \sum_{i=1}^m \max\{\mathbf{a}_i \mathbf{x} - b_i, 0\},$$

or, equivalently,

$$V(\mathbf{x}) = \begin{cases} \sum_{i \in J(\mathbf{x})} (\mathbf{a}_i \mathbf{x} - b_i) & \text{if } \mathbf{x} \notin \Omega, \\ 0 & \text{if } \mathbf{x} \in \Omega. \end{cases} \quad (6.38)$$

Observe that $V(\mathbf{x}) > 0$ for all $\mathbf{x} \notin \Omega$, and $V(\mathbf{x}) = 0$ for all $\mathbf{x} \in \Omega$. We can view $V(\mathbf{x})$ as a distance measure from the point \mathbf{x} to the feasible set Ω . We shall show that V is strictly decreasing when evaluated along any trajectory of system (6.34) outside Ω . Specifically, we will prove the existence of a positive constant, η , with the property that

$$\frac{d}{dt} V(\mathbf{x}(t)) \leq -\eta$$

when evaluated along any trajectory of system (6.34) outside Ω . In doing so, we use the following result:

Lemma 6.3 Suppose $\hat{\mathbf{x}} \notin \Omega$ and $\beta_1, \beta_2, \dots, \beta_m$ are nonnegative constants. If $\beta_j > 0$ for $j \in J(\hat{\mathbf{x}})$, then $\sum_{i \in I(\hat{\mathbf{x}}) \cup J(\hat{\mathbf{x}})} \beta_i \mathbf{a}_i \neq \mathbf{0}$.

Proof Let the function $q: \mathbb{R}^n \rightarrow \mathbb{R}$ be defined as

$$q(\mathbf{x}) = \sum_{i \in I(\hat{\mathbf{x}}) \cup J(\hat{\mathbf{x}})} \beta_i (\mathbf{a}_i \mathbf{x} - b_i).$$

Clearly, $q \in \mathcal{C}^1$ and

$$\nabla q^T(\mathbf{x}) = \sum_{i \in I(\hat{\mathbf{x}}) \cup J(\hat{\mathbf{x}})} \beta_i \mathbf{a}_i.$$

Because $\beta_1, \beta_2, \dots, \beta_m$ are all nonnegative, it follows that

1. the function q is convex over \mathbb{R}^n , and
2. $q(\mathbf{x}) \leq 0$ for all $\mathbf{x} \in \Omega$.

Furthermore, because $\beta_j > 0$ for $j \in J(\hat{\mathbf{x}})$, we also have $q(\hat{\mathbf{x}}) > 0$. Hence, we must have

$$\nabla q^T(\hat{\mathbf{x}}) = \sum_{i \in I(\hat{\mathbf{x}}) \cup J(\hat{\mathbf{x}})} \beta_i \mathbf{a}_i \neq \mathbf{0}.$$

We must consider two cases when analyzing the behavior of system (6.34). The first case is when the motion of the system is not confined to any of the surfaces h_1, h_2, \dots, h_m . The second case is when the system is in sliding mode on one or more of the surfaces h_1, h_2, \dots, h_m . We only need to consider these two cases because every trajectory of system (6.34) is composed of these two types of motion. That is, every trajectory $\mathbf{x}(t)$ of system (6.34) can be broken into a countable number of intervals, $(t_0, t_1), \dots, (t_{l-1}, t_l), \dots$, on each of which both index sets $I(\mathbf{x}(t))$ and $J(\mathbf{x}(t))$ are constant. If $I(\mathbf{x}(t))$ is not empty, then the system is in a sliding mode on that interval, and if $I(\mathbf{x}(t))$ is empty, then the system is not in a sliding mode.

Case 1 Let $\mathbf{x}(t)$ denote any particular trajectory of system (6.34), and let $t' \in (t_{l-1}, t_l)$. Suppose that on the interval (t_{l-1}, t_l) the trajectory $\mathbf{x}(t)$ does not intersect Ω or any of the surfaces h_1, h_2, \dots, h_m , that is,

$$\mathbf{x}(t) \notin \Omega \cup H \quad \text{for all } t \in (t_{l-1}, t_l). \quad (6.39)$$

Denoting the index set $J(\mathbf{x}(t'))$ by J' , it follows from (6.39) and the fact that $\mathbf{x}(t)$ is absolutely continuous that $I(\mathbf{x}(t)) = \emptyset$ and $J' \neq \emptyset$ for all $t \in (t_{l-1}, t_l)$. This in turn implies that on the interval (t_{l-1}, t_l) the trajectory $\mathbf{x}(t)$ must satisfy (6.34) in the usual sense. Hence, we have

$$\tau \dot{\mathbf{x}}(t) = -\mu \sum_{i \in J'} \mathbf{a}_i^T$$

almost everywhere on the interval (t_{l-1}, t_l) . Also, it follows from (6.38) that

$$V(\mathbf{x}(t)) = \sum_{i \in J'} (\mathbf{a}_i \mathbf{x}(t) - b_i) \quad (6.40)$$

for all $t \in (t_{l-1}, t_l)$. Applying the chain rule to (6.40), we obtain

$$\frac{d}{dt} V(\mathbf{x}(t)) = \nabla V^T(\mathbf{x}(t)) \dot{\mathbf{x}}(t) = -\frac{\mu}{\tau} \left\| \sum_{i \in J'} \mathbf{a}_i \right\|_2^2$$

almost everywhere on the interval (t_{l-1}, t_l) . However, Lemma 6.3 along with the fact that there is a finite number of constraints—that is, m is finite—implies the existence of a positive constant, η_1 , with the property that

$$\eta_1 \leq \left\| \sum_{i \in J(\mathbf{x})} \mathbf{a}_i \right\|_2^2$$

for all $\mathbf{x} \notin \Omega$. Hence,

$$\frac{d}{dt} V(\mathbf{x}(t)) \leq -\frac{\mu \eta_1}{\tau} < 0 \quad (6.41)$$

almost everywhere on the interval (t_{l-1}, t_l) , which concludes our analysis of the first case.

Case 2 Suppose now that on an interval (t_{l-1}, t_l) a particular trajectory $\mathbf{x}(t)$ is confined to one or more of the surfaces h_1, h_2, \dots, h_m and does not intersect Ω . Specifically, suppose that there exist nonempty index sets I' and J' with the property that

$$I(\mathbf{x}(t)) = I' \quad \text{for all } t \in (t_{l-1}, t_l), \quad (6.42)$$

and

$$J(\mathbf{x}(t)) = J' \quad \text{for all } t \in (t_{l-1}, t_l). \quad (6.43)$$

Let i_1, i_2, \dots, i_k denote the elements of I' , and let $S = \bigcap_{i \in I'} h_i$ denote the surface to which the trajectory $\mathbf{x}(t)$ is confined on the interval (t_{l-1}, t_l) , and let

$$\tilde{\mathbf{a}} = \sum_{i \in J'} \mathbf{a}_i^T.$$

It follows from (6.42) and (6.43) that for a sufficiently small $\delta > 0$ the surfaces h_{i_1}, \dots, h_{i_k} partition the δ -neighborhood of $\mathbf{x}(t')$ into $\ell \leq 2^k$ regions, $\mathcal{R}_1, \mathcal{R}_2, \dots, \mathcal{R}_\ell$, in each of which the function \mathbf{g} is continuous. Then, $G(\mathbf{x}(t'))$ is the set of all vectors \mathbf{v} of the form

$$\mathbf{v} = -\mu \sum_{j=1}^{\ell} \alpha_j \left(\tilde{\mathbf{a}} + \sum_{i \in I'} U(i, j) \mathbf{a}_i^T \right), \quad (6.44)$$

where

$$\alpha_j \geq 0, \quad j = 1, 2, \dots, \ell, \quad (6.45)$$

$$\sum_{j=1}^{\ell} \alpha_j = 1, \quad (6.46)$$

and

$$U(i, j) = \begin{cases} 1 & \text{if } \mathbf{a}_i \mathbf{x} > b_i \quad \text{for all } \mathbf{x} \in \mathcal{R}_j \\ 0 & \text{if } \mathbf{a}_i \mathbf{x} \leq b_i \quad \text{for all } \mathbf{x} \in \mathcal{R}_j. \end{cases}$$

Using (6.46), we rewrite (6.44) as

$$\mathbf{v} = -\mu \tilde{\mathbf{a}} - \mu \sum_{j=1}^{\ell} \alpha_j \left(\sum_{i \in I'} U(i, j) \mathbf{a}_i^T \right). \quad (6.47)$$

For notational simplicity, we let

$$\mathbf{r} = \sum_{j=1}^{\ell} \alpha_j \left(\sum_{i \in I'} U(i, j) \mathbf{a}_i^T \right).$$

Then, (6.47) becomes

$$\mathbf{v} = -\mu \tilde{\mathbf{a}} - \mu \mathbf{r}. \quad (6.48)$$

Let $T(\mathbf{x})$ denote the tangent plane to the surface S at the point \mathbf{x} . Observe that because the trajectory $\mathbf{x}(t)$ is confined to the surface S on the interval (t_{l-1}, t_l) , then $\tau \dot{\mathbf{x}}(t) \in T(\mathbf{x}(t))$ almost everywhere on the interval (t_{l-1}, t_l) . Thus, $\mathbf{x}(t)$ is an absolutely continuous function that satisfies

$$\tau \dot{\mathbf{x}}(t) \in G(\mathbf{x}(t)) \cap T(\mathbf{x}(t))$$

almost everywhere on the interval (t_{l-1}, t_l) . We note that because the trajectory $\mathbf{x}(t)$ is confined to the surface S on the interval (t_{l-1}, t_l) , the set $G(\mathbf{x}(t)) \cap T(\mathbf{x}(t))$ is nonempty for all (t_{l-1}, t_l) . In particular, $G(\mathbf{x}(t')) \cap T(\mathbf{x}(t'))$ is the set of all \mathbf{v} of the form (6.48) that also lie on the tangent plane $T(\mathbf{x}(t'))$. However,

$$\mathbf{r} \in \text{Span}\{\mathbf{a}_{i_1}^T, \mathbf{a}_{i_2}^T, \dots, \mathbf{a}_{i_k}^T\},$$

and thus \mathbf{r} is orthogonal to $T(\mathbf{x}(t'))$. Therefore, the set $G(\mathbf{x}(t')) \cap T(\mathbf{x}(t'))$ consists of exactly one element, the orthogonal projection of the vector $-\mu \tilde{\mathbf{a}}$ onto the tangent plane $T(\mathbf{x}(t'))$. Hence, for (t_{l-1}, t_l) the set $G(\mathbf{x}(t)) \cap T(\mathbf{x}(t))$ consists of exactly one element, and thus $\dot{\mathbf{x}}(t) = -(\mu/\tau) \mathbf{P} \tilde{\mathbf{a}}$ is uniquely determined almost everywhere on the interval (t_{l-1}, t_l) , where \mathbf{P} is the orthogonal

projector onto the tangent plane $T(\mathbf{x}(t))$. Note that

$$\mathbf{P} = \mathbf{I}_n - \mathbf{B}(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T, \quad (6.49)$$

where \mathbf{I}_n denotes the $n \times n$ identity matrix and

$$\mathbf{B} = [\mathbf{a}_{i_1}^T \quad \mathbf{a}_{i_2}^T \quad \cdots \quad \mathbf{a}_{i_k}^T].$$

By (6.38) and (6.43),

$$V(\mathbf{x}(t)) = \sum_{i \in J'} (\mathbf{a}_i \mathbf{x}(t) - b_i) \quad (6.50)$$

for all $t \in (t_{l-1}, t_l)$. Applying the chain rule to (6.50) yields

$$\frac{d}{dt} V(\mathbf{x}(t)) = \nabla V^T(\mathbf{x}(t)) \dot{\mathbf{x}}(t) = -\frac{\mu}{\tau} \tilde{\mathbf{a}}^T \mathbf{P} \tilde{\mathbf{a}} \quad (6.51)$$

almost everywhere on the interval (t_{l-1}, t_l) . Observing that $\mathbf{P} = \mathbf{P}^T = \mathbf{P}^2$, we rewrite (6.51) as

$$\frac{d}{dt} V(\mathbf{x}(t)) = -\frac{\mu}{\tau} \|\mathbf{P} \tilde{\mathbf{a}}\|_2^2. \quad (6.52)$$

We can show, using Lemma 6.3, that $\mathbf{P} \tilde{\mathbf{a}} \neq \mathbf{0}$ (see Exercise 6.17). Combining this with the fact that the number of constraints is finite, we conclude that there exists a positive constant, η_2 , such that

$$\eta_2 \leq \|\mathbf{P} \tilde{\mathbf{a}}\|_2^2 \quad \text{for all } \mathbf{x} \in H \setminus \Omega.$$

Hence,

$$\frac{d}{dt} V(\mathbf{x}(t)) \leq -\frac{\mu \eta_2}{\tau} < 0 \quad (6.53)$$

almost everywhere on the interval (t_{l-1}, t_l) , which concludes our analysis of the second case.

Using the above results, we can now state and prove the following theorem.

Theorem 6.3 Every trajectory of system (6.34) reaches Ω in finite time and is thereafter confined to Ω .

Proof Let $\mathbf{x}(t)$ be any particular trajectory of system (6.34) and let

$$\eta = \min \left\{ \frac{\mu \eta_1}{\tau}, \frac{\mu \eta_2}{\tau} \right\}.$$

It follows directly from (6.41) and (6.53) that

$$\frac{d}{dt} V(\mathbf{x}(t)) \leq -\eta$$

for almost all t when $\mathbf{x}(t) \notin \Omega$. Hence, $\mathbf{x}(t) \in \Omega$ for all $t \geq V(\mathbf{x}(0))/\eta$.

Convergence Phase We now analyze dynamical behavior of system (6.34) when its motion is confined to the feasible set Ω . We begin by introducing the function $F: \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$F(\mathbf{x}) = f(\mathbf{x}) - f^*,$$

where f^* is the optimal value of f for the optimization problem (6.32). Clearly, $F \in \mathcal{C}^1$. Furthermore, we see from Lemma 6.1 that $F(\mathbf{x}) = 0$ for all $\mathbf{x} \in \Gamma$ and $F(\mathbf{x}) > 0$ for all $\mathbf{x} \in \Omega \setminus \Gamma$.

Similar to the previous subsection, we consider two cases when analyzing the behavior of system (6.34) in the region Ω . We first consider the case when the motion of system (6.34) is confined to the interior of Ω . Then, we consider the case when the motion of system (6.34) is confined to the boundary of Ω —that is, when system (6.34) is in a sliding mode on one or more of the surfaces h_1, \dots, h_m .

Case 1 Let $\mathbf{x}(t)$ be any particular trajectory of system (6.34). Suppose that on an interval (t_{l-1}, t_l) the trajectory $\mathbf{x}(t)$ is confined to the interior of Ω . Then, $I(\mathbf{x}(t)) = \emptyset$ and $J(\mathbf{x}(t)) = \emptyset$ for all $t \in (t_{l-1}, t_l)$. This implies that on the interval (t_{l-1}, t_l) , the trajectory $\mathbf{x}(t)$ satisfies $\tau \dot{\mathbf{x}}(t) = -\nabla f(\mathbf{x}(t))$ almost everywhere on the interval (t_{l-1}, t_l) in the usual sense. Applying the chain rule to $F(\mathbf{x}(t))$, we obtain

$$\frac{d}{dt} F(\mathbf{x}(t)) = \nabla F^T(\mathbf{x}(t)) \dot{\mathbf{x}}(t) = -\frac{1}{\tau} \|\nabla f(\mathbf{x}(t))\|^2_2 \quad (6.54)$$

almost everywhere on the interval (t_{l-1}, t_l) . Using the fact that $I(\mathbf{x}(t)) = \emptyset$ for all $t \in (t_{l-1}, t_l)$, we see from Lemma 6.2 that on the interval (t_{l-1}, t_l) , we have $\nabla f(\mathbf{x}(t)) = \mathbf{0}$ if and only if $\mathbf{x}(t) \in \Gamma$. Hence, we have

$$\frac{d}{dt} F(\mathbf{x}(t)) = 0 \quad \text{if and only if } \mathbf{x}(t) \in \Gamma, \quad (6.55)$$

and

$$\frac{d}{dt} F(\mathbf{x}(t)) < 0 \quad \text{if and only if } \mathbf{x}(t) \notin \Gamma, \quad (6.56)$$

almost everywhere on the interval (t_{l-1}, t_l) . This concludes our analysis of the first case.

Case 2 Let $\mathbf{x}(t)$ denote any particular trajectory of system (6.34) and let $t \in (t_{l-1}, t_l)$. Suppose that on the interval (t_{l-1}, t_l) the trajectory $\mathbf{x}(t)$ is confined to the boundary of Ω . Thus:

1. There exists a nonempty index set I' with the property that $I(\mathbf{x}(t)) = I'$ for all $t \in (t_{l-1}, t_l)$.
2. $J(\mathbf{x}(t)) = \emptyset$ for all $t \in (t_{l-1}, t_l)$.

Let i_1, \dots, i_k denote the elements of I' , and let $S = \bigcap_{i \in I'} h_i$ denote the surface to which the trajectory $\mathbf{x}(t)$ is confined on the interval (t_{l-1}, t_l) . Let $\mathbf{B} \in \mathbb{R}^{n \times k}$ be the matrix with columns $\mathbf{a}_{i_1}^T, \dots, \mathbf{a}_{i_k}^T$. By assumption A3 the vectors $\mathbf{a}_{i_1}^T, \dots, \mathbf{a}_{i_k}^T$ are linearly independent, and hence the matrix \mathbf{B} has full rank.

It follows from the above that for a sufficiently small $\delta > 0$ the surfaces h_{i_1}, \dots, h_{i_k} partition the δ -neighborhood of $\mathbf{x}(t')$ into 2^k regions, in each of which the function g is continuous. Using a similar argument to the one in the previous subsection, we can show that the set $G(\mathbf{x}(t'))$ is

the set of all \mathbf{v} of the form

$$\mathbf{v} = -\alpha_0 \nabla f(\mathbf{x}(t')) - \mu \sum_{j=1}^{2^k-1} \alpha_j \mathbf{B} \mathbf{u}_j, \quad (6.57)$$

where

$$\alpha_j \geq 0, \quad j = 0, 1, \dots, 2^k - 1, \quad (6.58)$$

$$\sum_{j=0}^{2^k-1} \alpha_j = 1, \quad (6.59)$$

and $\mathbf{u}_j \in \mathbb{R}^k$, $j = 1, \dots, 2^k - 1$, are defined as

$$\begin{aligned} \mathbf{u}_1 &= [1 \ 0 \ 0 \ \cdots \ 0 \ 0 \ 0]^T, \\ \mathbf{u}_2 &= [0 \ 1 \ 0 \ \cdots \ 0 \ 0 \ 0]^T, \\ \mathbf{u}_3 &= [1 \ 1 \ 0 \ \cdots \ 0 \ 0 \ 0]^T, \\ &\vdots \\ \mathbf{u}_{2^{k-3}} &= [1 \ 1 \ 1 \ \cdots \ 1 \ 0 \ 1]^T, \\ \mathbf{u}_{2^{k-2}} &= [1 \ 1 \ 1 \ \cdots \ 1 \ 1 \ 0]^T, \\ \mathbf{u}_{2^{k-1}} &= [1 \ 1 \ 1 \ \cdots \ 1 \ 1 \ 1]^T. \end{aligned}$$

We illustrate this case in Figure 6.12 for a hypothetical situation of sliding occurring on the intersection of two hyperplanes. In the figure, the components of the vector \mathbf{v} as given by (6.57) are shown.

Let $T(\mathbf{x})$ denote the plane tangent to the surface S at the point \mathbf{x} . Observe that $\dot{\mathbf{x}}(t) \in T(\mathbf{x}(t))$ almost everywhere on the interval (t_{l-1}, t_l) because the trajectory $\mathbf{x}(t)$ is confined to the surface S on the interval (t_{l-1}, t_l) . Thus, $\mathbf{x}(t)$ is an absolutely continuous function that satisfies

$$\tau \dot{\mathbf{x}}(t) \in G(\mathbf{x}(t)) \cap T(\mathbf{x}(t))$$

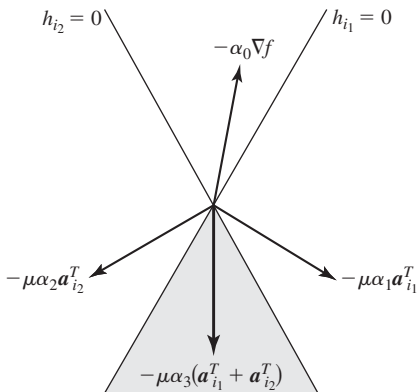


Figure 6.12 Components of the vector \mathbf{v} as given by (6.57) for sliding occurring on the intersection of two hyperplanes h_{i_1} and h_{i_2} .

almost everywhere on the interval (t_{l-1}, t_l) . We note that because the trajectory $\mathbf{x}(t)$ is confined to the surface S on the interval (t_{l-1}, t_l) , the set $G(\mathbf{x}(t)) \cap T(\mathbf{x}(t))$ is by assumption nonempty for all $t \in (t_{l-1}, t_l)$. In particular, we see that $G(\mathbf{x}(t')) \cap T(\mathbf{x}(t'))$ is the set of all \mathbf{v} that have the form (6.57) and also lie on the tangent plane $T(\mathbf{x}(t'))$. Let

$$\mathbf{r} = \sum_{j=1}^{2^k-1} \alpha_j \mathbf{B} \mathbf{u}_j.$$

Then, we represent (6.57) as

$$\mathbf{v} = -\alpha_0 \nabla f(\mathbf{x}(t')) - \mu \mathbf{r}. \quad (6.60)$$

Obviously, $\mathbf{r} \in \text{Span}\{\mathbf{a}_{i_1}^T, \dots, \mathbf{a}_{i_k}^T\}$, and thus \mathbf{r} is orthogonal to $T(\mathbf{x}(t'))$. However, $\mathbf{v} \in T(\mathbf{x}(t'))$ by assumption. Therefore, \mathbf{v} is the orthogonal projection of the vector $-\alpha_0 \nabla f(\mathbf{x}(t'))$ onto the tangent plane $T(\mathbf{x}(t'))$; that is,

$$\mathbf{v} = -\alpha_0 \mathbf{P} \nabla f(\mathbf{x}(t')), \quad (6.61)$$

where $\mathbf{P} = \mathbf{I}_n - \mathbf{B}(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T$.

Note that by virtue of assumptions A1 and A2, the set $\nabla f(\Omega)$ is compact. This implies the existence of a nonnegative constant, M , with the property that

$$M \geq \|\nabla f(\mathbf{x})\|_2 \quad \text{for all } \mathbf{x} \in \Omega.$$

Let

$$\gamma = \frac{\mu}{\mu + M\sqrt{k}\|(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T\|_2}.$$

Proposition 6.1

$$\gamma \leq \alpha_0 \leq 1.$$

Proof It follows from (6.58) and (6.59) that $\alpha_0 \leq 1$. Thus, it remains to show that $\alpha_0 \geq \gamma$. It follows from (6.60) and (6.61) that $\mu \mathbf{r}$ is the orthogonal projection of the vector $-\alpha_0 \nabla f(\mathbf{x}(t'))$ onto the subspace $T(\mathbf{x}(t'))^\perp$, where $T(\mathbf{x}(t'))^\perp$ is the orthogonal complement of $T(\mathbf{x}(t'))$ in \mathbb{R}^n . Therefore,

$$-\alpha_0(\mathbf{I}_n - \mathbf{P}) \nabla f(\mathbf{x}(t')) = \mu \sum_{j=1}^{2^k-1} \alpha_j \mathbf{B} \mathbf{u}_j. \quad (6.62)$$

Let $\mathbf{U} \in \mathbb{R}^{k \times (2^k-1)}$ be the matrix with columns $\mathbf{u}_1, \dots, \mathbf{u}_{2^k-1}$, and let $\boldsymbol{\alpha} \in \mathbb{R}^{2^k-1}$ be the column vector with the components $\alpha_1, \dots, \alpha_{2^k-1}$. Then, we can represent (6.62) as

$$-\frac{\alpha_0}{\mu}(\mathbf{I}_n - \mathbf{P}) \nabla f(\mathbf{x}(t')) = \mathbf{B} \mathbf{U} \boldsymbol{\alpha}. \quad (6.63)$$

Premultiplying both sides of (6.63) by $(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T$, we obtain

$$-\frac{\alpha_0}{\mu}(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T(\mathbf{I}_n - \mathbf{P}) \nabla f(\mathbf{x}(t')) = \mathbf{U} \boldsymbol{\alpha},$$

and hence

$$\frac{\alpha_0}{\mu} \|(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T (\mathbf{I}_n - \mathbf{P}) \nabla f(\mathbf{x}(t'))\|_2 = \|\mathbf{U}\alpha\|_2. \quad (6.64)$$

Applying the Cauchy–Schwarz inequality to (6.64) and using the fact that $\|\mathbf{I}_n - \mathbf{P}\| = 1$, we have

$$\frac{M\alpha_0}{\mu} \|(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T\|_2 \geq \|\mathbf{U}\alpha\|_2 \geq \frac{1}{\sqrt{k}} \|\mathbf{U}\alpha\|_1$$

since for $\mathbf{z} \in \mathbb{R}^k$, we have $\|\mathbf{z}\|_2 \geq \frac{1}{\sqrt{k}} \|\mathbf{z}\|_1$. However, using (6.58) and (6.59), we see that

$$\|\mathbf{U}\alpha\|_1 \geq \sum_{j=1}^{2^k-1} \alpha_j = 1 - \alpha_0.$$

Thus,

$$\frac{M\alpha_0\sqrt{k}}{\mu} \|(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T\|_2 \geq 1 - \alpha_0. \quad (6.65)$$

It follows from (6.65) that $\alpha_0 > 0$. Dividing both sides of (6.65) by α_0 yields

$$\frac{M\sqrt{k}}{\mu} \|(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T\|_2 \geq \frac{1 - \alpha_0}{\alpha_0} = \frac{1}{\alpha_0} - 1.$$

Hence,

$$\alpha_0 \geq \frac{1}{1 + \frac{M\sqrt{k}}{\mu} \|(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T\|_2} = \gamma.$$

The proof of the proposition is complete.

We have the following two corollaries.

Corollary 6.1

$$G(\mathbf{x}(t')) \cap T(\mathbf{x}(t')) \subset \{-\alpha \mathbf{P} \nabla f(\mathbf{x}(t')) : \gamma \leq \alpha \leq 1\}.$$

Proof This is directly from (6.61) and Proposition 6.1.

Corollary 6.2

$$\mathbf{P} \nabla f(\mathbf{x}(t')) = \mathbf{0} \quad \text{if and only if } \mathbf{0} \in G(\mathbf{x}(t')) \cap T(\mathbf{x}(t')).$$

Proof If $\mathbf{P} \nabla f(\mathbf{x}(t')) = \mathbf{0}$, then $\mathbf{0} \in G(\mathbf{x}(t')) \cap T(\mathbf{x}(t'))$ because, by Corollary 6.1,

$$\mathbf{P} \nabla f(\mathbf{x}(t')) \in G(\mathbf{x}(t')) \cap T(\mathbf{x}(t')).$$

Conversely, suppose that $\mathbf{0} \in G(\mathbf{x}(t')) \cap T(\mathbf{x}(t'))$. By Corollary 6.1, $-\alpha \mathbf{P} \nabla f(\mathbf{x}(t')) = \mathbf{0}$ for some $\alpha \in [\gamma, 1]$. Since $\gamma > 0$ for all $\mu > 0$, $-\alpha \mathbf{P} \nabla f(\mathbf{x}(t')) = \mathbf{0}$ is possible only if $\mathbf{P} \nabla f(\mathbf{x}(t')) = \mathbf{0}$, which completes the proof.

Proposition 6.2

$\mathbf{P} \nabla f(\mathbf{x}(t')) = \mathbf{0}$ if and only if $\mathbf{x}(t') \in \Gamma$.

Proof We first prove the sufficiency part of the proposition. Suppose $\mathbf{P} \nabla f(\mathbf{x}(t')) = \mathbf{0}$. It follows from (6.61) and (6.57) that

$$\alpha_0 \nabla f(\mathbf{x}(t')) + \mu \sum_{j=1}^{2^k-1} \alpha_j \mathbf{B} \mathbf{u}_j = \mathbf{0}. \quad (6.66)$$

Let

$$\mathbf{d} = \frac{\mu}{\alpha_0} \sum_{j=1}^{2^k-1} \alpha_j \mathbf{u}_j.$$

Note that because $\alpha_0 \neq 0$, \mathbf{d} is well-defined. It then follows from (6.66) that

$$\nabla f(\mathbf{x}(t')) + \mathbf{B} \mathbf{d} = \mathbf{0}.$$

We see from (6.58) and the definition of $\mathbf{u}_1, \dots, \mathbf{u}_{2^k-1}$ that each element of \mathbf{d} is non-negative. By Lemma 6.2, the point $\mathbf{x}(t')$ satisfies the Karush–Kuhn–Tucker conditions for problem (6.32) and hence $\mathbf{x}(t') \in \Gamma$. The sufficiency condition is thus proven.

We now prove the necessity condition. Suppose that $\mathbf{x}(t') \in \Gamma$. By Lemma 6.2, the point $\mathbf{x}(t')$ satisfies the Karush–Kuhn–Tucker optimality conditions for problem (6.32). This implies the existence of nonnegative constants μ_1, \dots, μ_k with the property that

$$\nabla f(\mathbf{x}(t')) + \mathbf{B} \boldsymbol{\mu} = \mathbf{0}, \quad (6.67)$$

where $\boldsymbol{\mu} = [\mu_1, \dots, \mu_k]^T$. Let

$$v = 1 + \frac{1}{\mu} \sum_{i=1}^k \mu_i, \quad (6.68)$$

$$\tilde{\alpha}_0 = \frac{1}{v}, \quad (6.69)$$

$$\tilde{\alpha}_{2^{i-1}} = \frac{\mu_i}{\mu v}, \quad i = 1, \dots, k, \quad (6.70)$$

$$\tilde{\alpha}_j = 0, \quad j = 1, \dots, 2^k - 1, \quad j \neq 2^{i-1} \quad \text{for } i = 1, \dots, k, \quad (6.71)$$

$$\tilde{\mathbf{v}} = -\tilde{\alpha}_0 \nabla f(\mathbf{x}(t')) - \mu \sum_{j=1}^{2^k-1} \tilde{\alpha}_j \mathbf{B} \mathbf{u}_j. \quad (6.72)$$

It follows from (6.68)–(6.71) that

$$\tilde{\alpha}_j \geq 0, \quad j = 0, 1, \dots, 2^k - 1 \quad \text{and} \quad \sum_{j=0}^{2^k-1} \tilde{\alpha}_j = 1.$$

Therefore, $\tilde{\mathbf{v}} \in G(\mathbf{x}(t'))$. Note that

$$\tilde{\mathbf{v}} = -\frac{\mu}{\mu + \sum_{i=1}^k \mu_i} (\nabla f(\mathbf{x}(t')) + \mathbf{B}\mu).$$

Using (6.67), we conclude that $\tilde{\mathbf{v}} = \mathbf{0}$, and thus $\tilde{\mathbf{v}} \in T(\mathbf{x}(t'))$. Hence,

$$\mathbf{0} \in G(\mathbf{x}(t')) \cap T(\mathbf{x}(t')).$$

Using Corollary 6.2, we get $\mathbf{P}\nabla f(\mathbf{x}(t')) = \mathbf{0}$, which completes the proof of the necessary condition and the proof of the proposition.

Applying now the chain rule to F , taking into account (6.61), Proposition 6.1, and the fact that $\mathbf{P} = \mathbf{P}^T = \mathbf{P}^2$, we obtain

$$\frac{d}{dt}F(\mathbf{x}(t)) \leq -\frac{\gamma}{\tau} \|\mathbf{P}\nabla f(\mathbf{x}(t))\|_2^2 \quad (6.73)$$

almost everywhere on the interval (t_{l-1}, t_l) . By Proposition 6.2,

$$\frac{d}{dt}F(\mathbf{x}(t)) = 0 \quad \text{if and only if } \mathbf{x}(t) \in \Gamma, \quad (6.74)$$

and

$$\frac{d}{dt}F(\mathbf{x}(t)) < 0 \quad \text{if and only if } \mathbf{x}(t) \notin \Gamma, \quad (6.75)$$

almost everywhere on the interval (t_{l-1}, t_l) , which concludes our analysis for the second case.

We now show that any trajectory of system (6.34) that is confined to Ω converges asymptotically to the region Γ .

Theorem 6.4 Let $\mathbf{x}(t)$ be any trajectory of system (6.34). If $\mathbf{x}(t) \in \Omega$ for all t , then the trajectory $\mathbf{x}(t)$ converges asymptotically to Γ .

Proof Recall that $F(\mathbf{x}) = 0$ if $\mathbf{x} \in \Gamma$ and $F(\mathbf{x}) > 0$ if $\mathbf{x} \in \Omega \setminus \Gamma$. Therefore, to prove the theorem it is enough to show that

$$\lim_{t \rightarrow \infty} F(\mathbf{x}(t)) = 0,$$

or equivalently, that for any $\varepsilon > 0$ there exists a finite time, $T(\mathbf{x}(0), \varepsilon)$, such that $F(\mathbf{x}(t)) < \varepsilon$ for all $t \geq T(\mathbf{x}(0), \varepsilon)$. We see from (6.56) and (6.75) that

$$\frac{d}{dt} F(\mathbf{x}(t)) < 0$$

for almost all t when $\mathbf{x}(t) \notin \Gamma$. Consequently $F(\mathbf{x}(t)) \leq F(\mathbf{x}(0))$ for all $t \geq 0$. Therefore, we only need to consider the case when $F(\mathbf{x}(0)) \geq \varepsilon$. Let

$$\theta(\varepsilon) = \left\{ \mathbf{x} : \mathbf{x} \in \Omega, F(\mathbf{x}) \geq \frac{\varepsilon}{2} \right\}.$$

Because $F(\mathbf{x}(0)) \geq \varepsilon$, we have $\mathbf{x}(0) \in \theta(\varepsilon)$ and $\theta(\varepsilon) \subset \Omega \setminus \Gamma$. Furthermore, because the function F is continuous, the set $\theta(\varepsilon)$ is compact. We will now show that there exists a positive constant, $\rho(\varepsilon)$, such that

$$\frac{d}{dt} F(\mathbf{x}(t)) \leq -\rho(\varepsilon)$$

for almost all t when $\mathbf{x}(t) \in \theta(\varepsilon)$, which guarantees that $F(\mathbf{x}(t)) \leq \varepsilon/2$ for all

$$t \geq \frac{F(\mathbf{x}(0)) - \varepsilon/2}{\rho(\varepsilon)}.$$

We analyze two cases. The first case is when the trajectory $\mathbf{x}(t)$ is in $\theta(\varepsilon)$ and is not in a sliding mode on any of the surfaces h_i . Then, we analyze the case when the trajectory is in $\theta(\varepsilon)$ in a sliding mode on a surface that is an intersection of some of the hyperplanes h_i .

Case 1 Because the function f is convex and $f \in C^1$, we have $\nabla f(\mathbf{x}) \neq \mathbf{0}$ for all $\mathbf{x} \in \Omega \setminus \Gamma$. Hence,

$$-\frac{1}{\tau} \|\nabla f(\mathbf{x})\|_2^2 < 0 \quad \text{for all } \mathbf{x} \in \theta(\varepsilon) \quad (6.76)$$

because $\theta \subset \Omega \setminus \Gamma$ and $\nabla f(\mathbf{x}) \neq \mathbf{0}$ in $\theta(\varepsilon)$. The fact that the function $-1/\tau \|\nabla f(\mathbf{x})\|_2^2$ is continuous on the compact set $\theta(\varepsilon)$ implies the existence of a positive constant, $\rho_1(\varepsilon)$, such that

$$-\frac{1}{\tau} \|\nabla f(\mathbf{x})\|_2^2 \leq -\rho_1(\varepsilon) \quad \text{for all } \mathbf{x} \in \theta(\varepsilon).$$

Combining the above with (6.54) yields

$$\frac{d}{dt} F(\mathbf{x}(t)) \leq -\rho_1(\varepsilon)$$

for almost all t when $\mathbf{x}(t) \in \theta(\varepsilon)$, which concludes our analysis of the first case.

Case 2 Let i_1, \dots, i_k denote the elements of I' , and let $S = \bigcap_{i \in I'} h_i$ denote the surface to which the trajectory $\mathbf{x}(t)$ is confined on the interval (t_{l-1}, t_l) . Suppose that $S \cap \theta(\varepsilon) \neq \emptyset$. Clearly,

$$S \cap \theta(\varepsilon) \subset \Omega \setminus \Gamma.$$

Using arguments similar to those of Glazos et al. [98, p. 692], we can show that

$$-\frac{1}{\tau} \|\mathbf{P} \nabla f(\mathbf{x})\|_2^2 < 0 \quad \text{for all } \mathbf{x} \in S \cap \theta(\varepsilon), \quad (6.77)$$

where \mathbf{P} is the orthogonal projector onto the tangent plane to the surface S . By the above and the fact that the function $-\frac{1}{\tau} \|\mathbf{P} \nabla f(\mathbf{x})\|_2^2$ is continuous on the compact set $S \cap \theta(\varepsilon)$, there exists a positive constant, $\hat{\rho}(\varepsilon)$, such that

$$-\frac{1}{\tau} \|\mathbf{P} \nabla f(\mathbf{x})\|_2^2 < -\hat{\rho}(\varepsilon) \quad \text{for all } \mathbf{x} \in S \cap \theta(\varepsilon).$$

It then follows directly from (6.73) that there exists a positive constant $\hat{\gamma}$ such that

$$\frac{d}{dt} F(\mathbf{x}(t)) \leq -\hat{\gamma} \hat{\rho}(\varepsilon)$$

for almost all t when $\mathbf{x}(t) \in S \cap \theta(\varepsilon)$. The above along with the fact that the number of constraints is finite implies the existence of a positive constant, $\rho_2(\varepsilon)$, such that

$$\frac{d}{dt} F(\mathbf{x}(t)) \leq -\rho_2(\varepsilon) \quad (6.78)$$

for almost all t when $\mathbf{x}(t) \in S \cap \theta(\varepsilon)$. This concludes our analysis for the second case.

Let

$$\rho(\varepsilon) = \min\{\rho_1(\varepsilon), \rho_2(\varepsilon)\}.$$

Then, using (6.77) and (6.78) yields

$$\frac{d}{dt} F(\mathbf{x}(t)) \leq -\rho(\varepsilon)$$

for almost all t when $\mathbf{x}(t) \in \theta(\varepsilon)$. Hence, $F(\mathbf{x}(t)) \leq \varepsilon/2 < \varepsilon$ for all

$$t \geq \frac{F(\mathbf{x}(0)) - \varepsilon/2}{\rho(\varepsilon)}.$$

Therefore,

$$T(\mathbf{x}(0), \varepsilon) = \begin{cases} 0 & \text{if } F(\mathbf{x}(0)) < \varepsilon, \\ \frac{F(\mathbf{x}(0)) - \varepsilon/2}{\rho(\varepsilon)} & \text{if } F(\mathbf{x}(0)) \geq \varepsilon. \end{cases}$$

The proof of the theorem is complete.

Theorem 6.5 Every trajectory of system (6.34), irrespective of the initial condition, $\mathbf{x}(0)$, converges to a solution of problem (6.32).

Proof This statement follows directly from Theorems 6.3 and 6.4.

◆ **Example 6.5** (Based on Glazos et al. [98])

Consider the quadratic programming problem

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^2 \\ & \text{subject to} \quad \mathbf{A} \mathbf{x} \leq \mathbf{b}, \end{aligned}$$

where

$$\mathbf{Q} = \begin{bmatrix} 4 & 1 \\ 1 & 2 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} -12 \\ -10 \end{bmatrix}, \quad \mathbf{A} = \begin{bmatrix} 5 & 1 \\ -5 & 1 \\ -1 & -2 \end{bmatrix}, \quad \text{and} \quad \mathbf{b} = \begin{bmatrix} 0 \\ 10 \\ 10 \end{bmatrix}.$$

The above optimization problem clearly satisfies assumptions A1–A4, and one can easily verify that the point $\mathbf{x}^* = [-19/22 \ 95/22]^T$ is a unique minimizer for the problem, that is, $\Gamma = \{\mathbf{x}^*\}$. Note that $\mathbf{x}^* \in h_1$. A phase-plane portrait for system (6.34) that solves the above optimization problem is shown in Figure 6.13. The feasible region is depicted by the shaded area. Also shown are level sets of the quadratic objective function.

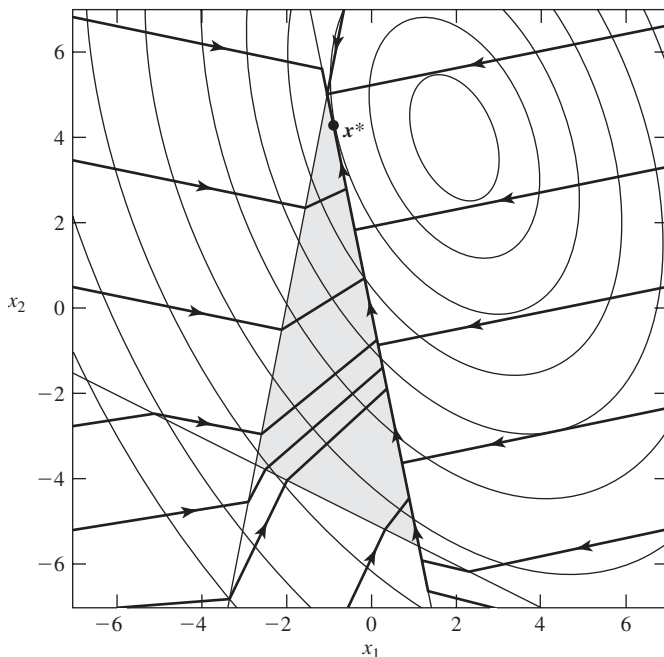


Figure 6.13 Phase-plane portrait for the system in Example 6.5.

Notes

The idea of controlling a dynamical system by varying its structure was originated in the 1950s. Early pioneers in the field of variable structure systems were S. V. Emelyanov and his

co-workers at the Institute of Control Sciences in Moscow in the former Soviet Union. A large body of literature was produced by these studies; unfortunately, most of the work was published in Russian only. The titles of books and papers on the subject of variable structure control systems published in the years 1957–1989 were collected by Emelyanov [77]. The first widely accessible books in English on variable structure sliding mode systems were by Itkis [140] and Utkin [285]. A survey paper by Utkin [284] refers to many early contributions to the field; see also a tutorial by DeCarlo, Žak, and Matthews [60] and a survey by Hung, Gao, and Hung [135]. At about the same time, A. F. Filippov was developing his theory of differential equations with discontinuous right-hand sides [83]. Filippov's approach is now accepted as the tool for the analysis of sliding mode systems. Sliding modes also found applications in the state estimation of dynamical uncertain systems as we show in Section 6.8. For further analysis of discontinuous state estimators, we refer the reader to Edwards and Spurgeon [71, 73]. For a stability analysis of uncertain closed-loop systems with the combined variable structure estimator–controller compensators we refer to Hui and Žak [131]. Output feedback sliding mode controllers are presented in references 64, 72, 133, 174, 310, and 311.

In 1967, Karpinskaya [150] identified the occurrence of sliding modes in an analog circuit for solving optimization problems. The idea of using analog circuits to solve constrained optimization problems first appeared in the literature in 1956 when Pyne [239] proposed a method for solving linear programming problems using an electronic analog computer. Soon after, other methods were proposed for solving various mathematical programming problems on an analog computer—see, for example, Rybashov [248–251] and Karpinskaya [150], among others. More recently, many new analog architectures have been proposed for solving constrained optimization problems. In particular, Chua and Lin [51] developed a canonical nonlinear programming circuit for solving general nonlinear programming problems. Tank and Hopfield [278] later proposed an analog circuit for solving linear programming problems. The results of Tank and Hopfield were then extended to nonlinear programming problems by Kennedy and Chua [155]. Rodríguez-Vázquez et al. [242] and Cichocki and Unbehauen [52] proposed analog networks that use switched-capacitor integrated circuit technology for solving certain constrained optimization problems. We analyzed dynamical behavior of a neural network of Rodríguez-Vázquez et al. [242] for solving a class of convex programming problems. In carrying out the analysis, we used concepts from the theory of differential equations with discontinuous right-hand sides and Lyapunov stability theory. Our method can also be used to analyze other classes of analog dynamical optimizers whose designs are based on exact penalty functions. An open problem is to extend the results obtained to a more general class of mathematical programming problems. Preliminary results in this direction can be found in Glazos et al. [98] and Chong et al. [49].

EXERCISES

- 6.1** Consider the following variable structure system that was analyzed in Hung et al. [135] (see also Utkin [284] for a similar system):

$$\begin{aligned}\dot{x}_1 &= x_2, \\ \dot{x}_2 &= -x_1 + 2x_2 + u,\end{aligned}\tag{6.79}$$

where

$$u = -4x_1 \operatorname{sign} \sigma(x_1, x_2), \quad \sigma(x_1, x_2) = x_1(0.5x_1 + x_2).$$

- (a) Sketch phase plane portraits of two linear subsystems that constitute the given variable structure system (6.79).
 (b) Sketch a phase-plane portrait of the variable structure system (6.79).

6.2 (a) For the dynamical system model,

$$\dot{\mathbf{x}} = \begin{bmatrix} 1 & 1 \\ 2 & 1 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 1 \\ -1 \end{bmatrix} u,$$

construct the switching surface

$$\sigma(\mathbf{x}) = \mathbf{s}\mathbf{x} = 0$$

such that the system in sliding along the surface has its pole at -2 .

- (b) Consider the following model of a dynamical system:

$$\dot{\mathbf{x}} = \begin{bmatrix} 2 & 0 & -1 \\ 0 & -2 & 2 \\ -1 & -2 & 3 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} u.$$

Design a switching surface $\sigma(\mathbf{x}) = \mathbf{s}\mathbf{x} = 0$ such that the system restricted to it has poles at $-4, -5$.

- (c) Given the following dynamical system model,

$$\dot{\mathbf{x}} = \begin{bmatrix} 2 & 0 & -1 \\ 0 & -2 & 2 \\ -1 & -2 & 3 \end{bmatrix} \mathbf{x} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} u,$$

along with the switching surface

$$\sigma(\mathbf{x}) = \mathbf{s}\mathbf{x} = [8 \quad 8 \quad -15]\mathbf{x} = 0.$$

Determine the reduced-order equations that describe the system behavior in a sliding mode.

6.3 Consider a dynamical system modeled by the equations

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ a_1(t) & a_2(t) \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u,$$

where

$$|a_1(t)| \leq 3, \quad |a_2(t)| \leq 4 \quad \text{for all } t.$$

- (a) Design a switching surface

$$\sigma(x_1, x_2) = s_1 x_1 + x_2 = 0$$

such that the system restricted to it has its pole at -2 .

(b) Design a sliding mode controller

$$u = k_1^\pm x_1 + k_2^\pm x_2,$$

so that the switching surface $\sigma = 0$ is globally attractive. The gains k_1 and k_2 are functions of $x_1\sigma$ and $x_2\sigma$, respectively.

6.4 Consider the following model of an armature controlled DC motor controlling an inverted pendulum:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} x_2 \\ 9.8 \sin x_1 + x_3 \\ -\frac{1}{L_a} x_2 - \frac{1}{L_a} x_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \frac{1}{L_a} \end{bmatrix} u$$

$$y = x_1,$$

where L_a is the armature inductance of the motor.

- (a) Derive the reduced-order model assuming L_a is negligible.
 (b) For the obtained reduced-order model from the first part, design the switching surface $\sigma = 0$ of the form

$$\sigma = \sigma(x_1, x_2) = [s_1 \quad 1] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$$

so that the system in sliding will have its pole at -3 .

- (c) Use the above sliding surface to design a variable structure controller of the form

$$u = -k_1 x_1 \operatorname{sign}(x_1 \sigma) - k_2 x_2 \operatorname{sign}(x_2 \sigma).$$

The smaller the gains k_1, k_2 , the better.

6.5 Consider a dynamical system model,

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} A_1 & A_2 \\ A_3 & A_4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ b \end{bmatrix} u, \quad (6.80)$$

$$y = [C_1 \quad C_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad (6.81)$$

where b and C_2 are nonzero scalars. Let

$$\sigma = [S_1 \quad S_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = F [C_1 \quad C_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = Fy, \quad (6.82)$$

where S_2 and F are nonzero scalars.

- (a) Find the reduced-order system that models (6.80) when restricted to the surface $\sigma = 0$.
 (b) Find the transfer function $H(s) = \mathcal{L}(y(t))/\mathcal{L}(u(t)) = Y(s)/U(s)$. What is the relation between the zeros of $H(s)$ and the poles of the reduced-order system found in part (a)?

6.6 (Based on Kwan [174]) Prove the following lemma.

Lemma 6.4 Let

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u},$$

where the matrix \mathbf{A} has negative real distinct eigenvalues, that is,

$$\text{eig}(\mathbf{A}) = \{-\lambda_1, -\lambda_2, \dots, -\lambda_n\}, \quad \lambda_i > 0, \lambda_i \neq \lambda_j \text{ for } i \neq j.$$

Let $\lambda_{\min} = \min\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ and K be a positive constant such that

$$\|e^{\mathbf{A}t}\| \leq K e^{-\lambda_{\min} t}$$

Then,

$$\|\mathbf{x}(t)\| \leq w(t),$$

where $w(t)$ is a solution of

$$\dot{w} = -\lambda_{\min} w + K \|\mathbf{B}\| \|\mathbf{u}\|, \quad \text{subject to } w(0) \geq K \|\mathbf{x}(0)\|.$$

6.7 (Based on Kwan [174]) Consider an n th-order model of an uncertain dynamical system,

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 \\ \mathbf{A}_3 & \mathbf{A}_4 \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} \mathbf{0} \\ b \end{bmatrix} (u + \xi), \quad (6.83)$$

$$y = [\mathbf{C}_1 \quad \mathbf{C}_2] \begin{bmatrix} \mathbf{x}_1 \\ x_2 \end{bmatrix},$$

where b and \mathbf{C}_2 are nonzero scalars and ξ models system's uncertainty, where

$$\xi = \mathbf{D}_1 \mathbf{x}_1 + D_2 x_2 + D_3 u + D_4(t).$$

In the above, \mathbf{D}_1 is a row vector whose size is equal to the dimension of $\mathbf{x}_1 \in \mathbb{R}^{n-1}$, and D_2 and D_3 are scalars, where $|D_3| < 1$ and D_4 is an unknown bounded function of time whose bound is known to us; that is, $|D_4(t)| \leq d_4$, where d_4 is known to us. Let $\sigma = Fy$ be as defined in Exercise 6.5. Assume that the switching surface,

$$\sigma = Fy = F[\mathbf{C}_1 \quad \mathbf{C}_2] \begin{bmatrix} \mathbf{x}_1 \\ x_2 \end{bmatrix} = 0,$$

was chosen in such a way that the nominal system (6.83) restricted to it has the eigenvalues

$$\{-\lambda_1, -\lambda_2, \dots, -\lambda_{n-1}\}, \quad \lambda_i > 0, \lambda_i \neq \lambda_j \text{ for } i \neq j.$$

(a) Evaluate $\dot{\sigma}$ on the trajectories of the uncertain system model (6.83).

(b) For the controller

$$u = (S_2 b)^{-1} (1 - |D_3|)^{-1} (-Hw(t) - G|\sigma| - \delta) \text{sign}(\sigma),$$

where $w(t)$ is a solution of a first-order differential equation such that $\|\mathbf{x}_1(t)\| \leq w(t)$, find lower bounds for H , G , and δ such that

$$\sigma \dot{\sigma} \leq -\eta |\sigma|$$

for some $\eta > 0$. Use the lemma of Exercise 6.6 to find the form of the differential equation whose solution satisfies the condition $\|\mathbf{x}_1(t)\| \leq w(t)$.

6.8 For the dynamical system model

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{b}(u + f(y)),$$

$$y = \mathbf{c}\mathbf{x},$$

where

$$\mathbf{A} = \begin{bmatrix} -2 & 2 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad \mathbf{c} = [0 \quad 1], \quad \text{and} \quad |f(y)| \leq 10y,$$

design a discontinuous state-estimator. The estimator poles should be placed at -2 and -3 .

Hint After constructing the required estimator gain \mathbf{L} , find a symmetric positive definite matrix \mathbf{Q} such that the solution, \mathbf{P} , to the Lyapunov matrix equation, $(\mathbf{A} - \mathbf{L}\mathbf{c})^T \mathbf{P} + \mathbf{P}(\mathbf{A} - \mathbf{L}\mathbf{c}) = -\mathbf{Q}$, satisfies the condition, $\mathbf{b}^T \mathbf{P} = \mathbf{c}$.

6.9 The notion of *feasible direction*, which we present next, plays an important role in the analysis of constrained optimization problems such as that given by (6.36). We say that a vector \mathbf{d} is a feasible direction at $\mathbf{x} \in \Omega$ if there is an $\bar{\alpha} > 0$ such that $\mathbf{x} + \alpha \mathbf{d} \in \Omega$ for all $\alpha \in [0, \bar{\alpha}]$. Suppose now that the minimizer \mathbf{x}^* of the optimization problem (6.36) lies on the boundary of Ω and there is a feasible direction \mathbf{d} at \mathbf{x}^* such that $\mathbf{d}^T \nabla f(\mathbf{x}^*) > 0$. Show, using contraposition, that if p is an exact penalty function for the problem, then p is not differentiable at \mathbf{x}^* (see [50]).

6.10 Consider an uncertain, or nonlinear, dynamical system model,

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{f}(t, \mathbf{x}(t)) + \mathbf{B}(\mathbf{u}(t) + \mathbf{h}(t, \mathbf{x}(t), \mathbf{u}(t))), \quad (6.84)$$

where $\mathbf{x}(t) \in \mathbb{R}^n$, $\mathbf{u}(t) \in \mathbb{R}^m$, the constant matrices $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times m}$ are known, and the function \mathbf{f} and \mathbf{h} model uncertainties, or nonlinearities, in the system. We refer to \mathbf{h} as the matched uncertainty, while we refer to \mathbf{f} as the unmatched uncertainty. The reason for calling \mathbf{h} the matched uncertainty is that it affects the system dynamics through the column space of the matrix \mathbf{B} in a similar fashion as the control \mathbf{u} does. In other words, \mathbf{h} matches the control action \mathbf{u} . We assume that

1. $\|\mathbf{f}(t, \mathbf{x})\| \leq \alpha_f \|\mathbf{x}\| + \beta_f$ and
2. $\|\mathbf{h}(t, \mathbf{x}, \mathbf{u})\| \leq \gamma_h \|\mathbf{u}\| + \alpha_h \|\mathbf{x}\| + \beta_h$,

where α_f , β_f , γ_h , α_h , and β_h are known nonnegative constants, and $\|\cdot\| = \|\cdot\|_2$ denotes the Euclidean norm for vectors and the induced spectral norm for matrices. We also assume that

3. the pair (\mathbf{A}, \mathbf{B}) is controllable and
4. $\text{rank } \mathbf{B} = m$.

Let

$$\tilde{\mathbf{x}} = \mathbf{M}\mathbf{x} = \begin{bmatrix} \mathbf{W}^g \\ \mathbf{S} \end{bmatrix} \mathbf{x},$$

where \mathbf{W}^g and \mathbf{S} are the matrices obtained from the Switching Surface Design Algorithm. Transform the uncertain system model (6.84) into the $\tilde{\mathbf{x}}$ coordinates.

- 6.11** The presence of unmatched uncertainties has a deteriorating effect on the stability properties of a closed-loop system. The above combined with the imposed hard bounds on the control action may make it impossible to globally asymptotically stabilize the closed-loop system. The best it can be achieved in such a situation may be UUB with some region of attraction $\Gamma \subset \mathbb{R}^n$ and uniformly stable set Σ containing the origin of \mathbb{R}^n . The goal of this exercise is to estimate the “size” of the set Γ and that of Σ .

It follows from Exercise 6.10 that the uncertain system (6.84) confined to the switching surface $\sigma = \mathbf{S}\mathbf{x} = \mathbf{0}$ is described, in the new coordinates, by

$$\dot{\mathbf{z}} = \mathbf{A}_{11}\mathbf{z} + \mathbf{W}^g \tilde{\mathbf{f}}(t, \mathbf{z}, \mathbf{0}), \quad (6.85)$$

where $\|\tilde{\mathbf{f}}(t, \mathbf{z}, \mathbf{0})\| \leq \alpha_f \|\mathbf{W}\|\|\mathbf{z}\| + \beta_f$. We refer to the above system as the *reduced-order system*. The matrix \mathbf{A}_{11} is asymptotically stable by construction. Therefore, for any symmetric and positive definite matrix $\mathbf{Q} \in \mathbb{R}^{(n-m) \times (n-m)}$ the solution \mathbf{P} to the Lyapunov equation

$$\mathbf{A}_{11}^T \mathbf{P} + \mathbf{P} \mathbf{A}_{11} = -2\mathbf{Q}$$

is (symmetric) positive definite. Consider now the positive definite function:

$$V_1 = (\mathbf{z}^T \mathbf{P} \mathbf{z})^{1/2}.$$

The matrix \mathbf{P} is symmetric positive definite, hence it can be represented as

$$\mathbf{P} = \mathbf{P}^{1/2} \mathbf{P}^{1/2},$$

where $\mathbf{P}^{1/2} \in \mathbb{R}^{(n-m) \times (n-m)}$ is symmetric and nonsingular. We refer to $\mathbf{P}^{1/2}$ as a square root of \mathbf{P} —see page 635 for a discussion of a matrix square root. Let

$$\boldsymbol{\zeta} = \mathbf{P}^{1/2} \mathbf{z}.$$

Then,

$$V_1 = (\mathbf{z}^T \mathbf{P} \mathbf{z})^{1/2} = (\boldsymbol{\zeta}^T \boldsymbol{\zeta})^{1/2} = \|\boldsymbol{\zeta}\|.$$

Evaluate the time derivative of V_1 on the trajectories of the reduced-order system (6.85). Specifically, let

$$\rho = \frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{P})} - \alpha_f \frac{\lambda_{\max}(\mathbf{P})}{\lambda_{\min}(\mathbf{P})} \|\mathbf{W}^g\| \|\mathbf{W}\|,$$

and

$$r = \frac{\beta_f}{\rho} \frac{\lambda_{\max}(\mathbf{P})}{\lambda_{\min}^{\frac{1}{2}}(\mathbf{P})} \|\mathbf{W}^g\|.$$

Assume that $\rho > 0$. Using the above notation, show that

$$\frac{d}{dt} \|\zeta\| \leq -\rho(\|\zeta\| - r).$$

6.12 Prove the following theorem concerning stability of the motion on the switching surface.

Theorem 6.6 The reduced-order system, analyzed in Exercise 6.11,

$$\dot{\mathbf{z}} = \mathbf{A}_{11} \mathbf{z} + \mathbf{W}^g \tilde{\mathbf{f}}(t, \mathbf{z}, \mathbf{0})$$

is globally UUB, with the uniformly stable set

$$\Sigma = \{ \mathbf{z}: \|\zeta\| = \|\mathbf{P}^{1/2} \mathbf{z}\| \leq r \} \subset \mathbb{R}^{n-m}.$$

In the proof of the above theorem, use the following lemma, which we recommend that you prove before using it.

Lemma 6.5 Suppose $\phi(t)$, $a, b \in \mathbb{R}$, and $a \neq 0$. If

$$\dot{\phi}(t) - a\phi(t) \leq -b,$$

then, for $t \geq t_0$,

$$\phi(t) \leq \frac{b}{a} + \left(\phi(t_0) - \frac{b}{a} \right) e^{(t-t_0)}.$$

6.13 Let $\sigma = \mathbf{S}\mathbf{x}$, $\|\cdot\|_p$ be the p norm of a vector, and let ∇_σ denote the gradient operation with respect to σ . Consider the following class of controllers:

$$\mathbf{u}^{(p)} = -\eta \nabla_\sigma \|\sigma\|_p, \quad 1 \leq p \leq \infty \quad (6.86)$$

where

$$\nabla_\sigma \|\sigma\|_p = \begin{bmatrix} \frac{\partial}{\partial \sigma_1} \|\sigma\|_p & \frac{\partial}{\partial \sigma_2} \|\sigma\|_p & \cdots & \frac{\partial}{\partial \sigma_m} \|\sigma\|_p \end{bmatrix}^T.$$

(a) Write down the expressions for the controllers when $p = 1, 2, \infty$.

(b) Find the bounds of $\mathbf{u}^{(p)}$ for $p = 1, 2, \infty$.

6.14 The goal of this exercise is to obtain certain estimates that can be used to investigate stability of the closed-loop system

$$\begin{aligned} \dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{f} + \mathbf{B}(\mathbf{u} + \mathbf{h}), \\ \mathbf{u} &= \mathbf{u}^{(p)} = -\eta \nabla_\sigma \|\sigma(\mathbf{x})\|_p \quad \text{for } p = 1, 2, \text{ or } \infty. \end{aligned} \quad (6.87)$$

Let

$$\begin{bmatrix} \zeta \\ \sigma \end{bmatrix} = \begin{bmatrix} \mathbf{P}^{1/2} \mathbf{W}^g \\ \mathbf{S} \end{bmatrix} \mathbf{x}.$$

Note that

$$\mathbf{x} = [\mathbf{W}\mathbf{P}^{-1/2} \quad \mathbf{B}] \begin{bmatrix} \boldsymbol{\zeta} \\ \boldsymbol{\sigma} \end{bmatrix}.$$

In the new coordinates the system takes the form

$$\begin{bmatrix} \dot{\boldsymbol{\zeta}} \\ \dot{\boldsymbol{\sigma}} \end{bmatrix} = \begin{bmatrix} \mathbf{P}^{1/2}\mathbf{A}_{11}\mathbf{P}^{-1/2} & \mathbf{P}^{1/2}\mathbf{A}_{12} \\ \mathbf{A}_{21}\mathbf{P}^{-1/2} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \boldsymbol{\zeta} \\ \boldsymbol{\sigma} \end{bmatrix} + \begin{bmatrix} \mathbf{O} \\ \mathbf{I}_m \end{bmatrix} (\mathbf{u} + \tilde{\mathbf{h}}) + \begin{bmatrix} \mathbf{P}^{1/2}\mathbf{W}^g\tilde{\mathbf{f}} \\ \mathbf{S}\tilde{\mathbf{f}} \end{bmatrix}. \quad (6.88)$$

Obtain estimates for $\frac{d}{dt}\|\boldsymbol{\sigma}\|$ ($\boldsymbol{\sigma} \neq \mathbf{0}$) and $\frac{d}{dt}\|\boldsymbol{\zeta}\|$ ($\boldsymbol{\zeta} \neq \mathbf{0}$) evaluated on any solution of the closed-loop system (6.88). Specifically, show that under some assumptions we obtain

$$\frac{d}{dt}\|\boldsymbol{\sigma}\| \leq a_{21}\|\boldsymbol{\zeta}\| + a_{22}\|\boldsymbol{\sigma}\| - \mu^{(p)}$$

and

$$\frac{d}{dt}\|\boldsymbol{\zeta}\| \leq -\rho\|\boldsymbol{\zeta}\| + a_{12}\|\boldsymbol{\sigma}\| + \tilde{\beta}_f$$

where

$$\begin{aligned} a_{21} &= \|\mathbf{A}_{21}\mathbf{P}^{-1/2}\| + \alpha_h\|\mathbf{W}\mathbf{P}^{-1/2}\| + \alpha_f\|\mathbf{S}\|\|\mathbf{W}\mathbf{P}^{-1/2}\|, \\ a_{22} &= \|\mathbf{A}_{22}\| + \alpha_h\|\mathbf{B}\| + \alpha_f\|\mathbf{S}\|\|\mathbf{B}\|, \\ a_{12} &= \|\mathbf{P}^{1/2}\mathbf{A}_{12}\| + \alpha_f\|\mathbf{P}^{1/2}\mathbf{W}^g\|\|\mathbf{B}\|, \\ \tilde{\beta}_f &= \beta_f\|\mathbf{P}^{1/2}\mathbf{W}^g\|, \\ \mu^{(1)} &= \eta(1 - \gamma_h\sqrt{m}) - \beta_h - \beta_f\|\mathbf{S}\|, \\ \mu^{(2)} &= \eta(1 - \gamma_h) - \beta_h - \beta_f\|\mathbf{S}\|, \\ \mu^{(\infty)} &= \frac{\eta}{\sqrt{m}}(1 - \gamma_h\sqrt{m}) - \beta_h - \beta_f\|\mathbf{S}\|. \end{aligned}$$

6.15 In the following lemma, a definition of sliding mode domain is used, which we provide below. This definition can be found, for example, in Utkin [284, p. 218].

Definition 6.4 A domain Δ in the manifold $\{\mathbf{x} : \boldsymbol{\sigma}(\mathbf{x}) = \mathbf{0}\}$ is a sliding mode domain if for each $\varepsilon > 0$ there exists a $\delta > 0$ such that any trajectory starting in the n -dimensional δ -neighborhood of Δ may leave the n -dimensional ε neighborhood of Δ only through the n -dimensional ε neighborhood of the boundaries of Δ .

The previous definition is illustrated in Figure 6.14. Prove the following lemma:

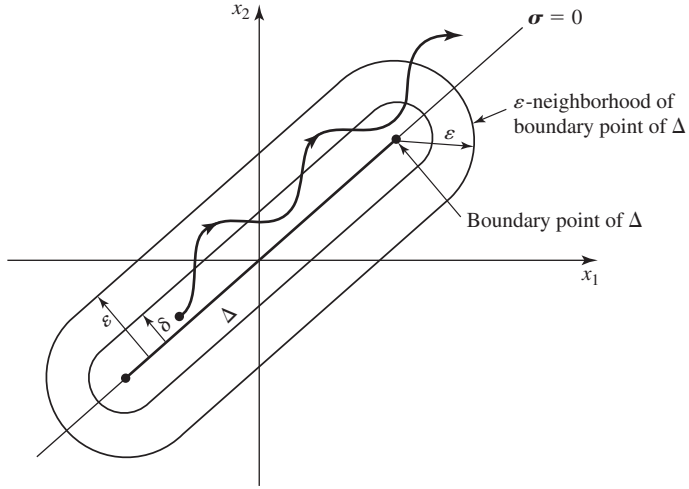


Figure 6.14 Two-dimensional illustration of a sliding mode domain.

Lemma 6.6 If $a_{21} \neq 0$, then the region

$$\Delta = \left\{ (\zeta, \sigma) : \|\zeta\| < \frac{\mu^{(p)}}{a_{21}}, \sigma = 0 \right\}$$

is a sliding mode domain for the closed-loop system (6.87) of Exercise 6.14.

The notation of the lemma is explained in Exercise 6.14.

6.16 Prove the following result concerning UUB of the closed-loop system (6.87) of Exercise 6.14. Consult Glazos and Žak [99] for a method of proving this result.

Theorem 6.7 Let $a_{21} \neq 0$ and

$$\Gamma = \{(\zeta, \sigma) : \alpha \|\zeta\| + \beta \|\sigma\| < \gamma\},$$

where

$$\alpha = \frac{\beta a_{21}}{\beta + \rho},$$

$$\beta = \frac{1}{2} (a_{22} - \rho + \sqrt{(a_{22} + \rho)^2 + 4a_{12}a_{21}}),$$

$$\gamma = \mu^{(p)} - \frac{\alpha \tilde{\beta}_f}{\beta}.$$

Then, the closed-loop system is UUB with respect to the uncertainty region

$$\Sigma \times \{\mathbf{0}\} = \{(\boldsymbol{\zeta}, \mathbf{0}) : \|\boldsymbol{\zeta}\| \leq r\}$$

with region of attraction Γ , where

$$r = \frac{\beta_f \lambda_{\max}(\mathbf{P}) \|\mathbf{W}^g\|}{\lambda_{\min}^{1/2}(\mathbf{P}) \rho}.$$

The set Σ is defined in Exercise 6.12.

For an illustration of the above mentioned region see Figure 6.15.

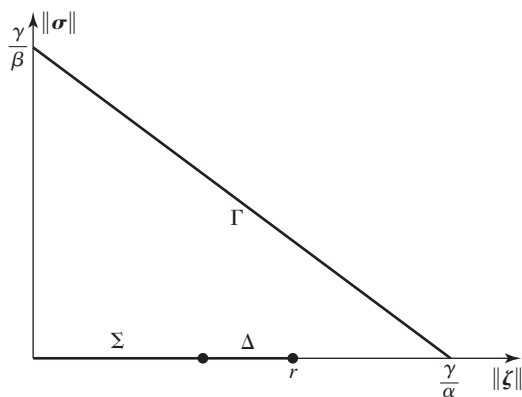


Figure 6.15 Illustration of the region of attraction and the uncertainty region of Theorem 6.7.

6.17 Show, using Lemma 6.3, that $\mathbf{P}\tilde{\mathbf{a}} \neq \mathbf{0}$, where \mathbf{P} and $\tilde{\mathbf{a}}$ are defined on page 347.



CHAPTER

7

Vector Field Methods

1. If your postage stamps get stuck together, place them in the freezer for an hour or two. They'll usually come apart and be usable.
2. To remove a stamp from an envelope for collecting purposes, put it in the freezer overnight; then take it out and slide a knife between it and the envelope.
3. Or cut off the corner of the envelope with the stamp and float—don't immerse—it in lukewarm water until the paper comes off. Dry the stamp on blotting paper, then flatten it between the pages of a book. Or lay a thin paper over the stamp and run moderately hot iron over it.

Caution: Never put an antique stamp in water because the ink might run. And don't use steam to loosen such a stamp: it might be ruined as a collector's item.

—*Reader's Digest Practical Problem Solver* [241]

In this chapter we discuss methods that use vector fields in the controller construction for a class of nonlinear dynamical systems. The control design process, as in the case of linear systems, involves three steps. The first step is the design of a state-feedback control law, the second step involves the design of a state estimator, and the third step combines the first two steps to obtain a combined controller–estimator compensator. We use the so-called *feedback linearization* method to construct a state-feedback control strategy in the first stage of the above design algorithm. The basic idea of the feedback linearization approach is to use a control consisting of two components: one component cancels out the plant's nonlinearities, and the other controls the resulting linear system. The method is limited to a class of dynamical systems, plants, whose models are sufficiently smooth; that is, the plants whose right-hand sides of the modeling differential equations are sufficiently many times differentiable. In the next section we describe the class of plants for which the feedback linearization method is applicable.

7.1 A Nonlinear Plant Model

Dynamical systems that we consider in this chapter are modeled by the equations

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u}, \quad (7.1)$$

$$\mathbf{y} = \mathbf{h}(\mathbf{x}), \quad (7.2)$$

where $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $\mathbf{G} : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times m}$, and the output map $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^p$. Note that $\mathbf{G}(\mathbf{x})$ may be written as $\mathbf{G}(\mathbf{x}) = [\mathbf{g}_1(\mathbf{x}) \ \mathbf{g}_2(\mathbf{x}) \ \cdots \ \mathbf{g}_m(\mathbf{x})]$, where $\mathbf{g}_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $i = 1, 2, \dots, m$. We write $\mathbf{f} \in \mathcal{C}^\infty$ if \mathbf{f} is in all \mathcal{C}^k for $k = 1, 2, \dots$. The functions \mathbf{f} , \mathbf{g}_i , and \mathbf{h} are assumed to be \mathcal{C}^∞ vector fields on \mathbb{R}^n . We also assume that $\mathbf{f}(\mathbf{0}) = \mathbf{0}$ and $\mathbf{h}(\mathbf{0}) = \mathbf{0}$. We note that the linear system model $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$, $\mathbf{y} = \mathbf{C}\mathbf{x}$ is a special case of the above more general model. Indeed, we have $\mathbf{f}(\mathbf{x}) = \mathbf{A}\mathbf{x}$, $\mathbf{G}(\mathbf{x}) = \mathbf{B}$, and $\mathbf{h}(\mathbf{x}) = \mathbf{C}\mathbf{x}$. We now give a simple example of a nonlinear dynamical system whose model fits into the above class of systems.

◆ Example 7.1

The example comes from Žak and MacCarley [313]. It is a simple model of a one-link robot manipulator whose reduced-order model we analyzed in Example 5.13. The motion of the robot arm is controlled by a DC motor via a gear as shown in Figure 5.11. The DC motor is armature-controlled, and its schematic is shown in Figure 5.12. We used Newton's second law to write the equation modeling the arm dynamics,

$$ml^2 \frac{d^2 \theta_p}{dt^2} = mgl \sin \theta_p + NK_m i_a,$$

where $g = 9.8 \text{ m/sec}^2$ is the gravitational constant. We next applied Kirchhoff's voltage law to the armature circuit to obtain

$$L_a \frac{di_a}{dt} + Ri_a + e_b = L_a \frac{di_a}{dt} + Ri_a + K_b N \frac{d\theta_p}{dt} = u,$$

where K_b is the back emf constant. We can now construct a third-order state-space model of the one-link robot. For this we choose the following state variables:

$$x_1 = \theta_p, \quad x_2 = \frac{d\theta_p}{dt} = \omega_p, \quad x_3 = i_a.$$

Then, the model in state-space format is

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} x_2 \\ \frac{g}{l} \sin x_1 + \frac{NK_m}{ml^2} x_3 \\ -\frac{K_b N}{L_a} x_2 - \frac{R_a}{L_a} x_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \frac{1}{L_a} \end{bmatrix} u, \quad (7.3)$$

Reasonable parameters of the robot are: $l = 1 \text{ m}$, $m = 1 \text{ kg}$, $N = 10$, $K_m = 0.1 \text{ Nm/A}$, $K_b = 0.1 \text{ V sec/rad}$, $R_a = 1 \ \Omega$, $L_a = 100 \text{ mH}$. With the above

parameters the robot model takes the form

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} x_2 \\ 9.8 \sin x_1 + x_3 \\ -10 x_2 - 10 x_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 10 \end{bmatrix} u. \quad (7.4)$$

We assume that the output, y , is

$$y = x_1. \quad (7.5)$$

7.2 Controller Form

One of the objectives of this chapter is to devise a method for constructing state-feedback stabilizing control law for a class of dynamical systems modeled by (7.1). To construct such a law it is advantageous to work with an equivalent system model rather than with the original one. We thus first discuss a method for reducing a nonlinear system model into an equivalent form that is a generalization of the controller form known from linear system theory. In our subsequent discussion, we will be using three types of *Lie derivatives*. They are as follows:

1. *Derivative of a vector field with respect to a vector field, also known as the **Lie bracket**.* Given the vector-valued functions $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$, where f and g are C^∞ vector fields, their Lie bracket is defined as

$$[f, g] = \frac{\partial f}{\partial x} g - \frac{\partial g}{\partial x} f, \quad (7.6)$$

where $\frac{\partial f}{\partial x}$ and $\frac{\partial g}{\partial x}$ are the Jacobian matrices of f and g , respectively. Sometimes in the literature the negative of the above definition is used. For convenience we will also be using an alternative notation of the form

$$[f, g] = (\text{ad}^1 f, g). \quad (7.7)$$

We define

$$(\text{ad}^k f, g) = [f, (\text{ad}^{k-1} f, g)],$$

where

$$(\text{ad}^0 f, g) = g.$$

For a single-input linear system model, $f(x) = Ax$ and $g(x) = b$. Hence,

$$(\text{ad}^0 f, g) = b,$$

$$(\text{ad}^1 f, g) = Ab,$$

$$\vdots$$

$$(\text{ad}^k f, g) = A^k b.$$

Thus, the controllability matrix of the linear system can be expressed in terms of Lie brackets as

$$[b \quad Ab \quad \cdots \quad A^{n-1}b] = [(\text{ad}^0 f, g) \quad (\text{ad}^1 f, g) \quad \cdots \quad (\text{ad}^{n-1} f, g)].$$

We will call the $n \times n$ matrix

$$[(\text{ad}^0 f, g) \quad (\text{ad}^1 f, g) \quad \cdots \quad (\text{ad}^{n-1} f, g)], \quad (7.8)$$

the controllability matrix of the single input nonlinear system modeled by

$$\dot{x} = f(x) + g(x)u. \quad (7.9)$$

2. *Derivative of a function with respect to a vector field.* Let $h: \mathbb{R}^n \rightarrow \mathbb{R}$ be a C^∞ function on \mathbb{R}^n . Let $Dh = \nabla h^T$, where ∇h is the gradient (a column vector) of h with respect to x . Then, the Lie derivative of the function h with respect to the vector field f , denoted $L_f h$ or $L_f(h)$, is defined as

$$L_f h = L_f(h) = \langle \nabla h, f \rangle = Dh f = \frac{\partial h}{\partial x_1} f_1 + \frac{\partial h}{\partial x_2} f_2 + \cdots + \frac{\partial h}{\partial x_n} f_n.$$

We use the following notation:

$$\begin{aligned} L_f^0 h &= h, \\ L_f^1 h &= L_f h, \\ L_f^2 h &= L_f L_f h, \\ &\vdots \\ L_f^k h &= L_f L_f^{k-1} h. \end{aligned}$$

3. *Derivative of Dh with respect to the vector field.* The Lie derivative of Dh with respect to the vector field f , denoted $L_f(Dh)$, is defined as

$$L_f(Dh) = \left(\frac{\partial \nabla h}{\partial x} f \right)^T + Dh \frac{\partial f}{\partial x}. \quad (7.10)$$

Note that

$$L_f(Dh) = DL_f h = \nabla L_f^T h. \quad (7.11)$$

The above three Lie derivatives satisfy the *Leibniz formula*,

$$L_{[f, g]} h = \langle \nabla h, [f, g] \rangle = L_g L_f h - L_f L_g h \quad (7.12)$$

Our goal now is to construct a C^∞ state variable transformation

$$z = T(x), \quad T(0) = 0,$$

for which there is a C^∞ inverse $\mathbf{x} = \mathbf{T}^{-1}(\mathbf{z})$, such that system model $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})u$ in the new coordinates has the form

$$\begin{bmatrix} \dot{z}_1 \\ \dot{z}_2 \\ \vdots \\ \dot{z}_{n-1} \\ \dot{z}_n \end{bmatrix} = \begin{bmatrix} z_2 \\ z_3 \\ \vdots \\ z_n \\ \tilde{f}_n(z_1, z_2, \dots, z_n) \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} u. \quad (7.13)$$

A C^∞ transformation $\mathbf{z} = \mathbf{T}(\mathbf{x})$ such that $\mathbf{T}(\mathbf{0}) = \mathbf{0}$ and for which there is a C^∞ inverse $\mathbf{x} = \mathbf{T}^{-1}(\mathbf{z})$ is called a *diffeomorphism*. We devise a method for constructing the desired transformation. Suppose that such a transformation exists. Differentiating $\mathbf{z}(t)$ with respect to time t gives

$$\dot{\mathbf{z}}(t) = \frac{\partial \mathbf{T}}{\partial \mathbf{x}} \dot{\mathbf{x}} = \frac{\partial \mathbf{T}}{\partial \mathbf{x}} (\mathbf{f}(\mathbf{x}) + \mathbf{g}(\mathbf{x})u). \quad (7.14)$$

Comparing the right-hand sides of (7.13) and (7.14), we obtain

$$\left. \frac{\partial \mathbf{T}}{\partial \mathbf{x}} \mathbf{f}(\mathbf{x}) \right|_{\mathbf{x}=\mathbf{T}^{-1}(\mathbf{z})} = \begin{bmatrix} z_2 \\ z_3 \\ \vdots \\ z_n \\ \tilde{f}_n(z_1, z_2, \dots, z_n) \end{bmatrix}, \quad (7.15)$$

and

$$\left. \frac{\partial \mathbf{T}}{\partial \mathbf{x}} \mathbf{g}(\mathbf{x}) \right|_{\mathbf{x}=\mathbf{T}^{-1}(\mathbf{z})} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}. \quad (7.16)$$

We first consider (7.15). We have

$$\begin{bmatrix} z_2 \\ z_3 \\ \vdots \\ z_n \\ \tilde{f}_n(z_1, z_2, \dots, z_n) \end{bmatrix} = \begin{bmatrix} T_2 \\ T_3 \\ \vdots \\ T_n \\ \tilde{f}_n(z_1, z_2, \dots, z_n) \end{bmatrix}$$

because $\mathbf{z} = \mathbf{T}(\mathbf{x})$. We conclude that

$$\frac{\partial T_i}{\partial \mathbf{x}} \mathbf{f}(\mathbf{x}) = T_{i+1}, \quad i = 1, 2, \dots, n-1,$$

or equivalently

$$T_{i+1} = L_f T_i, \quad i = 1, 2, \dots, n-1.$$

Therefore, the transformation $\mathbf{z} = \mathbf{T}(\mathbf{x})$ has the form

$$\mathbf{T} = \begin{bmatrix} T_1 \\ L_f T_1 \\ \vdots \\ L_f^{n-2} T_1 \\ L_f^{n-1} T_1 \end{bmatrix}. \quad (7.17)$$

The above means that the problem of constructing the desired transformation $\mathbf{z} = \mathbf{T}(\mathbf{x})$ is reduced to finding its first component T_1 . The remaining components of \mathbf{T} can be successively computed using T_1 . To specify T_1 , we use (7.16), which we represent componentwise as

$$L_g T_i = 0, \quad i = 1, 2, \dots, n-1,$$

and

$$L_g T_n = 1.$$

Note that $L_g T_1 = 0$ can be represented as

$$\frac{\partial T_1}{\partial \mathbf{x}} (\text{ad}^0 \mathbf{f}, \mathbf{g}) = 0.$$

Next, consider the relation $L_g T_2 = 0$. We will show that

$$L_g T_2 = \frac{\partial T_1}{\partial \mathbf{x}} (\text{ad}^1 \mathbf{f}, \mathbf{g}) = 0. \quad (7.18)$$

By the Leibniz formula, we obtain

$$\begin{aligned} L_g T_2 &= L_g L_f T_1 \\ &= L_f L_g T_1 + L_{[f, g]} T_1 \\ &= L_{[f, g]} T_1, \end{aligned} \quad (7.19)$$

because $L_g T_1 = 0$, and hence $L_f L_g T_1 = 0$. Using (7.19), we obtain (7.18). Continuing in this manner, we get

$$\frac{\partial T_1}{\partial \mathbf{x}} (\text{ad}^i \mathbf{f}, \mathbf{g}) = 0, \quad i = 0, 1, \dots, n-2,$$

and

$$\frac{\partial T_1}{\partial \mathbf{x}} (\text{ad}^{n-1} \mathbf{f}, \mathbf{g}) = 1.$$

We represent the above n equations in matrix form as

$$\begin{aligned} &\frac{\partial T_1}{\partial \mathbf{x}} [(\text{ad}^0 \mathbf{f}, \mathbf{g}) \quad (\text{ad}^1 \mathbf{f}, \mathbf{g}) \quad \cdots \quad (\text{ad}^{n-1} \mathbf{f}, \mathbf{g})] \\ &= [L_g T_1 \quad L_{(\text{ad}^1 \mathbf{f}, \mathbf{g})} T_1 \quad \cdots \quad L_{(\text{ad}^{n-1} \mathbf{f}, \mathbf{g})} T_1] \\ &= [0 \quad 0 \quad \cdots \quad 1]. \end{aligned}$$

Thus the derivative of T_1 —that is, the row vector $\frac{\partial T_1}{\partial \mathbf{x}}$ —is the last row of the inverse of the controllability matrix, provided the controllability matrix is invertible. We denote the last row of the inverse of the controllability matrix by $\mathbf{q}(\mathbf{x})$. Then, the problem we have to solve is to find $T_1 : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$\frac{\partial T_1(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{q}(\mathbf{x}), \quad T_1(\mathbf{0}) = 0, \quad (7.20)$$

or, equivalently,

$$\nabla T_1 = \mathbf{q}^T(\mathbf{x}), \quad T_1(\mathbf{0}) = 0. \quad (7.21)$$

The above problem can be viewed as finding a potential function T_1 for a given vector field \mathbf{q}^T . A necessary and sufficient condition for a given vector field to have a potential function is given by Theorem A.21 on page 672. We will show that if \mathbf{q}^T has a potential function, then the resulting state variable transformation is a diffeomorphism. Before we proceed with the proof, we first illustrate the results obtained so far with a numerical example.

◆ Example 7.2

Consider the nonlinear system model from Example 7.1. The controllability matrix of this system model is

$$[(\text{ad}^0 \mathbf{f}, \mathbf{g}) \quad (\text{ad}^1 \mathbf{f}, \mathbf{g}) \quad (\text{ad}^2 \mathbf{f}, \mathbf{g})] = \begin{bmatrix} 0 & 0 & 10 \\ 0 & 10 & -100 \\ 10 & -100 & 900 \end{bmatrix}.$$

The above controllability matrix is of full rank on \mathbb{R}^3 . The last row of its inverse is

$$\mathbf{q} = [0.1 \quad 0 \quad 0].$$

Hence,

$$T_1 = 0.1x_1.$$

Having obtained T_1 , we construct the desired transformation $\mathbf{z} = \mathbf{T}(\mathbf{x})$, where

$$\mathbf{T}(\mathbf{x}) = \begin{bmatrix} T_1 \\ L_{\mathbf{f}} T_1 \\ L_{\mathbf{f}}(L_{\mathbf{f}} T_1) \end{bmatrix} = \begin{bmatrix} 0.1x_1 \\ 0.1x_2 \\ 0.98 \sin x_1 + 0.1x_3 \end{bmatrix}.$$

Note that the inverse transformation $\mathbf{x} = \mathbf{T}^{-1}(\mathbf{z})$ exists and has the form

$$\mathbf{T}^{-1}(\mathbf{z}) = \begin{bmatrix} 10z_1 \\ 10z_2 \\ 10z_3 - 9.8 \sin(10z_1) \end{bmatrix}.$$

Furthermore,

$$\frac{\partial \mathbf{T}}{\partial \mathbf{x}} = \begin{bmatrix} 0.1 & 0 & 0 \\ 0 & 0.1 & 0 \\ 0.98 \cos x_1 & 0 & 0.1 \end{bmatrix}.$$

Applying the transformation $\mathbf{z} = \mathbf{T}(\mathbf{x})$ to the model of the one-link robot manipulator yields

$$\begin{aligned} \dot{\mathbf{z}} &= \left. \frac{\partial \mathbf{T}}{\partial \mathbf{x}} \mathbf{f}(\mathbf{x}) \right|_{\mathbf{x}=\mathbf{T}^{-1}(\mathbf{z})} + \left. \frac{\partial \mathbf{T}}{\partial \mathbf{x}} \mathbf{g}(\mathbf{x}) \right|_{\mathbf{x}=\mathbf{T}^{-1}(\mathbf{z})} u \\ &= \begin{bmatrix} z_2 \\ z_3 \\ 9.8 \sin(10z_1) + (9.8 \cos(10z_1) - 10)z_2 - 10z_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} u. \quad (7.22) \end{aligned}$$

We will now show that if there is a potential function for $\mathbf{q}^T(\mathbf{x})$, then the resulting transformation is indeed a diffeomorphism. By construction, \mathbf{T} is of \mathcal{C}^∞ class, and $\mathbf{T}(\mathbf{0}) = \mathbf{0}$. We now show that \mathbf{T} is locally \mathcal{C}^∞ -invertible. By the inverse function theorem, stated in Exercise 7.3, if the Jacobian matrix of \mathbf{T} evaluated at $\mathbf{0}$ is invertible, then \mathbf{T} is locally \mathcal{C}^∞ -invertible. The Jacobian matrix of \mathbf{T} is

$$D\mathbf{T} = \begin{bmatrix} DT_1 \\ DL_f T_1 \\ \vdots \\ DL_f^{n-2} T_1 \\ DL_f^{n-1} T_1 \end{bmatrix}.$$

We compute the product of the above Jacobian matrix and the controllability matrix

$$[(\text{ad}^0 \mathbf{f}, \mathbf{g}) \quad (\text{ad}^1 \mathbf{f}, \mathbf{g}) \quad \cdots \quad (\text{ad}^{n-1} \mathbf{f}, \mathbf{g})]$$

to get

$$\begin{bmatrix} DT_1 \\ DL_f T_1 \\ \vdots \\ DL_f^{n-2} T_1 \\ DL_f^{n-1} T_1 \end{bmatrix} [(\text{ad}^0 \mathbf{f}, \mathbf{g}) \quad (\text{ad}^1 \mathbf{f}, \mathbf{g}) \quad \cdots \quad (\text{ad}^{n-1} \mathbf{f}, \mathbf{g})].$$

Performing multiplications yields

$$\begin{aligned}
 DT[(\text{ad}^0 f, g) \quad (\text{ad}^1 f, g) \quad \cdots \quad (\text{ad}^{n-1} f, g)] \\
 &= \begin{bmatrix} L_g T_1 & L_{(\text{ad}^1 f, g)} T_1 & \cdots & L_{(\text{ad}^{n-1} f, g)} T_1 \\ L_g L_f T_1 & L_{(\text{ad}^1 f, g)} L_f T_1 & \cdots & L_{(\text{ad}^{n-1} f, g)} L_f T_1 \\ \vdots & \vdots & & \vdots \\ L_g L_f^{n-1} T_1 & L_{(\text{ad}^1 f, g)} L_f^{n-1} T_1 & \cdots & L_{(\text{ad}^{n-1} f, g)} L_f^{n-1} T_1 \end{bmatrix} \\
 &= \begin{bmatrix} L_g T_1 & L_{(\text{ad}^1 f, g)} T_1 & \cdots & L_{(\text{ad}^{n-1} f, g)} T_1 \\ L_g T_2 & L_{(\text{ad}^1 f, g)} T_2 & \cdots & L_{(\text{ad}^{n-1} f, g)} T_2 \\ \vdots & \vdots & & \vdots \\ L_g T_n & L_{(\text{ad}^1 f, g)} T_n & \cdots & L_{(\text{ad}^{n-1} f, g)} T_n \end{bmatrix} \\
 &= \begin{bmatrix} L_g T_1 & L_{(\text{ad}^1 f, g)} T_1 & \cdots & L_{(\text{ad}^{n-1} f, g)} T_1 \\ L_{(\text{ad}^1 f, g)} T_1 & L_{(\text{ad}^2 f, g)} T_1 & \cdots & L_{(\text{ad}^n f, g)} T_1 \\ \vdots & \vdots & & \vdots \\ L_{(\text{ad}^{n-1} f, g)} T_1 & L_{(\text{ad}^n f, g)} T_1 & \cdots & L_{(\text{ad}^{2n-2} f, g)} T_1 \end{bmatrix} \\
 &= \begin{bmatrix} 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & \cdots & 1 & x \\ \vdots & \vdots & & \vdots & \vdots \\ 1 & x & \cdots & x & x \end{bmatrix},
 \end{aligned}$$

where we use x to denote “don’t care” elements. It follows from the above that at $\mathbf{x} = \mathbf{0}$ the Jacobian matrix DT is invertible, that is,

$$\text{rank } DT(\mathbf{0}) = n.$$

In conclusion, a diffeomorphic state variable transformation exists, in a neighborhood of $\mathbf{x} = \mathbf{0}$, if

1. in a neighborhood of $\mathbf{x} = \mathbf{0}$,

$$\text{rank}[(\text{ad}^0 f, g) \quad (\text{ad}^1 f, g) \quad \cdots \quad (\text{ad}^{n-1} f, g)] = n,$$

2. the equation

$$\mathbf{q}(\mathbf{x}) = [0 \quad 0 \quad \cdots \quad 0 \quad 1][(\text{ad}^0 f, g) \quad (\text{ad}^1 f, g) \quad \cdots \quad (\text{ad}^{n-1} f, g)]^{-1}$$

has a potential function in a neighborhood of $\mathbf{x} = \mathbf{0}$; that is, $\mathbf{q}^T(\mathbf{x})$ is a conservative vector field in a neighborhood of $\mathbf{x} = \mathbf{0}$.

The above method of transforming a class of single-input nonlinear systems into the controller form can be generalized to multi-input systems in a way similar to that in the linear case discussed in Subsection 3.3.1. We will illustrate the method with a numerical example.

◆ Example 7.3

Consider the model of the permanent stepper motor given by (1.61). We first represent the model as

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{g}_1 u_1 + \mathbf{g}_2 u_2, \quad (7.23)$$

where

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} x_2 \\ -K_1 x_2 - K_2 x_3 \sin(K_3 x_1) + K_2 x_4 \cos(K_3 x_1) \\ K_4 x_2 \sin(K_3 x_1) - K_5 x_3 \\ -K_4 x_2 \cos(K_3 x_1) - K_5 x_4 \end{bmatrix},$$

$$\mathbf{g}_1 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \quad \text{and} \quad \mathbf{g}_2 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}.$$

We next form the controllability matrix for the above model to get

$$[\mathbf{g}_1 \ \mathbf{g}_2 \ (\text{ad}^1 \mathbf{f}, \mathbf{g}_1) \ (\text{ad}^1 \mathbf{f}, \mathbf{g}_2) \ (\text{ad}^2 \mathbf{f}, \mathbf{g}_1) \ (\text{ad}^2 \mathbf{f}, \mathbf{g}_2) \ (\text{ad}^3 \mathbf{f}, \mathbf{g}_1) \ (\text{ad}^3 \mathbf{f}, \mathbf{g}_2)], \quad (7.24)$$

where

$$(\text{ad}^1 \mathbf{f}, \mathbf{g}_1) = \begin{bmatrix} 0 \\ -K_2 \sin(K_3 x_1) \\ -K_5 \\ 0 \end{bmatrix}, \quad (\text{ad}^1 \mathbf{f}, \mathbf{g}_2) = \begin{bmatrix} 0 \\ K_2 \cos(K_3 x_1) \\ 0 \\ -K_5 \end{bmatrix},$$

$$(\text{ad}^2 \mathbf{f}, \mathbf{g}_1)$$

$$= [\mathbf{f}, (\text{ad}^1 \mathbf{f}, \mathbf{g}_1)] = \frac{\partial \mathbf{f}}{\partial \mathbf{x}} (\text{ad}^1 \mathbf{f}, \mathbf{g}_1) - \frac{\partial (\text{ad}^1 \mathbf{f}, \mathbf{g}_1)}{\partial \mathbf{x}} \mathbf{f}$$

$$= \begin{bmatrix} -K_2 \sin(K_3 x_1) \\ K_1 K_2 \sin(K_3 x_1) + K_2 K_5 \sin(K_3 x_1) + K_2 K_3 x_2 \cos(K_3 x_1) \\ -K_4 K_2 \sin^2(K_3 x_1) + K_5^2 \\ K_4 K_2 \cos(K_3 x_1) \sin(K_3 x_1) \end{bmatrix}, \quad (7.25)$$

and

$$\begin{aligned}
 (\text{ad}^2 \mathbf{f}, \mathbf{g}_2) &= [\mathbf{f}, (\text{ad}^1 \mathbf{f}, \mathbf{g}_2)] = \frac{\partial \mathbf{f}}{\partial \mathbf{x}} (\text{ad}^1 \mathbf{f}, \mathbf{g}_2) - \frac{\partial (\text{ad}^1 \mathbf{f}, \mathbf{g}_2)}{\partial \mathbf{x}} \mathbf{f} \\
 &= \begin{bmatrix} K_2 \cos(K_3 x_1) \\ -K_1 K_2 \cos(K_3 x_1) - K_2 K_5 \cos(K_3 x_1) + K_2 K_3 x_2 \sin(K_3 x_1) \\ K_4 K_2 \cos(K_3 x_1) \sin(K_3 x_1) \\ -K_4 K_2 \cos^2(K_3 x_1) + K_5^2 \end{bmatrix}.
 \end{aligned} \tag{7.26}$$

Observe that

$$(\text{ad}^1 \mathbf{f}, \mathbf{g}_2) = -K_5 \cot(K_3 x_1) \mathbf{g}_1 - K_5 \mathbf{g}_2 - \cot(K_3 x_1) (\text{ad}^1 \mathbf{f}, \mathbf{g}_1); \tag{7.27}$$

that is, $(\text{ad}^1 \mathbf{f}, \mathbf{g}_2)$ is a linear combination of \mathbf{g}_1 , \mathbf{g}_2 , and $(\text{ad}^1 \mathbf{f}, \mathbf{g}_1)$. We will use this fact when selecting columns to construct a transformation into a controller form. Proceeding from left to right and using (7.27), we select linearly independent columns from the controllability matrix (7.24). The PM stepper motor system is controllable because we have found four linearly independent columns— \mathbf{g}_1 , $(\text{ad}^1 \mathbf{f}, \mathbf{g}_1)$, $(\text{ad}^2 \mathbf{f}, \mathbf{g}_1)$, and \mathbf{g}_2 —in its controllability matrix. Thus, the controllability indices of the PM stepper motor system are $d_1 = 3$ and $d_2 = 1$. Therefore, the controller form of the PM stepper motor model is

$$\begin{bmatrix} \dot{z}_1 \\ \dot{z}_2 \\ \dot{z}_3 \\ \dot{z}_4 \end{bmatrix} = \begin{bmatrix} z_2 \\ z_3 \\ \tilde{f}_3(z_1, z_2, z_3, z_4) \\ \tilde{f}_4(z_1, z_2, z_3, z_4) \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ \tilde{g}_{31} & \tilde{g}_{32} \\ \tilde{g}_{41} & \tilde{g}_{42} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}, \tag{7.28}$$

where the elements \tilde{g}_{ij} may depend on z_k , $k = 1, 2, 3, 4$, and

$$\det \begin{bmatrix} \tilde{g}_{31} & \tilde{g}_{32} \\ \tilde{g}_{41} & \tilde{g}_{42} \end{bmatrix} \neq 0. \tag{7.29}$$

The transformation we are looking for has the form

$$\mathbf{z} = \mathbf{T}(\mathbf{x}) = \begin{bmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \end{bmatrix} = \begin{bmatrix} T_1 \\ L_{\mathbf{f}} T_1 \\ L_{\mathbf{f}}^2 T_1 \\ T_4 \end{bmatrix}, \tag{7.30}$$

where T_1 and T_4 are functions that we have to find. There is no unique solution for T_1 and T_4 . Differentiating $\mathbf{z} = \mathbf{T}(\mathbf{x})$ and using (7.23) yields

$$\begin{aligned}
 \dot{\mathbf{z}} &= \frac{\partial \mathbf{T}}{\partial \mathbf{x}} \dot{\mathbf{x}} \\
 &= \frac{\partial \mathbf{T}}{\partial \mathbf{x}} \mathbf{f}(\mathbf{x}) \Big|_{\mathbf{x}=\mathbf{T}^{-1}(\mathbf{z})} + u_1 \frac{\partial \mathbf{T}}{\partial \mathbf{x}} \mathbf{g}_1(\mathbf{x}) \Big|_{\mathbf{x}=\mathbf{T}^{-1}(\mathbf{z})} + u_2 \frac{\partial \mathbf{T}}{\partial \mathbf{x}} \mathbf{g}_2(\mathbf{x}) \Big|_{\mathbf{x}=\mathbf{T}^{-1}(\mathbf{z})},
 \end{aligned} \tag{7.31}$$

where

$$\frac{\partial \mathbf{T}}{\partial \mathbf{x}} = D\mathbf{T} = \begin{bmatrix} DT_1 \\ DT_2 \\ DT_3 \\ DT_4 \end{bmatrix}.$$

Comparing (7.28) and (7.31), we obtain

$$z_2 = T_2 = DT_1 \mathbf{f} = L_{\mathbf{f}} T_1, \quad z_3 = T_3 = DT_2 \mathbf{f} = L_{\mathbf{f}} T_2, \quad (7.32)$$

and

$$L_{\mathbf{g}_1} T_1 = 0, \quad L_{\mathbf{g}_1} T_2 = L_{\mathbf{g}_1} L_{\mathbf{f}} T_1 = 0, \quad (7.33)$$

$$L_{\mathbf{g}_2} T_1 = 0, \quad L_{\mathbf{g}_2} T_2 = L_{\mathbf{g}_2} L_{\mathbf{f}} T_1 = 0. \quad (7.34)$$

Applying the Leibniz formula (7.12) to the right-hand sides of (7.33) and (7.34) gives

$$L_{\mathbf{f}} L_{\mathbf{g}_1} T_1 = L_{[\mathbf{f}, \mathbf{g}_1]} T_1 \quad \text{and} \quad L_{\mathbf{f}} L_{\mathbf{g}_2} T_1 = L_{[\mathbf{f}, \mathbf{g}_2]} T_1.$$

Because by (7.33) and (7.34), $L_{\mathbf{g}_1} T_1 = 0$ and $L_{\mathbf{g}_2} T_1 = 0$, we can write

$$DT_1 [\mathbf{g}_1 \quad (\text{ad}^1 \mathbf{f}, \mathbf{g}_1) \quad \mathbf{g}_2 \quad (\text{ad}^1 \mathbf{f}, \mathbf{g}_2)] = \mathbf{0}^T, \quad (7.35)$$

that is,

$$\begin{bmatrix} \frac{\partial T_1}{\partial x_1} & \frac{\partial T_1}{\partial x_2} & \frac{\partial T_1}{\partial x_3} & \frac{\partial T_1}{\partial x_4} \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & -K_2 \sin(K_3 x_1) & 0 & K_2 \cos(K_3 x_1) \\ 1 & -K_5 & 0 & 0 \\ 0 & 0 & 1 & -K_5 \end{bmatrix} = [0 \ 0 \ 0 \ 0].$$

By inspection, a nontrivial solution to the above is

$$DT_1 = \begin{bmatrix} \frac{\partial T_1}{\partial x_1} & 0 & 0 & 0 \end{bmatrix}, \quad (7.36)$$

where the function T_1 may only be a function of x_1 for DT_1 to be integrable. We will specify T_1 later. To proceed, we note that by (7.30) and (7.31), we obtain

$$\begin{bmatrix} \tilde{g}_{31} & \tilde{g}_{32} \\ \tilde{g}_{41} & \tilde{g}_{42} \end{bmatrix} = \begin{bmatrix} (DT_3) \mathbf{g}_1 & (DT_3) \mathbf{g}_2 \\ (DT_4) \mathbf{g}_1 & (DT_4) \mathbf{g}_2 \end{bmatrix}, \quad (7.37)$$

Using the structure of the state transformation (7.30) and the Leibniz formula (7.12), we express the element $\tilde{g}_{31} = (DT_3) \mathbf{g}_1$ as

$$\begin{aligned} (DT_3) \mathbf{g}_1 &= L_{\mathbf{g}_1} L_{\mathbf{f}} T_2 \\ &= L_{[\mathbf{f}, \mathbf{g}_1]} T_2 + L_{\mathbf{f}} L_{\mathbf{g}_1} T_2 \\ &= L_{[\mathbf{f}, \mathbf{g}_1]} T_2, \end{aligned} \quad (7.38)$$

because $L_{\mathbf{g}_1} T_2 = 0$ and therefore $L_{\mathbf{f}} L_{\mathbf{g}_1} T_2 = 0$. Observing that $T_2 = L_{\mathbf{f}} T_1$, employing the Leibniz formula, and taking into account the fact that $L_{\mathbf{g}_1} T_1 = 0$,

we get

$$\begin{aligned}
 L_{[\mathbf{f}, \mathbf{g}_1]} T_2 &= L_{[\mathbf{f}, \mathbf{g}_1]} L_{\mathbf{f}} T_1 \\
 &= L_{[\mathbf{f}, [\mathbf{f}, \mathbf{g}_1]]} T_1 + L_{\mathbf{f}} L_{[\mathbf{f}, \mathbf{g}_1]} T_1 \\
 &= L_{[\mathbf{f}, [\mathbf{f}, \mathbf{g}_1]]} T_1 + L_{\mathbf{f}} (L_{\mathbf{g}_1} L_{\mathbf{f}} T_1 - L_{\mathbf{f}} L_{\mathbf{g}_1} T_1) \\
 &= L_{[\mathbf{f}, [\mathbf{f}, \mathbf{g}_1]]} T_1 \\
 &= DT_1(\text{ad}^2 \mathbf{f}, \mathbf{g}_1).
 \end{aligned} \tag{7.39}$$

Hence,

$$\tilde{g}_{31} = DT_1(\text{ad}^2 \mathbf{f}, \mathbf{g}_1). \tag{7.40}$$

Similarly, we find

$$\tilde{g}_{32} = DT_1(\text{ad}^2 \mathbf{f}, \mathbf{g}_2). \tag{7.41}$$

Substituting (7.25), (7.26), and (7.36) into (7.40) and (7.41), we find

$$\tilde{g}_{31} = -\frac{\partial T_1}{\partial x_1} K_2 \sin(K_3 x_1), \tag{7.42}$$

and

$$\tilde{g}_{32} = \frac{\partial T_1}{\partial x_1} K_2 \cos(K_3 x_1). \tag{7.43}$$

Substituting the above expressions into (7.29) yields

$$\begin{aligned}
 \tilde{g}_{31} \tilde{g}_{42} - \tilde{g}_{32} \tilde{g}_{41} &= (DT_1(\text{ad}^2 \mathbf{f}, \mathbf{g}_1))((DT_4) \mathbf{g}_2) - (DT_1(\text{ad}^2 \mathbf{f}, \mathbf{g}_2))((DT_4) \mathbf{g}_1) \\
 &= -\frac{\partial T_1}{\partial x_1} K_2 \sin(K_3 x_1) \frac{\partial T_4}{\partial x_4} - \frac{\partial T_1}{\partial x_1} K_2 \cos(K_3 x_1) \frac{\partial T_4}{\partial x_3} \\
 &\neq 0.
 \end{aligned} \tag{7.44}$$

Suppose that we choose

$$\frac{\partial T_1}{\partial x_1} = \frac{1}{K_2}, \quad \frac{\partial T_4}{\partial x_3} = \cos(K_3 x_1), \quad \text{and} \quad \frac{\partial T_4}{\partial x_4} = \sin(K_3 x_1). \tag{7.45}$$

Substituting (7.45) into (7.44) gives

$$\det \begin{bmatrix} \tilde{g}_{31} & \tilde{g}_{32} \\ \tilde{g}_{41} & \tilde{g}_{42} \end{bmatrix} = -1.$$

Using (7.36), (7.45), and (7.30), we obtain

$$\begin{aligned}
 T_1 &= x_1 / K_2, \\
 T_2 &= L_{\mathbf{f}} T_1 = x_2 / K_2, \\
 T_3 &= L_{\mathbf{f}} T_2 = -\frac{K_1}{K_2} x_2 - x_3 \sin(K_3 x_1) + x_4 \cos(K_3 x_1), \\
 T_4 &= x_3 \cos(K_3 x_1) + x_4 \sin(K_3 x_1).
 \end{aligned} \tag{7.46}$$

The Jacobian matrix of \mathbf{T} is

$$D\mathbf{T} = \begin{bmatrix} 1/K_2 & 0 & 0 & 0 \\ 0 & 1/K_2 & 0 & 0 \\ -K_3 x_3 \cos(K_3 x_1) - K_3 x_4 \sin(K_3 x_1) & -K_1/K_2 & -\sin(K_3 x_1) & \cos(K_3 x_1) \\ -K_3 x_3 \sin(K_3 x_1) + K_3 x_4 \cos(K_3 x_1) & 0 & \cos(K_3 x_1) & \sin(K_3 x_1) \end{bmatrix}. \quad (7.47)$$

We have

$$\mathbf{T}(\mathbf{0}) = \mathbf{0} \quad \text{and} \quad \det D\mathbf{T} = -1/K_2^2 \neq 0 \quad \text{for all } \mathbf{x} \in \mathbb{R}^4, \quad (7.48)$$

which means that we were able to find a global transformation into the controller form. The inverse, $\mathbf{x} = \mathbf{T}^{-1}(\mathbf{z})$, is

$$\begin{aligned} x_1 &= K_2 z_1, \\ x_2 &= K_2 z_2, \\ x_3 &= -K_1 z_2 \sin(K_3 K_2 z_1) - z_3 \sin(K_3 K_2 z_1) + z_4 \cos(K_3 K_2 z_1), \\ x_4 &= K_1 z_2 \cos(K_3 K_2 z_1) + z_3 \cos(K_3 K_2 z_1) + z_4 \sin(K_3 K_2 z_1). \end{aligned} \quad (7.49)$$

We transform the PM stepper motor model into the new coordinates using (7.30), (7.46), and (7.46) to obtain

$$\begin{aligned} \begin{bmatrix} \dot{z}_1 \\ \dot{z}_2 \\ \dot{z}_3 \\ \dot{z}_4 \end{bmatrix} &= \begin{bmatrix} z_2 \\ z_3 \\ -K_2 K_3 z_2 z_4 - K_1 z_3 - K_2 K_4 z_2 - K_5 z_3 - K_5 K_1 z_2 \\ K_2 K_3 z_2 z_3 + K_1 K_2 K_3 z_2^2 - K_5 z_4 \end{bmatrix} \\ &+ \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ -\sin(K_3 K_2 z_1) & \cos(K_3 K_2 z_1) \\ \cos(K_3 K_2 z_1) & \sin(K_3 K_2 z_1) \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}. \end{aligned} \quad (7.50)$$

7.3 Linearizing State-Feedback Control

Once the plant model is transformed into the controller form, we can construct a state-feedback controller in the new coordinates and then, using the inverse transformation, represent the controller in the original coordinates. While constructing the controller in the new coordinates, a part of the controller is used to cancel nonlinearities, thus resulting in a linear system in the new coordinates. This procedure, as we mentioned at the beginning of this chapter, is called feedback linearization. Then, we proceed to construct the other part of the controller. This part can be designed using linear control methods because the feedback linearized system is linear.

We illustrate the above method of state feedback design on a numerical example involving the one-link manipulator whose model was derived in Example 7.1. The state-space model of

the manipulator, given by (7.4), was transformed into the controller form in Example 7.2. The controller form of the one-link manipulator model is given by (7.22). It is easy to design a stabilizing state-feedback controller in the new coordinates. It takes the form

$$u = -\tilde{f}_3(z_1, z_2, z_3) - (k_1 z_1 + k_2 z_2 + k_3 z_3), \quad (7.51)$$

where $\tilde{f}_3(z_1, z_2, z_3) = 9.8 \sin(10z_1) + (9.8 \cos(10z_1) - 10)z_2 - 10z_3$. Suppose that the desired closed-loop poles of the feedback linearized one-link manipulator are to be located at $\{-3.5, -4, -5\}$. Then, the linear feedback gains k_i , $i = 1, 2, 3$, that shift the poles to these desired locations are

$$k_1 = 70, \quad k_2 = 51.5, \quad k_3 = 12.5.$$

Applying (7.51) with the above values of the linear feedback gains to the model (7.22) gives

$$\dot{z} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -70 & -51.5 & -12.5 \end{bmatrix} z.$$

Next, applying the inverse transformation yields the controller in the original coordinates,

$$u = -7x_1 - 4.15x_2 - 0.25x_3 - 0.98x_2 \cos(x_1) - 12.25 \sin(x_1). \quad (7.52)$$

In Figure 7.1, three plots of the manipulator's link angle, $x_1 = \theta$, versus time are presented. The initial conditions have the form $\mathbf{x}(0) = [\theta(0) \ 0 \ 0]^T$. The control law applied to the manipulator is given by (7.52).

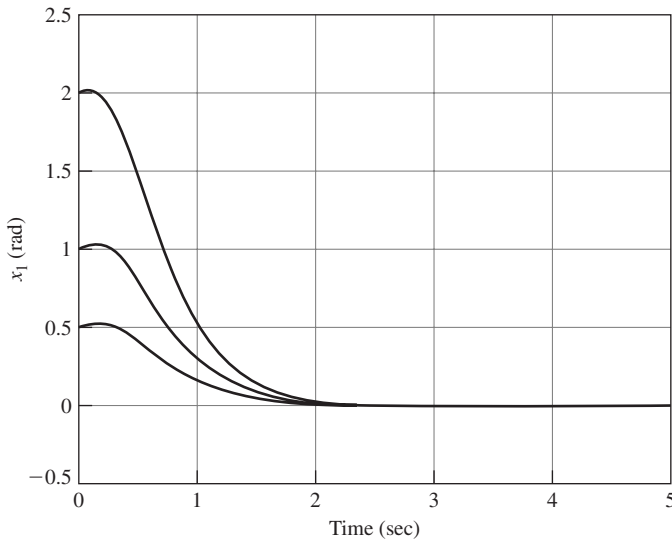


Figure 7.1 Plots of the one-link manipulator's link angle, $x_1 = \theta$, versus time for three different initial angles $\theta(0)$. The manipulator is driven by the state-feedback control law (7.52).

7.4 Observer Form

In this section we discuss the problem of transforming a class of nonlinear system models into an observer form, which is then used to construct state estimators for these systems. We consider a single output nonlinear system model of the form

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u}, \\ y &= h(\mathbf{x}).\end{aligned}\tag{7.53}$$

First, we analyze the case of unforced system when $\mathbf{u} = \mathbf{0}$. In such a case, (7.53) reduces to

$$\begin{aligned}\dot{\mathbf{x}} &= \mathbf{f}(\mathbf{x}), \\ y &= h(\mathbf{x}).\end{aligned}\tag{7.54}$$

We desire to find a state variable transformation, written as $\mathbf{x} = \mathbf{T}(\mathbf{z})$, such that the model (7.54) in the \mathbf{z} coordinates is

$$\begin{aligned}\dot{\mathbf{z}} &= \begin{bmatrix} 0 & 0 & \cdots & 0 & 0 \\ 1 & 0 & \cdots & 0 & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{bmatrix} \mathbf{z} + \begin{bmatrix} \alpha_0(y) \\ \alpha_1(y) \\ \vdots \\ \alpha_{n-1}(y) \end{bmatrix} = \mathbf{A}\mathbf{z} + \boldsymbol{\alpha}(y) = \tilde{\mathbf{f}}(\mathbf{z}), \\ y &= [0 \ 0 \ \cdots \ 0 \ 1]\mathbf{z} = \mathbf{c}\mathbf{z}.\end{aligned}\tag{7.55}$$

It is easy to verify that the state $\tilde{\mathbf{z}}$ of the dynamical system

$$\dot{\tilde{\mathbf{z}}} = (\mathbf{A} - \mathbf{l}\mathbf{c})\tilde{\mathbf{z}} + \boldsymbol{\alpha}(y) + \mathbf{l}y,\tag{7.56}$$

will asymptotically converge to the state \mathbf{z} of (7.55) if the matrix $\mathbf{A} - \mathbf{l}\mathbf{c}$ is asymptotically stable. This is because the dynamics of the error, $\mathbf{e} = \tilde{\mathbf{z}} - \mathbf{z}$, are described by

$$\dot{\mathbf{e}} = (\mathbf{A} - \mathbf{l}\mathbf{c})\mathbf{e}, \quad \mathbf{e}(t_0) = \tilde{\mathbf{z}}(t_0) - \mathbf{z}(t_0).\tag{7.57}$$

We will now derive a sufficiency condition for the existence of a state variable transformation that brings (7.54) into the observer form (7.55).

Taking the derivative of $\mathbf{x} = \mathbf{T}(\mathbf{z})$ with respect to time yields

$$\dot{\mathbf{x}} = \frac{\partial \mathbf{T}}{\partial \mathbf{z}} \dot{\mathbf{z}} = \frac{\partial \mathbf{T}}{\partial \mathbf{z}} \tilde{\mathbf{f}}(\mathbf{z}),\tag{7.58}$$

where we represent the Jacobian matrix of $\mathbf{T}(\mathbf{z})$, denoted $\frac{\partial \mathbf{T}}{\partial \mathbf{z}}$, as

$$\frac{\partial \mathbf{T}}{\partial \mathbf{z}} = D\mathbf{T} = \begin{bmatrix} \frac{\partial \mathbf{T}}{\partial z_1} & \cdots & \frac{\partial \mathbf{T}}{\partial z_n} \end{bmatrix}.\tag{7.59}$$

Equating the right-hand sides of $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x})$ and (7.58) yields

$$\mathbf{f}(\mathbf{x}) = \mathbf{f}(\mathbf{T}(\mathbf{z})) = \frac{\partial \mathbf{T}}{\partial \mathbf{z}} \tilde{\mathbf{f}}(\mathbf{z}).\tag{7.60}$$

Taking the partial derivative of (7.60) with respect to z_k yields

$$\boxed{\frac{\partial f(T(z))}{\partial z_k} = \frac{\partial \left(\frac{\partial T}{\partial z} \right)}{\partial z_k} \tilde{f}(z) + \frac{\partial T}{\partial z} \frac{\partial \tilde{f}(z)}{\partial z_k}} \quad (7.61)$$

On the other hand, using the chain rule gives

$$\frac{\partial f(T(z))}{\partial z_k} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial T(z)}{\partial z_k}. \quad (7.62)$$

Using (7.60) and the fact that $\mathbf{x} = T(z)$, the first term on the right-hand side of (7.61) evaluates to

$$\begin{aligned} \frac{\partial \left(\frac{\partial T}{\partial z} \right)}{\partial z_k} \tilde{f}(z) &= \frac{\partial \left(\frac{\partial T}{\partial z_k} \right)}{\partial z} \tilde{f}(z) = \frac{\partial \left(\frac{\partial T}{\partial z_k} \right)}{\partial z} \frac{\partial \mathbf{x}}{\partial \mathbf{x}} \tilde{f}(z) = \frac{\partial \left(\frac{\partial T}{\partial z_k} \right)}{\partial \mathbf{x}} \frac{\partial T(z)}{\partial z} \tilde{f}(z) \\ &= \frac{\partial \left(\frac{\partial T}{\partial z_k} \right)}{\partial \mathbf{x}} f(\mathbf{x}). \end{aligned} \quad (7.63)$$

Employing (7.55) and (7.59), we evaluate the second term on the right-hand side of (7.61) to get

$$\frac{\partial T}{\partial z} \frac{\partial \tilde{f}(z)}{\partial z_k} = \frac{\partial T}{\partial z_{k+1}} \quad \text{for } k = 1, 2, \dots, n-1, \quad (7.64)$$

Using the above and substituting (7.62) and (7.63) into (7.61) gives

$$\boxed{\frac{\partial T}{\partial z_{k+1}} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial T(z)}{\partial z_k} - \frac{\partial \left(\frac{\partial T}{\partial z_k} \right)}{\partial \mathbf{x}} f(\mathbf{x}), \quad k = 1, 2, \dots, n-1} \quad (7.65)$$

With the help of the Lie-derivative notation introduced in Section 7.2, we express (7.65) in an alternative way. Specifically, using (7.6) and (7.7), we represent (7.65) as

$$\frac{\partial T}{\partial z_{k+1}} = \left[f, \frac{\partial T}{\partial z_k} \right] = \left(\text{ad}^1 f, \frac{\partial T}{\partial z_k} \right), \quad k = 1, 2, \dots, n-1. \quad (7.66)$$

Observe that

$$\begin{aligned} \frac{\partial T}{\partial z_2} &= \left(\text{ad}^1 f, \frac{\partial T}{\partial z_1} \right), \\ \frac{\partial T}{\partial z_3} &= \left(\text{ad}^1 f, \frac{\partial T}{\partial z_2} \right) = \left(\text{ad}^2 f, \frac{\partial T}{\partial z_1} \right), \\ &\vdots \quad \quad \quad \vdots \end{aligned}$$

and in general

$$\boxed{\frac{\partial \mathbf{T}}{\partial z_k} = \left(\text{ad}^{k-1} \mathbf{f}, \frac{\partial \mathbf{T}}{\partial z_1} \right), \quad k = 1, 2, \dots, n.} \quad (7.67)$$

With the aid of the above, we express all the columns of the Jacobian matrix (7.59) of $\mathbf{T}(\mathbf{z})$ in terms of the starting vector $\partial \mathbf{T} / \partial z_1$, that is,

$$\frac{\partial \mathbf{T}}{\partial \mathbf{z}} = \left[\left(\text{ad}^0 \mathbf{f}, \frac{\partial \mathbf{T}}{\partial z_1} \right) \quad \left(\text{ad}^1 \mathbf{f}, \frac{\partial \mathbf{T}}{\partial z_1} \right) \quad \cdots \quad \left(\text{ad}^{n-1} \mathbf{f}, \frac{\partial \mathbf{T}}{\partial z_1} \right) \right] \quad (7.68)$$

To obtain an expression for the starting vector $\partial \mathbf{T} / \partial z_1$, we use the output equation

$$\mathbf{y} = h(\mathbf{x}) = z_n. \quad (7.69)$$

We use the chain rule when taking the partial derivative of (7.69) with respect to \mathbf{z} to get

$$\frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} \frac{\partial \mathbf{T}}{\partial \mathbf{z}} = [0 \quad 0 \quad \cdots \quad 0 \quad 1]. \quad (7.70)$$

We utilize (7.10) to express the first component on the left-hand side of the above as

$$\frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} \frac{\partial \mathbf{T}}{\partial z_1} = L_f^0(Dh) \frac{\partial \mathbf{T}}{\partial z_1} = 0. \quad (7.71)$$

Using the Leibniz formula, we represent the second component of the left-hand side of (7.70) as

$$\begin{aligned} 0 &= \frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} \frac{\partial \mathbf{T}}{\partial z_2} \\ &= \left\langle \nabla h, \left[\mathbf{f}, \frac{\partial \mathbf{T}}{\partial z_1} \right] \right\rangle \\ &= L_{\frac{\partial \mathbf{T}}{\partial z_1}} L_f h - L_f L_{\frac{\partial \mathbf{T}}{\partial z_1}} h \\ &= L_{\frac{\partial \mathbf{T}}{\partial z_1}} L_f h, \end{aligned} \quad (7.72)$$

because $\frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} \frac{\partial \mathbf{T}}{\partial z_1} = 0$ by (7.71). We use (7.11) to represent (7.72) as

$$\begin{aligned} 0 &= L_{\frac{\partial \mathbf{T}}{\partial z_1}} L_f h \\ &= DL_f h \frac{\partial \mathbf{T}}{\partial z_1} \\ &= L_f^1(Dh) \frac{\partial \mathbf{T}}{\partial z_1}. \end{aligned} \quad (7.73)$$

By repeated application of the Leibniz formula and (7.68), we express (7.70) equivalently as

$$\begin{bmatrix} L_f^0(Dh) \\ L_f^1(Dh) \\ \vdots \\ L_f^{n-2}(Dh) \\ L_f^{n-1}(Dh) \end{bmatrix} \frac{\partial \mathbf{T}}{\partial z_1} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}. \quad (7.74)$$

The matrix

$$\begin{bmatrix} L_f^0(Dh) \\ L_f^1(Dh) \\ \vdots \\ L_f^{n-2}(Dh) \\ L_f^{n-1}(Dh) \end{bmatrix} \quad (7.75)$$

is called the observability matrix of the nonlinear system model (7.54). It follows from (7.74) that the starting vector $\partial \mathbf{T} / \partial z_1$ is equal to the last column of the inverse of the observability matrix given by (7.75).

We now discuss the case of transforming a nonlinear system model with inputs, described by (7.53), into its observer form. We proceed as recommended by Krener and Respondek [170, Section 6]. First transform the unforced system into the observer form (7.55). Then, check if the vector field $\mathbf{G}(\mathbf{x})\mathbf{u}$ is transformed into the form

$$\left(\frac{\partial \mathbf{T}}{\partial \mathbf{z}} \right)^{-1} \mathbf{G}(\mathbf{x})|_{\mathbf{x}=\mathbf{T}(\mathbf{z})} \mathbf{u} = \boldsymbol{\beta}(\mathbf{y}, \mathbf{u}). \quad (7.76)$$

If this is the case, then set

$$\boldsymbol{\gamma}(\mathbf{y}, \mathbf{u}) = \boldsymbol{\alpha}(\mathbf{y}) + \boldsymbol{\beta}(\mathbf{y}, \mathbf{u}).$$

The above means that the system with inputs is transformed into the observer form

$$\begin{aligned} \dot{\mathbf{z}} &= \mathbf{A}\mathbf{z} + \boldsymbol{\gamma}(\mathbf{y}, \mathbf{u}), \\ y &= \mathbf{c}\mathbf{z}, \end{aligned} \quad (7.77)$$

where the matrices \mathbf{A} and \mathbf{c} are as in (7.55). If the unforced system cannot be transformed into the observer form or $\mathbf{G}(\mathbf{x})\mathbf{u}$ does not transform into (7.76), then the original system (7.53) cannot be transformed into the observer form (7.77).

◆ Example 7.4

Consider the system consisting of an inverted pendulum controlled by an armature-controlled DC motor that we analyzed in Example 7.1. The observability matrix

for the system is

$$\begin{bmatrix} L_{\mathbf{f}}^0(Dh) \\ L_{\mathbf{f}}^1(Dh) \\ L_{\mathbf{f}}^2(Dh) \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 9.8 \cos(x_1) & 0 & 1 \end{bmatrix}. \quad (7.78)$$

The starting vector $\partial \mathbf{T} / \partial z_1$ is the last column of the inverse of the above observability matrix,

$$\frac{\partial \mathbf{T}}{\partial z_1} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}. \quad (7.79)$$

We will now find the Jacobian matrix of the desired transformation:

$$\frac{\partial \mathbf{T}}{\partial \mathbf{z}} = \left[\left(\text{ad}^0 \mathbf{f}, \frac{\partial \mathbf{T}}{\partial z_1} \right) \quad \left(\text{ad}^1 \mathbf{f}, \frac{\partial \mathbf{T}}{\partial z_1} \right) \quad \left(\text{ad}^2 \mathbf{f}, \frac{\partial \mathbf{T}}{\partial z_1} \right) \right] = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & -10 \\ 1 & -10 & 90 \end{bmatrix}. \quad (7.80)$$

Therefore, we can take

$$\mathbf{x} = \mathbf{T}(\mathbf{z}) = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & -10 \\ 1 & -10 & 90 \end{bmatrix} \mathbf{z}. \quad (7.81)$$

Note that $\mathbf{T}(\mathbf{0}) = \mathbf{0}$. The inverse transformation is

$$\mathbf{z} = \mathbf{T}^{-1} \mathbf{x} = \begin{bmatrix} 10 & 10 & 1 \\ 10 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \mathbf{x}. \quad (7.82)$$

The system in the new coordinates has the form

$$\begin{aligned} \dot{\mathbf{z}} &= \left(\frac{\partial \mathbf{T}}{\partial \mathbf{z}} \right)^{-1} \mathbf{f}(\mathbf{x}) \Big|_{\mathbf{x}=\mathbf{T}(\mathbf{z})} + \left(\frac{\partial \mathbf{T}}{\partial \mathbf{z}} \right)^{-1} \mathbf{g}(\mathbf{x}) \Big|_{\mathbf{x}=\mathbf{T}(\mathbf{z})} u = \mathbf{A}\mathbf{z} + \mathbf{y}(y, u), \\ y &= \mathbf{c}\mathbf{z} \end{aligned}$$

Performing simple manipulations yields

$$\dot{\mathbf{z}} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \mathbf{z} + \begin{bmatrix} 98 \sin(y) + 10u \\ 9.8 \sin(y) - 10y \\ -10y \end{bmatrix}, \quad (7.83)$$

$$y = [0 \quad 0 \quad 1] \mathbf{z} = z_3. \quad (7.84)$$

7.5 Asymptotic State Estimator

Once the nonlinear system model (7.53) is in the observer form, we can proceed with the construction of an asymptotic state estimator using the technique from the linear systems. The state estimator dynamics in the new coordinates are

$$\dot{\tilde{z}} = (A - lc)\tilde{z} + ly + \gamma(y, u), \quad (7.85)$$

where the estimator gain vector l is chosen so that the matrix $(A - lc)$ has its eigenvalues in the desired locations in the open left-hand complex plane. It is easy to verify that the state estimation error, $z - \tilde{z}$, satisfies the linear differential equation

$$\frac{d}{dt}(z - \tilde{z}) = (A - lc)(z - \tilde{z}).$$

Applying the inverse transformation to (7.85), we obtain the state estimator dynamics in the original coordinates.

We illustrate the above method of constructing a state estimator using the model of the one-link manipulator from Example 7.1. We first construct a state estimator for the one-link manipulator model in the observer form given by (7.83) and (7.84). Suppose that we wish the eigenvalues of the state estimator matrix, $A - lc$, to be located at $\{-9, -10, -11\}$. Then, the estimator gain vector is $l = [990 \ 299 \ 30]^T$. With the above choice of the state estimator design parameters, its dynamics in the new coordinates become

$$\dot{\tilde{z}} = \begin{bmatrix} 0 & 0 & -990 \\ 1 & 0 & -299 \\ 0 & 1 & -30 \end{bmatrix} \tilde{z} + \begin{bmatrix} 990 \\ 299 \\ 30 \end{bmatrix} y + \begin{bmatrix} 98 \sin(y) + 10u \\ 9.8 \sin(y) - 10y \\ -10y \end{bmatrix}, \quad (7.86)$$

$$y = [0 \ 0 \ 1] z = z_3. \quad (7.87)$$

Using transformations (7.81) and (7.82), we represent the state estimator dynamics in the original coordinates,

$$\dot{\tilde{x}} = \begin{bmatrix} -20 & 1 & 0 \\ -89 & 0 & 1 \\ 100 & -10 & -10 \end{bmatrix} \tilde{x} + \begin{bmatrix} 30 \\ -1 \\ 700 \end{bmatrix} y + \begin{bmatrix} -10y \\ 9.8 \sin(y) + 90y \\ -800y + 10u \end{bmatrix}, \quad (7.88)$$

$$y = [1 \ 0 \ 0] \tilde{x} = \tilde{x}_1. \quad (7.89)$$

Performing simple manipulations, we represent the above state estimator as

$$\dot{\tilde{x}} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -10 & -10 \end{bmatrix} \tilde{x} + \begin{bmatrix} 0 \\ 9.8 \sin(y) \\ 0 \end{bmatrix} - \begin{bmatrix} 20 \\ 89 \\ -100 \end{bmatrix} (\tilde{y} - y) + \begin{bmatrix} 0 \\ 0 \\ 10 \end{bmatrix} u, \quad (7.90)$$

$$y = [1 \ 0 \ 0] \tilde{x} = \tilde{x}_1, \quad (7.91)$$

$$\tilde{y} = [1 \ 0 \ 0] \tilde{x} = \tilde{x}_1. \quad (7.92)$$

7.6 Combined Controller–Estimator Compensator

Having constructed a state estimator, we use the state estimates, rather than the states themselves, in the state-feedback control law implementation. The question now is how the incorporation of the state estimator affects the stability of the closed-loop system with the combined estimator–controller compensator in the loop. Recall that in the case of linear systems, we were able to show that the closed-loop system was asymptotically stable. In the present case, one has to be careful when analyzing the stability of the closed-loop system with the combined estimator–controller compensator in the loop. The plant’s nonlinearities may cause difficulties in stabilizing the closed-loop system using the combined controller–estimator compensator. However, for certain “nice” nonlinearities—for example, the ones for which the Lipschitz condition is satisfied—one should be able to find stability conditions in terms of the nonlinearity’s Lipschitz constant and the location of the estimator’s poles rather easily.

We now illustrate the implementation of the combined estimator–controller compensator on the one-link manipulator model. In this example, the controller is

$$u = -7\tilde{x}_1 - 4.15\tilde{x}_2 - 0.25\tilde{x}_3 - 0.98\tilde{x}_2 \cos(\tilde{x}_1) - 12.25 \sin(\tilde{x}_1). \quad (7.93)$$

In Figure 7.2, plots of the one-link manipulator output for three different initial conditions are shown. The manipulator is driven by the combined controller–estimator compensator consisting of the state estimator (7.90), (7.91), and (7.92) and the control law (7.93). The initial conditions of the estimator were set to zero while the initial conditions of the plant were the same as when we implemented the state-feedback control alone. Comparing plots of Figures 7.1 and 7.2, we can see that the presence of the estimator alters the system output in a noticeable way.

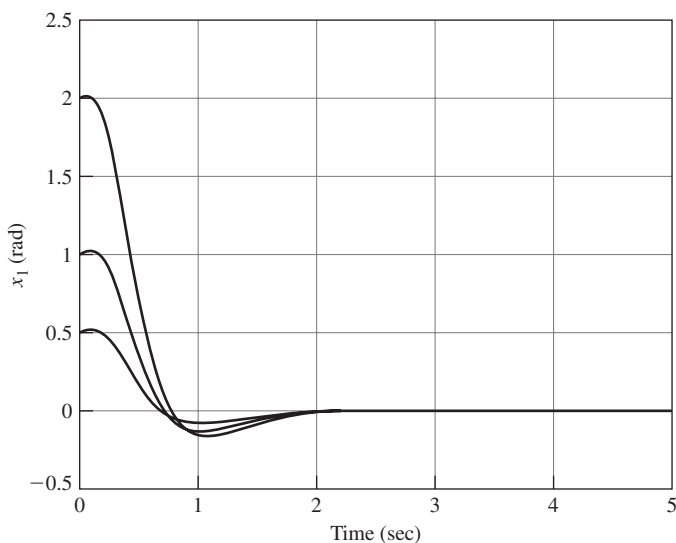


Figure 7.2 Plots of the one-link manipulator’s link angle, $x_1 = \theta$, versus time for three different initial angles $\theta(0)$. The manipulator is driven by the combined controller–estimator compensator.

Notes

Friedland [91, p. 123] has the following remarks to offer on the method of feedback linearization: “The idea underlying feedback linearization is fairly obvious. The method probably was used for many years in practical control system design without anyone having reported it. Concurrently, a number of mathematical investigations on the use of differential geometry in dynamical systems were conducted without much notice by control system designers.”

Our exposition of feedback linearization is based on the papers of Sommer [268] and Žak and MacCarley [313]. The research in applying differential geometric methods to the control design intensified after the papers by Su [270] and by Hunt, Su, and Meyer [136] appeared. Bakker and Nijmeijer [17] showed the equivalence of the Sommer’s approach and that of Hunt, Su, and Meyer. Our treatment of the subject of state estimation follows the approach of Bestle and Zeitz [27], which in turn was presented, using the differential geometry notation, in the paper by Walcott, Corless, and Žak [291]. Zeitz [316] offers a generalization of the method proposed by Walcott et al. [291]. For further study of differential geometric techniques in nonlinear feedback design, we recommend books by Isidori [139] and by Marino and Tomei [196].

EXERCISES

- 7.1** Verify the Leibniz formula (7.12).
- 7.2** Consider the one-link robot manipulator of Example 7.1 whose model is given by (7.4) and (7.5). Treat this model as a simulation model. In order to obtain a design model, eliminate the armature inductance from the design model. In other words, neglect high-frequency parasitics represented by L_a in the simulation model. For this set $L_a = 0$, compute x_3 from the equation obtained by setting $L_a = 0$ and substitute the result into the second of the state equations. The result will be a second-order design model. Represent the design model in the form

$$\dot{\mathbf{x}} = \mathbf{Ax} + \mathbf{b}(u + \xi),$$

where $\xi = 9.8 \sin(x_1)$ represents the model nonlinearity. Use the design model to do the following:

- Construct the sliding surface, $\sigma = \mathbf{s}\mathbf{x} = 0$, so that the design model restricted to this surface will have a pole at -3 .
- Construct a control law of the form

$$u = -k_1 x_1 \operatorname{sign}(x_1 \sigma) - k_2 x_2 \operatorname{sign}(x_2 \sigma)$$

so that $\sigma \dot{\sigma} < 0$ for all $\mathbf{x} \in \mathbb{R}^2$ for which $\sigma(\mathbf{x}) \neq 0$. While constructing the controller u use the fact that

$$-1 \leq \frac{\sin(x_1)}{x_1} \leq 1.$$

- Apply the above control law to the simulation model in which $L_a = 10$ mH. Plot x_1 , x_2 , and x_3 versus time and compare these plots with the plots of x_1 and x_2 versus time when the above controller was applied to the design model. Investigate

a destabilizing effect of increasing L_a on the behavior of the closed-loop system composed of the simulation model and the control law from part (b).

7.3 Let S be an open set in \mathbb{R}^n and let $f: S \rightarrow \mathbb{R}^n$, where

$$f(x) = [f_1(x) \quad f_2(x) \quad \cdots \quad f_n(x)]^T.$$

Recall that if all partial derivatives of all functions f_i exist and are continuous in S , then we say that f is a C^1 function, map, in S . We say that f is a C^p function, map, in S if each component f_i is continuous and has continuous partial derivatives in S for all $k = 1, 2, \dots, p$.

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a C^p map, $p > 0$, on an open subset S of \mathbb{R}^n , and let a be a point of S . We say that f is locally C^p -invertible at a if there is an open set $\mathcal{U} \subset S$ containing a and an open set \mathcal{V} containing $f(a)$ such that $f: \mathcal{U} \rightarrow \mathcal{V}$ has a C^p inverse $h: \mathcal{V} \rightarrow \mathcal{U}$ so that for all $x \in \mathcal{U}$,

$$h(f(x)) = x.$$

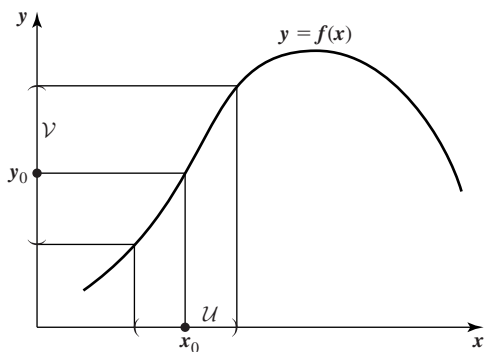


Figure 7.3 An illustration of the sets \mathcal{U} and \mathcal{V} that appear in the inverse function theorem.

Figure 7.3 illustrates the sets \mathcal{U} and \mathcal{V} for one variable case. The Jacobian matrix of $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$, denoted $J(x)$, is defined as

$$J(x) = Df(x) = \begin{bmatrix} \nabla f_1^T(x) \\ \nabla f_2^T(x) \\ \vdots \\ \nabla f_n^T(x) \end{bmatrix} = \begin{bmatrix} Df_1(x) \\ Df_2(x) \\ \vdots \\ Df_n(x) \end{bmatrix}.$$

Prove the inverse function theorem for maps:

Theorem 7.1 Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a C^p map, $p > 0$, on an open subset S of \mathbb{R}^n , and let a be a point of S such that the Jacobian matrix $J(a)$ is invertible, and $b = f(a)$. Then, there exist open sets \mathcal{U} and \mathcal{V} in \mathbb{R}^n such that $a \in \mathcal{U}$, $b \in \mathcal{V}$, f is one-to-one on \mathcal{U} , $f(\mathcal{U}) = \mathcal{V}$, and f is locally C^p -invertible at a .

7.4 For the dynamical system model

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_1 x_2 + x_2 \\ x_1 + x_2 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} u,$$

do the following:

- Construct a smooth state-space transformation, $\mathbf{z} = \mathbf{T}(\mathbf{x})$, that transforms the system into the controller form.
- Transform the system into the controller form.
- In the \mathbf{z} coordinates, design a linearizing state-feedback controller so that the linear closed-loop system has poles at $-2 \pm 2j$.
- Represent the control law in the \mathbf{x} coordinates.

7.5 Consider the following model of a dynamical system:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_1^2 + x_2 \\ x_1^3 + x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u.$$

- Transform the system into the controller form.
- Design the linearizing state-feedback controller such that the closed-loop system in the new coordinates has its poles at -1 and -2 .
- Linearize the closed-loop system about $\mathbf{x} = \mathbf{0}$ and find the poles of the linearized model.
- Perform numerical simulations of the dynamical behavior of the closed-loop system for a few initial conditions. Plot x_1 versus time as well as x_2 versus time.

7.6 Consider the following model of a dynamical system:

$$\dot{\mathbf{x}} = \begin{bmatrix} x_2 + 2x_2^2 \\ 2x_2 \end{bmatrix} + \begin{bmatrix} x_2 \\ 1 \end{bmatrix} u.$$

- Construct a smooth state-space transformation $\mathbf{z} = \mathbf{T}(\mathbf{x})$ that transforms the system into the controller form, and then transform the system into the controller form.
- In the \mathbf{z} coordinates, design a linearizing state-feedback controller so that the linear closed-loop system has poles at $-2 \pm 2j$, and then represent the control law in the \mathbf{x} coordinates.
- Perform numerical simulations of the dynamical behavior of the closed-loop system for a few initial conditions. Plot x_1 versus time as well as x_2 versus time.

7.7 Consider a dynamical system modeled by the equations

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} -x_1 \\ x_1 - x_2 - x_2^3 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} u,$$

$$y = x_2.$$

- Find a smooth state-space transformation $\mathbf{z} = \mathbf{T}(\mathbf{x})$ that brings the above system model into the controller form.

- (b) Find the system's equations in the new coordinates.
- (c) Find a smooth state-space transformation $\mathbf{x} = \mathbf{T}(\mathbf{z})$ that brings the above system model into the observer form.
- (d) Find the system's observer form.

7.8 Consider a dynamical system modeled by the following equations:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_1^2 + x_2 \\ x_1^3 + x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u,$$

$$y = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}.$$

Note that this is the same plant as in Exercise 7.5 plus the output map.

- (a) Transform the above model into the observer form.
- (b) Design a state estimator so that the poles of the linearized error dynamics are located at -3 and -4 .
- (c) Use the control law developed in Exercise 7.5 and the state estimator from the previous part of this exercise to construct a combined controller-estimator compensator. Simulate the behavior of the closed-loop system for a few different initial conditions. Note that because the state variables of the state estimator are available to you, you can always set zero initial conditions on the estimator if you wish to do so. Plot the output of the plant versus time.



CHAPTER 8

Fuzzy Systems

It is in the response to mistakes, of technique or of judgment, that one sees most clearly the human side of science's practitioners: a mix of pride, stubbornness, humility, honesty, dismay, and regret that one could expect to encounter in any section of the population. There are the inevitable diehards, whose pride forces them to concoct ever-more-elaborate explanations for the facts in order to avoid confession of error. There are the crusaders, determined at all costs to expose the foolishness. And there are the "told-you-so's," generally members of the scientific elite, who make little contribution while the debate is raging but are ready to heap scorn on the perpetrators afterward for their "obvious" mistakes. All this, of course, is the antithesis of objective science, but follows as inevitably as night follows day from the fact that science is practiced by people and not by automatons.

—*A Biography of Water: Life's Matrix* [18, p. 291]

8.1 Motivation and Basic Definitions

Conventional control algorithms require a mathematical model of the dynamical system to be controlled. The mathematical model is then used to construct a controller. In many practical situations, however, it is not always feasible to obtain an accurate mathematical model of the controlled system. *Fuzzy logic control* offers a way of dealing with modeling problems by implementing linguistic, nonformally expressed control laws derived from expert knowledge. This approach is especially successful in the case of processes such as washing machines, camera focusing, or drying processes that are being controlled well by human operators. The control of such processes can be captured in the form of rules

IF (process state) THEN (control action).

Thus, by a rule we will understand an IF-THEN statement that defines the set of facts that must be true (the IF part) before a set of actions (the THEN part) can be executed. The basic notion of fuzzy logic is a *fuzzy set*. An ordinary set, or simply a set, is a collection of objects. Ordinary sets have two membership states: inclusion or exclusion. An object is either in the set or it is

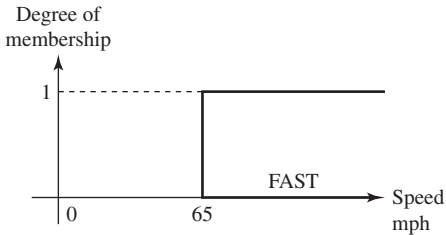


Figure 8.1 An example of a crisp set—fast cars.

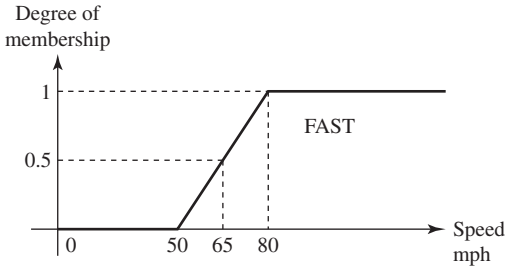


Figure 8.2 A fuzzy set FAST.

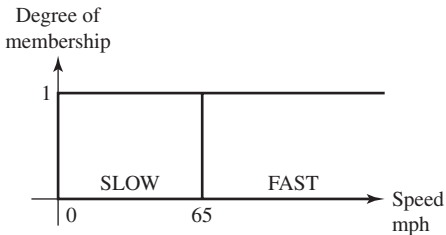


Figure 8.3 Crisp sets for SLOW and FAST.

not. Thus the degree to which an object is a member of a set is either 0 or 1. For example, a car is fast if it is driven faster than 65 mph, and it is slow if it is driven slower than 65 mph (see Figure 8.1). This is the reason we sometimes refer to ordinary sets as *crisp* sets. Fuzzy sets, on the other hand, allow for partial membership in a set. The degree to which an object is a member of a fuzzy set can vary between 0 and 1. With a fuzzy set, we have a gradual transition from membership to nonmembership. For example, a car moving with velocity of 65 mph is a member of the fuzzy set FAST to degree 0.5. A car being driven 57.5 mph is fast to degree 0.25 (see Figure 8.2). We next illustrate the crisp and fuzzy versions of the sets SLOW and FAST. Observe the sudden change of membership in the crisp SLOW and FAST sets at a speed of 65 mph in Figure 8.3. The fuzzy SLOW and FAST sets allow for a gradual change from SLOW to FAST for a car moving between 50 and 80 mph, as illustrated in Figure 8.4. Notice that a car moving at a speed between 50 and 80 mph has partial membership in SLOW and partial membership in FAST at the same time. In the above example, the variable SPEED can have two possible linguistic values: SLOW and FAST. For this reason, we refer to the variable SPEED as the *linguistic variable* because its possible values are the linguistic values: in this example they are SLOW and FAST.

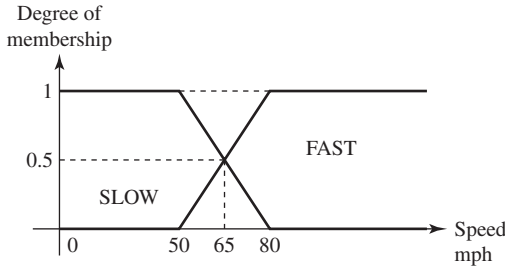


Figure 8.4 Fuzzy sets for SLOW and FAST.

With the above examples in mind, we now formally define a fuzzy set. For this, we need to introduce some notation. Let U be a set. We use $u \in U$ to denote that u is an element of U . We call U the *universe of discourse*, or just the *universe*. We first recall the definition of a function. A function A from a set X to a set Y is a set of ordered pairs in $X \times Y$ with the property that if (x, y) and (x, z) are elements of A , then $y = z$.

Definition 8.1 (Fuzzy Set) A fuzzy set A on U is a function that maps U into the interval $[0, 1]$. Equivalently,

$$A = \{(u, \mu_A(u)) : u \in U\},$$

where μ_A is commonly called the membership function of the fuzzy set A .

Note that the concept of a fuzzy set generalizes that of a crisp set in the following way. Let $V \subset U$. Define $\mu_V : U \rightarrow [0, 1]$ by $\mu_V(u) = 1$ if $u \in V$ and $\mu_V(u) = 0$ otherwise. Then V can be identified with the fuzzy set associated with μ_V . In other words, the subsets of U are identified with the functions from U to the two point set $\{0, 1\}$.

◆ Example 8.1

Consider the following universe $U = \{4, 5, 6, 7, 8\}$. Let

$$A = \{(4, 0.0), (5, 0.5), (6, 1.0), (7, 0.5), (8, 0.0)\}.$$

Then A is a fuzzy set. One interpretation of A is closeness to 6.

Let A and B be fuzzy sets on U . We define the *union*, *intersection*, and *complement* using the following membership functions:

Union: $\mu_{A \cup B}(u) = \max\{\mu_A(u), \mu_B(u)\}$ for all $u \in U$;

Intersection: $\mu_{A \cap B}(u) = \min\{\mu_A(u), \mu_B(u)\}$ for all $u \in U$;

Complement: $\mu_{\bar{A}}(u) = 1 - \mu_A(u)$ for all $u \in U$.

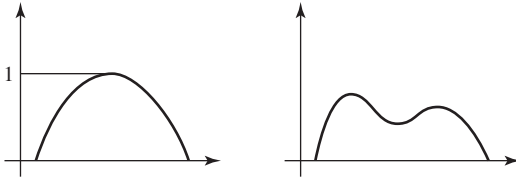


Figure 8.5 The first fuzzy set is a fuzzy number, while the second one is not.

◆ Example 8.2

Consider another fuzzy set B , “number close to 4,” given by

$$B = \{(4, 1.0), (5, 0.8), (6, 0.5), (7, 0.2), (8, 0.0)\}.$$

Then,

$$A \cup B = \{(4, 1.0), (5, 0.8), (6, 1.0), (7, 0.5), (8, 0.0)\}.$$

The intersection of the two fuzzy sets A and B is

$$A \cap B = \{(4, 0.0), (5, 0.5), (6, 0.5), (7, 0.2), (8, 0.0)\}.$$

The complement of the fuzzy set A is

$$\bar{A} = \{(4, 1.0), (5, 0.5), (6, 0.0), (7, 0.5), (8, 1.0)\}.$$

Definition 8.2 (Fuzzy Number) A fuzzy number is a fuzzy set F on a universe $U \subset \mathbb{R}$ such that

$$\max_{u \in U} \mu_F(u) = 1$$

and

$$\mu_F(\lambda u_1 + (1 - \lambda)u_2) \geq \min(\mu_F(u_1), \mu_F(u_2))$$

for all $u_1, u_2 \in U$ and all $\lambda \in [0, 1]$.

A fuzzy set that satisfies the conditions $\max_{u \in U} \mu_F(u) = 1$ and $\mu_F(\lambda u_1 + (1 - \lambda)u_2) \geq \min(\mu_F(u_1), \mu_F(u_2))$ is called *normal* and *convex*, respectively. Thus, a fuzzy number is a fuzzy set that is both normal and convex. Examples of convex and nonconvex fuzzy sets are given in Figure 8.5.

8.2 Fuzzy Arithmetic and Fuzzy Relations

8.2.1 Interval Arithmetic

We begin our discussion by introducing the concept of the *interval of confidence*. For this, imagine an experiment where a measurement was performed, resulting in a data point a . The measurement was performed in nonideal conditions. As a result of the presence of the noise

measurement in the experiment, we can say that the data point a lies in the interval $[a_1, a_2]$, that is,

$$a \in [a_1, a_2].$$

We refer to the closed interval $A = [a_1, a_2]$ as the interval of confidence. A fuzzy number can be viewed as an extension of the concept of the interval of confidence. This extension can be used to devise useful interpretations to various operations on fuzzy numbers. We now discuss operations that we can perform on the intervals of confidence. We start with the operation of addition.

Let $A = [a_1, a_2]$ and $B = [b_1, b_2]$ be two intervals of confidence. If $x \in [a_1, a_2]$ and $y \in [b_1, b_2]$, then $x + y \in [a_1 + b_1, a_2 + b_2]$. We write this symbolically as

$$\begin{aligned} A (+) B &= [a_1, a_2] (+) [b_1, b_2] \\ &= [a_1 + b_1, a_2 + b_2]. \end{aligned}$$

To prove the above, note that if $a_1 \leq x$ and $b_1 \leq y$, then $a_1 + b_1 \leq x + y$. Similarly, if $x \leq a_2$ and $y \leq b_2$, then $x + y \leq a_2 + b_2$.

We now define the operation of subtraction on the intervals of confidence. Let $A = [a_1, a_2]$ and $B = [b_1, b_2]$ be two intervals of confidence. If $x \in [a_1, a_2]$ and $y \in [b_1, b_2]$, then $x - y \in [a_1 - b_2, a_2 - b_1]$. We write this symbolically as

$$\begin{aligned} A (-) B &= [a_1, a_2] (-) [b_1, b_2] \\ &= [a_1 - b_2, a_2 - b_1]. \end{aligned}$$

We can also multiply the given intervals of confidence. We first discuss the simple case when the intervals of confidence belong to the set of nonnegative real numbers denoted \mathbb{R}^+ . Let $A = [a_1, a_2]$ and $B = [b_1, b_2]$ be two intervals of confidence such that $A, B \subset \mathbb{R}^+$. If $x \in [a_1, a_2]$, and $y \in [b_1, b_2]$, then $xy \in [a_1b_1, a_2b_2]$, written symbolically as

$$\begin{aligned} A (\cdot) B &= [a_1, a_2] (\cdot) [b_1, b_2] \\ &= [a_1b_1, a_2b_2]. \end{aligned}$$

When the intervals of confidence are arbitrary intervals of the set of real numbers \mathbb{R} , we have nine possible combinations that are discussed in Kaufmann and Gupta [151]. We list the above-mentioned nine cases in Table 8.1.

Table 8.1 Multiplication of the Intervals of Confidence in \mathbb{R}

A	B	$A (\cdot) B$
$0 \leq a_1 \leq a_2$	$0 \leq b_1 \leq b_2$	$[a_1b_1, a_2b_2]$
$0 \leq a_1 \leq a_2$	$b_1 \leq 0 \leq b_2$	$[a_2b_1, a_2b_2]$
$0 \leq a_1 \leq a_2$	$b_1 \leq b_2 \leq 0$	$[a_2b_1, a_1b_2]$
$a_1 \leq 0 \leq a_2$	$0 \leq b_1 \leq b_2$	$[a_1b_1, a_2b_2]$
$a_1 \leq 0 \leq a_2$	$b_1 \leq 0 \leq b_2$	$[\min(a_1b_2, a_2b_1), \max(a_1b_1, a_2b_2)]$
$a_1 \leq 0 \leq a_2$	$b_1 \leq b_2 \leq 0$	$[a_2b_1, a_1b_1]$
$a_1 \leq a_2 \leq 0$	$0 \leq b_1 \leq b_2$	$[a_1b_2, a_2b_1]$
$a_1 \leq a_2 \leq 0$	$b_1 \leq 0 \leq b_2$	$[a_1b_2, a_1b_1]$
$a_1 \leq a_2 \leq 0$	$b_1 \leq b_2 \leq 0$	$[a_2b_2, a_1b_1]$

Observe that the multiplication of an interval of confidence $A \subset \mathbb{R}^+$ by a real number $k \in \mathbb{R}$ can be easily deduced from the operation of multiplication of two intervals of confidence. First, assume that $k \in \mathbb{R}^+$ and represent it as

$$k = [k, \quad k].$$

Then,

$$\begin{aligned} k \cdot A &= k \cdot [a_1, \quad a_2] \\ &= [k, \quad k] (\cdot) [a_1, \quad a_2] \\ &= [ka_1, \quad ka_2]. \end{aligned}$$

If $k < 0$, then

$$k \cdot [a_1, \quad a_2] = [ka_2, \quad ka_1].$$

We now discuss the operations of maximum and minimum of two intervals of confidence. First, we explain the meaning of the symbols we are going to use. Consider two real numbers a and b . By definition, we have

$$\begin{aligned} a \wedge b &= \min(a, b) \\ &= \begin{cases} a & \text{if } a \leq b, \\ b & \text{if } b < a, \end{cases} \end{aligned}$$

and

$$\begin{aligned} a \vee b &= \max(a, b) \\ &= \begin{cases} b & \text{if } a \leq b, \\ a & \text{if } b < a. \end{cases} \end{aligned}$$

Consider now two intervals of confidence $A, B \subset \mathbb{R}$, where $A = [a_1, \quad a_2]$ and $B = [b_1, \quad b_2]$. We introduce two operations (\wedge) and (\vee) on the intervals of confidence. We define

$$\begin{aligned} A (\wedge) B &= [a_1, \quad a_2] (\wedge) [b_1, \quad b_2] \\ &= [a_1 \wedge b_1, \quad a_2 \wedge b_2] \end{aligned}$$

and

$$\begin{aligned} A (\vee) B &= [a_1, \quad a_2] (\vee) [b_1, \quad b_2] \\ &= [a_1 \vee b_1, \quad a_2 \vee b_2]. \end{aligned}$$

The above operations on the intervals will be used in the next section to define the corresponding operations on fuzzy numbers.

8.2.2 Manipulating Fuzzy Numbers

By relating the concept of the interval of confidence with the new concept of the *level of presumption*, we arrive at the concept of fuzzy number. Recall that a fuzzy number is a fuzzy set that is convex and normal. The fuzzy set $A \subset \mathbb{R} \times \mathbb{R}$ is convex if for any $x, y \in \mathbb{R}$ and $\lambda \in [0, 1]$

$$\mu_A(\lambda x + (1 - \lambda)y) \geq \mu_A(x) \wedge \mu_A(y).$$

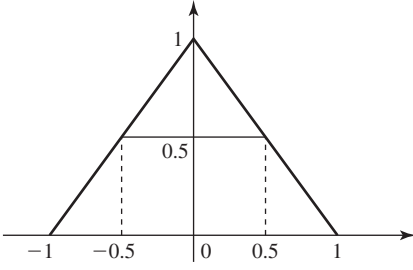


Figure 8.6 A fuzzy zero.

The fuzzy set $A \subset \mathbb{R} \times \mathbb{R}$ is normal if for all $x \in \mathbb{R}$

$$\bigvee_x \mu_A(x) = 1,$$

where \bigvee denotes the operation of taking a maximum. Thus, the fuzzy set $A \subset \mathbb{R} \times \mathbb{R}$ is normal if the largest value of $\mu_A(x)$ is equal to 1.

We consider now a fuzzy zero as depicted in Figure 8.6. The interval of confidence corresponding to the level of presumption of 0 is $[-1, 1]$. The interval of confidence corresponding to the level of presumption of 0.5 is $[-0.5, 0.5]$, and so on. Thus the level of presumption $\alpha \in [0, 1]$ gives an interval of confidence

$$A_\alpha = [a_1^{(\alpha)}, a_2^{(\alpha)}].$$

Note that if $\alpha_1 < \alpha_2$, then

$$[a_1^{(\alpha_2)}, a_2^{(\alpha_2)}] \subseteq [a_1^{(\alpha_1)}, a_2^{(\alpha_1)}],$$

or, equivalently, $A_{\alpha_2} \subseteq A_{\alpha_1}$. Thus, if the presumption level α is increasing, then the corresponding interval of confidence is nonincreasing; that is, A_α is a nonincreasing function of α . In the previous section we discussed how to add intervals of confidence. The addition of fuzzy numbers can be performed following the same process, level by level. For example, let A and B be two fuzzy numbers, and let A_α and B_α be their intervals of confidence for the level of presumption $\alpha \in [0, 1]$; that is, $A_\alpha = \{x : \mu_A(x) \geq \alpha\}$ and $B_\alpha = \{x : \mu_B(x) \geq \alpha\}$. Then,

$$\begin{aligned} A_\alpha (+) B_\alpha &= [a_1^{(\alpha)}, a_2^{(\alpha)}] (+) [b_1^{(\alpha)}, b_2^{(\alpha)}] \\ &= [a_1^{(\alpha)} + b_1^{(\alpha)}, a_2^{(\alpha)} + b_2^{(\alpha)}]. \end{aligned}$$

Similarly, we can subtract two given fuzzy numbers,

$$\begin{aligned} A_\alpha (-) B_\alpha &= [a_1^{(\alpha)}, a_2^{(\alpha)}] (-) [b_1^{(\alpha)}, b_2^{(\alpha)}] \\ &= [a_1^{(\alpha)} - b_2^{(\alpha)}, a_2^{(\alpha)} - b_1^{(\alpha)}]. \end{aligned}$$

◆ Example 8.3

Consider two triangular fuzzy numbers whose membership functions are

$$\mu_A(x) = \begin{cases} 0 & \text{if } x \leq 0.5, \\ 2/5x - 1/5 & \text{if } 0.5 < x < 3, \\ -1/2x + 5/2 & \text{if } 3 \leq x < 5, \\ 0 & \text{if } x \geq 5 \end{cases}$$

and

$$\mu_B(x) = \begin{cases} 0 & \text{if } x \leq 1, \\ x - 1 & \text{if } 1 < x < 2, \\ -1/5x + 7/5 & \text{if } 2 \leq x < 7, \\ 0 & \text{if } x \geq 7. \end{cases}$$

We now compute the intervals of confidence for the above fuzzy numbers for each level of presumption $\alpha \in [0, 1]$. We begin with the fuzzy number A . In a sense, our goal is to express the triangular shape of the fuzzy number A in the form $A_\alpha = [a_1^{(\alpha)}, a_2^{(\alpha)}]$ in terms of α . From the description of A above, we obtain

$$\alpha = 2/5a_1^{(\alpha)} - 1/5$$

and

$$\alpha = -1/2a_2^{(\alpha)} + 5/2.$$

Hence, the interval of confidence at the level α for the fuzzy number A is

$$\begin{aligned} A_\alpha &= [a_1^{(\alpha)}, a_2^{(\alpha)}] \\ &= [5/2\alpha + 1/2, -2\alpha + 5]. \end{aligned}$$

Proceeding in a similar manner, we express the triangular shape of the fuzzy number B in the form $B_\alpha = [b_1^{(\alpha)}, b_2^{(\alpha)}]$ in terms of α . From the description of B above, we obtain

$$\alpha = b_1^{(\alpha)} - 1,$$

and

$$\alpha = -1/5b_2^{(\alpha)} + 7/5.$$

To find the interval of confidence at the level α for the fuzzy number B , we perform simple manipulations to get

$$\begin{aligned} B_\alpha &= [b_1^{(\alpha)}, b_2^{(\alpha)}] \\ &= [\alpha + 1, -5\alpha + 7]. \end{aligned}$$

We next compute the sum of the above intervals of confidence,

$$\begin{aligned} A_\alpha (+) B_\alpha &= [a_1^{(\alpha)} + b_1^{(\alpha)}, a_2^{(\alpha)} + b_2^{(\alpha)}] \\ &= [7/2\alpha + 3/2, -7\alpha + 12]. \end{aligned}$$

From the above, we can find analytic description of the fuzzy number $A(+)B$ in terms of its membership function $\mu_{A(+)B}$. Because $a_1^{(\alpha)} + b_1^{(\alpha)} = 7/2\alpha + 3/2$, and $a_2^{(\alpha)} + b_2^{(\alpha)} = -7\alpha + 12$, solving for α yields

$$\mu_{A(+)B}(x) = \begin{cases} 0 & \text{if } x \leq 3/2, \\ 2/7x - 3/7 & \text{if } 3/2 < x \leq 5, \\ -1/7x + 12/7 & \text{if } 5 < x \leq 12, \\ 0 & \text{if } x > 12. \end{cases}$$

We now compute the fuzzy number $A(-)B$. First, we find the expression for the interval of confidence $A_\alpha(-)B_\alpha$ corresponding to the level α :

$$\begin{aligned} A_\alpha(-)B_\alpha &= [a_1^{(\alpha)} - b_2^{(\alpha)}, a_2^{(\alpha)} - b_1^{(\alpha)}] \\ &= [15/2\alpha - 13/2, -3\alpha + 4]. \end{aligned}$$

Thus,

$$a_1^{(\alpha)} - b_2^{(\alpha)} = 15/2\alpha - 13/2$$

and

$$a_2^{(\alpha)} - b_1^{(\alpha)} = -3\alpha + 4.$$

We use the above two equations to solve for α , which then allows us to find analytic expression for the membership function $\mu_{A(-)B}$. Performing simple manipulations yields

$$\mu_{A(-)B}(x) = \begin{cases} 0 & \text{if } x \leq -13/2, \\ 2/15x + 13/15 & \text{if } -13/2 < x \leq 1, \\ -1/3x + 4/3 & \text{if } 1 < x \leq 4, \\ 0 & \text{if } x > 4. \end{cases}$$

As far as the multiplication is concerned, we consider only the case when for $\alpha \in [0, 1]$ the intervals of confidence $A_\alpha, B_\alpha \subset \mathbb{R}^+$. In such a case the multiplication of the two fuzzy numbers A and B can be performed level by level on the intervals of confidence following the process of multiplication of the intervals of confidence in \mathbb{R}^+ , where

$$\begin{aligned} A_\alpha(\cdot)B_\alpha &= [a_1^{(\alpha)}, a_2^{(\alpha)}](\cdot)[b_1^{(\alpha)}, b_2^{(\alpha)}] \\ &= [a_1^{(\alpha)}b_1^{(\alpha)}, a_2^{(\alpha)}b_2^{(\alpha)}]. \end{aligned}$$

From the above we can deduce the way we can multiply a fuzzy number by a nonnegative real number k . We have

$$\begin{aligned} k \cdot A_\alpha &= [k, k](\cdot)[a_1^{(\alpha)}, a_2^{(\alpha)}] \\ &= [ka_1^{(\alpha)}, ka_2^{(\alpha)}]. \end{aligned}$$

The multiplication of fuzzy numbers in \mathbb{R} is discussed in detail in Kaufmann and Gupta [151, pp. 325–334].

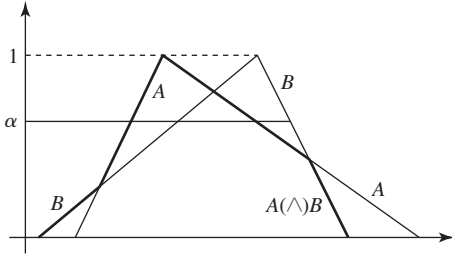


Figure 8.7 Calculating the fuzzy minimum.

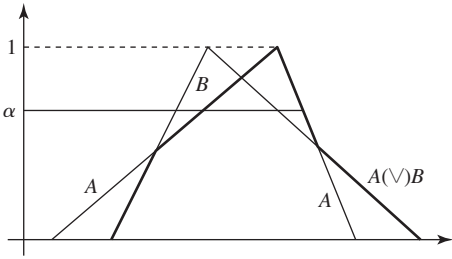


Figure 8.8 Calculating the fuzzy maximum of two fuzzy numbers.

The fuzzy minimum of two fuzzy numbers A and B is the fuzzy number that can be obtained by performing the following level by level operation:

$$\begin{aligned} A_{\alpha} (\wedge) B_{\alpha} &= [a_1^{(\alpha)}, a_2^{(\alpha)}] (\wedge) [b_1^{(\alpha)}, b_2^{(\alpha)}] \\ &= [a_1^{(\alpha)} \wedge b_1^{(\alpha)}, a_2^{(\alpha)} \wedge b_2^{(\alpha)}]. \end{aligned}$$

We illustrate a graphical procedure for calculating the fuzzy minimum in Figure 8.7.

The fuzzy maximum of the two fuzzy numbers is a fuzzy number that can be obtained by performing the following level-by-level operation:

$$\begin{aligned} A_{\alpha} (\vee) B_{\alpha} &= [a_1^{(\alpha)}, a_2^{(\alpha)}] (\vee) [b_1^{(\alpha)}, b_2^{(\alpha)}] \\ &= [a_1^{(\alpha)} \vee b_1^{(\alpha)}, a_2^{(\alpha)} \vee b_2^{(\alpha)}]. \end{aligned}$$

The process of calculating the fuzzy maximum is illustrated in Figure 8.8. We now present a numerical example to illustrate analytical method for calculating the fuzzy minimum and maximum of two fuzzy numbers.

◆ Example 8.4

Consider the two triangular fuzzy numbers A and B from Example 8.3. Recall that

$$A_{\alpha} = [a_1^{(\alpha)}, a_2^{(\alpha)}] = [5/2\alpha + 1/2, -2\alpha + 5]$$

and

$$B_{\alpha} = [b_1^{(\alpha)}, b_2^{(\alpha)}] = [\alpha + 1, -5\alpha + 7].$$

Therefore, for the fuzzy minimum

$$A_\alpha (\wedge) B_\alpha = [(5/2\alpha + 1/2) \wedge (\alpha + 1), \quad (-2\alpha + 5) \wedge (-5\alpha + 7)]$$

we find that

$$A_\alpha (\wedge) B_\alpha = \begin{cases} [5/2\alpha + 1/2, -2\alpha + 5] & \text{if } 0 \leq \alpha \leq 1/3, \\ [\alpha + 1, -2\alpha + 5] & \text{if } 1/3 < \alpha \leq 2/3, \\ [\alpha + 1, -5\alpha + 7] & \text{if } 2/3 < \alpha \leq 1. \end{cases}$$

From the above we can find an analytic expression for the membership function of the fuzzy number $A(\wedge) B$:

$$\mu_{A(\wedge) B}(x) = \begin{cases} 0 & \text{if } x \leq 1/2, \\ 2/5x - 1/5 & \text{if } 1/2 < x \leq 4/3, \\ x - 1 & \text{if } 4/3 < x \leq 2, \\ -1/5x + 7/5 & \text{if } 2 < x \leq 11/3, \\ -1/2x + 5/2 & \text{if } 11/3 < x \leq 5, \\ 0 & \text{if } 5 < x. \end{cases}$$

We next find the fuzzy maximum of the fuzzy numbers A and B :

$$\begin{aligned} A_\alpha (\vee) B_\alpha &= [a_1^{(\alpha)} \vee b_1^{(\alpha)}, \quad a_2^{(\alpha)} \vee b_2^{(\alpha)}] \\ &= [(5/2\alpha + 1/2) \vee (\alpha + 1), \quad (-2\alpha + 5) \vee (-5\alpha + 7)]. \end{aligned}$$

Performing simple manipulations yields

$$A_\alpha (\vee) B_\alpha = \begin{cases} [\alpha + 1, -5\alpha + 7] & \text{if } 0 \leq \alpha \leq 1/3, \\ [5/2\alpha + 1/2, -5\alpha + 7] & \text{if } 1/3 < \alpha \leq 2/3, \\ [5/2\alpha + 1/2, -2\alpha + 5] & \text{if } 2/3 < \alpha \leq 1. \end{cases}$$

Using the above, we compute

$$\mu_{A(\vee) B}(x) = \begin{cases} 0 & \text{if } x \leq 1, \\ x - 1 & \text{if } 1 < x \leq 4/3, \\ 2/5x - 1/5 & \text{if } 4/3 < x \leq 3, \\ -1/2x + 5/2 & \text{if } 3 < x \leq 11/3, \\ -1/5x + 7/5 & \text{if } 11/3 < x \leq 7, \\ 0 & \text{if } 7 < x. \end{cases}$$

8.2.3 Fuzzy Relations

In this subsection we discuss fuzzy relations. We begin with the definition of a fuzzy relation.

Definition 8.3 Let U and V be two universes of discourse. A fuzzy relation R between U and V is a fuzzy set in the product space

$$U \times V = \{(u, v) : u \in U, v \in V\}$$

characterized by a membership function

$$\mu_R : U \times V \rightarrow [0, 1].$$

Suppose we are given two universes of discourse:

$$U = \{u_1, u_2, \dots, u_m\} \quad \text{and} \quad V = \{v_1, v_2, \dots, v_n\}.$$

A fuzzy relation in $U \times V$ can be represented by an $m \times n$ matrix:

$$\mathbf{R} = \begin{bmatrix} \mu_R(u_1, v_1) & \mu_R(u_1, v_2) & \cdots & \mu_R(u_1, v_n) \\ \mu_R(u_2, v_1) & \mu_R(u_2, v_2) & \cdots & \mu_R(u_2, v_n) \\ \vdots & \vdots & \ddots & \vdots \\ \mu_R(u_m, v_1) & \mu_R(u_m, v_2) & \cdots & \mu_R(u_m, v_n) \end{bmatrix}.$$

Terano, Asai, and Sugeno [280] refer to the matrix \mathbf{R} as a fuzzy matrix. Because $\mu_R : U \times V \rightarrow [0, 1]$, the values of the elements of the fuzzy matrix \mathbf{R} are in the interval $[0, 1]$.

◆ Example 8.5

Let $U = \{-1, 0, 1\}$ and $V = \{1, 2, 3, 4\}$. An example of a fuzzy matrix corresponding to a fuzzy relation in $U \times V$ is

$$\mathbf{R} = \begin{bmatrix} 0.2 & 0.6 & 0.7 & 0.9 \\ 0.0 & 0.3 & 0.6 & 1 \\ 0.1 & 0.4 & 0.5 & 0.8 \end{bmatrix}.$$

Definition 8.3 is concerned with a binary fuzzy relation. We can generalize this definition to the n -ary fuzzy relation R in $X_1 \times X_2 \times \cdots \times X_n$ as a fuzzy set in the product space $X_1 \times X_2 \times \cdots \times X_n$ with a membership function:

$$\mu_R : X_1 \times X_2 \times \cdots \times X_n \rightarrow [0, 1].$$

A special kind of a fuzzy relation is a fuzzy implication that we discuss next. Let A and B be fuzzy sets in U and V , respectively. A fuzzy implication from A to B , denoted $A \rightarrow B$, is a fuzzy relation in $U \times V$ that is characterized by a specific membership function. Among different types of membership functions that can characterize a fuzzy implication, we list two classes:

- A fuzzy conjunction denoted $\mu_A(u) \star \mu_B(v)$
- A fuzzy disjunction denoted $\mu_A(u) \dot{+} \mu_B(v)$

A fuzzy conjunction, also called a t -norm, represents a generalized intersection operator, while a fuzzy disjunction, also called a t -conorm, represents a generalized union operator. Examples of a t -norm are as follows:

- Fuzzy intersection = $\min\{\mu_A(u), \mu_B(v)\}$
- Algebraic product = $\mu_A(u)\mu_B(v)$

- Bounded product = $\max\{0, \mu_A + \mu_B - 1\}$
- Drastic product =
$$\begin{cases} \mu_A & \text{if } \mu_B = 1, \\ \mu_B & \text{if } \mu_A = 1, \\ 0 & \text{if } \mu_A, \mu_B < 1 \end{cases}$$

Examples of a t -conorm are as follows:

- Fuzzy union = $\max\{\mu_A(u), \mu_B(v)\}$
- Algebraic sum = $\mu_A(u) + \mu_B(v) - \mu_A(u)\mu_B(v)$
- Bounded sum = $\min\{1, \mu_A + \mu_B\}$
- Drastic sum =
$$\begin{cases} \mu_A & \text{if } \mu_B = 0 \\ \mu_B & \text{if } \mu_A = 0 \\ 1 & \text{if } \mu_A, \mu_B > 0 \end{cases}$$

A fuzzy implication $A \rightarrow B$ can be viewed as a fuzzy IF–THEN rule of the form

$$\text{IF } x \text{ is } A, \quad \text{THEN } y \text{ is } B,$$

where x and y are linguistic variables and A and B are the corresponding values of the linguistic variables.

8.2.4 Composition of Fuzzy Relations

Suppose we are given a fuzzy relation in $U \times V$ labeled R , along with a fuzzy relation in $V \times W$ labeled S . The max-star composition of R and S is a fuzzy relation, denoted $R \circ S$, that is characterized by the membership function

$$\mu_{R \circ S}(u, w) = \bigvee_{v \in V} \{\mu_R(u, v) \star \mu_S(v, w)\},$$

where \bigvee denotes the maximum operator, and \star refers to any operator in the class of t norms. Note that $R \circ S$ is a fuzzy set in $U \times W$.

◆ Example 8.6

Let $U = \{u_1, u_2\}$, $V = \{v_1, v_2, v_3\}$, and $W = \{w_1, w_2\}$ be the given three universes of discourse. Let

$$\mathbf{R} = \begin{bmatrix} 0.5 & 0.7 & 0.0 \\ 0.8 & 1.0 & 0.2 \end{bmatrix}$$

be the fuzzy matrix that represents the relation R in $U \times V$, and let

$$\mathbf{S} = \begin{bmatrix} 0.4 & 0.6 \\ 0.2 & 0.9 \\ 0.0 & 0.5 \end{bmatrix}$$

be the given fuzzy matrix representing the relation S in $V \times W$. Let the \star operator be the minimum operator \wedge . Then the fuzzy matrix representing the max–min

fuzzy relation $R \circ S$ has the form

$$\begin{aligned}
 & \mathbf{R} \circ \mathbf{S} \\
 &= \begin{bmatrix} 0.5 & 0.7 & 0.0 \\ 0.8 & 1.0 & 0.2 \end{bmatrix} \circ \begin{bmatrix} 0.4 & 0.6 \\ 0.2 & 0.9 \\ 0.1 & 0.5 \end{bmatrix} \\
 &= \begin{bmatrix} (0.5 \wedge 0.4) \vee (0.7 \wedge 0.2) \vee (0.0 \wedge 0.1) & (0.5 \wedge 0.6) \vee (0.7 \wedge 0.9) \vee (0.0 \wedge 0.5) \\ (0.8 \wedge 0.4) \vee (1.0 \wedge 0.2) \vee (0.2 \wedge 0.1) & (0.8 \wedge 0.6) \vee (1.0 \wedge 0.9) \vee (0.2 \wedge 0.5) \end{bmatrix} \\
 &= \begin{bmatrix} 0.4 & 0.7 \\ 0.4 & 0.9 \end{bmatrix}.
 \end{aligned}$$

Other kinds of composition of fuzzy relations are possible. In particular, we can define the dual composition to the $\max - \star$ as the $\min - \dot{+}$. The $\min - \dot{+}$ composition of the fuzzy relations R and S is characterized by the membership function

$$\mu_{R \circ S}(u, w) = \bigwedge_{v \in V} \{\mu_R(u, v) \dot{+} \mu_S(v, w)\},$$

where \bigwedge refers to the minimum operator, and $\dot{+}$ can be any operator from the class of t conorm that we discussed in the previous section.

A particular type of the composition of fuzzy relations is the $\max - \star$ composition when S is just a fuzzy set A in V , rather than a fuzzy relation in the product space $V \times W$. In this case the fuzzy matrix representing S is a column vector. This means that the fuzzy matrix representing the composition $R \circ A$ is also a column vector. A general form of the membership function that characterizes such a composition is

$$\mu_{R \circ A}(w) = \bigvee_{v \in V} \{\mu_R(u, v) \star \mu_A(v)\}.$$

Note that the composition of the fuzzy relation R and the fuzzy set A is a fuzzy set in U .

◆ Example 8.7

Let $U = \{u_1, u_2, u_3\}$ and $V = \{v_1, v_2, v_3\}$ be given two universes of discourse. Let A be a fuzzy set in V described by the fuzzy column vector $\mathbf{A} = [0.2 \ 0.8 \ 1.0]^T$. The fuzzy matrix \mathbf{R} that describes the fuzzy relation R in $U \times V$ is

$$\mathbf{R} = \begin{bmatrix} 0.5 & 0.4 & 0.8 \\ 0.8 & 0.6 & 0.4 \\ 0.9 & 0.7 & 0.1 \end{bmatrix}.$$

Assume the \wedge operator for the \star operator—that is, the membership function of the composition of R and A —has the form

$$\mu_{R \circ A}(w) = \bigvee_{v \in V} \{\mu_R(u, v) \wedge \mu_A(v)\}.$$

The fuzzy column vector corresponding to the composition $R \circ A$ is

$$\begin{aligned}
 R \circ A &= \begin{bmatrix} 0.5 & 0.4 & 0.8 \\ 0.8 & 0.6 & 0.4 \\ 0.9 & 0.7 & 0.1 \end{bmatrix} \circ \begin{bmatrix} 0.2 \\ 0.8 \\ 1.0 \end{bmatrix} \\
 &= \begin{bmatrix} (0.5 \wedge 0.2) \vee (0.4 \wedge 0.8) \vee (0.8 \wedge 1.0) \\ (0.8 \wedge 0.2) \vee (0.6 \wedge 0.8) \vee (0.4 \wedge 1.0) \\ (0.9 \wedge 0.2) \vee (0.7 \wedge 0.8) \vee (0.1 \wedge 1.0) \end{bmatrix} \\
 &= \begin{bmatrix} 0.8 \\ 0.6 \\ 0.7 \end{bmatrix}.
 \end{aligned}$$

We can think of A as the fuzzy input to the fuzzy system modeled by the fuzzy relation R . Let $B = R \circ A$. Then the output of the fuzzy system is B , which is a fuzzy set. In other words, the fuzzy system modeled by a fuzzy relation R yields a fuzzy output in response to a fuzzy input. We illustrate the above in Figure 8.9.

8.3 Standard Additive Model

Let X and Y be nonempty sets and let I and J be nonempty index sets. Suppose we are given collections of fuzzy sets $\{A_\alpha : \alpha \in I\}$ and $\{B_\beta : \beta \in J\}$ on X and Y , respectively, and rules of the form

$$\text{IF } A_\alpha, \quad \text{THEN } B_\beta.$$

The IF–THEN rules along with the collections of fuzzy sets form a function R from $\{A_\alpha : \alpha \in I\}$ to $\{B_\beta : \beta \in J\}$, which we call the *fuzzy system* associated with the rules.

The *support* of a function into \mathbb{R} is the subset of the domain where the function is nonzero. The support of the fuzzy sets A_α and B_β in the IF–THEN rules describe subsets, or patches, in the Cartesian product space $X \times Y$, as illustrated in Figure 8.10. Suppose the support of each A_α and each B_β is a one-point set, say $x_\alpha \in X$ and $y_\beta \in Y$, respectively. If $X = \bigcup \{x_\alpha : \alpha \in I\}$, then R can be identified with a function \tilde{R} from X to Y by

$$\tilde{R}(x_\alpha) = y_\alpha$$

if

$$R(A_\alpha) = B_\alpha.$$

One can think of this case as having complete knowledge and the patches are just points.

A fuzzy system can be thought of as a “fuzzified” function from X to Y with fuzzy sets replacing points in the domain. The size of a patch is usually inversely proportional to the

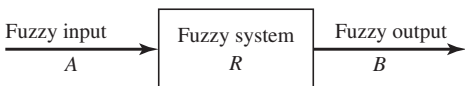


Figure 8.9 A fuzzy system produces a fuzzy output in response to a fuzzy input.

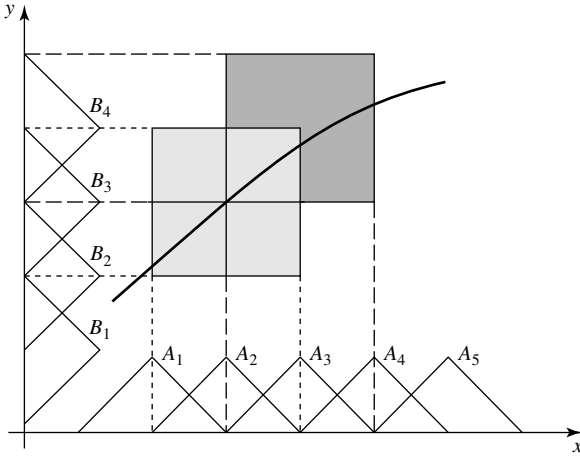


Figure 8.10 Function approximation using IF-THEN rules.

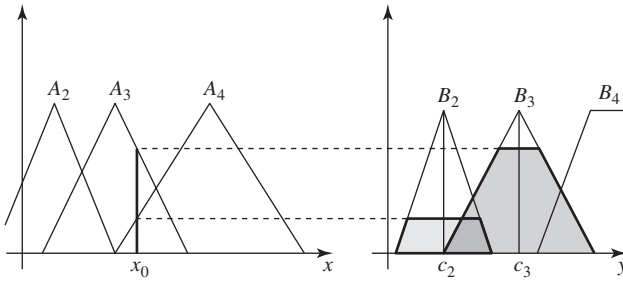


Figure 8.11 Generating output fuzzy sets from THEN-part sets of IF-THEN rules.

knowledge or confidence of the situation. A point $x \in X$ is usually in the support of more than one A_α with confidence represented by the values of the membership functions μ_α . The fuzzy set $R(A_\alpha)$ represents the desired action corresponding to A_α . This is illustrated in Figure 8.11. Suppose we are given the following two rules:

IF A_3 , THEN B_3 ,

IF A_4 , THEN B_2 .

The point x_0 is contained in the support of A_3 and A_4 . Thus, according to the given rules, both B_3 and B_2 need to be taken into account. The process by which we combine these actions is called *defuzzification*. The process of defuzzification is also known as fuzzy reasoning in the literature. The outcome of defuzzification is a numerical value $y \in Y$. An example of a defuzzification algorithm is a centroidal defuzzifier of the form

$$y = \frac{\sum_\alpha \mu_{A_\alpha}(x) \xi(x, A_\alpha, R(A_\alpha))}{\sum_\alpha \mu_{A_\alpha}(x)}, \quad (8.1)$$

where $\xi(x, A_\alpha, R(A_\alpha))$ is a function into the real numbers. We now give an example of ξ .

Let $\xi(x, A_\alpha, R(A_\alpha))$ be the centroid of the support of $R(A_\alpha)$. For example, the centroids of the support of $R(A_3) = B_3$ is c_3 , and that of $R(A_4) = B_2$ is c_2 . A fuzzy system with the centroidal

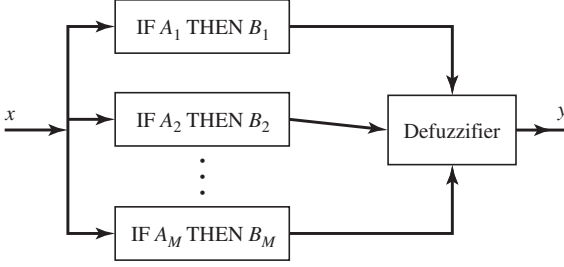


Figure 8.12 Single-input, single-output standard additive model schematic representation.

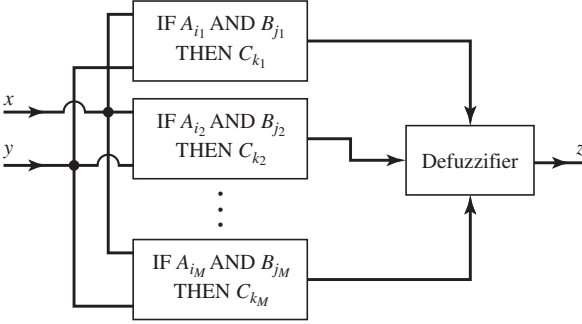


Figure 8.13 Two-input, single-output standard additive model representation.

defuzzifier is called the *standard additive model* (SAM). The term standard additive model was coined by Kosko [165]. A schematic representation of a SAM is shown in Figure 8.12. In the above description of the SAM, we were concerned with a single-input, single-output case. The above construction can be extended to a many-inputs case. Suppose we are given IF–THEN rules of the form

$$\text{IF } A_i \text{ AND } B_i, \quad \text{THEN } C_i.$$

The SAM corresponding to this case is illustrated in Figure 8.13. We could use a defuzzification algorithm of the form

$$z = \frac{\sum_{l=1}^M \mu_{A_l}(x) \mu_{B_l}(y) c_l}{\sum_{l=1}^M \mu_{A_l}(x) \mu_{B_l}(y)}, \quad (8.2)$$

where M is the number of IF–THEN rules, and c_l is the centroid of the fuzzy set C_l , $l = 1, 2, \dots, M$. This type of defuzzifier is called the *center average defuzzifier*. Another type of defuzzifier that we could use is

$$z = \frac{\sum_{l=1}^M \min\{\mu_{A_l}(x), \mu_{B_l}(y)\} c_l}{\sum_{l=1}^M \min\{\mu_{A_l}(x) \mu_{B_l}(y)\}}. \quad (8.3)$$

An illustration of fuzzy reasoning using defuzzification scheme (8.3) is given in Figure 8.14. In general, for q inputs the center average defuzzifier takes the form

$$z = \frac{\sum_{l=1}^M \prod_{j=1}^q \mu_{A_j^l}(x_j) c_l}{\sum_{l=1}^M \prod_{j=1}^q \mu_{A_j^l}(x_j)}. \quad (8.4)$$

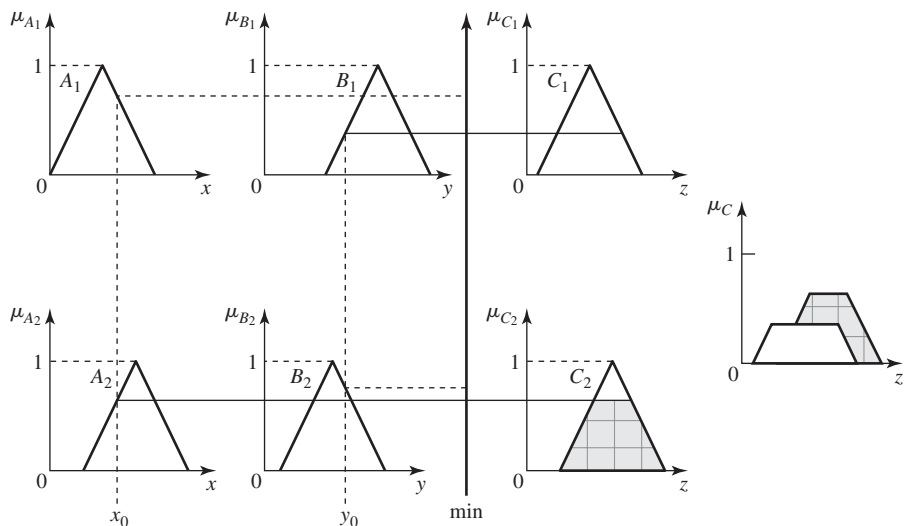


Figure 8.14 An illustration of fuzzy reasoning in a two-input fuzzy system.

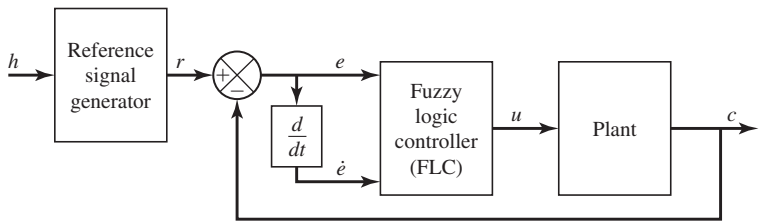


Figure 8.15 A diagram of a closed-loop system with a fuzzy logic controller.

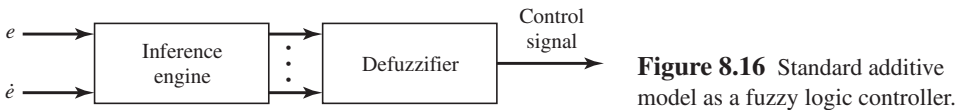


Figure 8.16 Standard additive model as a fuzzy logic controller.

8.4 Fuzzy Logic Control

A block diagram of a closed-loop system with a fuzzy logic controller (FLC) is shown in Figure 8.15. A typical fuzzy logic control architecture is depicted in Figure 8.16. The controller uses the control error e and the rate of change in the control error, \dot{e} or Δe , as its inputs. A typical fuzzy logic control architecture has the structure of a SAM; hence, it is composed of two basic components: a decision-making logic module, also referred to as an inference engine, and a defuzzifier. The inference engine processes the inputs using a knowledge base. The outputs of the inference engine are converted into crisp values, defuzzified, by a defuzzifier. Designing a fuzzy logic controller can be summarized in the following six-step algorithm:

- Identify ranges of the controller inputs.
- Identify ranges of the controller outputs.
- Create fuzzy sets for each input and output.

Table 8.2 The Rules Matrix

Error \ Change-in-error	LN	N	ZE	P	LP
LN	LN	LN	LN	N	ZE
N	LN	LN	N	ZE	P
ZE	N	N	ZE	P	P
P	N	ZE	P	LP	LP
LP	ZE	P	LP	LP	LP

- Translate the interaction of the inputs and outputs into IF–THEN rules, and then form the rules matrix.
- Decide on the inference engine and the defuzzifier, and then apply them to crisp inputs to produce crisp outputs.
- Implement the controller, test it, and modify it if necessary.

As an example, suppose we use five membership functions for the controller’s inputs and its output. We then generate the rules. A possible set of rules is shown in Table 8.2, where LN denotes large negative, N means negative, ZE is zero, P is positive, and LP refers to large positive. The next step involves choosing an inference engine and a defuzzifier. We can use the product inference rule and the modified center average defuzzifier to get

$$u = \frac{\sum_{k=1}^{25} m_k^o p_k / s_k^o}{\sum_{k=1}^{25} p_k / s_k^o},$$

where for $k = 1, 2, \dots, 25$, following the rules shown in Table 8.2, we compute

$$p_k = \mu_{i_k}^e \mu_{j_k}^{ce}, \quad i, j = 1, 2, \dots, 5.$$

The functions μ_i^e , μ_j^{ce} are the membership functions corresponding to the linguistic variables ERROR and CHANGE-IN-ERROR, respectively. The parameters m_k^o and s_k^o are the m and s parameters of the controller output membership functions—see figures below for their definition. The modified center average defuzzifier takes into account the shape of the membership functions, and in particular their “spread.” As pointed out by Wang [295], the sharper the shape of μ , the stronger is our belief that the value should be closer to its m . This observation is reflected in the formula for the defuzzifier we use, where we divide the weights p_k by s_k^o —the spread of the corresponding controller output membership function. We can use triangular membership functions whose equations are given next. The expression for the function shown in Figure 8.17 is

$$\mu(x, m, s) = \max\left(0, 1 - \frac{|x - m|}{s}\right).$$

The expression for the function shown in Figure 8.18 is

$$\mu(x, m, s) = \min\left(1, \max\left(0, 1 + \frac{x - m}{s}\right)\right).$$

Finally, the expression for the function shown in Figure 8.19 is

$$\mu(x, m, s) = \min\left(1, \max\left(0, 1 - \frac{x - m}{s}\right)\right).$$

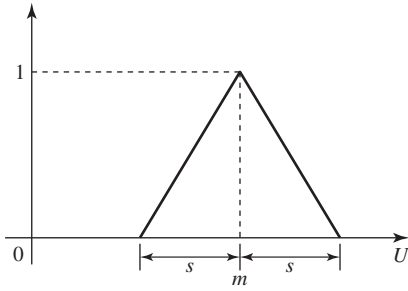


Figure 8.17 A triangular shaped fuzzy set.

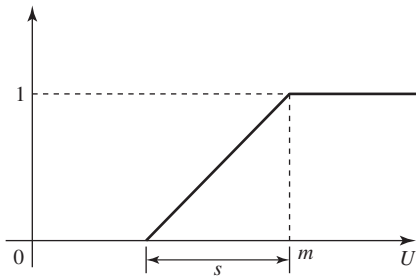


Figure 8.18 A fuzzy set at the right end of a universe of discourse.

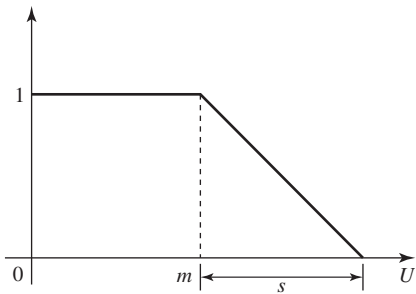


Figure 8.19 A fuzzy set at the left end of a universe of discourse.

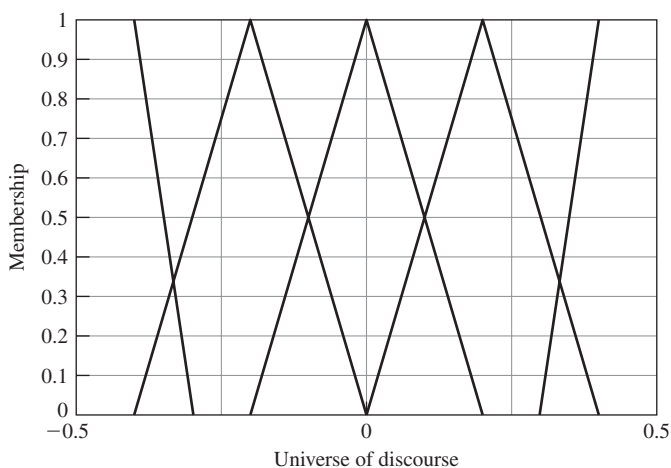
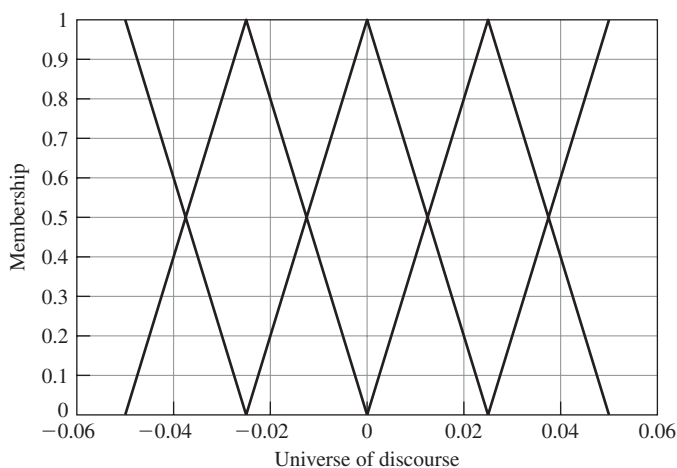
◆ Example 8.8

Consider the vehicle model given in Subsection 2.4.2. Let the input signal be $u = \delta_f$ (the front-wheel steer angle) and $\delta_r = 0$. We will design a simple fuzzy logic controller for the above vehicle model. The range for the error is $[-0.5, 0.5]$. The range for the change-in-error is $[-0.06, 0.06]$. The range for the controller output is $[-3^\circ, 3^\circ]$. The rules matrix is as in Table 8.2. The parameters of the membership functions are given in Table 8.3.

Plots of the error, the change-in-error, and controller's output membership functions are given in Figures 8.20, 8.21, and 8.22, respectively. We created a nonlinear mapping from the controller's input space into its output space and formed a look-up table that implements the obtained nonlinear mapping representing fuzzy logic controller. We used SIMULINK's 2-D Look-Up Table block to implement the

Table 8.3 The m and s Parameters of the Membership Functions

ERROR		CHANGE-IN-ERROR		OUTPUT	
m	s	m	s	m	s
$m_1 = -0.4$	$s_1 = 0.1$	$m_1 = -0.05$	$s_1 = 0.025$	$m_1 = -2.5$	$s_1 = 1$
$m_2 = -0.2$	$s_2 = 0.2$	$m_2 = -0.025$	$s_2 = 0.025$	$m_2 = -1$	$s_2 = 1$
$m_3 = 0$	$s_3 = 0.2$	$m_3 = 0$	$s_3 = 0.025$	$m_3 = 0$	$s_3 = 1$
$m_4 = 0.2$	$s_4 = 0.2$	$m_4 = 0.025$	$s_4 = 0.025$	$m_4 = 1$	$s_4 = 1$
$m_5 = 0.4$	$s_5 = 0.1$	$m_5 = 0.05$	$s_5 = 0.025$	$m_5 = 2.5$	$s_5 = 1$

**Figure 8.20** Fuzzy numbers of the linguistic variable ERROR in Example 8.8.**Figure 8.21** Fuzzy numbers of the linguistic variable CHANGE-IN-ERROR in Example 8.8.

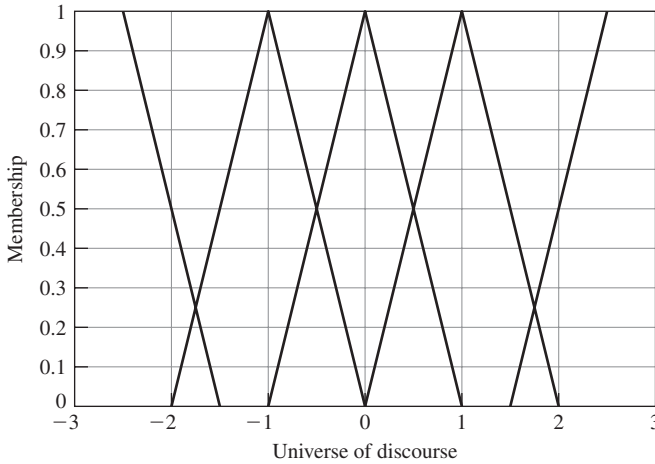


Figure 8.22 Fuzzy numbers of the linguistic variable CONTROL in Example 8.8.

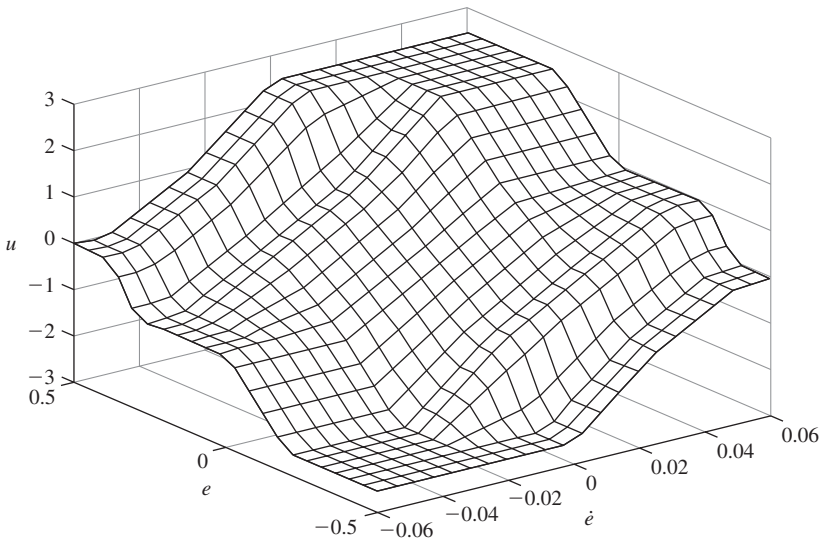


Figure 8.23 A plot of the control surface $u = f(\dot{e}, e)$ in Example 8.8.

controller. The look-up table size was 22×22 . The output of the look-up table was scaled by the factor $-2\pi/180$ before using it as the vehicle's steering input. The negative sign is to take into account the opposite directions of the lateral axes of the vehicle coordinate system and the fixed system. The factor $\pi/180$ converts degrees into radians. We multiplied the output of the look-up table by 2 to increase the control authority. This allows us to use the same controller to track step inputs with large amplitudes. A plot of the control surface is given in Figure 8.23. Next, a reference signal generator was constructed. We selected the following transfer function of the signal generator:

$$\frac{R(s)}{H(s)} = \frac{1}{s^3 + 1.75s^2 + 2.15s + 1}'$$

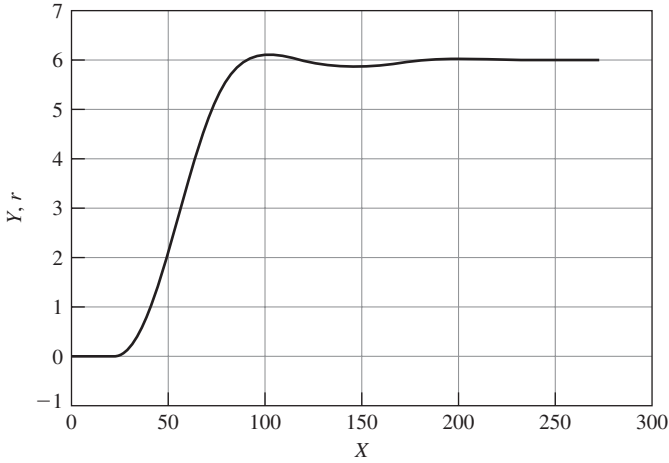


Figure 8.24 Plots of Y and r versus X in Example 8.8.

where $R(s) = \mathcal{L}(r(t))$, $H(s) = \mathcal{L}(6h(t))$, and $h(t)$ is the unit step function. The closed-loop system was simulated using SIMULINK. Plots of Y and r versus X are shown in Figure 8.24.

◆ Example 8.9

This example is based on Example 6.1 found in Driankov, Hellendoorn, and Reinfrank [67, pp. 266–272]. We consider a nonlinear dynamical system modeled by the differential equation

$$\ddot{x} + a_1(1 - x^2)\dot{x} + a_2x = bu, \quad (8.5)$$

where $|u| \leq 15$, $a_1 = 2.2165$, and $a_2 = b = 12.7388$. Let $x_1 = x$ and $x_2 = \dot{x}$. Then, the above system model in state-space format has the form

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ -a_2x_1 - a_1(1 - x_1^2)x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ b \end{bmatrix} u.$$

A phase-plane portrait of the uncontrolled system is shown in Figure 8.25. Plots of x_1 and x_2 versus time for the initial conditions $x_1(0) = -1$ and $x_2(0) = 0$ are shown in Figure 8.26. The relevant membership functions are shown in Figure 8.27, 8.28, 8.29, and 8.30. The following MATLAB function can be used to generate all the relevant membership functions.

```
function y = trap_a(x,a1,a2,c,s)
%TRAP_A is an asymmetric trapezoid type membership
% function,
% also called a 'pi' function,
% where x is the function domain, a1 is the length of the
```

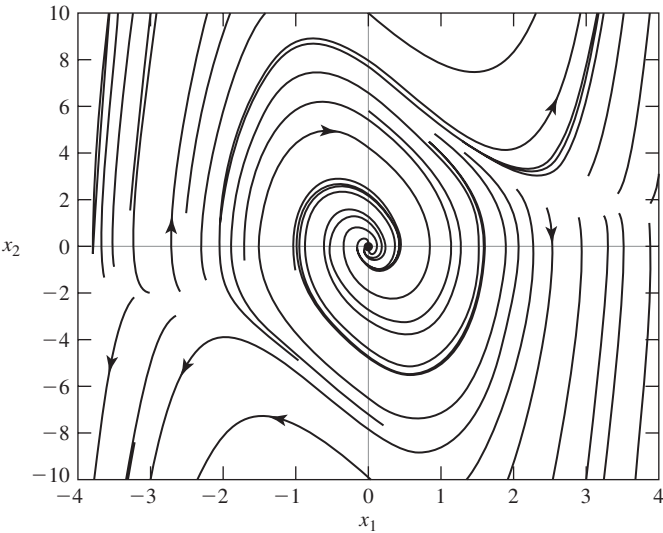


Figure 8.25 A phase-plane portrait of the uncontrolled system in Example 8.9.

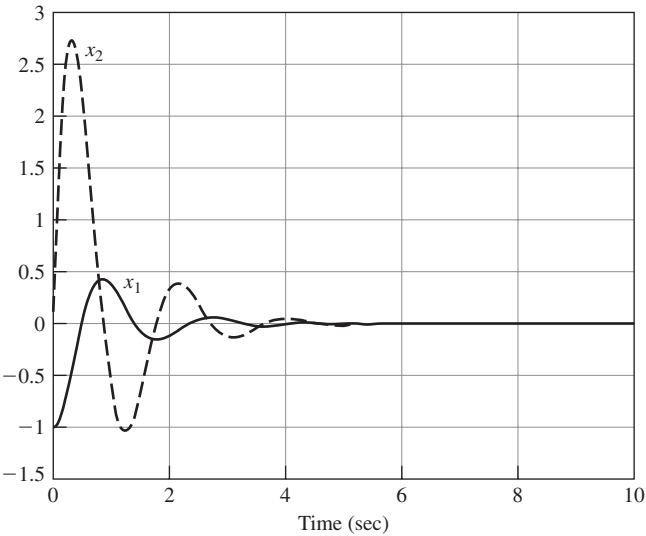


Figure 8.26 Plots of x_1 and x_2 versus time of the uncontrolled system of Example 8.9.

```
% left trapezoid's side
% at the base, a2 is the length of the right side at the
% base,
% c is the trapezoid's center, and s is a half of the
% trapezoid upper side. Thus the length of the trapezoid
% upper
```

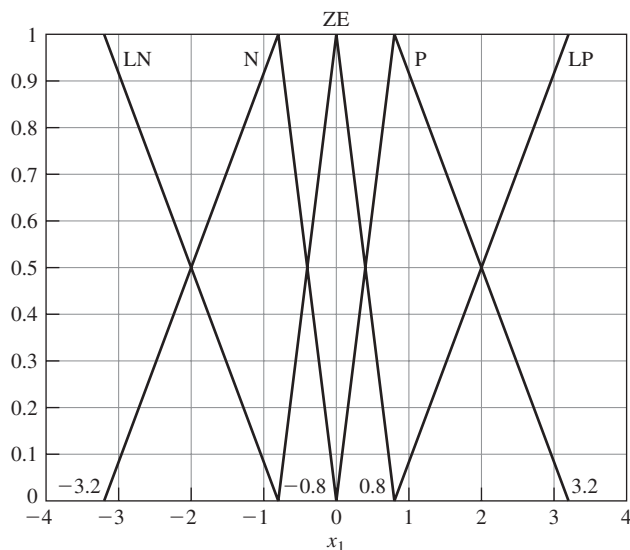


Figure 8.27 Membership functions of x_1 of Example 8.9.

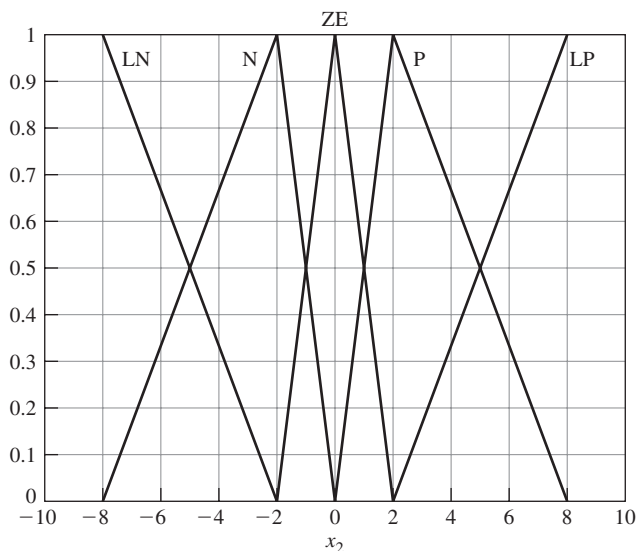


Figure 8.28 Membership functions of x_2 of Example 8.9.

```
% side is 2*s. The length of the trapezoid base is then
% 2*s+a1+a2.
% Set s=0 to get a triangular membership function.
```

```
y1=min(0.5,(max((-0.5),(0.5 +((x-c+s)./a1)))));
y2=min(0.5,(max((-0.5),(0.5 -((x-c-s)./a2)))));
y=y1+y2;
```

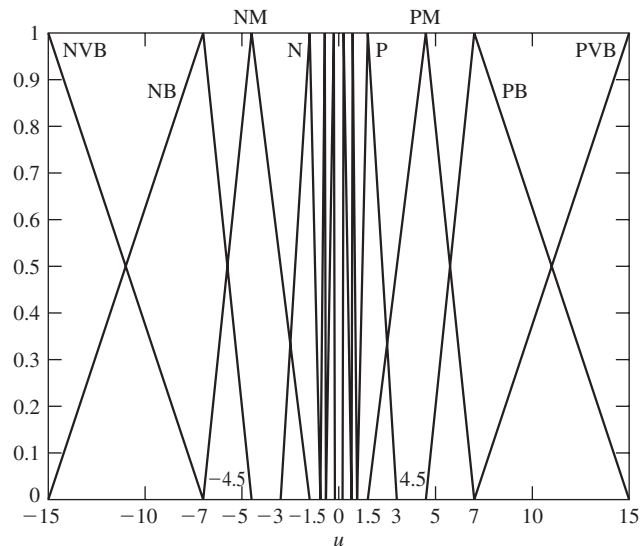


Figure 8.29 Membership functions of u of Example 8.9.

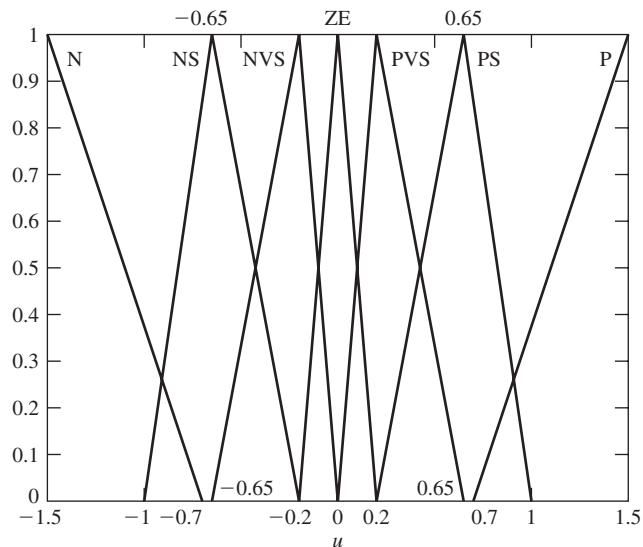
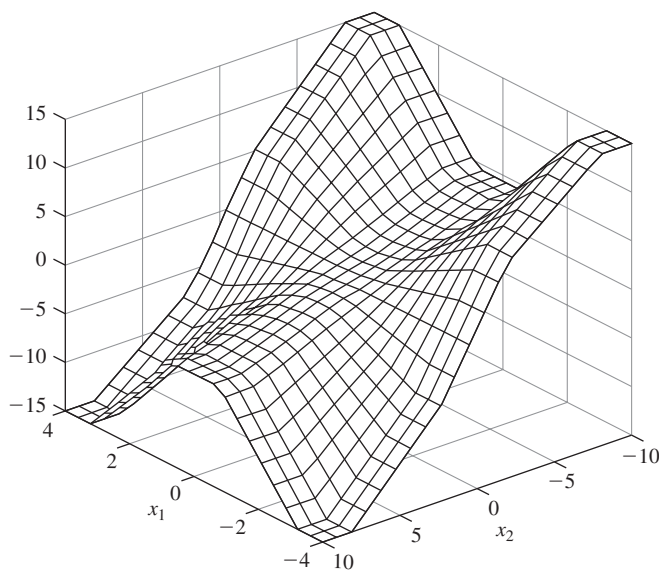
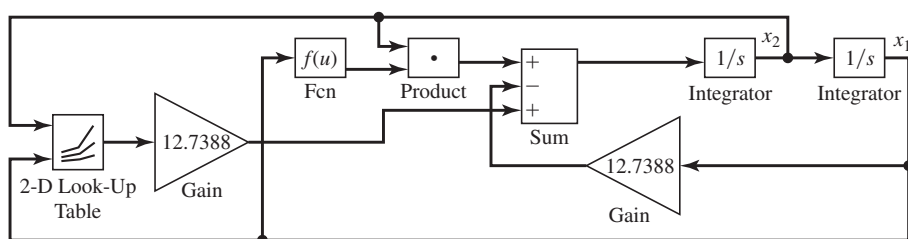


Figure 8.30 Magnification of the membership functions of u on the interval $[-1.5, 1.5]$ of Example 8.9.

We used the rules matrix given in Table 8.4 to produce a nonlinear mapping $u = \Phi(x_1, x_2)$ and to construct a look-up table that implements the obtained mapping representing a fuzzy logic controller. The 2-D Look-Up Table block to implement the controller was used. The look-up table size was 22×22 . A plot

Table 8.4 The Rules Matrix of Example 8.9

$x_1 \backslash x_2$	LN	N	ZE	P	LP
LN	PVB	PM	PM	PM	PVB
N	PB	P	PS	PS	PM
ZE	P	PVS	ZE	NVS	N
P	NM	NS	NS	N	NB
LP	NVB	NM	NM	NM	NVB

**Figure 8.31** A plot of $u = \Phi(x_1, x_2)$ of Example 8.9.**Figure 8.32** A SIMULINK diagram of the closed-loop system of Example 8.9.

of $u = \Phi(x_1, x_2)$ is shown in Figure 8.31. A SIMULINK diagram of the closed-loop system is given in Figure 8.32. A phase portrait of the closed-loop system is depicted in Figure 8.33. Plots of x_1 and x_2 versus time for $x_1(0) = -1$ and $x_2(0) = 0$ for the closed-loop system are shown in Figure 8.34.

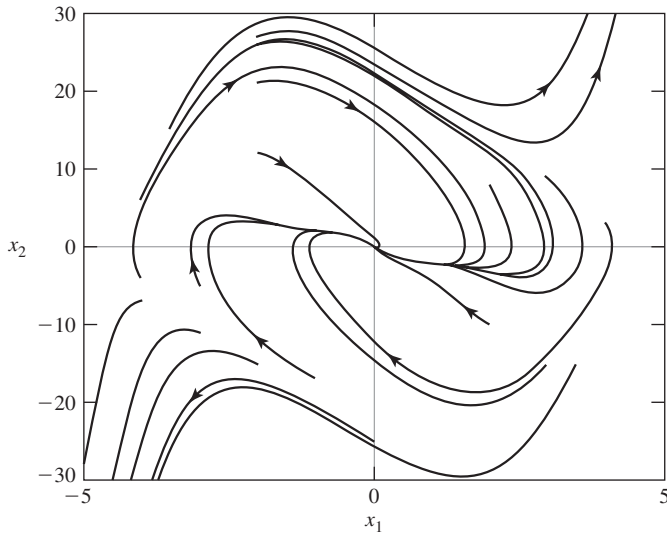


Figure 8.33 A phase-plane portrait of the closed-loop system of Example 8.9.

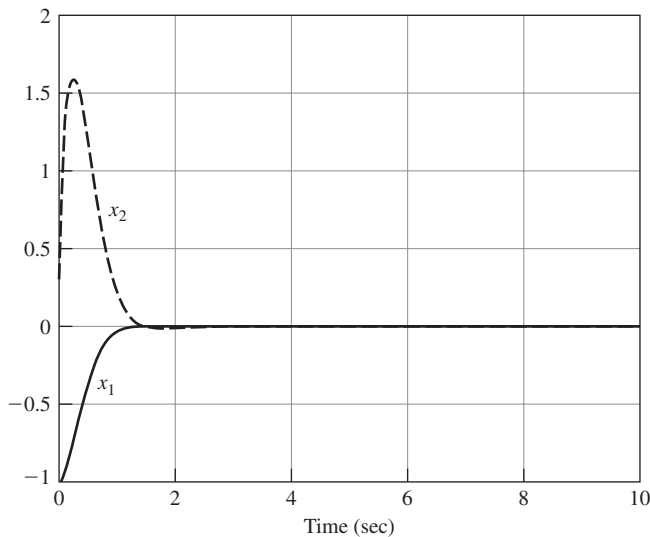


Figure 8.34 Plots of x_1 and x_2 versus time for $x_1(0) = -1$ and $x_2(0) = 0$ for the closed-loop system in Example 8.9.

8.5 Stabilization Using Fuzzy Models

In this section a systematic Lyapunov-based stabilizing control design method for uncertain nonlinear dynamical systems using fuzzy models is presented. The controller is constructed using a design model of the dynamical process to be controlled. Recall that the dynamical process to be controlled is also called the plant. We obtain the design model from the truth model using a fuzzy modeling approach.

8.5.1 Fuzzy Modeling

In many situations there may be human experts that can provide a linguistic description, in terms of IF–THEN rules, of the plant dynamical behavior. Combining available mathematical description of the plant with its linguistic description results in a fuzzy system model. Such an approach to modeling of continuous-time dynamical processes was proposed by Takagi and Sugeno [275] and further developed by Sugeno and Kang [271]. This type of model is called the Takagi–Sugeno (T–S) or Takagi–Sugeno–Kang (TSK) fuzzy model. We mention here that the discrete-time version of the TSK model was analyzed by Tanaka and Sugeno [277], and Tanaka [276], while Wang, Tanaka, and Griffin [293], Cao, Rees, and Feng [42], and Feng et al. [82] used TSK fuzzy models to design stabilizing controllers for nonlinear systems. One can construct a TSK model if local description of the plant is available in terms of local linear models,

$$\dot{\mathbf{x}}(t) = \mathbf{A}_i \mathbf{x}(t) + \mathbf{B}_i \mathbf{u}(t), \quad i = 1, 2, \dots, r,$$

where the state vector $\mathbf{x}(t) \in \mathbb{R}^n$, the control input $\mathbf{u}(t) \in \mathbb{R}^m$, and the matrices \mathbf{A}_i and \mathbf{B}_i are of appropriate dimensions. The above information is then fused with the available IF–THEN rules, where the i th rule can have the following form:

$$\begin{aligned} \textbf{Rule } i: \quad & \text{IF } x_1(t) \text{ is } F_1^i \text{ AND } \dots \text{ AND } x_n(t) \text{ is } F_n^i, \\ & \text{THEN } \dot{\mathbf{x}}(t) = \mathbf{A}_i \mathbf{x}(t) + \mathbf{B}_i \mathbf{u}(t), \end{aligned}$$

where F_j^i , $j = 1, 2, \dots, n$, is the j th fuzzy set of the i th rule. Let $\mu_j^i(x_j)$ be the membership function of the fuzzy set F_j^i , and let

$$w^i = w^i(\mathbf{x}) = \prod_{j=1}^n \mu_j^i(x_j).$$

Then, the resulting fuzzy system model is constructed as the weighted average of the local models and has the form

$$\begin{aligned} \dot{\mathbf{x}} &= \frac{\sum_{i=1}^r w^i (\mathbf{A}_i \mathbf{x} + \mathbf{B}_i \mathbf{u})}{\sum_{i=1}^r w^i} \\ &= \sum_{i=1}^r \alpha_i (\mathbf{A}_i \mathbf{x} + \mathbf{B}_i \mathbf{u}) \\ &= \left(\sum_{i=1}^r \alpha_i \mathbf{A}_i \right) \mathbf{x} + \left(\sum_{i=1}^r \alpha_i \mathbf{B}_i \right) \mathbf{u} \\ &= \mathbf{A}(\boldsymbol{\alpha}) \mathbf{x} + \mathbf{B}(\boldsymbol{\alpha}) \mathbf{u}, \end{aligned} \tag{8.6}$$

where for $i = 1, 2, \dots, r$,

$$\alpha_i = \frac{w^i}{\sum_{i=1}^r w^i}. \tag{8.7}$$

Note that for $i = 1, 2, \dots, r$, we have

$$\alpha_i \geq 0, \quad \sum_{i=1}^r \alpha_i = 1, \quad \text{and} \quad \boldsymbol{\alpha} = [\alpha_1 \quad \alpha_2 \quad \dots \quad \alpha_r]^T \in [0, 1]^r. \tag{8.8}$$

A system model given by (8.6) and (8.8) is also referred to as in the literature as the *polytopic system*—see, for example, Takaba [274].

8.5.2 Constructing a Fuzzy Design Model Using a Nonlinear Model

In the stabilization problem of a nonlinear plant, we are concerned with constructing a controller so that starting from an arbitrary point, in some neighborhood of the operating point, the controller forces the closed-loop system trajectory to converge to this operating point. On the other hand, if the starting point coincides with the operating point, the closed-loop system trajectory is expected to stay at this point for all subsequent time. When designing a controller, we use a TSK fuzzy design model obtained by combining a mathematical description of the plant in terms of its truth model and a linguistic description of the plant. The essential components of the TSK model are linear local models, which describe the plant dynamical behavior at its different operating points. We now show that using a common linearization approach to construct local models may result in a fuzzy design model different from what the designer had in his or her mind. For example, suppose that the truth model of a plant has the form

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u}. \quad (8.9)$$

For simplicity of notation, let

$$\mathbf{F}(\mathbf{x}, \mathbf{u}) = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u}.$$

Then, we can represent (8.9) as

$$\dot{\mathbf{x}} = \mathbf{F}(\mathbf{x}, \mathbf{u}). \quad (8.10)$$

Expanding \mathbf{F} into a Taylor series around $(\mathbf{x}_0, \mathbf{u}_0)$ yields

$$\dot{\mathbf{x}} = \mathbf{F}(\mathbf{x}_0, \mathbf{u}_0) + \left. \frac{\partial \mathbf{F}}{\partial \mathbf{x}} \right|_{\substack{\mathbf{x}=\mathbf{x}_0 \\ \mathbf{u}=\mathbf{u}_0}} (\mathbf{x} - \mathbf{x}_0) + \left. \frac{\partial \mathbf{F}}{\partial \mathbf{u}} \right|_{\substack{\mathbf{x}=\mathbf{x}_0 \\ \mathbf{u}=\mathbf{u}_0}} (\mathbf{u} - \mathbf{u}_0) + \text{higher-order terms}, \quad (8.11)$$

where

$$\mathbf{F}(\mathbf{x}, \mathbf{u}) = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u} = \mathbf{f}(\mathbf{x}) + \sum_{k=1}^m u_k \mathbf{g}_k(\mathbf{x}), \quad (8.12)$$

where \mathbf{g}_k is the k th column of \mathbf{G} . To write expressions for the second and the third terms of (8.11) as functions of \mathbf{f} and \mathbf{G} , let g_{ij} be the (i, j) th element of the matrix \mathbf{G} . Then,

$$\left. \frac{\partial \mathbf{F}}{\partial \mathbf{x}} \right|_{\substack{\mathbf{x}=\mathbf{x}_0 \\ \mathbf{u}=\mathbf{u}_0}} = \left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{x}_0} + \sum_{k=1}^m u_k \left. \frac{\partial \mathbf{g}_k}{\partial \mathbf{x}} \right|_{\substack{\mathbf{x}=\mathbf{x}_0 \\ \mathbf{u}=\mathbf{u}_0}} = \left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{x}_0} + \mathbf{H}(\mathbf{x}_0, \mathbf{u}_0), \quad (8.13)$$

where the (i, j) th element of the $n \times n$ matrix \mathbf{H} has the form

$$\sum_{k=1}^m u_k \left. \frac{\partial g_{ik}(\mathbf{x})}{\partial x_j} \right|_{\substack{\mathbf{x}=\mathbf{x}_0 \\ \mathbf{u}=\mathbf{u}_0}}.$$

Finally,

$$\left. \frac{\partial \mathbf{F}}{\partial \mathbf{u}} \right|_{\substack{\mathbf{x}=\mathbf{x}_0 \\ \mathbf{u}=\mathbf{u}_0}} = \mathbf{G}(\mathbf{x}_0). \quad (8.14)$$

A point $(\mathbf{x}_0^T, \mathbf{u}_0^T)^T \in \mathbb{R}^{n+m}$ is an equilibrium point of (8.10) if $\mathbf{F}(\mathbf{x}_0, \mathbf{u}_0) = \mathbf{0}$ —that is, if at $(\mathbf{x}_0, \mathbf{u}_0)$ we have $\dot{\mathbf{x}} = \mathbf{0}$. Let $\Delta\mathbf{x} = \mathbf{x} - \mathbf{x}_0$ and $\Delta\mathbf{u} = \mathbf{u} - \mathbf{u}_0$. Note that

$$\frac{d\mathbf{x}_0}{dt} = \mathbf{0}.$$

Then, the linearized model about the equilibrium $(\mathbf{x}_0, \mathbf{u}_0)$ is obtained by neglecting higher-order terms and observing that at the equilibrium point $\mathbf{F}(\mathbf{x}_0, \mathbf{u}_0) = \mathbf{0}$. The linearized model has the form

$$\frac{d}{dt}\Delta\mathbf{x} = \mathbf{A}\Delta\mathbf{x} + \mathbf{B}\Delta\mathbf{u}, \quad (8.15)$$

where

$$\mathbf{A} = \left. \frac{\partial \mathbf{F}}{\partial \mathbf{x}} \right|_{\substack{\mathbf{x}=\mathbf{x}_0 \\ \mathbf{u}=\mathbf{u}_0}} \quad \text{and} \quad \mathbf{B} = \left. \frac{\partial \mathbf{F}}{\partial \mathbf{u}} \right|_{\substack{\mathbf{x}=\mathbf{x}_0 \\ \mathbf{u}=\mathbf{u}_0}}.$$

Suppose now that we wish to construct a stabilizing controller for the truth model $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u}$, where $\mathbf{f}(\mathbf{0}) = \mathbf{0}$, so that the pair $(\mathbf{x} = \mathbf{0}, \mathbf{u} = \mathbf{0})$ constitutes an asymptotically stable equilibrium point of the closed-loop system. Our first step is to construct a fuzzy design model. We start with generating linear local models describing plant's behavior at selected operating points. It is natural to first construct a local linear model that describes plant's behavior about the equilibrium state $\mathbf{x} = \mathbf{0}$. We can obtain this first local linear model using the above-described linearization technique. The resulting model is given by

$$\dot{\mathbf{x}} = \mathbf{A}_1\mathbf{x} + \mathbf{B}_1\mathbf{u}.$$

Next, we construct local linear models describing plant's behavior at the remaining operating points. Suppose that $\mathbf{x} = \mathbf{x}_j$ is the next operating state of interest. The result of Taylor's linearization of a nonlinear model about an operating nonequilibrium point is an affine, rather than linear, model. Even if the operating point is an equilibrium point, Taylor series linearization, in general, will not yield a local model that is linear in \mathbf{x} and \mathbf{u} . Indeed, suppose that the operating point $(\mathbf{x}_j, \mathbf{u}_j)$ is an equilibrium point, that is,

$$\mathbf{f}(\mathbf{x}_j) + \mathbf{G}(\mathbf{x}_j)\mathbf{u}_j = \mathbf{0}. \quad (8.16)$$

By (8.11), the resulting linearized model is

$$\frac{d}{dt}(\mathbf{x} - \mathbf{x}_j) = \mathbf{f}(\mathbf{x}_j) + \mathbf{G}(\mathbf{x}_j)\mathbf{u}_j + \mathbf{A}_j(\mathbf{x} - \mathbf{x}_j) + \mathbf{B}_j(\mathbf{u} - \mathbf{u}_j) \quad (8.17)$$

$$= \mathbf{A}_j(\mathbf{x} - \mathbf{x}_j) + \mathbf{B}_j(\mathbf{u} - \mathbf{u}_j). \quad (8.18)$$

We can represent model (8.18) in the form

$$\dot{\mathbf{x}} = \mathbf{A}_j\mathbf{x} + \mathbf{B}_j\mathbf{u} - (\mathbf{A}_j\mathbf{x}_j + \mathbf{B}_j\mathbf{u}_j). \quad (8.19)$$

The term $(\mathbf{A}_j\mathbf{x}_j + \mathbf{B}_j\mathbf{u}_j)$ does not have to be equal to zero, and hence (8.19) is not a linear model in \mathbf{x} and \mathbf{u} , but rather an affine model. We illustrate, on a numerical example, that the term $(\mathbf{A}_j\mathbf{x}_j + \mathbf{B}_j\mathbf{u}_j)$ is indeed, in general, not equal to zero.

◆ Example 8.10

Consider a subsystem of model (1.78) of the form

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ \frac{g \sin(x_1) - mla x_2^2 \sin(2x_1)/2}{4l/3 - mla \cos^2(x_1)} \end{bmatrix} + \begin{bmatrix} 0 \\ -\frac{a \cos(x_1)}{4l/3 - mla \cos^2(x_1)} \end{bmatrix} u, \quad (8.20)$$

where $g = 9.8 \text{ m/sec}^2$, $m = 2 \text{ kg}$, $M = 8 \text{ kg}$, $a = 1/(m + M)$, and $l = 0.5 \text{ m}$. We linearize (8.20) about $x_1 = \pi/4$. For $x_1 = \pi/4$ to be a component of an equilibrium state, we have to have $x_2 = 0$. Hence, the equilibrium state is

$$\mathbf{x}_e = [\pi/4 \quad 0]^T.$$

We next compute $u = u_e$ so that equation (8.16) holds. Performing simple calculations, we get $u_e = 98$. Then, using (8.13) we obtain

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 22.4745 & 0 \end{bmatrix}. \quad (8.21)$$

Next, using (8.14) gives

$$\mathbf{B} = \begin{bmatrix} 0 \\ -0.1147 \end{bmatrix}. \quad (8.22)$$

Note that

$$\mathbf{A}\mathbf{x}_e + \mathbf{B}u_e = \begin{bmatrix} 0 \\ 6.4108 \end{bmatrix} \neq \mathbf{0}.$$

The above example illustrates the fact that the system being linear in $\Delta \mathbf{x} = \mathbf{x} - \mathbf{x}_0$ and $\Delta \mathbf{u} = \mathbf{u} - \mathbf{u}_0$ is, in general, only affine and not linear in \mathbf{x} and \mathbf{u} because of the presence of the offset term $-(\mathbf{A}\mathbf{x}_0 + \mathbf{B}u_0)$. The designer must then be careful when constructing local models using Taylor series linearization because this approach, in general, will not yield a local model that is linear in \mathbf{x} and \mathbf{u} . We note that Taylor's linearization will always yield a linear system in \mathbf{x} and \mathbf{u} if the equilibrium point is $(\mathbf{x}_0^T, u_0^T) = (\mathbf{0}^T, 0^T)$.

We now describe our method of constructing local models using the truth model of the given plant. We assume that the truth model has the form of (8.9). Suppose that we are given an operating state \mathbf{x}_0 that does not have to be an equilibrium state of (8.9). Our goal is to construct a linear model in \mathbf{x} and \mathbf{u} that approximates the behavior of (8.9) in the vicinity of the operating state \mathbf{x}_0 ; that is, we wish to find constant matrices \mathbf{A} and \mathbf{B} such that in a neighborhood of \mathbf{x}_0 ,

$$\mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u} \approx \mathbf{A}\mathbf{x} + \mathbf{B}u \quad (8.23)$$

and

$$\mathbf{f}(\mathbf{x}_0) + \mathbf{G}(\mathbf{x}_0)\mathbf{u} = \mathbf{A}\mathbf{x}_0 + \mathbf{B}u \quad \text{for all } u. \quad (8.24)$$

Because \mathbf{u} is arbitrary, we have to have

$$\mathbf{G}(\mathbf{x}_0) = \mathbf{B}. \quad (8.25)$$

Thus, we are left with finding a constant matrix \mathbf{A} such that in a neighborhood of \mathbf{x}_0 ,

$$\mathbf{f}(\mathbf{x}) \approx \mathbf{A}\mathbf{x} \quad (8.26)$$

and

$$\mathbf{f}(\mathbf{x}_0) = \mathbf{A}\mathbf{x}_0. \quad (8.27)$$

Let \mathbf{a}_i^T denote the i th row of the matrix \mathbf{A} . Then, for the purpose of further analysis, we represent condition (8.26) as

$$f_i(\mathbf{x}) \approx \mathbf{a}_i^T \mathbf{x}, \quad i = 1, 2, \dots, n, \quad (8.28)$$

and we represent (8.27) as

$$f_i(\mathbf{x}_0) = \mathbf{a}_i^T \mathbf{x}_0, \quad i = 1, 2, \dots, n, \quad (8.29)$$

where the i th component of \mathbf{f} is $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$. Expanding the left-hand side of (8.28) about \mathbf{x}_0 and neglecting second- and higher-order terms yields

$$f_i(\mathbf{x}_0) + \nabla^T f_i(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) \approx \mathbf{a}_i^T \mathbf{x}, \quad (8.30)$$

where $\nabla f_i(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is the gradient (a column vector) of f_i at \mathbf{x} . Now, using (8.29), we represent (8.30) as

$$\nabla^T f_i(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) \approx \mathbf{a}_i^T (\mathbf{x} - \mathbf{x}_0), \quad (8.31)$$

where \mathbf{x} is arbitrary but “close” to \mathbf{x}_0 . Our task now is to determine a constant vector \mathbf{a}_i that is as “close as possible” to $\nabla f_i(\mathbf{x}_0)$ and satisfies the constraint $\mathbf{a}_i^T \mathbf{x}_0 = f_i(\mathbf{x}_0)$. Let

$$E = \frac{1}{2} \|\nabla f_i(\mathbf{x}_0) - \mathbf{a}_i\|_2^2.$$

Then, we can formulate our task as a constrained optimization problem of the form

$$\begin{aligned} & \underset{\mathbf{a}_i}{\text{minimize}} \quad E \\ & \text{subject to} \quad \mathbf{a}_i^T \mathbf{x}_0 = f_i(\mathbf{x}_0). \end{aligned} \quad (8.32)$$

We note that (8.32) is a convex constrained optimization problem. This means that the first-order necessary condition for a minimum of E is also sufficient; see, for example, Luenberger [190, Section 6.5] or Chong and Żak [50, Chapter 21]. The first-order conditions for the optimization problem given by (8.32) are

$$\nabla_{\mathbf{a}_i} E + \lambda \nabla_{\mathbf{a}_i} (\mathbf{a}_i^T \mathbf{x}_0 - f_i(\mathbf{x}_0)) = \mathbf{0}, \quad (8.33)$$

$$\mathbf{a}_i^T \mathbf{x}_0 = f_i(\mathbf{x}_0), \quad (8.34)$$

where λ is the Lagrange multiplier and the subscript \mathbf{a}_i in $\nabla_{\mathbf{a}_i}$ indicates that the gradient, ∇ is computed with respect to \mathbf{a}_i . Performing the required differentiation in (8.33) yields

$$\mathbf{a}_i - \nabla f_i(\mathbf{x}_0) + \lambda \mathbf{x}_0 = \mathbf{0}, \quad (8.35)$$

$$\mathbf{a}_i^T \mathbf{x}_0 = f_i(\mathbf{x}_0). \quad (8.36)$$

Recall that we consider now the case when $\mathbf{x}_0 \neq \mathbf{0}$. Premultiplying (8.35) by \mathbf{x}_0^T and substituting (8.36) into the resulting equation yields

$$\lambda = \frac{\mathbf{x}_0^T \nabla f_i(\mathbf{x}_0) - f_i(\mathbf{x}_0)}{\|\mathbf{x}_0\|^2}. \quad (8.37)$$

Substituting λ given by (8.37) into (8.35), we obtain

$$\boxed{\mathbf{a}_i = \nabla f_i(\mathbf{x}_0) + \frac{f_i(\mathbf{x}_0) - \mathbf{x}_0^T \nabla f_i(\mathbf{x}_0)}{\|\mathbf{x}_0\|^2} \mathbf{x}_0, \quad \mathbf{x}_0 \neq \mathbf{0}} \quad (8.38)$$

◆ Example 8.11

We illustrate the above procedure for constructing local linear models on the nonlinear model (8.20) that we analyzed in Example 8.10. This model was also considered by Wang et al. [293, p. 20]. Suppose that, as in Wang et al. [293], we wish to construct two local models; one describing the plant operation about $x_1 = 0$ and the other local model when x_1 is about $\pm\pi/2$. In both models, we assume that $x_2 = 0$ and u is arbitrary. Note that the nonlinear model (8.20) is in the format given by (8.9). We first construct a local model about the operating state $x_1 = x_2 = 0$. The local model for this case can be obtained using Taylor's linearization about the equilibrium point $(\mathbf{x}_0, \mathbf{u}_0) = (\mathbf{0}, \mathbf{0})$.

We next construct the second local model that corresponds to $x_1 = 88^\circ\pi/180^\circ$ and $x_2 = 0$. Let $\beta = \cos(88^\circ)$. Then, the input matrix \mathbf{B} can be obtained from (8.20) using (8.25):

$$\mathbf{B} = \begin{bmatrix} 0 \\ -\frac{a\beta}{4l/3 - mla\beta^2} \end{bmatrix}.$$

We then compute the matrix \mathbf{A} of the second local model using (8.38). We first compute the first row \mathbf{A} . We note that $f_1(\mathbf{x}) = x_2$, and hence

$$\nabla f_1 = [0 \quad 1]^T.$$

The operating state is $\mathbf{x}_0 = [88^\circ\pi/180^\circ \ 0]^T$. We compute \mathbf{a}_1^T , the first row of \mathbf{A} , using (8.38) to get

$$\mathbf{a}_1^T = \nabla^T f_1(\mathbf{x}_0) + [0 \quad 0] = [0 \quad 1].$$

We again use (8.38) to compute \mathbf{a}_2^T , the second row of the matrix \mathbf{A} , of the second local model:

$$\mathbf{a}_2^T = [f_2(\mathbf{x}_0)/(88^\circ\pi/180^\circ) \quad 0] \approx \left[\frac{g}{4l/3 - mla\beta^2} \left(\frac{2}{\pi} \right) \quad 0 \right].$$

The above-obtained local models are the same as the ones that were proposed by Wang et al. [293] using a heuristic approach.

8.5.3 Stabilizability of Fuzzy Models

In this subsection we first present sufficient conditions for asymptotic stability in the large of the unforced continuous-time fuzzy model

$$\dot{\mathbf{x}} = \sum_{i=1}^r \alpha_i \mathbf{A}_i \mathbf{x}, \quad (8.39)$$

where $\alpha_i = \alpha_i(\mathbf{x}(t)) \geq 0$ for $i = 1, 2, \dots, r$, and $\sum_{i=1}^r \alpha_i = 1$. The above model can be viewed as an example of unforced continuous-time polytopic model. We next discuss the stabilizability problem of forced fuzzy models using fuzzy state-feedback control strategies.

It is straightforward to show that a sufficient condition for asymptotic stability in the large of the equilibrium state $\mathbf{x} = \mathbf{0}$ of system (8.39) is that there exists a symmetric positive definite matrix \mathbf{P} such that for $i = 1, 2, \dots, r$,

$$\mathbf{A}_i^T \mathbf{P} + \mathbf{P} \mathbf{A}_i < 0. \quad (8.40)$$

It is obvious that a necessary condition for the existence of a common symmetric positive definite \mathbf{P} satisfying (8.40) is that each \mathbf{A}_i be asymptotically stable—that is, the eigenvalues of each \mathbf{A}_i be in the open left-hand complex plane. We now give another necessary condition for the existence of a common \mathbf{P} that satisfies (8.40).

Theorem 8.1 Suppose that each \mathbf{A}_i , $i = 1, 2, \dots, r$, in the model (8.39) is asymptotically stable and there exists a symmetric positive definite \mathbf{P} such that conditions (8.40) hold. Then, the matrices

$$\sum_{k=1}^s \mathbf{A}_{i_k},$$

where $i_k \in \{1, 2, \dots, r\}$ and $s = 1, 2, 3, \dots, r$ are asymptotically stable.

Proof Suppose that conditions (8.40) hold. Let

$$\mathbf{A}_{i_k}^T \mathbf{P} + \mathbf{P} \mathbf{A}_{i_k} = -\mathbf{Q}_{i_k}, \quad (8.41)$$

where $\mathbf{Q}_{i_k} = \mathbf{Q}_{i_k}^T > 0$ for $i_k \in \{1, 2, \dots, r\}$. Summing both sides of (8.41) yields

$$\left(\sum_{k=1}^s \mathbf{A}_{i_k}^T \right) \mathbf{P} + \mathbf{P} \left(\sum_{k=1}^s \mathbf{A}_{i_k} \right) = - \sum_{k=1}^s \mathbf{Q}_{i_k}. \quad (8.42)$$

By assumption $\mathbf{P} = \mathbf{P}^T > 0$ and because $\mathbf{Q}_{i_k} = \mathbf{Q}_{i_k}^T > 0$ for $i_k \in \{1, 2, \dots, r\}$, the symmetric matrices

$$\sum_{k=1}^s \mathbf{Q}_{i_k}, \quad s = 2, 3, \dots, r,$$

are also positive definite. Hence, by the Lyapunov theorem, the matrices $\sum_{k=1}^s \mathbf{A}_{i_k}$, $s = 2, 3, \dots, r$, are asymptotically stable.

The following corollary can be viewed as a continuous-time counterpart of Theorem 4.3 of Tanaka and Sugeno [277, pp. 146–147].

Corollary 8.1 Suppose that each A_i , $i = 1, 2, \dots, r$ in fuzzy model (8.39) is asymptotically stable and there exists a symmetric positive definite P such that conditions (8.40) hold. Then, the matrices

$$A_i + A_j, \quad i, j = 1, 2, \dots, r,$$

are asymptotically stable.

We mention that Theorem 4.3 in Tanaka and Sugeno [277] can be generalized to correspond to Theorem 8.1 rather than to the above corollary.

The above corollary can be restated in an equivalent way as follows. If there exist i and j such that a matrix $A_i + A_j$ is not asymptotically stable, then there is no positive definite matrix P such that $A_i^T P + P A_i < 0$ for $i = 1, 2, \dots, r$. The following example illustrates this fact.

◆ Example 8.12

Given the unforced fuzzy system model described by (8.39), where $r = 2$ and

$$A_1 = \begin{bmatrix} -1 & 4 \\ 0 & -2 \end{bmatrix}, \quad A_2 = \begin{bmatrix} -1 & 0 \\ 4 & -2 \end{bmatrix}.$$

Obviously each A_i , $i = 1, 2$, is asymptotically stable. However, there is no common positive definite P such that $A_i^T P + P A_i < 0$ for $i = 1, 2$ because the matrix

$$A_1 + A_2 = \begin{bmatrix} -1 & 4 \\ 0 & -2 \end{bmatrix} + \begin{bmatrix} -1 & 0 \\ 4 & -2 \end{bmatrix} = \begin{bmatrix} -2 & 4 \\ 4 & -4 \end{bmatrix}$$

is unstable since its eigenvalues are

$$\{1.1231, -7.1231\}.$$

Theorem 8.1 can be used to check if there does not exist a positive definite matrix P satisfying conditions (8.40) because this theorem gives us a sufficient condition for nonexistence of a common P that satisfies (8.40).

We now present a sufficiency condition for asymptotic stabilizability in the large of the model given by (8.6) and (8.8) using fuzzy state feedback control law:

$$u = - \sum_{j=1}^r \alpha_j K_j x. \quad (8.43)$$

The closed-loop system has the form

$$\dot{x} = \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (A_i - B_i K_j) x. \quad (8.44)$$

To obtain the above model, we used the fact that $\sum_{j=1}^r \alpha_j = 1$. Suppose that the gain matrices, K_i , were chosen so that the matrices

$$A_i - B_i K_i, \quad i = 1, 2, \dots, r,$$

have their eigenvalues in the open left-hand complex plane; that is, they are asymptotically stable. Let

$$\mathbf{G}_{ij} = (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) + (\mathbf{A}_j - \mathbf{B}_j \mathbf{K}_i) \quad \text{for } i < j \leq r.$$

We now present sufficient conditions of Wang et al. [293] for the closed-loop system modeled by (8.44) to be asymptotically stable in the large.

Theorem 8.2 If for some symmetric positive definite \mathbf{P} the conditions

$$(\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_i)^T \mathbf{P} + \mathbf{P}(\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_i) < 0, \quad i = 1, 2, \dots, r, \quad (8.45)$$

and

$$\mathbf{G}_{ij}^T \mathbf{P} + \mathbf{P} \mathbf{G}_{ij} < 0, \quad i < j \leq r, \quad (8.46)$$

are satisfied, then the closed-loop system modeled by (8.44) is asymptotically stable in the large.

Proof Let \mathbf{P} be such that (8.45) and (8.46) are satisfied. We evaluate the time derivative of $\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)$ on the trajectories of (8.44) to get

$$\begin{aligned} \frac{d}{dt}(\mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)) &= 2\mathbf{x}^T \mathbf{P} \dot{\mathbf{x}} \\ &= 2\mathbf{x}^T \mathbf{P} \left(\sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) \mathbf{x} \right) \\ &= \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j \mathbf{x}^T ((\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j)^T \mathbf{P} + \mathbf{P}(\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j)) \mathbf{x} \\ &= \sum_{i=1}^r \alpha_i^2 \mathbf{x}^T ((\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_i)^T \mathbf{P} + \mathbf{P}(\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_i)) \mathbf{x} \\ &\quad + \sum_{i=1}^r \sum_{j>i}^r \alpha_i \alpha_j \mathbf{x}^T (\mathbf{G}_{ij}^T \mathbf{P} + \mathbf{P} \mathbf{G}_{ij}) \mathbf{x}. \end{aligned} \quad (8.47)$$

Hence, when (8.45) and (8.46) are satisfied, the positive definite quadratic function $\mathbf{x}^T \mathbf{P} \mathbf{x}$ becomes a Lyapunov function for the system (8.44), which implies asymptotic stability in the large of (8.44). The proof of the theorem is complete.

We now give a sufficient condition for the stability of system (8.44) for the case when conditions (8.46) are not satisfied. To proceed, we need the following notation. Suppose that there is a positive definite \mathbf{P} such that conditions (8.45) are satisfied. Represent conditions (8.45) as

$$(\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_i)^T \mathbf{P} + \mathbf{P}(\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_i) = -\mathbf{Q}_i, \quad i = 1, 2, \dots, r, \quad (8.48)$$

where each \mathbf{Q}_i is symmetric positive definite. Let λ_i denote the minimum eigenvalue of \mathbf{Q}_i . We have

$$\lambda_i > 0, \quad i = 1, 2, \dots, r, \quad (8.49)$$

because $\mathbf{Q}_i = \mathbf{Q}_i^T > 0$. Let

$$\mathbf{G}_{ij}^T \mathbf{P} + \mathbf{P} \mathbf{G}_{ij} = -\mathbf{Q}_{ij}, \quad i < j \leq r, \quad (8.50)$$

where now we do not require \mathbf{Q}_{ij} to be positive definite. Note that $\mathbf{Q}_{ij} = \mathbf{Q}_{ij}^T$, and hence the eigenvalues of \mathbf{Q}_{ij} are real. Let λ_{ij} denote the minimum eigenvalue of \mathbf{Q}_{ij} . We are ready to state and prove the following theorem.

Theorem 8.3 Suppose that each $(\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_i)$, $i = 1, 2, \dots, r$, in (8.44) is asymptotically stable and there exists a symmetric positive definite \mathbf{P} such that conditions (8.45) hold. Then, the closed-loop fuzzy model (8.44) is asymptotically stable in the large if the matrix

$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & \lambda_{12}/2 & \cdots & \lambda_{1r}/2 \\ \lambda_{12}/2 & \lambda_2 & & \lambda_{2r}/2 \\ \vdots & & \ddots & \vdots \\ \lambda_{1r}/2 & \lambda_{2r}/2 & \cdots & \lambda_r \end{bmatrix} \quad (8.51)$$

is positive definite.

Proof Let

$$W = W(\mathbf{x}(t)) = \mathbf{x}^T(t) \mathbf{P} \mathbf{x}(t)$$

be the Lyapunov function candidate for the closed-loop system modeled by (8.44). We then evaluate the time derivative of W on the trajectories of (8.44) to get

$$\begin{aligned} \dot{W} &= 2 \mathbf{x}^T \mathbf{P} \dot{\mathbf{x}} \\ &= 2 \mathbf{x}^T \mathbf{P} \left(\sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) \mathbf{x} \right) \\ &= \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j \mathbf{x}^T ((\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j)^T \mathbf{P} + \mathbf{P}(\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j)) \mathbf{x} \\ &= - \sum_{i=1}^r \alpha_i^2 \mathbf{x}^T \mathbf{Q}_i \mathbf{x} - \sum_{i=1}^r \sum_{j>i}^r \alpha_i \alpha_j \mathbf{x}^T \mathbf{Q}_{ij} \mathbf{x}. \end{aligned} \quad (8.52)$$

Recall that for any symmetric matrix $\mathbf{M} = \mathbf{M}^T$, we have

$$\lambda_{\min}(\mathbf{M}) \|\mathbf{x}\|^2 \leq \mathbf{x}^T \mathbf{M} \mathbf{x},$$

where $\lambda_{\min}(\mathbf{M})$ is the smallest eigenvalue of \mathbf{M} . Using this fact in (8.52) yields

$$\begin{aligned} \dot{W} &\leq - \left(\sum_{i=1}^r \alpha_i^2 \lambda_i + \sum_{i=1}^r \sum_{j>i}^r \alpha_i \alpha_j \lambda_{ij} \right) \|\mathbf{x}\|^2 \\ &= - \left(\begin{bmatrix} \alpha_1 & \alpha_2 & \cdots & \alpha_r \end{bmatrix} \begin{bmatrix} \lambda_1 & \lambda_{12}/2 & \cdots & \lambda_{1r}/2 \\ \lambda_{12}/2 & \lambda_2 & & \lambda_{2r}/2 \\ \vdots & & \ddots & \vdots \\ \lambda_{1r}/2 & \lambda_{2r}/2 & \cdots & \lambda_r \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_r \end{bmatrix} \right) \|\mathbf{x}\|^2 \\ &= -(\boldsymbol{\alpha}^T \mathbf{\Lambda} \boldsymbol{\alpha}) \|\mathbf{x}\|^2. \end{aligned} \quad (8.53)$$

The Lyapunov derivative is negative definite if the matrix $\mathbf{\Lambda}$ is positive definite, which completes the proof.

8.5.4 A Lyapunov-Based Stabilizer

In this subsection we apply a Lyapunov approach to the fuzzy control design for uncertain dynamical systems of the form

$$\dot{\mathbf{x}} = \mathbf{A}(\boldsymbol{\alpha})\mathbf{x} + \mathbf{B}(\boldsymbol{\alpha})(\mathbf{u} + \mathbf{h}(t, \mathbf{x})), \quad (8.54)$$

where $\boldsymbol{\alpha}$ satisfies (8.8) and the vector-valued function \mathbf{h} represents modeling uncertainties or disturbances. The only information about uncertainties available to the designer are their bounds. Specifically, we assume that the uncertainty modeling function is bounded by a nonnegative function $\eta = \eta(t, \mathbf{x})$; that is,

$$\|\mathbf{h}(t, \mathbf{x})\| \leq \eta(t, \mathbf{x}),$$

where $\|\cdot\|$ denotes the vector p -norm. Observe that \mathbf{h} affects the system dynamics via the input matrix $\mathbf{B}(\boldsymbol{\alpha})$ in a similar fashion as the control input \mathbf{u} does. For this reason, we say that the uncertainty \mathbf{h} satisfies the matching condition.

Some of the Lyapunov-based control strategies, using nonfuzzy controllers, for uncertain systems are documented in Zinober [318]. The method of constructing stabilizing controllers, presented in this section, requires that each of the matrices \mathbf{A}_i , $i = 1, 2, \dots, r$, is already asymptotically stable and that there exists a common symmetric positive definite $\mathbf{P} = \mathbf{P}^T > 0$ such that, for $i = 1, 2, \dots, r$, we have

$$\mathbf{A}_i^T \mathbf{P} + \mathbf{P} \mathbf{A}_i < 0. \quad (8.55)$$

If this is not the case, we break the design of stabilizing control law, \mathbf{u} , into two parts $\mathbf{u} = \mathbf{u}_1 + \mathbf{u}_2$. We first design \mathbf{u}_1 of the form

$$\mathbf{u}_1 = - \sum_{j=1}^r \alpha_j \mathbf{K}_j \mathbf{x} \quad (8.56)$$

so that the system

$$\dot{\mathbf{x}} = \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) \mathbf{x}$$

is asymptotically stable in the large and there exists a common \mathbf{P} such that conditions (8.45) and (8.46) or conditions of Theorem 8.3 are satisfied. Note that the system modeled by (8.54) and driven by (8.56) has the form

$$\begin{aligned} \dot{\mathbf{x}} &= \sum_{i=1}^r \alpha_i \mathbf{A}_i \mathbf{x} - \sum_{i=1}^r \alpha_i \mathbf{B}_i \left(\sum_{j=1}^r \alpha_j \mathbf{K}_j \mathbf{x} \right) + \sum_{i=1}^r \alpha_i \mathbf{B}_i (\mathbf{u}_2 + \mathbf{h}) \\ &= \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) \mathbf{x} + \sum_{i=1}^r \alpha_i \mathbf{B}_i (\mathbf{u}_2 + \mathbf{h}). \end{aligned} \quad (8.57)$$

We next proceed with the design of \mathbf{u}_2 . The role of \mathbf{u}_2 is to reject the uncertainty \mathbf{h} . Then, $\mathbf{u} = \mathbf{u}_1 + \mathbf{u}_2$ will asymptotically stabilize in the large the uncertain system (8.54). We propose a family of fuzzy controllers \mathbf{u}_2 that reject the uncertainty \mathbf{h} . These controllers are obtained by

“fuzzifying” the controllers proposed in Hui and Žak [133]. To proceed, we need the following notation. Let

$$\boldsymbol{\mu} = [\mu_1 \quad \mu_2 \quad \dots \quad \mu_m]^T = \sum_{i=1}^r \alpha_i \mathbf{B}_i^T \mathbf{P} \mathbf{x}. \quad (8.58)$$

Let $\|\cdot\|_p$ denote the p -norm of a vector; that is, if a vector $\mathbf{v} \in \mathbb{R}^m$, then its p -norm is defined as

$$\|\mathbf{v}\|_p = (|v_1|^p + |v_2|^p + \dots + |v_m|^p)^{1/p}, \quad p \geq 1.$$

A well-known result concerning p -norms is the Hölder inequality,

$$|\mathbf{v}^T \mathbf{w}| \leq \|\mathbf{v}\|_p \|\mathbf{w}\|_q, \quad \frac{1}{p} + \frac{1}{q} = 1.$$

All norms on \mathbb{R}^m are equivalent in the sense that if $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$ are norms on \mathbb{R}^m , then there exist positive constants C_1 and C_2 such that

$$C_1 \|\mathbf{v}\|_\alpha \leq \|\mathbf{v}\|_\beta \leq C_2 \|\mathbf{v}\|_\alpha \quad \text{for all } \mathbf{v} \in \mathbb{R}^m. \quad (8.59)$$

Finally, the gradient of $\|\boldsymbol{\mu}\|_p$, where it is defined, has the form

$$\nabla \|\boldsymbol{\mu}\|_p = \left[\frac{\partial \|\boldsymbol{\mu}\|_p}{\partial \mu_1} \quad \frac{\partial \|\boldsymbol{\mu}\|_p}{\partial \mu_2} \quad \dots \quad \frac{\partial \|\boldsymbol{\mu}\|_p}{\partial \mu_m} \right]^T, \quad p \geq 1. \quad (8.60)$$

Applying the chain rule, we compute

$$\nabla \|\boldsymbol{\mu}\|_p = \frac{1}{\|\boldsymbol{\mu}\|_p^{p-1}} \begin{bmatrix} |\mu_1|^{p-1} \text{sign}(\mu_1) \\ |\mu_2|^{p-1} \text{sign}(\mu_2) \\ \vdots \\ |\mu_m|^{p-1} \text{sign}(\mu_m) \end{bmatrix}. \quad (8.61)$$

Let $\|\mathbf{h}(t, \mathbf{x})\|_q \leq \eta_q$, where $\eta_q \geq 0$ is a constant. We then show that when

$$\mathbf{u}_2 = -\eta_q \nabla \|\boldsymbol{\mu}\|_p, \quad p \geq 1, \quad \text{and} \quad \frac{1}{p} + \frac{1}{q} = 1, \quad (8.62)$$

then the closed-loop system driven by $\mathbf{u} = \mathbf{u}_1 + \mathbf{u}_2$ is asymptotically stable in the large for any uncertainty \mathbf{h} such that $\|\mathbf{h}\|_q \leq \eta_q$.

Theorem 8.4 Suppose that \mathbf{u}_1 , given by (8.56), is such that there exists a common \mathbf{P} satisfying conditions (8.45) and (8.46) or conditions of Theorem 8.3. Then, the system (8.54) driven by $\mathbf{u} = \mathbf{u}_1 + \mathbf{u}_2$, where \mathbf{u}_2 is given by (8.62), is asymptotically stable in the large for any uncertainty \mathbf{h} such that $\|\mathbf{h}\|_q \leq \eta_q$.

Proof Suppose that \mathbf{u}_1 , given by (8.56), was constructed such that (8.45) and (8.46) or the conditions of Theorem 8.3 hold. Then, (8.54) driven by \mathbf{u}_1 takes the form (8.57); that is,

$$\dot{\mathbf{x}} = \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) \mathbf{x} + \sum_{i=1}^r \alpha_i \mathbf{B}_i (\mathbf{u}_2 + \mathbf{h}).$$

We now show that

$$W = \mathbf{x}^T \mathbf{P} \mathbf{x}, \quad (8.63)$$

is the Lyapunov function for the closed-loop system, where $\mathbf{P} = \mathbf{P}^T > 0$ satisfies (8.45) and (8.46) or the conditions of Theorem 8.3. For this, we evaluate the Lyapunov derivative of W —that is, we evaluate $\frac{d}{dt} W$ on the trajectories of (8.57)—and show that the Lyapunov derivative of W is negative definite. Before we proceed with evaluating $\frac{d}{dt} W$, observe that using the fact that

$$\text{sign}(\mu_j) = \frac{|\mu_j|}{\mu_j}$$

and using (8.61), we obtain

$$\boldsymbol{\mu}^T \nabla \|\boldsymbol{\mu}\|_p = \|\boldsymbol{\mu}\|_p.$$

Hence,

$$\begin{aligned} \frac{d}{dt} W &= 2\mathbf{x}^T \mathbf{P} \left(\sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) \mathbf{x} + \sum_{i=1}^r \alpha_i \mathbf{B}_i (\mathbf{u}_2 + \mathbf{h}) \right) \\ &< -2\eta_q \boldsymbol{\mu}^T \nabla \|\boldsymbol{\mu}\|_p + 2\boldsymbol{\mu}^T \mathbf{h} \\ &\leq -2\eta_q \|\boldsymbol{\mu}\|_p + 2|\boldsymbol{\mu}^T \mathbf{h}|. \end{aligned} \quad (8.64)$$

Applying Hölder's inequality to the above yields

$$\begin{aligned} \frac{d}{dt} W &< -2\eta_q \|\boldsymbol{\mu}\|_p + 2\|\boldsymbol{\mu}\|_p \|\mathbf{h}\|_q \\ &\leq -2\eta_q \|\boldsymbol{\mu}\|_p + 2\|\boldsymbol{\mu}\|_p \eta_q. \end{aligned} \quad (8.65)$$

Therefore $\frac{d}{dt} W < 0$, and by Lyapunov's theorem the closed loop-system is asymptotically stable in the large.

◆ Example 8.13

Using system model (1.79), we propose the following two rules describing the plant dynamics:

Rule 1: IF $x_1(t)$ is about 0,

THEN $\dot{\mathbf{x}}(t) = \mathbf{A}_1 \mathbf{x}(t) + \mathbf{B}_1(u(t) + h)$

and

Rule 2: IF $x_1(t)$ is about $\pm \pi/4$,
THEN $\dot{\mathbf{x}}(t) = \mathbf{A}_2 \mathbf{x}(t) + \mathbf{B}_2(u(t) + h)$,

where $h = -f_c + m/x_2^2 \sin(x_1)$. Local linear models are constructed using the procedure described in the previous section, and in particular relations (8.25) and (8.38). We have

$$\mathbf{A}_1 = \begin{bmatrix} 0 & 1 & 0 & 0 \\ \frac{g}{4l/3 - aml} & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ -\frac{amg}{4/3 - am} & 0 & 0 & 0 \end{bmatrix}, \quad \mathbf{B}_1 = \begin{bmatrix} 0 \\ a \\ -\frac{4l/3 - aml}{0} \\ \frac{4a/3}{4/3 - ma} \end{bmatrix}.$$

The matrices \mathbf{A}_2 and \mathbf{B}_2 of the local model corresponding to the second rule are:

$$\mathbf{A}_2 = \begin{bmatrix} 0 & 1 & 0 & 0 \\ \frac{g \sin(x_1)}{4l/3 - mla \cos^2(x_1)} \left(\frac{1}{x_1} \right) & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ \frac{-mag \sin(2x_1)/2}{4/3 - ma \cos^2(x_1)} \left(\frac{1}{x_1} \right) & 0 & 0 & 0 \end{bmatrix}_{|x_1 = \pm \pi/4}$$

and

$$\mathbf{B}_2 = \begin{bmatrix} 0 \\ \frac{-a \cos(x_1)}{4l/3 - mla \cos^2(x_1)} \\ 0 \\ \frac{4a/3}{4/3 - ma \cos^2(x_1)} \end{bmatrix}_{|x_1 = \pm \pi/4}.$$

The parameter values are the same as in Example 8.10. We use weighting functions of the form:

$$w_1(x_1) = \frac{1 - 1/(1 + \exp(-14(x_1 - \pi/8)))}{1 + \exp(-14(x_1 + \pi/8))} \quad \text{and} \quad w_2(x_1) = 1 - w_1(x_1).$$

The corresponding fuzzy system model is

$$\dot{\mathbf{x}} = (\alpha_1 \mathbf{A}_1 + \alpha_2 \mathbf{A}_2) \mathbf{x} + (\alpha_1 \mathbf{B}_1 + \alpha_2 \mathbf{B}_2) (u - \mu_c \text{sign}(x_4) + m/x_2^2 \sin(x_1)), \quad (8.66)$$

where $\alpha_i = w_i$, $i = 1, 2$. Note that $w_1 + w_2 = 1$ for all t . The friction coefficient $\mu_c = 0.005$. We then proceed to construct a stabilizing controller. We first construct

$$u_1 = -(\alpha_1 \mathbf{K}_1 + \alpha_2 \mathbf{K}_2) \mathbf{x}$$

such that the system

$$\dot{\mathbf{x}} = \sum_{i=1}^2 \sum_{j=1}^2 \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_j \mathbf{K}_j) \mathbf{x} \quad (8.67)$$

is asymptotically stable in the large. To test our stabilization algorithm, we randomly chose the set of desired poles for $(\mathbf{A}_1 - \mathbf{B}_1 \mathbf{K}_1)$ to be in the interval $[-4, -1]$. The set of desired poles is

$$\{-1.0970, -2.1263, -2.5553, -3.9090\}.$$

The desired poles for $(\mathbf{A}_2 - \mathbf{B}_2 \mathbf{K}_2)$ were also chosen in random fashion to be in the interval $[-4, -1]$. In this case, we generated the set of poles of the form

$$\{-1.5794, -1.6857, -2.7908, -3.6895\}.$$

For this choice of the closed-loop poles of the local models, we obtain the following feedback gain matrices:

$$\mathbf{K}_1 = [-294.8766 \quad -73.1208 \quad -13.4726 \quad -27.3362]$$

and

$$\mathbf{K}_2 = [-440.3915 \quad -118.5144 \quad -19.1611 \quad -35.5575].$$

Let

$$\mathbf{G}_{12} = \mathbf{A}_1 - \mathbf{B}_1 \mathbf{K}_2 + \mathbf{A}_2 - \mathbf{B}_2 \mathbf{K}_1$$

and

$$\mathbf{P} = \begin{bmatrix} 54.9580 & 15.6219 & 6.9389 & 12.1165 \\ 15.6219 & 4.5429 & 2.1011 & 3.5488 \\ 6.9389 & 2.1011 & 1.3972 & 1.7978 \\ 12.1165 & 3.5488 & 1.7978 & 2.9375 \end{bmatrix}.$$

Then, it is easy to check that sufficient conditions, given by (8.45) and (8.46), for asymptotic stability in the large of system (8.67) are satisfied. This completes the design of u_1 .

We now construct u_2 . We first use (8.58) to construct

$$\mu = (\alpha_1 \mathbf{B}_1^T \mathbf{P} + \alpha_2 \mathbf{B}_2^T \mathbf{P}) \mathbf{x}.$$

In this example, the number of inputs is $m = 1$, and hence (8.62) reduces to

$$u_2 = -\eta_q \text{sign}(\mu),$$

where we can take $\eta_q = |h|$. We combine u_1 and u_2 to obtain the following stabilizing controller:

$$u = u_1 + u_2 = -(\alpha_1 \mathbf{K}_1 + \alpha_2 \mathbf{K}_2) \mathbf{x} - \eta_q \text{sign}(\mu). \quad (8.68)$$

The above controller globally stabilizes the design model. In our simulations we connected the controller to the truth model given by (1.79). Typical responses are shown in Figures 8.35 and 8.36. The controller locally stabilizes the truth model. We observed that when $|x_1(0)| > 1.3$ rads and $x_2(0) = x_3(0) = x_4(0) = 0$, then the equilibrium state $\mathbf{x} = \mathbf{0}$ of the closed-loop system consisting of the truth model driven by (8.68) becomes unstable.

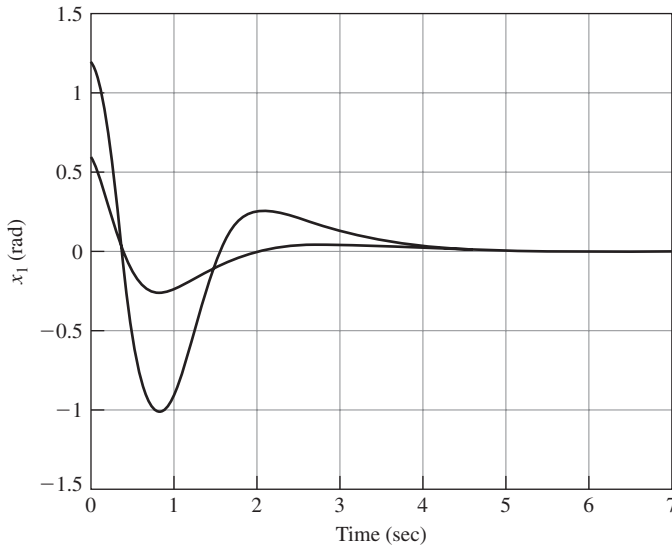


Figure 8.35 Plots of x_1 versus time for two different initial values of x_1 : $x_1(0) = 1.2$ and $x_1(0) = 0.6$. In both cases, $x_2(0) = x_3(0) = x_4(0) = 0$.

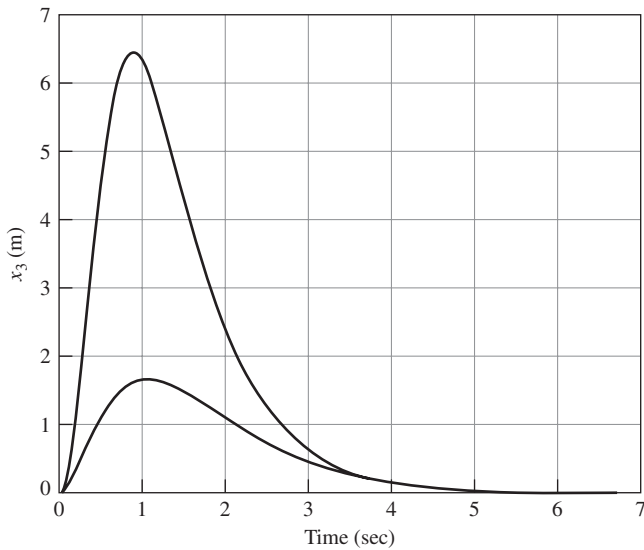


Figure 8.36 Plots of x_3 versus time for two different initial values of x_1 : $x_1(0) = 1.2$ and $x_1(0) = 0.6$. In both cases, $x_2(0) = x_3(0) = x_4(0) = 0$.

We summarize this section as follows. A common approach to the control system design is to use a design model to construct a controller. The design model is a simplified version of the truth model and usually neglects high-frequency dynamics of the plant. Thus, the designer has to cope with an incomplete mathematical model of the process to be controlled. To account for

modeling uncertainties, we used a two-part controller. The first part of the controller stabilized the nominal plant—that is, the plant model that does not include modeling uncertainties. The role of the second part of the controller was to reject modeling uncertainties. In both cases, fuzzy state-feedback control strategies were used. In order to be able to apply the proposed controllers in real-life systems, we have to account for limited authority of actuators, unmatched uncertainties, and inaccessibility of state variables.

8.6 Stability of Discrete Fuzzy Models

We analyze the model of the form

$$\mathbf{x}(k+1) = \sum_{i=1}^r w_i(\mathbf{x}(k)) \mathbf{A}_i \mathbf{x}(k), \quad (8.69)$$

where

$$w_i(\mathbf{x}(k)) \geq 0 \quad \text{and} \quad \sum_{i=1}^r w_i(\mathbf{x}(k)) = 1.$$

To proceed, we need the following notation. Given two $n \times n$ matrices \mathbf{A} and \mathbf{B} . Then, $\mathbf{A} < \mathbf{B}$ means that

$$\mathbf{x}^T \mathbf{A} \mathbf{x} < \mathbf{x}^T \mathbf{B} \mathbf{x} \quad \text{for all } \mathbf{x} \in \mathbb{R}^n, \mathbf{x} \neq \mathbf{0};$$

that is, the quadratic form $\mathbf{x}^T (\mathbf{A} - \mathbf{B}) \mathbf{x}$ is negative definite.

We now prove the following matrix inequality lemma that is useful in the stability analysis of discrete fuzzy models.

Lemma 8.1 Let $\mathbf{A}, \mathbf{B}, \mathbf{P} \in \mathbb{R}^{n \times n}$. If $\mathbf{P} = \mathbf{P}^T > 0$ such that

$$\mathbf{A}^T \mathbf{P} \mathbf{A} - \mathbf{P} < 0 \quad \text{and} \quad \mathbf{B}^T \mathbf{P} \mathbf{B} - \mathbf{P} < 0,$$

then

$$\mathbf{A}^T \mathbf{P} \mathbf{B} + \mathbf{B}^T \mathbf{P} \mathbf{A} - 2\mathbf{P} < 0.$$

Proof We first represent $\mathbf{A}^T \mathbf{P} \mathbf{B} + \mathbf{B}^T \mathbf{P} \mathbf{A} - 2\mathbf{P}$ as

$$\begin{aligned} \mathbf{A}^T \mathbf{P} \mathbf{B} + \mathbf{B}^T \mathbf{P} \mathbf{A} - 2\mathbf{P} &= -(\mathbf{A} - \mathbf{B})^T \mathbf{P} (\mathbf{A} - \mathbf{B}) + \mathbf{A}^T \mathbf{P} \mathbf{A} + \mathbf{B}^T \mathbf{P} \mathbf{B} - 2\mathbf{P} \\ &= -(\mathbf{A} - \mathbf{B})^T \mathbf{P} (\mathbf{A} - \mathbf{B}) + (\mathbf{A}^T \mathbf{P} \mathbf{A} - \mathbf{P}) \\ &\quad + (\mathbf{B}^T \mathbf{P} \mathbf{B} - \mathbf{P}). \end{aligned}$$

The first term, $-(\mathbf{A} - \mathbf{B})^T \mathbf{P} (\mathbf{A} - \mathbf{B})$, is nonpositive definite because $\mathbf{P} = \mathbf{P}^T > 0$ by assumption, while the second and the third terms are negative definite by assumption. Hence,

$$\mathbf{A}^T \mathbf{P} \mathbf{B} + \mathbf{B}^T \mathbf{P} \mathbf{A} - 2\mathbf{P} < 0,$$

and the lemma is proven.

Using Lemma 8.1, we now prove a sufficient condition for global asymptotic stability of model (8.69).

Theorem 8.5 The equilibrium state, $\mathbf{x}_e = \mathbf{0}$, of model (8.69) is globally asymptotically stable if there exists a positive definite matrix $\mathbf{P} = \mathbf{P}^T > 0$ such that for $i = 1, 2, \dots, r$,

$$\mathbf{A}_i^T \mathbf{P} \mathbf{A}_i - \mathbf{P} < 0. \quad (8.70)$$

Proof Suppose that there exists a symmetric positive definite matrix $\mathbf{P} = \mathbf{P}^T > 0$ such that for $i = 1, 2, \dots, r$, inequalities given by (8.70) are satisfied. Consider then the Lyapunov function candidate

$$V(\mathbf{x}(k)) = \mathbf{x}^T(k) \mathbf{P} \mathbf{x}(k). \quad (8.71)$$

To prove the theorem, it is enough to show that

$$\Delta V(\mathbf{x}(k)) < 0.$$

We have

$$\begin{aligned} \Delta V(\mathbf{x}(k)) &= V(\mathbf{x}(k+1)) - V(\mathbf{x}(k)) \\ &= \mathbf{x}^T(k+1) \mathbf{P} \mathbf{x}(k+1) - \mathbf{x}^T(k) \mathbf{P} \mathbf{x}(k) \\ &= \left(\sum_{i=1}^r w_i \mathbf{x}^T \mathbf{A}_i^T \right) \mathbf{P} \left(\sum_{i=1}^r w_i \mathbf{A}_i \mathbf{x} \right) - \mathbf{x}^T \mathbf{P} \mathbf{x} \\ &= \mathbf{x}^T \left(\left(\sum_{i=1}^r w_i \mathbf{A}_i^T \right) \mathbf{P} \left(\sum_{i=1}^r w_i \mathbf{A}_i \right) - \mathbf{P} \right) \mathbf{x}. \end{aligned} \quad (8.72)$$

To proceed, we represent (8.72) in an equivalent form as

$$\Delta V(\mathbf{x}) = \mathbf{x}^T \left(\left(\sum_{i=1}^r w_i \mathbf{A}_i^T \right) \mathbf{P} \left(\sum_{j=1}^r w_j \mathbf{A}_j \right) - \mathbf{P} \right) \mathbf{x}. \quad (8.73)$$

By assumption we have

$$\sum_{i=1}^r w_i = 1,$$

and hence

$$\sum_{i=1}^r w_i \sum_{j=1}^r w_j = \sum_{i=1}^r \sum_{j=1}^r w_i w_j = 1. \quad (8.74)$$

Substituting (8.74) into (8.73) yields

$$\begin{aligned} \Delta V(\mathbf{x}) &= \mathbf{x}^T \left(\left(\sum_{i=1}^r w_i \mathbf{A}_i^T \right) \mathbf{P} \left(\sum_{j=1}^r w_j \mathbf{A}_j \right) - \sum_{i=1}^r \sum_{j=1}^r w_i w_j \mathbf{P} \right) \mathbf{x} \\ &= \mathbf{x}^T \left(\sum_{i=1}^r \sum_{j=1}^r w_i w_j (\mathbf{A}_i^T \mathbf{P} \mathbf{A}_j - \mathbf{P}) \right) \mathbf{x}. \end{aligned} \quad (8.75)$$

We rearrange (8.75) to get

$$\begin{aligned}
 \Delta V(\mathbf{x}(k)) &= \sum_{i=1}^r w_i^2 \mathbf{x}^T (\mathbf{A}_i^T \mathbf{P} \mathbf{A}_i - \mathbf{P}) \mathbf{x} + \sum_{i=1}^r \sum_{\substack{j=1 \\ j \neq i}}^r w_i w_j \mathbf{x}^T (\mathbf{A}_i^T \mathbf{P} \mathbf{A}_j - \mathbf{P}) \mathbf{x} \\
 &= \sum_{i=1}^r w_i^2 \mathbf{x}^T (\mathbf{A}_i^T \mathbf{P} \mathbf{A}_i - \mathbf{P}) \mathbf{x} \\
 &\quad + \sum_{i=1}^r \sum_{j>i}^r w_i w_j \mathbf{x}^T (\mathbf{A}_i^T \mathbf{P} \mathbf{A}_j + \mathbf{A}_j^T \mathbf{P} \mathbf{A}_i - 2\mathbf{P}) \mathbf{x}. \tag{8.76}
 \end{aligned}$$

The first term in (8.76) is negative definite because of the assumption and the fact that at least one of the weights, w_i , is positive because $w_i \geq 0$ and $\sum_{i=1}^r w_i = 1$. The second term in (8.76) is nonpositive definite by virtue of Lemma 8.1. Hence, $V = \mathbf{x}^T \mathbf{P} \mathbf{x}$ is a Lyapunov function for the fuzzy model (8.69), and by the Lyapunov theorem the equilibrium state $\mathbf{x}_e = \mathbf{0}$ of the model is globally asymptotically stable.

The above theorem gives no method for finding a common positive definite matrix \mathbf{P} that satisfies r inequalities (8.70). All of the matrices \mathbf{A}_i , $i = 1, 2, \dots, r$, are asymptotically stable if there exists a common positive definite \mathbf{P} that satisfies (8.70). However, the asymptotic stability of all the matrices \mathbf{A}_i is only necessary but not sufficient for the existence of a common symmetric positive definite \mathbf{P} that satisfies (8.70). The following theorem provides a necessary, but not sufficient, condition for the existence of positive definite \mathbf{P} satisfying (8.70).

Theorem 8.6 Given r square $n \times n$ matrices \mathbf{A}_i that are asymptotically stable and nonsingular. If there exists a symmetric positive definite matrix $\mathbf{P} = \mathbf{P}^T > 0$ such that for $i = 1, 2, \dots, r$,

$$\mathbf{A}_i^T \mathbf{P} \mathbf{A}_i - \mathbf{P} < 0, \tag{8.77}$$

then for $i, j = 1, 2, \dots, r$, the matrices $\mathbf{A}_i \mathbf{A}_j$ are asymptotically stable.

Proof By assumption, there exists a symmetric positive definite matrix \mathbf{P} such that for $j = 1, 2, \dots, r$, we have

$$\mathbf{A}_j^T \mathbf{P} \mathbf{A}_j - \mathbf{P} < 0. \tag{8.78}$$

Premultiplying the above by $(\mathbf{A}_j^T)^{-1}$ and postmultiplying by \mathbf{A}_j^{-1} yields

$$\mathbf{P} - (\mathbf{A}_j^T)^{-1} \mathbf{P} \mathbf{A}_j^{-1} < 0. \tag{8.79}$$

Hence,

$$\mathbf{P} < (\mathbf{A}_j^T)^{-1} \mathbf{P} \mathbf{A}_j^{-1}. \tag{8.80}$$

On the other hand, using (8.77) we get

$$\mathbf{A}_i^T \mathbf{P} \mathbf{A}_i < \mathbf{P}. \tag{8.81}$$

Combining (8.80) and (8.81), we obtain

$$\mathbf{A}_i^T \mathbf{P} \mathbf{A}_i < (\mathbf{A}_j^T)^{-1} \mathbf{P} \mathbf{A}_j^{-1}. \tag{8.82}$$

Premultiplying both sides of (8.82) by \mathbf{A}_j^T and postmultiplying by \mathbf{A}_j , we get, for $i, j = 1, 2, \dots, r$,

$$(\mathbf{A}_i \mathbf{A}_j)^T \mathbf{P} (\mathbf{A}_i \mathbf{A}_j) - \mathbf{P} < 0. \quad (8.83)$$

Therefore, for $i, j = 1, 2, \dots, r$, the matrices $\mathbf{A}_i \mathbf{A}_j$ are asymptotically stable, which completes the proof.

The above theorem can be stated in an equivalent way as follows. Assume that for $i = 1, 2, \dots, r$, the matrices \mathbf{A}_i are asymptotically stable and nonsingular. If for some $i, j = 1, 2, \dots, r$, the matrix $\mathbf{A}_i \mathbf{A}_j$ is unstable, then there does not exist a common positive definite \mathbf{P} such that

$$\mathbf{A}_i^T \mathbf{P} \mathbf{A}_i - \mathbf{P} < 0 \quad \text{for } i = 1, 2, \dots, r.$$

◆ Example 8.14

Let

$$\mathbf{A}_1 = \begin{bmatrix} 1 & -0.5 \\ 1 & 0 \end{bmatrix} \quad \text{and} \quad \mathbf{A}_2 = \begin{bmatrix} -1 & -0.5 \\ 1 & 0 \end{bmatrix}.$$

Both matrices, \mathbf{A}_1 and \mathbf{A}_2 , are asymptotically stable. Indeed, the eigenvalues of \mathbf{A}_1 are located at $\{0.5 + j0.5, 0.5 - j0.5\}$, while the eigenvalues of \mathbf{A}_2 are located at $\{-0.5 + j0.5, -0.5 - j0.5\}$. However, the matrix $\mathbf{A}_1 \mathbf{A}_2$ has its eigenvalues located at $\{-1.866, -0.134\}$, and hence is unstable. Thus, it follows from Theorem 8.6, that there is no positive definite \mathbf{P} such that

$$\mathbf{A}_i^T \mathbf{P} \mathbf{A}_i - \mathbf{P} < 0 \quad \text{for } i = 1, 2. \quad (8.84)$$

8.7 Fuzzy Estimator

In this section we are concerned with the problem of designing a state estimator using a TSK model of the plant. We first discuss the issue of constructing a TSK model using input–output data. To construct a TSK model, we need local description of the plant in terms of linear models,

$$\begin{aligned} \dot{\mathbf{x}}(t) &= \mathbf{A}_i \mathbf{x}(t) + \mathbf{B}_i \mathbf{u}(t), \\ \mathbf{y}(t) &= \mathbf{C}_i \mathbf{x}(t), \end{aligned}$$

where $i = 1, 2, \dots, r$, the state vector $\mathbf{x}(t) \in \mathbb{R}^n$, the control input $\mathbf{u}(t) \in \mathbb{R}^m$, the output $\mathbf{y}(t) \in \mathbb{R}^p$, and the matrices \mathbf{A}_i , \mathbf{B}_i , and \mathbf{C}_i are of appropriate dimensions. The above information is then fused with the available IF–THEN rules, where the i th rule can have the following form:

$$\begin{aligned} \textbf{Rule } i: \quad & \text{IF } z_1(t) \text{ is } F_1^i \text{ AND } \dots \text{ AND } z_q(t) \text{ is } F_q^i, \\ & \text{THEN } \begin{cases} \dot{\mathbf{x}}(t) = \mathbf{A}_i \mathbf{x}(t) + \mathbf{B}_i \mathbf{u}(t), \\ \mathbf{y} = \mathbf{C}_i \mathbf{x}(t), \end{cases} \end{aligned}$$

where $F_j^i, j = 1, 2, \dots, q$, is the j th fuzzy set of the i th rule, and $z_1(t), \dots, z_q(t)$ are measurable system variables. Let $\mu_j^i(z_j)$ be the membership function of the fuzzy set F_j^i , and let

$$w^i = w^i(z) = \prod_{j=1}^q \mu_j^i(z_j), \quad i = 1, 2, \dots, r,$$

where r is the number of fuzzy rules. The resulting fuzzy system model is the weighted average of the local models and has the form of (8.6), that is,

$$\begin{aligned} \dot{\mathbf{x}} &= \frac{\sum_{i=1}^r w^i (\mathbf{A}_i \mathbf{x} + \mathbf{B}_i \mathbf{u})}{\sum_{i=1}^r w^i} \\ &= \sum_{i=1}^r \alpha_i (\mathbf{A}_i \mathbf{x} + \mathbf{B}_i \mathbf{u}) \\ &= \left(\sum_{i=1}^r \alpha_i \mathbf{A}_i \right) \mathbf{x} + \left(\sum_{i=1}^r \alpha_i \mathbf{B}_i \right) \mathbf{u} \\ &= \mathbf{A}(\boldsymbol{\alpha}) \mathbf{x} + \mathbf{B}(\boldsymbol{\alpha}) \mathbf{u}, \end{aligned}$$

where for $i = 1, 2, \dots, r$ we have

$$\alpha_i = \frac{w^i}{\sum_{i=1}^r w^i}.$$

Similarly,

$$\begin{aligned} \mathbf{y} &= \frac{\sum_{i=1}^r w^i \mathbf{C}_i \mathbf{x}}{\sum_{i=1}^r w^i} \\ &= \left(\sum_{i=1}^r \alpha_i \mathbf{C}_i \right) \mathbf{x} \\ &= \mathbf{C}(\boldsymbol{\alpha}) \mathbf{x}. \end{aligned} \tag{8.85}$$

The coefficients α_i satisfy (8.8); that is,

$$\alpha_i \geq 0, \quad \sum_{i=1}^r \alpha_i = 1 \quad \text{and} \quad \boldsymbol{\alpha} = [\alpha_1 \quad \alpha_2 \quad \dots \quad \alpha_r]^T \in [0, 1]^r.$$

A fuzzy state estimator for the system modeled by (8.6) and (8.85) can have the form

$$\dot{\tilde{\mathbf{x}}} = \mathbf{A}(\boldsymbol{\alpha}) \tilde{\mathbf{x}} + \mathbf{B}(\boldsymbol{\alpha}) \mathbf{u} + \mathbf{L}(\boldsymbol{\alpha}) (\mathbf{y} - \tilde{\mathbf{y}}), \tag{8.86}$$

where $\tilde{\mathbf{x}}$ is the estimate of the plant state \mathbf{x} and

$$\mathbf{L}(\boldsymbol{\alpha}) = \sum_{i=1}^r \alpha_i \mathbf{L}_i.$$

The state estimator (8.86) can be represented equivalently as

$$\boxed{\dot{\tilde{\mathbf{x}}} = \sum_{i=1}^r \alpha_i \mathbf{A}_i \tilde{\mathbf{x}} + \sum_{i=1}^r \alpha_i \mathbf{B}_i \mathbf{u} + \sum_{i=1}^r \alpha_i \mathbf{L}_i (\mathbf{y} - \tilde{\mathbf{y}})} \tag{8.87}$$

We note that the state estimator (8.87) can be viewed as a result of implementing the following r fuzzy rules:

$$\begin{aligned} \textbf{Rule i:} \quad & \text{IF } z_1(t) \text{ is } F_1^i \text{ AND } \cdots \text{ AND } z_q(t) \text{ is } F_q^i, \\ & \text{THEN } \begin{cases} \dot{\tilde{\mathbf{x}}} = \mathbf{A}_i \tilde{\mathbf{x}} + \mathbf{B}_i \mathbf{u} + \mathbf{L}_i (\mathbf{y} - \tilde{\mathbf{y}}), \\ \tilde{\mathbf{y}} = \mathbf{C}_i \tilde{\mathbf{x}}. \end{cases} \end{aligned}$$

We use the state estimate in the controller implementation to get

$$\mathbf{u} = - \sum_{j=1}^r \alpha_j \mathbf{K}_j \tilde{\mathbf{x}} + \mathbf{v}. \quad (8.88)$$

Substituting the above into the plant model yields

$$\dot{\mathbf{x}} = \sum_{i=1}^r \alpha_i \mathbf{A}_i \mathbf{x} - \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j \mathbf{B}_i \mathbf{K}_j \tilde{\mathbf{x}} + \sum_{i=1}^r \alpha_i \mathbf{B}_i \mathbf{v}. \quad (8.89)$$

Substituting (8.85) and (8.88) into the estimator equation (8.87), we obtain

$$\begin{aligned} \dot{\tilde{\mathbf{x}}} &= \sum_{i=1}^r \alpha_i \mathbf{A}_i \tilde{\mathbf{x}} - \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j \mathbf{B}_i \mathbf{K}_j \tilde{\mathbf{x}} + \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j \mathbf{L}_i \mathbf{C}_j (\mathbf{x} - \tilde{\mathbf{x}}) + \sum_{i=1}^r \alpha_i \mathbf{B}_i \mathbf{v} \\ &= \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) \tilde{\mathbf{x}} + \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j \mathbf{L}_i \mathbf{C}_j (\mathbf{x} - \tilde{\mathbf{x}}) + \sum_{i=1}^r \alpha_i \mathbf{B}_i \mathbf{v}. \end{aligned} \quad (8.90)$$

Equations (8.89) and (8.90) constitute the closed-loop system driven by the combined fuzzy controller–estimator compensator. We represent the above two equations modeling the closed-loop system in matrix form as

$$\begin{aligned} \begin{bmatrix} \dot{\mathbf{x}} \\ \dot{\tilde{\mathbf{x}}} \end{bmatrix} &= \begin{bmatrix} \sum_{i=1}^r \alpha_i \mathbf{A}_i & - \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j \mathbf{B}_i \mathbf{K}_j \\ \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j \mathbf{L}_i \mathbf{C}_j & \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j - \mathbf{L}_i \mathbf{C}_j) \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \tilde{\mathbf{x}} \end{bmatrix} \\ &+ \begin{bmatrix} \sum_{i=1}^r \alpha_i \mathbf{B}_i \\ \sum_{i=1}^r \alpha_i \mathbf{B}_i \end{bmatrix} \mathbf{v} \end{aligned} \quad (8.91)$$

$$\mathbf{y} = \begin{bmatrix} \sum_{i=1}^r \alpha_i \mathbf{C}_i & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \tilde{\mathbf{x}} \end{bmatrix}. \quad (8.92)$$

8.7.1 The Comparison Method for Linear Systems

We will now present the comparison method for linear systems. This method will then be used to investigate stability properties of the closed-loop system, modeled by (8.91) and (8.92), consisting of a fuzzy model of a given plant and a fuzzy controller–estimator compensator.

Let \mathbf{A} be a constant $n \times n$ matrix and let $e^{\mathbf{A}t}$ denote the matrix exponential. We will now give a necessary and sufficient condition for the matrix exponential to have all its elements nonnegative for $t \geq 0$. This result appears in Bellman [24, p. 171].

Lemma 8.2 A necessary and sufficient condition that all elements of e^{At} be nonnegative for $t \geq 0$ is that

$$a_{ij} \geq 0, \quad i \neq j. \quad (8.93)$$

Proof We follow the arguments of Bellman [24, p. 172]; see also Beckenbach and Bellman [23]. We first prove the necessary condition. Because

$$e^{At} = I_n + At + \frac{1}{2}A^2t^2 + \cdots$$

it is clear that (8.93) is necessary for the result to hold for small t , which establishes the necessity condition.

To prove the sufficiency part, let c_1 be a scalar so that all the elements of $A + c_1 I_n$ are nonnegative. Then, obviously, all the elements of $e^{(A+c_1 I_n)t}$ are nonnegative. Also, the elements of $e^{-c_1 I_n t}$ are nonnegative because exponentials are always nonnegative. The matrices $A + c_1 I_n$ and $-c_1 I_n$ commute. Hence,

$$\begin{aligned} e^{At} &= e^{(A+c_1 I_n)t - c_1 I_n t} \\ &= e^{(A+c_1 I_n)t} e^{-c_1 I_n t}, \end{aligned}$$

which establishes the desired nonnegativity condition and the sufficiency condition.

Using the above lemma, we can prove the following lemma due to Bellman [25]. We use the following notation: $\mathbf{x} \leq \mathbf{y}$ implies component-by-component majorization; that is, $x_i \leq y_i$ for $i = 1, 2, \dots, n$, where x_i and y_i are the i th components of \mathbf{x} and \mathbf{y} , respectively.

Lemma 8.3 If $a_{ij} \geq 0$, $i \neq j$, then for $t \geq 0$,

$$\frac{d\mathbf{x}}{dt} \leq A\mathbf{x}, \quad \mathbf{x}(0) = \mathbf{c}, \quad (8.94)$$

implies $\mathbf{x} \leq \mathbf{y}$, where

$$\frac{d\mathbf{y}}{dt} = A\mathbf{y}, \quad \mathbf{y}(0) = \mathbf{c}. \quad (8.95)$$

Proof We can represent (8.94) as

$$\frac{d\mathbf{x}(t)}{dt} = A\mathbf{x}(t) - \mathbf{f}(t), \quad \mathbf{x}(0) = \mathbf{c}, \quad (8.96)$$

where $\mathbf{f}(t)$ is a vector-valued function whose components are nonnegative for all $t \geq 0$. The solution of (8.96) can be represented as

$$\begin{aligned} \mathbf{x}(t) &= e^{At} \mathbf{c} - \int_0^t e^{A(t-s)} \mathbf{f}(s) ds \\ &= \mathbf{y}(t) - \int_0^t e^{A(t-s)} \mathbf{f}(s) ds. \end{aligned} \quad (8.97)$$

In the above, s is the integration variable and $t \geq 0$. Hence, $t \geq s \geq 0$. Using Lemma 8.2, we conclude that the components of the integrand of the second term in (8.97) are all nonnegative for $t \geq s \geq 0$. Therefore,

$$\int_0^t e^{A(t-s)} \mathbf{f}(s) ds \geq \mathbf{0} \quad \text{for } t \geq 0,$$

which implies

$$\mathbf{x} \leq \mathbf{y}.$$

The proof is complete.

Equation (8.95) is referred to as the *equation of comparison* for (8.94).

8.7.2 Stability Analysis of the Closed-Loop System

We now perform a stability analysis of the closed-loop system given by (8.91) and (8.92). We first apply a state variable transformation to bring the system into a form convenient for further analysis. The transformation has the form

$$\begin{bmatrix} \mathbf{x} \\ \tilde{\mathbf{x}} - \mathbf{x} \end{bmatrix} = \begin{bmatrix} \mathbf{I}_n & \mathbf{O} \\ -\mathbf{I}_n & \mathbf{I}_n \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \tilde{\mathbf{x}} \end{bmatrix}. \quad (8.98)$$

The closed-loop system model in the new coordinates is

$$\begin{aligned} \begin{bmatrix} \dot{\mathbf{x}} \\ \dot{\tilde{\mathbf{x}}} - \dot{\mathbf{x}} \end{bmatrix} &= \begin{bmatrix} \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) & -\sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j \mathbf{B}_i \mathbf{K}_j \\ \mathbf{O} & \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{L}_i \mathbf{C}_j) \end{bmatrix} \\ &\quad \times \begin{bmatrix} \mathbf{x} \\ \tilde{\mathbf{x}} - \mathbf{x} \end{bmatrix} + \begin{bmatrix} \sum_{i=1}^r \alpha_i \mathbf{B}_i \\ \mathbf{O} \end{bmatrix} \mathbf{v}, \end{aligned} \quad (8.99)$$

$$\mathbf{y} = \begin{bmatrix} \sum_{i=1}^r \alpha_i \mathbf{C}_i & \mathbf{O} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \tilde{\mathbf{x}} - \mathbf{x} \end{bmatrix}. \quad (8.100)$$

Let $\mathbf{e} = \tilde{\mathbf{x}} - \mathbf{x}$. We now present a sufficient condition for the unforced closed-loop system to be globally uniformly asymptotically stable; that is, we will give a condition for uniform asymptotic stability in the large of the system

$$\begin{bmatrix} \dot{\mathbf{x}} \\ \dot{\mathbf{e}} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) & -\sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j \mathbf{B}_i \mathbf{K}_j \\ \mathbf{O} & \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{L}_i \mathbf{C}_j) \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{e} \end{bmatrix}. \quad (8.101)$$

This condition was proposed by Ma, Sun, and He [191].

Theorem 8.7 If there exist two scalar functions $V(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ and $\tilde{V}(\mathbf{e}) : \mathbb{R}^n \rightarrow \mathbb{R}$ and positive real numbers $\gamma_i, \tilde{\gamma}_i, i = 1, \dots, 4$, such that

$$\gamma_1 \|\mathbf{x}\|^2 \leq V(\mathbf{x}) \leq \gamma_2 \|\mathbf{x}\|^2, \quad \tilde{\gamma}_1 \|\mathbf{e}\|^2 \leq \tilde{V}(\mathbf{e}) \leq \tilde{\gamma}_2 \|\mathbf{e}\|^2, \quad (8.102)$$

$$\begin{aligned} \frac{\partial V(\mathbf{x})}{\partial \mathbf{x}} \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) \mathbf{x} &\leq -\gamma_3 \|\mathbf{x}\|^2, \\ \frac{\partial \tilde{V}(\mathbf{e})}{\partial \mathbf{e}} \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{L}_i \mathbf{C}_j) \mathbf{e} &\leq -\tilde{\gamma}_3 \|\mathbf{e}\|^2, \end{aligned} \quad (8.103)$$

and

$$\left\| \frac{\partial V(\mathbf{x})}{\partial \mathbf{x}} \right\| \leq \gamma_4 \|\mathbf{x}\|, \quad \left\| \frac{\partial \tilde{V}(\mathbf{e})}{\partial \mathbf{e}} \right\| \leq \tilde{\gamma}_4 \|\mathbf{e}\|, \quad (8.104)$$

then the system (8.101) is globally uniformly asymptotically stable.

Proof We evaluate the time derivative of V on the trajectories of the system (8.101),

$$\begin{aligned} \dot{V}(\mathbf{x}) &= \frac{\partial V(\mathbf{x})}{\partial \mathbf{x}} \dot{\mathbf{x}} \\ &= \frac{\partial V(\mathbf{x})}{\partial \mathbf{x}} \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{B}_i \mathbf{K}_j) \mathbf{x} - \frac{\partial V(\mathbf{x})}{\partial \mathbf{x}} \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j \mathbf{B}_i \mathbf{K}_j \mathbf{e} \\ &\leq -\gamma_3 \|\mathbf{x}\|^2 + \left\| \frac{\partial V(\mathbf{x})}{\partial \mathbf{x}} \right\| \sum_{i=1}^r \sum_{j=1}^r \|\mathbf{B}_i \mathbf{K}_j\| \|\mathbf{e}\| \\ &\leq -\frac{\gamma_3}{2} \|\mathbf{x}\|^2 - \frac{\gamma_3}{2} \|\mathbf{x}\|^2 + \gamma_4 \sum_{i=1}^r \sum_{j=1}^r \|\mathbf{B}_i \mathbf{K}_j\| \|\mathbf{e}\| \|\mathbf{x}\|. \end{aligned} \quad (8.105)$$

To proceed, we observe that

$$\frac{\gamma_3}{2} \|\mathbf{x}\|^2 - \gamma_4 \sum_{i=1}^r \sum_{j=1}^r \|\mathbf{B}_i \mathbf{K}_j\| \|\mathbf{e}\| \|\mathbf{x}\| + \frac{\gamma_4^2}{2\gamma_3} \left(\sum_{i=1}^r \sum_{j=1}^r \|\mathbf{B}_i \mathbf{K}_j\| \right)^2 \|\mathbf{e}\|^2 \geq 0. \quad (8.106)$$

Indeed, we can represent the left-hand side of the above as

$$\begin{bmatrix} \|\mathbf{x}\| & \|\mathbf{e}\| \end{bmatrix} \begin{bmatrix} \frac{\gamma_3}{2} & -\frac{\gamma_4}{2} \sum_{i=1}^r \sum_{j=1}^r \|\mathbf{B}_i \mathbf{K}_j\| \\ -\frac{\gamma_4}{2} \sum_{i=1}^r \sum_{j=1}^r \|\mathbf{B}_i \mathbf{K}_j\| & \frac{\gamma_4^2}{2\gamma_3} \left(\sum_{i=1}^r \sum_{j=1}^r \|\mathbf{B}_i \mathbf{K}_j\| \right)^2 \end{bmatrix} \begin{bmatrix} \|\mathbf{x}\| \\ \|\mathbf{e}\| \end{bmatrix}. \quad (8.107)$$

It is easy to check that the above quadratic form is positive semidefinite—all three principle minors are nonnegative. We use this fact and (8.102) to further evaluate (8.105) as

$$\begin{aligned}\dot{V}(\mathbf{x}) &\leq -\frac{\gamma_3}{2}\|\mathbf{x}\|^2 + \frac{\gamma_4^2}{2\gamma_3} \left(\sum_{i=1}^r \sum_{j=1}^r \|\mathbf{B}_i \mathbf{K}_j\| \right)^2 \|\mathbf{e}\|^2 \\ &\leq -\frac{\gamma_3}{2\gamma_2} V(\mathbf{x}) + \frac{\gamma_4^2}{2\gamma_3 \tilde{\gamma}_1} \left(\sum_{i=1}^r \sum_{j=1}^r \|\mathbf{B}_i \mathbf{K}_j\| \right)^2 \tilde{V}(\mathbf{e}).\end{aligned}\quad (8.108)$$

Let $a = \left(\sum_{i=1}^r \sum_{j=1}^r \|\mathbf{B}_i \mathbf{K}_j\| \right)^2$. Then, (8.108) can be written as

$$\dot{V}(\mathbf{x}) \leq -\frac{\gamma_3}{2\gamma_2} V(\mathbf{x}) + \frac{a\gamma_4^2}{2\gamma_3 \tilde{\gamma}_1} \tilde{V}(\mathbf{e}).\quad (8.109)$$

We now evaluate the time derivative of $\tilde{V}(\mathbf{e})$ on the trajectories of (8.101),

$$\begin{aligned}\dot{\tilde{V}}(\mathbf{e}) &= \frac{\partial \tilde{V}(\mathbf{e})}{\partial \mathbf{e}} \dot{\mathbf{e}} \\ &\leq \frac{\partial \tilde{V}(\mathbf{e})}{\partial \mathbf{e}} \sum_{i=1}^r \sum_{j=1}^r \alpha_i \alpha_j (\mathbf{A}_i - \mathbf{L}_i \mathbf{C}_j) \mathbf{e} \\ &\leq -\tilde{\gamma}_3 \|\mathbf{e}\|^2 \\ &\leq -\frac{\tilde{\gamma}_3}{\tilde{\gamma}_2} \tilde{V}(\mathbf{e}).\end{aligned}\quad (8.110)$$

We combine (8.109) and (8.110) into a single matrix differential inequality of the form

$$\begin{aligned}\begin{bmatrix} \dot{V}(\mathbf{x}) \\ \dot{\tilde{V}}(\mathbf{e}) \end{bmatrix} &\leq \begin{bmatrix} -\frac{\gamma_3}{2\gamma_2} & \frac{a\gamma_4^2}{2\gamma_3 \tilde{\gamma}_1} \\ 0 & -\frac{\tilde{\gamma}_3}{\tilde{\gamma}_2} \end{bmatrix} \begin{bmatrix} V(\mathbf{x}) \\ \tilde{V}(\mathbf{e}) \end{bmatrix} \\ &= \mathbf{A} \begin{bmatrix} V(\mathbf{x}) \\ \tilde{V}(\mathbf{e}) \end{bmatrix}.\end{aligned}\quad (8.111)$$

We then refer to Lemma 8.3 to conclude that the system (8.101) is globally uniformly asymptotically stable. Indeed, because the equation of comparison for (8.111) is globally uniformly asymptotically stable, we have to have

$$V(\mathbf{x}) \rightarrow 0 \quad \text{and} \quad \tilde{V}(\mathbf{e}) \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

Because V and \tilde{V} are both positive definite decrescent and radially unbounded functions, the above can hold provided that $\mathbf{x} \rightarrow \mathbf{0}$ and $\mathbf{e} \rightarrow \mathbf{0}$ uniformly with respect to t and irrespective of the initial conditions $\mathbf{x}(0)$ and $\mathbf{e}(0)$. Therefore, the closed-loop system is globally uniformly asymptotically stable.

8.8 Adaptive Fuzzy Control

In this section we analyze adaptive fuzzy robust tracking controllers for a class of uncertain dynamical systems. The controllers' construction and their analysis involves sliding modes.

8.8.1 Plant Model and Control Objective

We consider a class of plants modeled by the differential equation

$$\dot{x}^{(n)} = f(\mathbf{x}) + g(\mathbf{x})u + \eta(t, \mathbf{x}), \quad (8.112)$$

where $\mathbf{x} = [x \ \dot{x} \ \cdots \ x^{(n-1)}]^T$, and the functions $f = f(\mathbf{x})$ and $g = g(\mathbf{x})$ are unknown to us. The function $\eta = \eta(t, \mathbf{x})$ models the plant's disturbances. We assume that

$$|\eta(t, \mathbf{x})| \leq d, \quad (8.113)$$

where $d > 0$ is known to us. We further assume that the function $g(\mathbf{x})$ has a positive lower bound; that is, there exists a positive constant \underline{g} such that $g(\mathbf{x}) \geq \underline{g} > 0$. Let

$$\mathbf{x}_d = \mathbf{x}_d(t) = [x_d(t) \ \dot{x}_d(t) \ \cdots \ x_d^{(n-1)}(t)]^T$$

denote the desired state trajectory and let

$$\begin{aligned} \mathbf{e} &= \mathbf{x} - \mathbf{x}_d \\ &= [x - x_d \ \dot{x} - \dot{x}_d \ \cdots \ x^{(n-1)} - x_d^{(n-1)}]^T \\ &= [e \ \dot{e} \ \cdots \ e^{(n-1)}]^T. \end{aligned}$$

We wish to construct a controller u such that

$$\lim_{t \rightarrow \infty} \mathbf{e}(t) = \mathbf{0}.$$

We analyze two different fuzzy adaptive tracking controllers and show that the plant (8.112) driven by these controllers is stable and the adaptation parameters are bounded.

In the next subsection we present background results used in the adaptive fuzzy controllers' construction.

8.8.2 Background Results

Fuzzy Logic System Recall, from Section 8.3, that a fuzzy logic controller with q inputs, v_1, v_2, \dots, v_q , and the center average defuzzifier can be represented as

$$u = \frac{\sum_{l=1}^m \prod_{j=1}^q \mu_{F_j^l}(v_j) \theta_l}{\sum_{l=1}^m \prod_{j=1}^q \mu_{F_j^l}(v_j)}, \quad (8.114)$$

where m is the number of fuzzy rules used to construct the controller, F_j^l is the j th fuzzy set corresponding to the l th fuzzy rule, and θ_l can be taken as a centroid of the l th fuzzy set corresponding to the controller's output, u . To proceed, we represent (8.114) as

$$u = \sum_{l=1}^m \left(\frac{\prod_{j=1}^q \mu_{F_j^l}(v_j)}{\sum_{l=1}^m \prod_{j=1}^q \mu_{F_j^l}(v_j)} \right) \theta_l. \quad (8.115)$$

Let

$$\boldsymbol{\theta} = [\theta_1 \quad \theta_2 \quad \cdots \quad \theta_m]^T$$

and

$$\boldsymbol{\xi} = [\xi_1 \quad \xi_2 \quad \cdots \quad \xi_m]^T = \left[\frac{\prod_{j=1}^q \mu_{F_j^1}(v_j)}{\sum_{l=1}^m \prod_{j=1}^q \mu_{F_j^l}(v_j)} \quad \cdots \quad \frac{\prod_{j=1}^q \mu_{F_j^m}(v_j)}{\sum_{l=1}^m \prod_{j=1}^q \mu_{F_j^l}(v_j)} \right]^T.$$

Then, we can represent (8.115) as

$$u = \boldsymbol{\theta}^T \boldsymbol{\xi}, \quad (8.116)$$

where the vector $\boldsymbol{\xi}$ is referred to as the *regressor* vector in the control literature.

Projection Operator Consider a vector-valued function of time,

$$\boldsymbol{\theta}(t) = [\theta_1(t) \quad \cdots \quad \theta_m(t)]^T \in \mathbb{R}^m,$$

where $\dot{\theta}_i(t) = \alpha_i(t)$, $i = 1, \dots, m$. We wish the components $\theta_i(t)$ to lie between $\underline{\theta}_i$ and $\bar{\theta}_i$, that is, $\underline{\theta}_i \leq \theta_i(t) \leq \bar{\theta}_i$. One way to achieve this goal is by setting $\dot{\theta}_i(t) = 0$ when $\theta_i(t)$ reaches either of the bounds and tends to go beyond the bound, that is,

$$\dot{\theta}_i(t) = \begin{cases} 0 & \text{if } \theta_i = \underline{\theta}_i \text{ and } \alpha_i(t) < 0, \\ 0 & \text{if } \theta_i = \bar{\theta}_i \text{ and } \alpha_i(t) > 0, \\ \alpha_i(t) & \text{otherwise.} \end{cases} \quad (8.117)$$

For compactness, we use the notation *Proj* for the right-hand side of (8.117). Then, we write

$$\dot{\theta}_i(t) = \text{Proj}_{\theta_i}(\alpha_i), \quad i = 1, \dots, m, \quad (8.118)$$

or in vector notation,

$$\dot{\boldsymbol{\theta}}(t) = \text{Proj}_{\boldsymbol{\theta}}(\boldsymbol{\alpha}). \quad (8.119)$$

Sliding Modes Let

$$\mathbf{s} = [s_1 \quad s_2 \quad \cdots \quad s_{n-1} \quad 1] \in \mathbb{R}^{1 \times n}$$

and $\sigma(\mathbf{e}) = \mathbf{s}\mathbf{e}$. Then, $\{\mathbf{e} : \sigma(\mathbf{e}) = 0\}$ represents a sliding surface in the tracking error, \mathbf{e} , space. The coefficients s_i of the sliding surface are chosen so that the system (8.112) restricted to the sliding surface is asymptotically stable—see Section 6.7 for a sliding surface design algorithm. Our objective is to construct fuzzy tracking controllers that force the tracking error to approach the sliding surface, $\{\sigma(\mathbf{e}) = 0\}$, and then stay in some neighborhood of the sliding surface for all subsequent time. If f and g were known exactly to us and $\eta = 0$, then the controller

$$u^* = \frac{1}{g}(-f - \mu\sigma + x_d^{(n)} - \mathbf{k}^T \mathbf{e}), \quad (8.120)$$

where $\mathbf{k} = [0 \quad s_1 \quad \cdots \quad s_{n-1}]^T$ and $\mu > 0$ is a design parameter, would force the tracking error, \mathbf{e} , to behave as desired. We can verify this using standard arguments from sliding mode control practice—see Section 6.2. We consider the function $V = 0.5\sigma^2$. Note that V is positive definite with respect to the sliding surface $\{\sigma = 0\}$. The function V can be viewed as a distance measure to the sliding surface $\{\sigma = 0\}$. The time derivative of V evaluated on the trajectories of the

closed-loop system (8.112), (8.120) is

$$\begin{aligned}
 \dot{V} &= \sigma \dot{\sigma} \\
 &= \sigma \mathbf{s} \dot{\mathbf{e}} \\
 &= \sigma (e^{(n)} + s_{n-1}e^{(n-1)} + \cdots + s_1 \dot{e}) \\
 &= \sigma (x^{(n)} - x_d^{(n)} + \mathbf{k}^T \mathbf{e}) \\
 &= \sigma (f + gu - x_d^{(n)} + \mathbf{k}^T \mathbf{e}).
 \end{aligned} \tag{8.121}$$

Substituting into the equation above the control law given by (8.120) yields

$$\dot{V} = -\mu \sigma^2. \tag{8.122}$$

Thus, the sliding surface, $\{\sigma = 0\}$, is asymptotically attractive and the system restricted to the sliding surface can be made asymptotically stable with respect to the origin by an appropriate choice of the parameters, s_i , of the sliding surface.

The control law (8.120) requires perfect knowledge of the plant's model. However, the plant model components f and g are unknown to us. In the following subsection, we describe how to use adaptive fuzzy logic controllers to approximate, as closely as possible, the performance of (8.120). In the following, we consider two cases. In the first case, the function $f = f(\mathbf{x})$ is unknown to us while $g = g(\mathbf{x}) > 0$ is known to us. In the second case, both f and g are unknown to us. In both cases, we consider the presence of a bounded disturbance modeled by η .

8.8.3 Controllers

Unknown $f(\mathbf{x})$ When in our plant model (8.112) only $g(\mathbf{x})$ is known, $\eta \neq 0$, and $f(\mathbf{x})$ is unknown, we propose to approximate $f(\mathbf{x})$ with fuzzy logic system of the form $\boldsymbol{\theta}_f^T \boldsymbol{\xi}_f(\mathbf{x})$, where the functions $\boldsymbol{\theta}_f$ and $\boldsymbol{\xi}_f$ will be described after defining relevant symbols. We use $\boldsymbol{\theta}_f^*$ to denote an “optimal” vector such that

$$\boldsymbol{\theta}_f^* = \operatorname{argmin}_{\boldsymbol{\theta}_f} \sup_{\mathbf{x} \in \Omega} |f(\mathbf{x}) - \boldsymbol{\theta}_f^T \boldsymbol{\xi}_f(\mathbf{x})|, \tag{8.123}$$

where $\Omega \subseteq \mathbb{R}^n$ is a region to which the state \mathbf{x} is constrained to reside. We assume that

$$|f(\mathbf{x}) - \boldsymbol{\theta}_f^{*T} \boldsymbol{\xi}_f(\mathbf{x})| \leq d_f \quad \text{for all } \mathbf{x} \in \Omega, \tag{8.124}$$

where $d_f > 0$ and each element of $\boldsymbol{\theta}_f^*$ is constant and bounded below and above as follows:

$$\underline{\theta}_{fi} \leq \theta_{fi}^* \leq \bar{\theta}_{fi} \quad \text{for all } i = 1, \dots, m, \tag{8.125}$$

or in vector notation,

$$\underline{\boldsymbol{\theta}}_f \leq \boldsymbol{\theta}_f^* \leq \bar{\boldsymbol{\theta}}_f. \tag{8.126}$$

We define the adaptation parameter error as

$$\boldsymbol{\phi}_f = \boldsymbol{\theta}_f - \boldsymbol{\theta}_f^*. \tag{8.127}$$

We will use the adaptation law

$$\dot{\boldsymbol{\phi}}_f(t) = \dot{\boldsymbol{\theta}}_f(t) = \operatorname{Proj}_{\boldsymbol{\theta}_f}(\gamma_f \sigma \boldsymbol{\xi}_f(t)), \tag{8.128}$$

where $\gamma_f > 0$ is a design parameter. It is easy to verify, using the definition of $Proj$ and the fact that $\underline{\theta}_{fi} \leq \theta_{fi}^* \leq \bar{\theta}_{fi}$ for each i , that

$$\boldsymbol{\phi}_f^T \left(\frac{1}{\gamma_f} Proj_{\theta_f}(\gamma_f \sigma \boldsymbol{\xi}_f) - \sigma \boldsymbol{\xi}_f \right) \leq 0. \quad (8.129)$$

We now analyze the following fuzzy adaptive control law,

$$u = \frac{1}{g} (-\boldsymbol{\theta}_f^T \boldsymbol{\xi}_f - \mu \sigma + x_d^{(n)} - \mathbf{k}^T \mathbf{e} + u_s) \quad (8.130)$$

with u_s satisfying

$$\sigma(f - \boldsymbol{\theta}_f^T \boldsymbol{\xi}_f + u_s + \eta) \leq \epsilon, \quad (8.131)$$

$$\sigma u_s \leq 0, \quad (8.132)$$

where $\epsilon > 0$ is a design parameter. We state and prove a theorem concerning the dynamical behavior of the closed-loop system driven by the control law (8.130).

Theorem 8.8 For the closed-loop system,

$$\begin{aligned} \dot{\mathbf{x}}^{(n)} &= \mathbf{f} + g\mathbf{u} + \boldsymbol{\eta}, \\ u &= \frac{1}{g} (-\boldsymbol{\theta}_f^T \boldsymbol{\xi}_f - \mu \sigma + x_d^{(n)} - \mathbf{k}^T \mathbf{e} + u_s), \\ \dot{\boldsymbol{\theta}}_f &= Proj_{\theta_f}(\gamma_f \sigma \boldsymbol{\xi}_f), \end{aligned} \quad (8.133)$$

where u_s satisfies (8.131) and (8.132), we have the following:

1.

$$\sigma^2(t) \leq e^{-2\mu t} \sigma^2(0) + \frac{\epsilon}{\mu}. \quad (8.134)$$

2. If $\eta = 0$ and there exists $\boldsymbol{\theta}_f^*$ such that $\mathbf{f}(\mathbf{x}) = \boldsymbol{\theta}_f^{*T} \boldsymbol{\xi}_f(\mathbf{x})$, then the origin of the $(\sigma, \boldsymbol{\phi}_f)$ space is stable, and hence $\sigma(t)$ and $\boldsymbol{\phi}_f(t)$ are bounded for all $t \geq 0$ and $\epsilon(t) \rightarrow 0$ as $t \rightarrow \infty$.

Proof To prove 1, consider $V = 0.5\sigma^2$. The time derivative of V evaluated on the solutions of the closed-loop system (8.112), (8.130) is

$$\begin{aligned} \dot{V} &= \sigma \dot{\sigma} \\ &= -\mu \sigma^2 + \sigma(f - \boldsymbol{\theta}_f^T \boldsymbol{\xi}_f + u_s + \eta). \end{aligned}$$

Taking into account (8.131) yields

$$\dot{V} \leq -\mu \sigma^2 + \epsilon = -2\mu V + \epsilon.$$

Applying to the above the comparison theorem, presented in Subsection A.14.2, we obtain

$$\begin{aligned} V(t) &\leq e^{-2\mu t} V(0) + \frac{\epsilon}{2\mu} (1 - e^{-2\mu t}) \\ &\leq e^{-2\mu t} V(0) + \frac{\epsilon}{2\mu}. \end{aligned}$$

Noting in the above that $V = 0.5\sigma^2$, we obtain (8.134), which proves the first part of the theorem.

To prove 2, consider the following Lyapunov function candidate:

$$V = \frac{1}{2} \left(\sigma^2 + \frac{1}{\gamma_f} \boldsymbol{\phi}_f^T \boldsymbol{\phi}_f \right). \quad (8.135)$$

We evaluate the time derivative of the above Lyapunov function candidate on the solutions of the closed-loop system (8.133), where $\eta = 0$. We obtain

$$\dot{V} = \sigma(-\mu\sigma + f - \boldsymbol{\theta}_f^T \boldsymbol{\xi}_f + u_s) + \frac{1}{\gamma_f} \boldsymbol{\phi}_f^T \dot{\boldsymbol{\phi}}_f.$$

Taking into account the assumption that $f = \boldsymbol{\theta}_f^{*T} \boldsymbol{\xi}_f$ and using (8.127), (8.128), (8.129), and (8.132) gives

$$\begin{aligned} \dot{V} &= -\mu\sigma^2 + \sigma u_s - (\boldsymbol{\theta}_f^T \boldsymbol{\xi}_f - \boldsymbol{\theta}_f^{*T} \boldsymbol{\xi}_f) \sigma + \frac{1}{\gamma_f} \boldsymbol{\phi}_f^T \dot{\boldsymbol{\phi}}_f \\ &= -\mu\sigma^2 + \sigma u_s - \boldsymbol{\phi}_f^T \boldsymbol{\xi}_f \sigma + \boldsymbol{\phi}_f^T \frac{1}{\gamma_f} \text{Proj}_{\boldsymbol{\theta}_f}(\gamma_f \sigma \boldsymbol{\xi}_f) \\ &\leq -\mu\sigma^2 + \sigma u_s \\ &\leq -\mu\sigma^2. \end{aligned} \quad (8.136)$$

Thus, $\dot{V} \leq 0$ in the $[\sigma \ \boldsymbol{\phi}_f]^T$ space, which implies that the zero solution of the closed-loop system (8.133) is stable; therefore, $\sigma(t)$ and $\boldsymbol{\phi}_f(t)$ are bounded for $t \geq 0$.

To show that $\sigma(t) \rightarrow 0$ as $t \rightarrow \infty$, we integrate (8.136) from 0 to t to obtain

$$\int_0^t \dot{V}(\tau) d\tau = V(t) - V(0) \leq \int_0^t (-\mu\sigma^2) d\tau, \quad (8.137)$$

or, equivalently,

$$\int_0^t \mu\sigma^2 d\tau \leq V(0) - V(t) \leq V(0). \quad (8.138)$$

Thus, $\lim_{t \rightarrow \infty} \int_0^t \sigma^2(\tau) d\tau$ exists. It follows from the first part of the theorem that σ^2 is bounded and because

$$\frac{d}{dt} \sigma^2 = 2\sigma\dot{\sigma} \leq -\mu\sigma^2 + \epsilon, \quad (8.139)$$

we conclude that $\frac{d}{dt} \sigma^2$ is also bounded. Therefore, σ^2 is uniformly continuous. By Barbalat's lemma (see Section 4.12), $\sigma^2 \rightarrow 0$ as $t \rightarrow \infty$. This implies the following for tracking error: $e(t) \rightarrow 0$ as $t \rightarrow \infty$, which concludes the proof of the theorem.

Selecting u_s in Control Law (8.130) We now discuss a method that can be used to select the component u_s while implementing the control law (8.130). Let h be a constant such that

$$h \geq \|\boldsymbol{\xi}_f\| \|\bar{\boldsymbol{\theta}}_f - \boldsymbol{\theta}_f\|.$$

Given the design parameter ϵ , choose positive numbers ϵ_1 , ϵ_2 , and ϵ_3 such that $\epsilon = \epsilon_1 + \epsilon_2 + \epsilon_3$. Then, u_s can be chosen as

$$u_s = -k_s \sigma, \quad (8.140)$$

where

$$k_s \geq \frac{d_f^2}{4\epsilon_1} + \frac{h^2}{4\epsilon_2} + \frac{d^2}{4\epsilon_3}. \quad (8.141)$$

We will now verify that the above choice of u_s guarantees that (8.140) satisfies (8.131) and (8.132). We first consider the left-hand side of (8.131). Adding and subtracting $\theta_f^{*T} \xi_f$ and using (8.127) yields

$$\begin{aligned} \sigma(f - \theta_f^T \xi_f + u_s + \eta) &= \sigma(f - \theta_f^{*T} \xi_f + \theta_f^{*T} \xi_f - \theta_f^T \xi_f + u_s + \eta) \\ &= \sigma(f - \theta_f^{*T} \xi_f - \phi_f^T \xi_f - k_s \sigma + \eta). \end{aligned} \quad (8.142)$$

We now use (8.141) and group the terms to obtain

$$\begin{aligned} \sigma(f - \theta_f^T \xi_f + u_s + \eta) &\leq \sigma\left(f - \theta_f^{*T} \xi_f - \frac{d_f^2}{4\epsilon_1} \sigma\right) \\ &\quad - \sigma\left(\phi_f^T \xi_f + \frac{h^2}{4\epsilon_2} \sigma\right) + \sigma\left(\eta - \frac{d^2}{4\epsilon_3} \sigma\right). \end{aligned} \quad (8.143)$$

Next, we consider each term on the right-hand side of the above inequality separately. We start with the first term. Completing the squares, we obtain

$$\sigma\left(f - \theta_f^{*T} \xi_f - \frac{d_f^2}{4\epsilon_1} \sigma\right) = -\left(\frac{d_f}{2\sqrt{\epsilon_1}} \sigma - \frac{f - \theta_f^{*T} \xi_f}{d_f/\sqrt{\epsilon_1}}\right)^2 + \left(\frac{f - \theta_f^{*T} \xi_f}{d_f/\sqrt{\epsilon_1}}\right)^2.$$

Neglecting the first term in the above and using the assumption (8.124), we get

$$\sigma\left(f - \theta_f^{*T} \xi_f - \frac{d_f^2}{4\epsilon_1} \sigma\right) \leq \frac{d_f^2}{d_f^2/\epsilon_1} \leq \epsilon_1. \quad (8.144)$$

Completing the squares in the second term on the right-hand side of (8.143) gives

$$\begin{aligned} -\sigma\left(\phi_f^T \xi_f + \frac{h^2}{4\epsilon_2} \sigma\right) &= -\left(\frac{h}{2\sqrt{\epsilon_2}} \sigma + \frac{\phi_f^T \xi_f}{h/\sqrt{\epsilon_2}}\right)^2 + \left(\frac{\phi_f^T \xi_f}{h/\sqrt{\epsilon_2}}\right)^2 \\ &\leq \left(\frac{\phi_f^T \xi_f}{h/\sqrt{\epsilon_2}}\right)^2 \\ &= \left(\frac{(\theta_f - \theta_f^*)^T \xi_f}{h}\right)^2 \epsilon_2 \\ &\leq \frac{\|\theta_f - \theta_f^*\|^2 \|\xi_f\|^2}{h^2} \epsilon_2 \\ &\leq \frac{\|\theta_f - \theta_f^*\|^2 \|\xi_f\|^2}{\|\xi_f\|^2 \|\theta_f - \theta_f^*\|^2} \epsilon_2 \\ &\leq \epsilon_2. \end{aligned} \quad (8.145)$$

$$\leq \epsilon_2. \quad (8.146)$$

Completing the squares in the third term on the right-hand side of (8.143) yields

$$\begin{aligned}\sigma \left(\eta - \frac{d^2}{4\epsilon_3} \sigma \right) &= - \left(\frac{d}{2\sqrt{\epsilon_3}} \sigma - \frac{\eta}{d/\sqrt{\epsilon_3}} \right)^2 + \frac{\eta^2}{d^2} \epsilon_3 \\ &\leq \epsilon_3.\end{aligned}$$

Thus,

$$\sigma(f - \theta_f^T \xi_f + u_s + \eta) \leq \epsilon_1 + \epsilon_2 + \epsilon_3 = \epsilon.$$

Hence (8.131) is satisfied.

Also, $u_s = -k_s \sigma$, $k_s > 0$, satisfies (8.132).

Unknown $f(\mathbf{x})$ and $g(\mathbf{x})$ We now assume that $g(\mathbf{x})$ and $f(\mathbf{x})$ are unknown to us and $\eta \neq 0$. We approximate $f(\mathbf{x})$ and $g(\mathbf{x})$ with fuzzy logic systems, $\theta_f^T \xi_f(\mathbf{x})$ and $\theta_g^T \xi_g(\mathbf{x})$, respectively. The conditions on $\theta_f^T \xi_f(\mathbf{x})$ are the same as in the previous case, and the conditions on $\theta_g^T \xi_g(\mathbf{x})$ are similar; that is, the optimal parameter vector, θ_g^* , satisfies

$$\theta_g^* = \operatorname{argmin}_{\theta_g} \sup_{\mathbf{x} \in \Omega} |g(\mathbf{x}) - \theta_g^T \xi_g(\mathbf{x})|. \quad (8.147)$$

We assume that

$$|g(\mathbf{x}) - \theta_g^{*T} \xi_g(\mathbf{x})| \leq d_g \quad \text{for all } \mathbf{x} \in \Omega, \quad (8.148)$$

where

$$\underline{\theta}_{gi} \leq \theta_{gi}^* \leq \bar{\theta}_{gi} \quad \text{for all } i = 1, \dots, m. \quad (8.149)$$

The parameter error, $\phi_g = \theta_g - \theta_g^*$, is adapted according to the law

$$\dot{\phi}_g = \dot{\theta}_g = \operatorname{Proj}_{\theta_g}(\gamma_g \sigma \xi_g u_a),$$

where

$$u_a = \frac{1}{\theta_g^T \xi_g} (-\theta_f^T \xi_f + x_d^{(n)} - \mathbf{k}^T \mathbf{e}). \quad (8.150)$$

The initial values of the components of θ_g are chosen to satisfy (8.149), where for each $i = 1, \dots, m$ we have $\theta_{gi} > 0$. This ensures that $\theta_g^T \xi_g > 0$ because we choose fuzzy sets so that at any t at least one rule fires and so at least one component of ξ is nonzero and in fact positive. Using the definition of Proj , one can verify that

$$\phi_g^T \left(\frac{1}{\gamma_g} \operatorname{Proj}_{\theta_g}(\gamma_g \sigma \xi_g u_a) - \sigma \xi_g u_a \right) \leq 0. \quad (8.151)$$

We now analyze the following fuzzy adaptive control law,

$$u = \frac{1}{\theta_g^T \xi_g} (-\theta_f^T \xi_f + x_d^{(n)} - \mathbf{k}^T \mathbf{e}) - \frac{1}{\underline{g}} \mu \sigma + u_s \quad (8.152)$$

with u_s satisfying the following conditions:

$$\sigma(f - x_d^{(n)} + \mathbf{k}^T \mathbf{e} + g u_a + g u_s + \eta) \leq \epsilon, \quad (8.153)$$

$$\sigma u_s \leq 0, \quad (8.154)$$

where $\epsilon > 0$ is a design parameter. Using the definition of u_a , given by (8.150), we represent (8.152) as

$$u = u_a - \frac{1}{\underline{g}}\mu\sigma + u_s. \quad (8.155)$$

Theorem 8.9 For the closed-loop system,

$$\begin{aligned} x^{(n)} &= f + gu + \eta, \\ u &= \frac{1}{\theta_g^T \xi_g} \left(-\theta_f^T \xi_f + x_d^{(n)} - \mathbf{k}^T \mathbf{e} \right) - \frac{1}{\underline{g}}\mu\sigma + u_s, \\ \dot{\theta}_f &= \text{Proj}_{\theta_f}(\gamma_f \sigma \xi_f), \\ \dot{\theta}_g &= \text{Proj}_{\theta_g}(\gamma_g \sigma \xi_g u_a), \end{aligned} \quad (8.156)$$

where u_s satisfies (8.153) and (8.154), we have the following:

1.

$$\sigma^2(t) \leq e^{-2\mu t} \sigma^2(0) + \frac{\epsilon}{\mu}. \quad (8.157)$$

2. If $\eta=0$ and there exist θ_f^* and θ_g^* such that $f(\mathbf{x}) = \theta_f^{*T} \xi_f(\mathbf{x})$ and $g(\mathbf{x}) = \theta_g^{*T} \xi_g(\mathbf{x})$, then the origin of the $[\sigma \ \phi_f \ \phi_g]^T$ -space is stable and hence $\sigma(t)$, $\phi_f(t)$, and $\phi_g(t)$ are bounded and $e(t) \rightarrow 0$ as $t \rightarrow \infty$.

Proof To prove 1, consider $V = 0.5\sigma^2$. The time derivative of V evaluated on the solutions of the closed-loop system (8.112), (8.152) is

$$\begin{aligned} \dot{V} &= \sigma \dot{\sigma} \\ &= \sigma (f + gu + \eta - x_d^{(n)} + \mathbf{k}^T \mathbf{e}). \end{aligned}$$

Substituting into the above the alternative expression for the control law given by (8.155) and applying (8.153) gives

$$\begin{aligned} \dot{V} &= \sigma \left(f + gu_a - \frac{g}{\underline{g}}\mu\sigma + gu_s + \eta - x_d^{(n)} + \mathbf{k}^T \mathbf{e} \right) \\ &\leq -\frac{g}{\underline{g}}\mu\sigma^2 + \epsilon \\ &= -2\mu V + \epsilon. \end{aligned}$$

Invoking now the comparison theorem, we obtain (8.157).

To prove 2, consider the following Lyapunov function candidate:

$$V = \frac{1}{2} \left(\sigma^2 + \frac{1}{\gamma_f} \phi_f^T \phi_f + \frac{1}{\gamma_g} \phi_g^T \phi_g \right). \quad (8.158)$$

By assumption, $f = \theta_f^{*T} \xi_f$, $g = \theta_g^{*T} \xi_g$, $\eta = 0$ as well as conditions (8.129), (8.151), and (8.154) hold. In addition, it follows from (8.150) that

$$x_d^{(n)} - \mathbf{k}^T \mathbf{e} = \theta_f^T \xi_f + \theta_g^T \xi_g u_a.$$

Taking into account all of the above in evaluating the time derivative of V on the trajectories of the closed-loop system (8.156), we obtain

$$\begin{aligned}
 \dot{V} &= \sigma \left(-\frac{g}{\underline{g}} \mu \sigma + f - \theta_f^T \xi_f + (g - \theta_g^T \xi_g) u_a + g u_s \right) + \frac{1}{\gamma_f} \phi_f^T \dot{\phi}_f + \frac{1}{\gamma_g} \phi_g^T \dot{\phi}_g \\
 &= \sigma \left(f - \theta_f^T \xi_f + (g - \theta_g^T \xi_g) u_a + g u_s \right) - \frac{g}{\underline{g}} \mu \sigma^2 \\
 &\quad + \phi_f^T \frac{1}{\gamma_f} \text{Proj}_{\theta_f}(\gamma_f \sigma \xi_f) + \phi_g^T \frac{1}{\gamma_g} \text{Proj}_{\theta_g}(\gamma_g \sigma \xi_g u_a) \\
 &\leq -\frac{g}{\underline{g}} \mu \sigma^2 + \sigma g u_s - \phi_f^T \xi_f \sigma + \phi_f^T \frac{1}{\gamma_f} \text{Proj}_{\theta_f}(\gamma_f \sigma \xi_f) \\
 &\quad - \sigma \phi_g^T \xi_g u_a + \phi_g^T \frac{1}{\gamma_g} \text{Proj}_{\theta_g}(\gamma_g \sigma \xi_g u_a) \\
 &\leq -\mu \sigma^2 + g \sigma u_s \\
 &\leq -\mu \sigma^2.
 \end{aligned}$$

It follows from the above that the closed-loop system (8.156) is stable; therefore, $\sigma(t)$, $\phi_f(t)$, and $\phi_g(t)$ are bounded for $t \geq 0$.

Applying now Barbalat's lemma, similar to the way we did in the proof of Theorem 8.8, we conclude that $\sigma(t) \rightarrow 0$ as $t \rightarrow \infty$, which completes the proof of the theorem.

Selecting u_s in Control Law (8.152) The component u_s appearing in the control law (8.152) can be chosen similarly as in the previous case. We now give a lower bound on the gain k_s . Let

$$h_f \geq \|\xi_f\| \|\bar{\theta}_f - \underline{\theta}_f\| \quad \text{and} \quad h_g \geq \|\xi_g\| \|\bar{\theta}_g - \underline{\theta}_g\| |u_a|.$$

Given the design parameter ϵ , select positive numbers ϵ_1 , ϵ_2 , ϵ_3 , and ϵ_4 such that $\epsilon = \epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_4$. Then, k_s should satisfy the following condition:

$$k_s \geq \frac{1}{\underline{g}} \left(\frac{d_f^2 + |u_a|^2 d_g^2}{2\epsilon_1} + \frac{h_f^2}{4\epsilon_2} + \frac{h_g^2}{4\epsilon_3} + \frac{d^2}{4\epsilon_4} \right). \quad (8.159)$$

When the gain k_s satisfies the above condition, $u_s = -k_s \sigma$ guarantees that (8.153) and (8.154) are satisfied. To show that this is indeed the case, consider the left-hand side of (8.153). Use $\theta_g^T \xi_g u_a + \theta_f^T \xi_f = x_d^{(n)} - \mathbf{k}^T \mathbf{e}$ from (8.150), and then add and subtract $\theta_f^{*T} \xi_f$ and $\theta_g^{*T} \xi_g u_a$ to obtain

$$\begin{aligned}
 \sigma (f - x_d^{(n)} + \mathbf{k}^T \mathbf{e} + g u_a + g u_s + \eta) &= \sigma (f - \theta_f^T \xi_f - \theta_g^T \xi_g u_a + g u_a + g u_s + \eta) \\
 &= \sigma (f - \theta_f^{*T} \xi_f - \phi_f^T \xi_f + (g - \theta_g^{*T} \xi_g) u_a \\
 &\quad - \phi_g^T \xi_g u_a + g u_s + \eta).
 \end{aligned}$$

Taking into account (8.159) and rearranging terms gives

$$\begin{aligned} & \sigma(f - x_d^{(n)} + \mathbf{k}^T \mathbf{e} + gu_a + gu_s + \eta) \\ & \leq \sigma\left(f - \boldsymbol{\theta}_f^T \boldsymbol{\xi}_f - \frac{d_f^2}{2\epsilon_1} \sigma\right) + \sigma\left((g - \boldsymbol{\theta}_g^{*T} \boldsymbol{\xi}_g)u_a - \frac{|u_a|^2 d_g^2}{2\epsilon_1} \sigma\right) \\ & \quad - \sigma\left(\boldsymbol{\phi}_f^T \boldsymbol{\xi}_f + \frac{h_f^2}{4\epsilon_2} \sigma\right) - \sigma\left(\boldsymbol{\phi}_g^T \boldsymbol{\xi}_g u_a + \frac{h_g^2}{4\epsilon_3} \sigma\right) + \sigma\left(\eta - \frac{d^2}{4\epsilon_4} \sigma\right). \end{aligned}$$

Completing the squares, similarly as in the previous case, yields

$$\begin{aligned} & \sigma(f - x_d^{(n)} + \mathbf{k}^T \mathbf{e} + gu_a + gu_s + \eta) \\ & \leq \left(\frac{f - \boldsymbol{\theta}_f^T \boldsymbol{\xi}_f}{\sqrt{2}d_f/\sqrt{\epsilon_1}}\right)^2 + \left(\frac{(g - \boldsymbol{\theta}_g^{*T} \boldsymbol{\xi}_g)u_a}{\sqrt{2}|u_a|d_g/\sqrt{\epsilon_1}}\right)^2 + \left(\frac{\boldsymbol{\phi}_f^T \boldsymbol{\xi}_f}{h_f^2/\sqrt{\epsilon_2}}\right)^2 + \left(\frac{\boldsymbol{\phi}_g^T \boldsymbol{\xi}_g u_a}{h_g^2/\sqrt{\epsilon_3}}\right)^2 + \left(\frac{\eta}{d/\sqrt{\epsilon_4}}\right)^2 \\ & \leq \frac{1}{2}\epsilon_1 + \frac{1}{2}\epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_4 \\ & = \epsilon. \end{aligned}$$

Hence, (8.153) is satisfied, and so is (8.154) by the fact that $u_s = -k_s \sigma$, where $k_s > 0$.

8.8.4 Examples

In this subsection we present simulation results illustrating the performance of the control strategies discussed above. The controllers are tested on a simplified model of an electric drive system used by Fischle and Schröder [85] to test their fuzzy adaptive controllers. The plant is modeled by the differential equation

$$\ddot{x} = -\frac{5}{J} \arctan(5\dot{x}) + \frac{c_T}{J} u = f(x, \dot{x}) + gu, \quad (8.160)$$

where x is the angular position, $c_T = 10$ Nm/A, and $J = 0.1$ kg·m². We generate the reference signal x_d using the reference system whose transfer function is given by $\frac{400}{s^2 + 40s + 400}$. The input signal, $w(t)$, to the reference system is changing its value randomly between -1.5 and 1.5 every 0.5 seconds.

◆ Example 8.15

In this example, we assume that we know g , where $g = 100$, and that there are no disturbances affecting the system; that is, $\eta = 0$. We use the control law given by (8.130), where $f(\mathbf{x})$ is approximated by a fuzzy logic system. We use fuzzy sets for x and \dot{x} as in Fischle and Schröder [85]. They are shown in Figure 8.37. There are two fuzzy sets for x and six fuzzy sets for \dot{x} . Thus, we can have twelve fuzzy rules possible. They have the following form:

Rule 1: IF x is N AND \dot{x} is LN, THEN $y = \theta_1$

⋮

Rule 12: IF x is P AND \dot{x} is LP, THEN $y = \theta_{12}$,

where θ_{fi} , $i = 1, \dots, 12$ are the adaptation parameters. The bounds on the components of the adaptation parameter vector, θ_f , are chosen to be $\bar{\theta}_{fi} = 200$ and $\underline{\theta}_{fi} = -200$. All the initial values of the components of θ_f are set to zero. We chose $d_f = 50$, $d_g = 10$, and $\gamma_f = 5000$. The values of the remaining design parameters are: $\mu = 250$, $\epsilon_1 = 50$, $\epsilon_2 = 50$. As $\eta(t) = 0$, we do not need to worry about ϵ_3 .

SIMULINK is used to simulate the closed-loop system. Its block diagram is depicted in Figure 8.38. The regressor generator block produces ξ using the input x and the parameter adaptation block updates θ_f . It is made up of sets of integrators as shown in Figure 8.39. Adaptation rates of the parameters θ_{fi} are computed with the MATLAB function block, depicted in Figure 8.39, using the m-file shown in Figure 8.40. The vector s that defines the sliding surface, $\{e : se = 0\}$, is chosen as $s = [40 \ 1]$. The simulation results are shown in Figures 8.41 and 8.42.

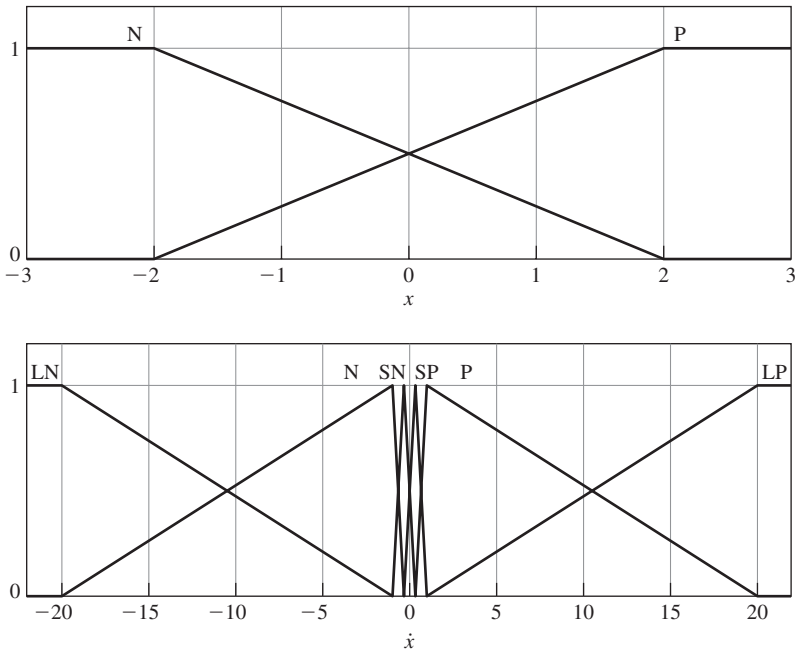


Figure 8.37 Fuzzy sets for x and \dot{x} .

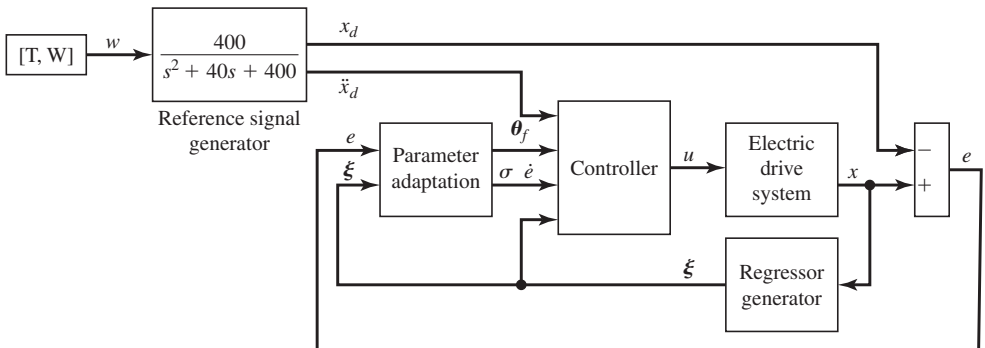


Figure 8.38 SIMULINK block diagram of fuzzy adaptive robust control system of Example 8.15.

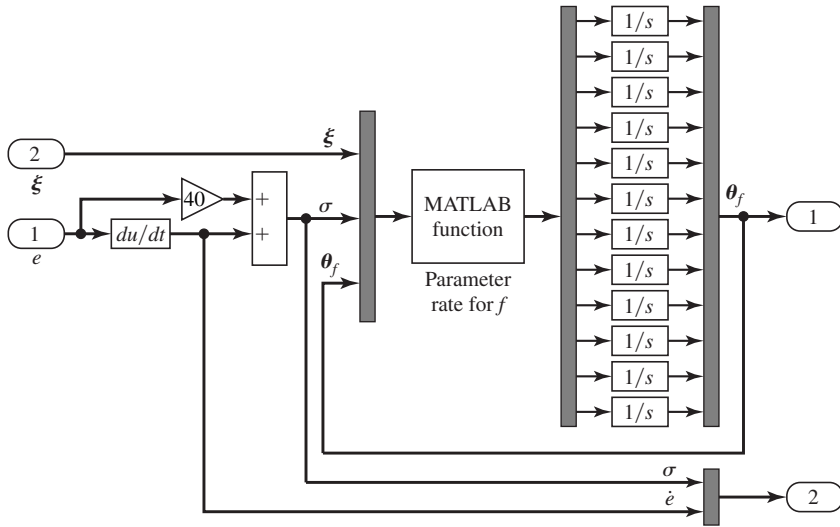


Figure 8.39 Parameter adaptation block of Example 8.15.

```
%FARC1f_UPDATE.M calculates adaptation rates
% (needs 'farcl.mdl' 'farcl_fuzzifier.m', 'af_input_gen.m')

function rate = farclf_update(x)
% x(1)--x(12) are components of xi,
% x(13) is sigma, x(14)--x(25) are components of theta
rate = zeros(12,1);
r = 5000; %adaptation gain(learning rate)
theta_lb = -200;
theta_ub = 200;

for i = 1:12;
    alpha = r*x(13)*x(i);
    rate(i) = alpha*((x(i+13)>theta_lb)|(alpha>0))&((x(i+13)<...
        theta_ub)|(alpha<0)));
end;
```

Figure 8.40 An m-file used in the block diagram of Figure 8.39 to compute adaptation rates.

The output $x(t)$, the reference signal $x_d(t)$, and the tracking error $e(t)$ are shown in Figure 8.41. The tracking error, $e(t)$, is so small that one cannot distinguish between the actual state $x(t)$ and the desired state $x_d(t)$. The plots of the control signal u and the components of θ_f versus time are depicted in Figure 8.42. As the fuzzy logic system adapts the parameter θ_f , the error, $e(t)$, gets smaller and smaller. If we used smaller ϵ , we would achieve even smaller tracking error, at the expense of higher control effort.

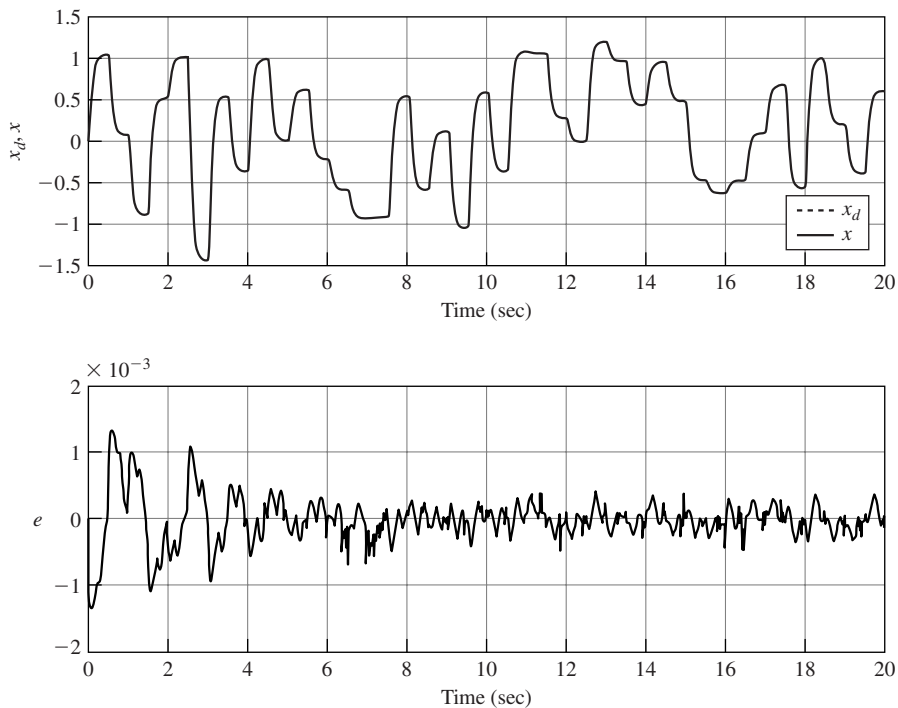


Figure 8.41 Plots of $x_d(t)$ and error $e(t)$ in Example 8.15.

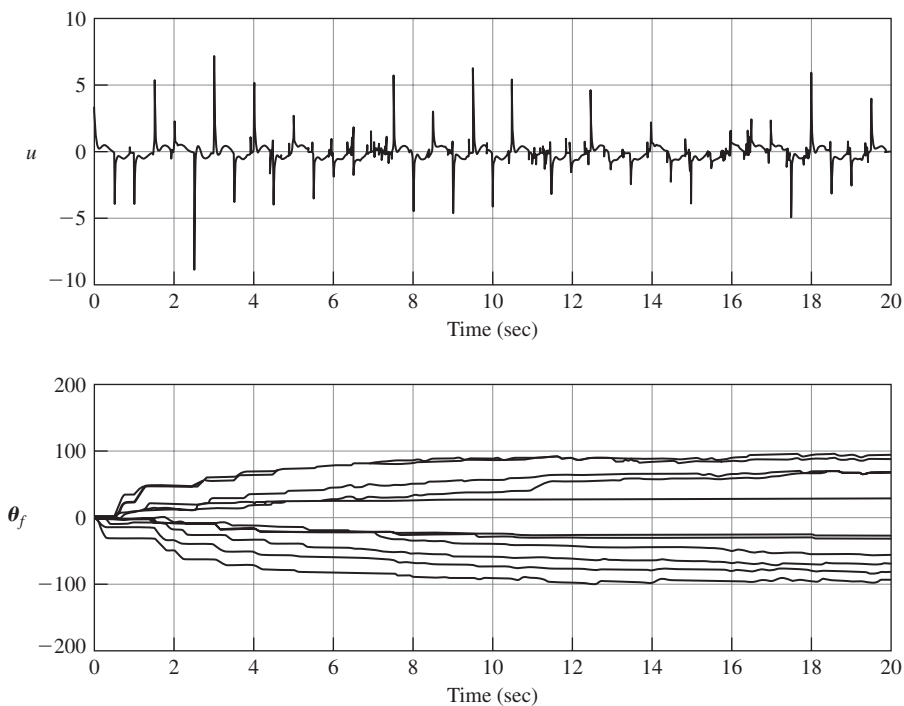


Figure 8.42 Plots of control effort $u(t)$ and the components of the parameter vector θ_f versus time in Example 8.15.

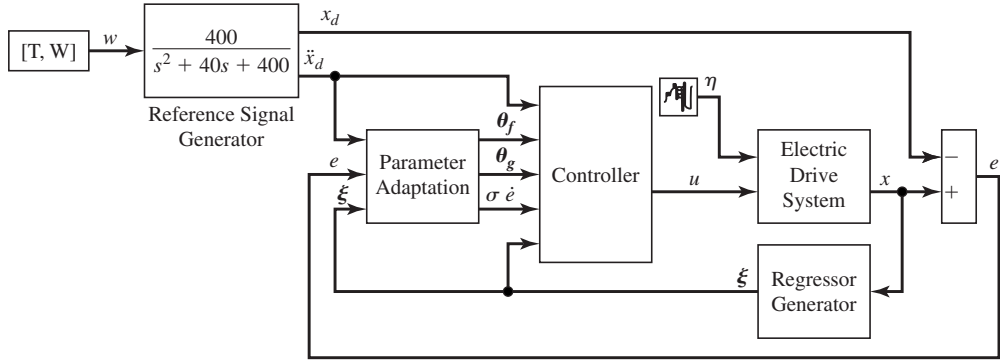


Figure 8.43 SIMULINK block diagram of the fuzzy adaptive robust control system of Example 8.16.

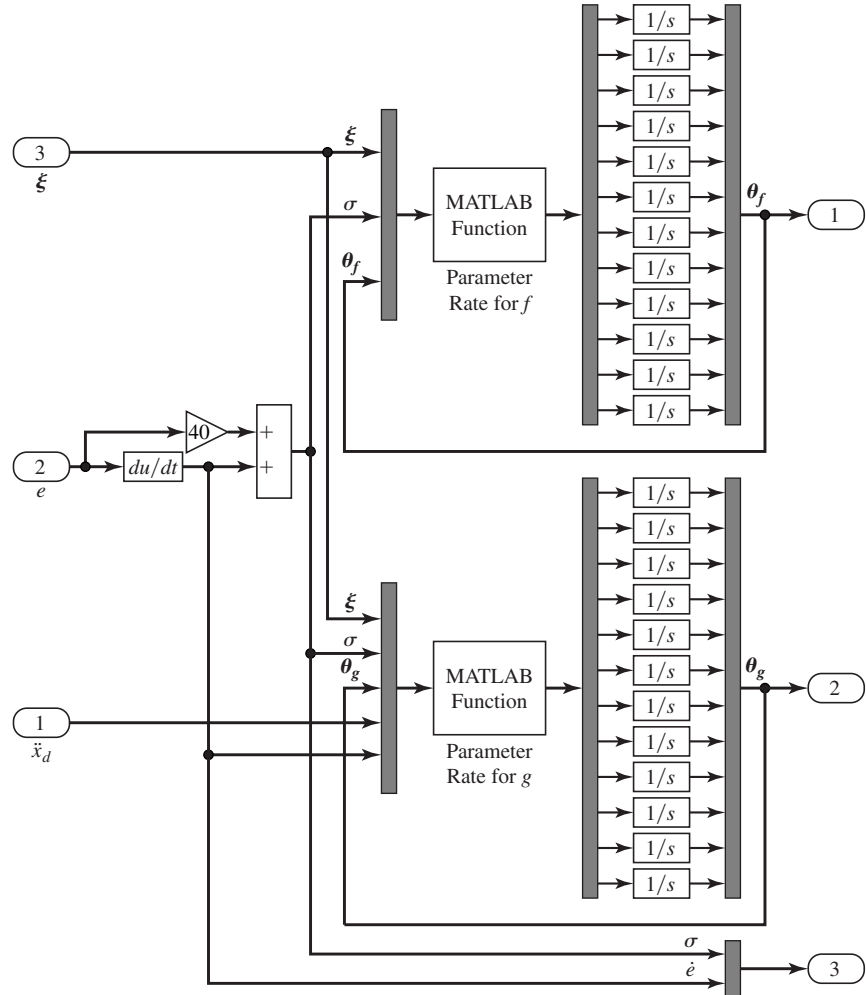


Figure 8.44 Parameter adaptation block of Example 8.16.

◆ Example 8.16

In this example, we assume that $f(\mathbf{x})$ and g are unknown to us and we employ the control law (8.152). Unknown $f(\mathbf{x})$ and g are approximated by two separate fuzzy logic systems using fuzzy sets shown in Figure 8.37 for both fuzzy logic systems.

The bounds on the adaptation parameters are chosen as follows: $\bar{\theta}_{fi} = 200$, $\underline{\theta}_{fi} = -200$, $\bar{\theta}_{gi} = 150$, and $\underline{\theta}_{gi} = 50$. All the initial values of the components of the adaptation parameter vector θ_f are set to zero. The initial values of θ_g are set to 150 to avoid excessively large control signal in the early stages of the simulation run. The selected gains are $\gamma_f = 5000$ and $\gamma_g = 1000$. The other design parameters are $\epsilon_1 = 50$, $\epsilon_2 = 100$, $\epsilon_3 = 50$, and $\epsilon_4 = 50$. The remaining parameters are the same as in the previous example.

SIMULINK block diagrams of the closed-loop system and the parameter adaptation law are depicted in Figure 8.43 and Figure 8.44, respectively. The disturbance, η , is a random signal whose plot versus time is shown in Figure 8.45. We used $d = 100$ so that the condition $|\eta(t)| \leq d$ is satisfied. We chose, as before, $\mathbf{s} = [40 \ 1]$. The simulation results are shown in Figures 8.45, 8.46, and 8.47. As can be seen in Figure 8.46, the tracking error, $e(t)$, remains very small even in the presence of the disturbance, η . Plots of the time history of the adaptation parameters are shown in Figure 8.47.

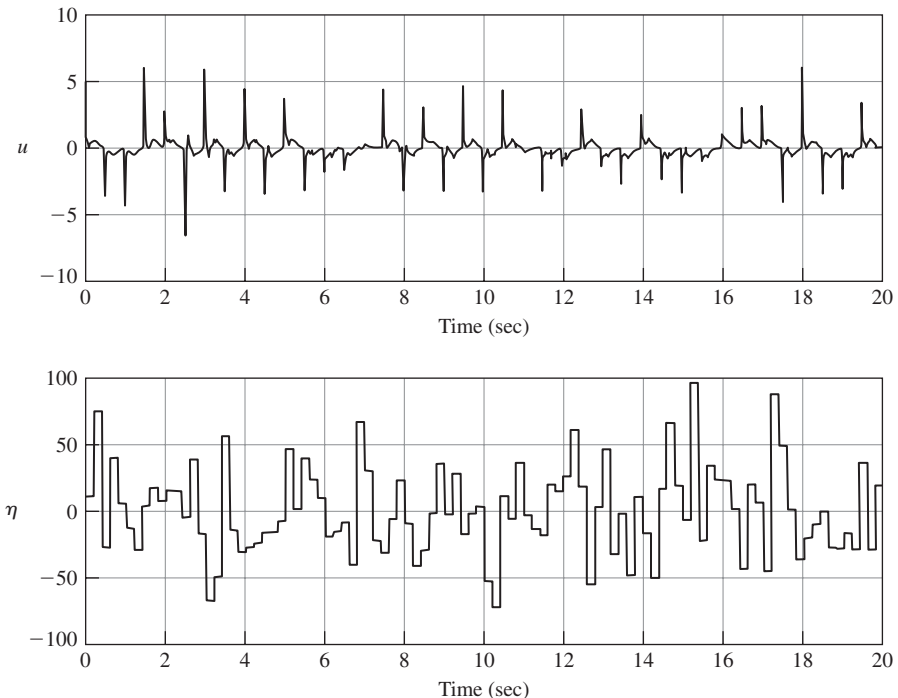


Figure 8.45 Plots of the control effort $u(t)$ and disturbance $\eta(t)$ versus time in Example 8.16.

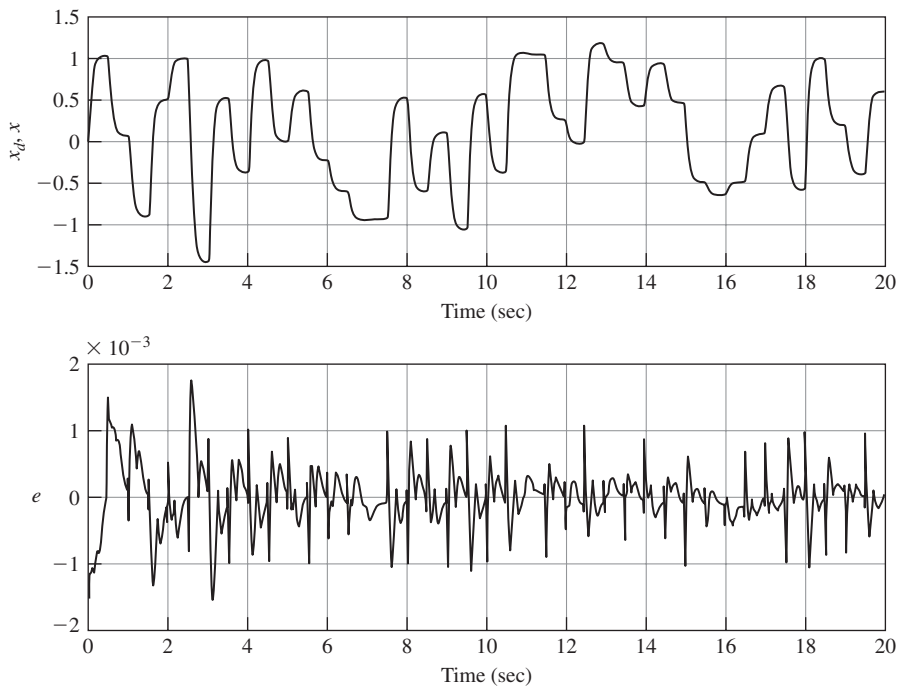


Figure 8.46 Plots of $x_d(t)$ and error $e(t)$ versus time in Example 8.16.

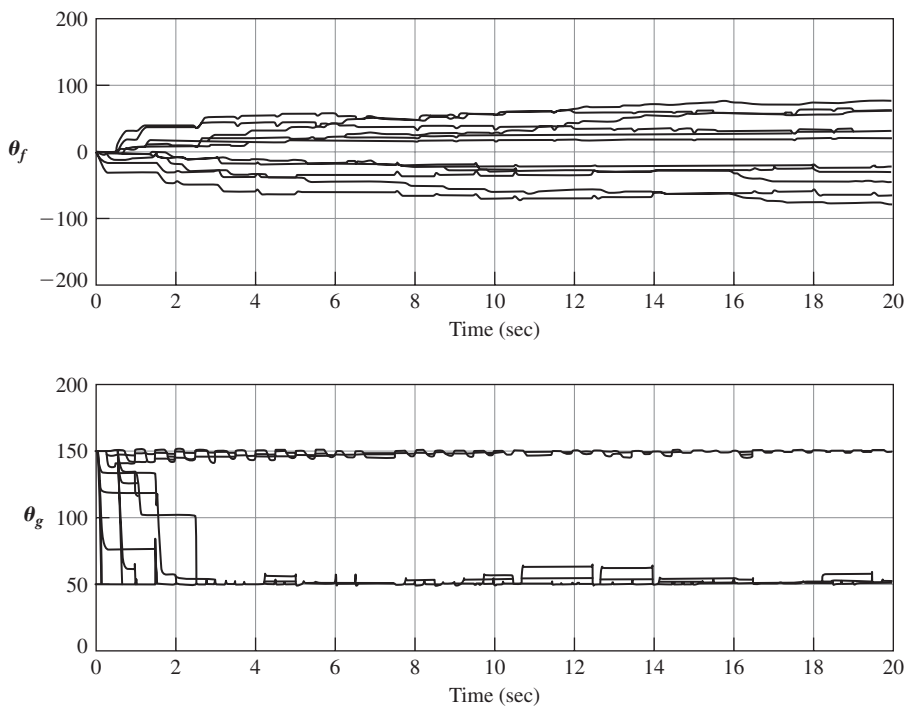


Figure 8.47 Plots of the components of the parameter vectors θ_f and θ_g versus time in Example 8.16.

Notes

The reader is encouraged to peruse a text on classical logic to become familiar with basic concepts of logic before undertaking his or her further study of fuzzy logic. An excellent introduction to logic is Copi [54]. For an introduction to fuzzy set theory the reader is referred to Klir, St.Clair, and Yuan [161], while for an introduction to fuzzy control, see Driankov, Hellendoorn, and Reinfrank [67]. Kosko [165] is the first book on fuzzy engineering. Graduate texts on fuzzy systems and control are Pedrycz [225] and Wang [296]. Paradigms and practice of fuzzy modeling are the topics of the edited book by Pedrycz [226]. Lin and Lee [185], and Jang, Sun, and Mizutani [142] cover applications of fuzzy systems as well as issues of merging fuzzy logic and neural networks. Mendel [201], on the other hand, presents extensions of the original fuzzy logic to deal with uncertainties.

The first chapter of Kosko's book [165] contains a thrilling account of the beginnings of fuzzy logic and fuzzy engineering. Kosko [165, p. 5] observes that "Fuzziness began as vagueness in the late nineteenth century." A concept is vague when it has blurred boundaries. For example, the concept *mountain* is vague because we do not know where a mountain ends and a hill begins. It was Bertrand Russell who in 1923 published a paper in which he described the concept of vagueness in terms of symbolic logic. In the 1920s, Polish mathematician Jan Lukasiewicz invented the multivalued logic. Max Black applied multivalued logic to sets of objects in the 1930s. In the 1950s, a number of logicians worked out a vague set algebra. In his 1965 paper, Zadeh [309] used the term *fuzzy* to mean *vague*. According to Kosko [165, p. 8], "Zadeh's 1965 paper applied Lukasiewicz's logic to each object in a set to work out a complete fuzzy set algebra and to extend the convex separation theorem to pattern recognition." The beginning of fuzzy engineering, on the other hand, can be marked by the work of Mamdani and Assilian [194], who first used fuzzy calculus to construct a fuzzy logic controller.

The method of constructing fuzzy design models described in Section 8.5 as well as the stabilizability conditions of fuzzy models along with a Lyapunov-based stabilizer were first reported by Teixeira and Žak in [279]. The controllers analyzed in Section 8.8 come from Lee et al. [178]. Other fuzzy adaptive control strategies were developed by Wang [294] and by Fischle and Schröder [85]. Fischle and Schröder [85, p. 38] observed that a major limitation of all existing stable adaptive fuzzy controllers, including the ones analyzed by us, is that they were developed for systems with unlimited actuators' authority—that is, for systems with no constraints on the effort level of the actuators. If we are serious about real-life applications of fuzzy adaptive controllers, then the case of the actuators with limited authority has to be rigorously analyzed. Another issue of practical importance is the convergence of the adaptation parameters to their "correct" values. Friedland [91, p. 345] states: "One may argue that the failure of a parameter estimate to converge to its correct value is a subsidiary issue as long as the adaptive control works (even with the incorrect parameter estimate). This reasoning might be valid if you can ascertain that the algorithm works for every input to which the system might be subjected. But in most applications the performance of the algorithm is tested only with a limited set of inputs. Satisfactory performance for this set of inputs is not a reliable indicator that its performance will be satisfactory for inputs outside the test set." In the current state-of-the-art adaptive control algorithms, the convergence of the adaptation parameters is usually ensured through the so-called *persistently exciting* inputs to the adaptation algorithms. However, the injection of the persistently exciting signals into the adaptation algorithms may not be feasible in real life applications. Thus, an alternative approach to the issues of the convergence of the adaptation parameters is needed.

EXERCISES

- 8.1** Consider a process whose input u is a function of two variables e and \dot{e} . The membership functions for e , \dot{e} , and u are given in Table 8.5. The control rules are:

Table 8.5 Membership Matrix Table for e , \dot{e} , and u for Exercise 8.1

Fuzzy Numbers	UNIVERSE OF DISCOURSE						
	−3	−2	−1	0	1	2	3
LN	1	0.5	0.25	0	0	0	0
SN	0.5	1	0.5	0.25	0	0	0
ZE	0	0.25	0.5	1	0.5	0.25	0
SP	0	0	0	0.25	0.5	1	0.5
LP	0	0	0	0	0.25	0.5	1

1. If e is ZE and \dot{e} is SP, then u is SN.
2. If e is ZE and \dot{e} is ZE, then u is ZE.
3. If e is SN and \dot{e} is SN, then u is LP.
4. If e is SP and \dot{e} is ZE, then u is SN.

Determine the value of u , for $e = 1$ and $\dot{e} = -2$, using

- (a) the algebraic product defuzzifier;
- (b) the center-of-gravity defuzzifier.

- 8.2** Use Lemma 8.3 to prove the following result due to Bailey [16]:

Let \mathbf{B} be a matrix with negative diagonal elements and nonnegative off-diagonal elements. If $\mathbf{x}(t; \mathbf{x}_0, t_0)$ and $\mathbf{y}(t; \mathbf{y}_0, t_0)$ are solutions of $\dot{\mathbf{x}} \leq \mathbf{B}\mathbf{x}$ and $\dot{\mathbf{y}} = \mathbf{B}\mathbf{y}$, respectively, and $\mathbf{x}_0 = \mathbf{y}_0$, then $\mathbf{x}(t; \mathbf{x}_0, t_0) \leq \mathbf{y}(t; \mathbf{y}_0, t_0)$ for all $t \geq t_0$.

- 8.3** Consider the one-link robot manipulator of Exercise 7.2, whose model is given by (7.4) and (7.5). Let

$$\mathbf{P} = \mathbf{s}^T (\mathbf{s}\mathbf{s}^T)^{-1} \mathbf{s},$$

where $\mathbf{s} \in \mathbb{R}^{1 \times 2}$ is the row vector that was found in Exercise 7.2. Then, $(\mathbf{I}_n - \mathbf{P})$ is the orthogonal projector of \mathbf{x} onto the sliding surface $\{\mathbf{x} : \mathbf{s}\mathbf{x} = 0\}$. The Euclidean distance of the point \mathbf{x} to the surface $\{\mathbf{x} : \mathbf{s}\mathbf{x} = 0\}$ is

$$\|\mathbf{P}\mathbf{x}\|_2.$$

Design a fuzzy sliding mode controller whose inputs are

$$\|\mathbf{P}\mathbf{x}(t)\|_2 \text{sign}(\sigma(\mathbf{x}(t))) \quad \text{and} \quad \Delta \|\mathbf{P}\mathbf{x}(t)\|_2 = \|\mathbf{P}\mathbf{x}(t)\|_2 - \|\mathbf{P}\mathbf{x}(t - T_d)\|_2,$$

where T_d is the time delay that is a design parameter. Use the rules matrix given in Table 8.6.

Table 8.6 The Rules Matrix for Exercise 8.3

Error \ Change-in-error	LN	N	ZE	P	LP
LN	LP	LP	LP	LP	LP
N	P	P	P	LP	LP
ZE	ZE	ZE	ZE	ZE	ZE
P	LN	LN	N	N	N
LP	LN	LN	LN	LN	LN

8.4 Consider a dynamical system modeled by (8.6)–(8.8). Let

$$\alpha_j(t) = 1 - \sum_{\substack{i=1 \\ i \neq j}}^r \alpha_i(t). \quad (8.161)$$

Represent (8.6) as

$$\begin{aligned} \dot{\mathbf{x}} &= \left(\alpha_j \mathbf{A}_j + \sum_{\substack{i=1 \\ i \neq j}}^r \alpha_i \mathbf{A}_i \right) \mathbf{x} + \left(\alpha_j \mathbf{B}_j + \sum_{\substack{i=1 \\ i \neq j}}^r \alpha_i \mathbf{B}_i \right) \mathbf{u} \\ &= \left(\mathbf{A}_j + \sum_{\substack{i=1 \\ i \neq j}}^r \alpha_i (\mathbf{A}_i - \mathbf{A}_j) \right) \mathbf{x} + \left(\mathbf{B}_j + \sum_{\substack{i=1 \\ i \neq j}}^r \alpha_i (\mathbf{B}_i - \mathbf{B}_j) \right) \mathbf{u}. \end{aligned} \quad (8.162)$$

Let

$$\boldsymbol{\alpha}^j(t) = \boldsymbol{\alpha}^j = [\alpha_1 \quad \cdots \quad \alpha_{j-1} \quad \alpha_{j+1} \quad \cdots \quad \alpha_r]^T.$$

Observe that $\boldsymbol{\alpha}^j \in [0, 1]^{r-1}$; that is, the values of the vector-valued function $\boldsymbol{\alpha}^j(\cdot)$ are in the $(r - 1)$ -dimensional unit hypercube. Next, let

$$\Delta \mathbf{A}(\boldsymbol{\alpha}^j) = \sum_{\substack{i=1 \\ i \neq j}}^r \alpha_i (\mathbf{A}_i - \mathbf{A}_j) \quad \text{and} \quad \Delta \mathbf{B}(\boldsymbol{\alpha}^j) = \sum_{\substack{i=1 \\ i \neq j}}^r \alpha_i (\mathbf{B}_i - \mathbf{B}_j).$$

Using the above notation, represent model (8.162), in an equivalent way, as

$$\begin{aligned} \dot{\mathbf{x}} &= (\mathbf{A}_j + \Delta \mathbf{A}(\boldsymbol{\alpha}^j)) \mathbf{x} + (\mathbf{B}_j + \Delta \mathbf{B}(\boldsymbol{\alpha}^j)) \mathbf{u} \\ &= \mathbf{A}_j \mathbf{x} + \mathbf{B}_j \mathbf{u} + \mathbf{F}(\boldsymbol{\alpha}^j, \mathbf{x}, \mathbf{u}), \end{aligned} \quad (8.163)$$

where $\mathbf{F}(\boldsymbol{\alpha}^j, \mathbf{x}, \mathbf{u}) = \Delta \mathbf{A}(\boldsymbol{\alpha}^j) \mathbf{x} + \Delta \mathbf{B}(\boldsymbol{\alpha}^j) \mathbf{u}$. Let

$$\mathbf{F}(\cdot) = \mathbf{f}(\cdot) + \mathbf{B}_j \mathbf{h}(\cdot). \quad (8.164)$$

We refer to \mathbf{f} as the *unmatched uncertainty*, while \mathbf{h} is referred to as the *matched uncertainty*. Hence, (8.163) takes the form

$$\dot{\mathbf{x}} = \mathbf{A}_j \mathbf{x} + \mathbf{B}_j \mathbf{u} + \mathbf{f} + \mathbf{B}_j \mathbf{h}. \quad (8.165)$$

Assume that

$$\begin{aligned} \|\mathbf{f}\| &= \|\mathbf{f}(\boldsymbol{\alpha}^j, \mathbf{x})\| \leq \alpha_f \|\mathbf{x}\|, \\ \|\mathbf{h}\| &= \|\mathbf{h}(\boldsymbol{\alpha}^j, \mathbf{x}, \mathbf{u})\| \leq \beta_x \|\mathbf{x}\| + \beta_u \|\mathbf{u}\|. \end{aligned} \quad (8.166)$$

Prove the following theorem:

Theorem 8.10 Suppose that A_j is asymptotically stable, and $P = P^T > 0$ is the solution to the Lyapunov matrix equation $A_j^T P + P A_j = -2Q$ for some $Q = Q^T > 0$. Suppose also that

$$\alpha_f < \frac{\lambda_{\min}(Q)}{\lambda_{\max}(P)} \quad \text{and} \quad \beta_u < 1.$$

Then, the state feedback controller $u = -\gamma B_j^T P x$, where

$$\gamma > \frac{\beta_x^2}{4(\lambda_{\min}(Q) - \alpha_f \lambda_{\max}(P))(1 - \beta_u)}$$

globally asymptotically stabilizes the uncertain system for arbitrary f and h that satisfy the norm bounds (8.166).

8.5 In this exercise, the following definition is used:

Definition 8.4 We say that system (8.163) satisfies the matching conditions if there exist continuous matrix functions $D(\cdot) : \mathbb{R}^{r-1} \rightarrow \mathbb{R}^{m \times n}$ and $E(\cdot) : \mathbb{R}^{r-1} \rightarrow \mathbb{R}^{m \times m}$ such that

$$\begin{aligned} \Delta A(\alpha^j) &= B_j D(\alpha^j), \\ \Delta B(\alpha^j) &= B_j E(\alpha^j) \end{aligned} \tag{8.167}$$

for all $\alpha^j \in [0, 1]^{r-1}$.

Consider the following algorithm for linear controller design for fuzzy system models of the form (8.162) that is the subject of Exercise 8.4. The algorithm is a consequence of Theorem 8.10 of Exercise 8.4.

LINEAR CONTROLLER DESIGN ALGORITHM FOR FUZZY SYSTEMS

Given fuzzy model (8.162):

STEP 1. Set $j = 1$

STEP 2. Check matching conditions (8.167). If the matching conditions are satisfied, calculate matrices D and E , and check if $\alpha_u < 1$. If the matching conditions are not satisfied, express the uncertain element as $F = f + B_j h$, calculate β_u , and check if $\beta_u < 1$. If $\beta_u > 1$, set $j = j + 1$ and repeat the above computations.

STEP 3. Stabilize A_j , if necessary. Calculate $u_1 = -K_1 x$ so that the matrix $A_j - B_j K_1$ has its eigenvalues in the desired locations. Update the matched uncertainty involving the state vector; that is, replace D with $D - E K_1$.

STEP 4. Solve the Lyapunov equation $A_j^T P + P A_j = -2Q$ for some $Q = Q^T > 0$. Choosing $Q = I_n$ maximizes the ratio $\lambda_{\min}(Q)/\lambda_{\max}(P)$. Check if $\alpha_f < \lambda_{\min}(Q)/\lambda_{\max}(P)$. If α_f does not satisfy the condition, set $j = j + 1$ and go to step 2.

STEP 5. Construct $u_2 = -\gamma B_j^T P x$.

STEP 6. Construct the resulting state feedback linear control law $u = u_1 + u_2$.

Table 8.7 Parameter Numerical Values of the Model of Exercise 8.5

m = mass of pendulum	2.0 kg
M = mass of cart	8.0 kg
$2l$ = length of pendulum	1.0 m
$a = 1/(m + M)$	

Illustrate the above algorithm on a reduced-order model of an inverted pendulum mounted on a cart. The equations of motion for the pendulum are

$$\begin{cases} \dot{x}_1 = x_2, \\ \dot{x}_2 = \frac{g \sin(x_1) - amlx_2^2 \sin(2x_1)/2 - a \cos(x_1)u}{4l/3 - aml \cos^2(x_1)}, \end{cases} \quad (8.168)$$

where x_1 is the angle of the pendulum from the vertical line, x_2 is the angular velocity of the pendulum, u is the control force applied to the cart, and $g = 9.8 \text{ m/sec}^2$ is the magnitude of the acceleration due to gravity. Numerical values of the parameters are given in Table 8.7. Consider modeling equations (8.168) as a truth model. Use the truth model in your simulations to evaluate the performance of the design. Design a linear stabilizing controller using the above algorithm utilizing a design model. Assume that, by using insight and experience, you obtained the following two rules describing the plant dynamics:

Rule 1: IF $x_1(t)$ is about 0,
THEN $\dot{\mathbf{x}}(t) = \mathbf{A}_1 \mathbf{x}(t) + \mathbf{b}_1 u(t)$,

and

Rule 2: IF $x_1(t)$ is about $\pm\pi/2$,
THEN $\dot{\mathbf{x}}(t) = \mathbf{A}_2 \mathbf{x}(t) + \mathbf{b}_2 u(t)$,

where

$$\mathbf{A}_1 = \begin{bmatrix} 0 & 1 \\ \frac{g}{4l/3 - aml} & 0 \end{bmatrix}, \quad \mathbf{b}_1 = \begin{bmatrix} 0 \\ a \end{bmatrix},$$

and

$$\mathbf{A}_2 = \begin{bmatrix} 0 & 1 \\ \frac{2g}{\pi(4l/3 - amlb^2)} & 0 \end{bmatrix}, \quad \mathbf{b}_2 = \begin{bmatrix} 0 \\ ab \end{bmatrix},$$

where $b = \cos(88^\circ)$. Use the following membership functions:

$$\mu_1(x_1) = \frac{1 - 1/(1 + \exp(-7(x_1 - \pi/4)))}{1 + \exp(-7(x_1 + \pi/4))} \quad \text{and} \quad \mu_2(x_1) = 1 - \mu_1(x_1).$$

The above membership functions are the same as the ones employed by Cao et al. [42]. The corresponding fuzzy system model is

$$\dot{\mathbf{x}} = (\alpha_1 \mathbf{A}_1 + \alpha_2 \mathbf{A}_2) \mathbf{x} + (\alpha_1 \mathbf{b}_1 + \alpha_2 \mathbf{b}_2) u, \quad (8.169)$$

where $\alpha_i = \mu_i$, $i = 1, 2$, because $\mu_1 + \mu_2 = 1$ for all t . Proceed to construct a stabilizing linear controller according to the above algorithm. First, represent model (8.169) as

$$\dot{\mathbf{x}} = (\mathbf{A}_1 + \alpha_2 (\mathbf{A}_2 - \mathbf{A}_1)) \mathbf{x} + (\mathbf{b}_1 + \alpha_2 (\mathbf{b}_2 - \mathbf{b}_1)) u. \quad (8.170)$$

Check that the uncertain elements of system model (8.170) satisfy the matching conditions. Suppose that the desired eigenvalues of $\mathbf{A}_1 - \mathbf{b}_1 \mathbf{k}_1$ are $\{-2, -2\}$. Construct the controller, and test its performance on the truth model. Generate plots of x_1 versus time for different initial angular displacements $x_1(0)$. You can set $x_2(0) = 0$.

- 8.6** The objective of this exercise is to understand better the dynamical behavior of the nonlinear projector used by us in adaptive fuzzy control and defined by (8.117). For example, take $\alpha(t) = \sin(0.5\pi t)$, $\underline{\theta} = -0.3$, and $\bar{\theta} = 0.5$. Compare solutions of $\dot{\theta} = \text{Proj}_{\theta}(\alpha)$ and $\dot{\theta} = \alpha$ for a number of initial conditions. Below is an example of MATLAB's function that implements the differential equation, $\dot{\theta} = \text{Proj}_{\theta}(\alpha)$, where $\text{low} = \underline{\theta}$ and $\text{upper} = \bar{\theta}$.

```
function theta_dot = Proj(t, theta)
low=-0.3;
upper=0.5;
alpha=sin(0.5*pi*t);
if (theta <= low)&(alpha<0)
    theta_dot=0;
elseif (theta >= upper)&(alpha>0)
    theta_dot=0;
else
    theta_dot=alpha;
end
```

- 8.7** For the fuzzy model the stick balancer given in Example 8.13 by equation (8.66), where

$$\mathbf{y} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix},$$

construct a fuzzy state estimator. (Note that we have $\mathbf{C}_1 = \mathbf{C}_2 = \mathbf{C}(\alpha)$.) Implement the combined controller–estimator compensator using first the fuzzy model and then the nonlinear model of the stick balancer. Write MATLAB scripts that animate the dynamical behavior of the closed-loop systems along with animation of the angle θ and the input u data.



CHAPTER 9

Neural Networks

We tend to regard forgetting as a failure of memory, the result of an inadequacy in the machinery for remembering, but a moment's reflection shows that that cannot always be right. If we were to remember every telephone number we have ever used, every scene we have ever looked at, every conversation we have ever heard, every sentence we have ever read, we should be overwhelmed with memories. In fact we transfer to our long-term memory stores only a tiny fraction of the contents of our short-term stores. We retain only what is likely to be useful.

—*An Anatomy of Thought* [101, p. 331]

This chapter is concerned with neurocomputing, which is also called brainlike computation. The motivation behind studying artificial neural networks, commonly referred to as neural networks, or just neural nets, is to build computers whose construction would mimic the organization of the brain. The brain is the most complex natural information processing system. It is capable of organizing its building blocks, called neurons, to perform computations in a very efficient way. A drawing of a typical biological neuron is given in Figure 9.1. The neuron collects signals from other neurons through its dendrites and receptor sites. The signals have the form of electrical pulses. The pulses received by the dendrites travel to the neuron's cell body where they are being processed. The surface of the cell body of the neuron is dotted with numerous receptor sites that can receive messages from other neurons in a fashion similar to the dendrites. The neuron may produce its own pulses in response to excitation from other neurons. The neuron's own pulses travel through a long, thin strand called an axon, which splits into many end branches. The neuron's message to the other neurons is delivered via the end branches. The key to the transmission of the neuron's message is the synapse, shown in Figure 9.2. The end branches of the neuron's axon terminate in synaptic knobs. The synaptic knob of the first neuron is separated from the receptor site of another neuron by a microscopic distance. The cell body of the first neuron produces chemical substances called neurotransmitters, which are delivered down to the synaptic vesicles. The neurotransmitters are stored in the synaptic vesicles until the neuron fires and a burst of neurotransmitters is released by the vesicles. The neurotransmitters then flow across the synaptic cleft and act on the second neuron. The neurotransmitters of the first neuron may stimulate or inhibit the second neuron activity. The first neuron has many other synaptic

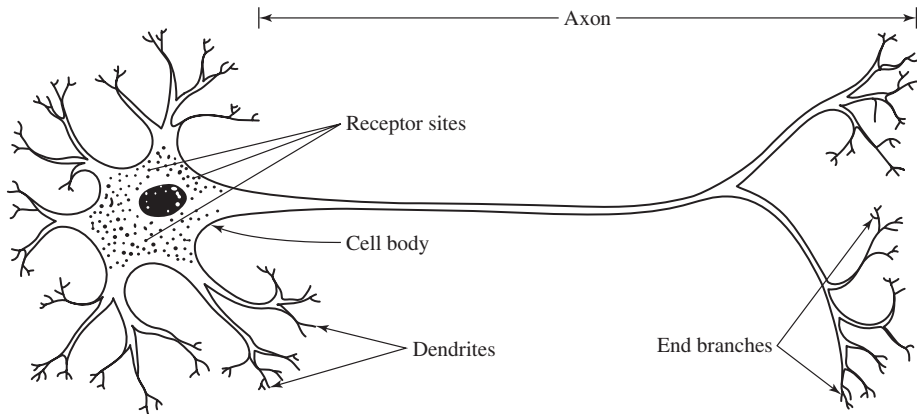


Figure 9.1 A typical biological neuron.

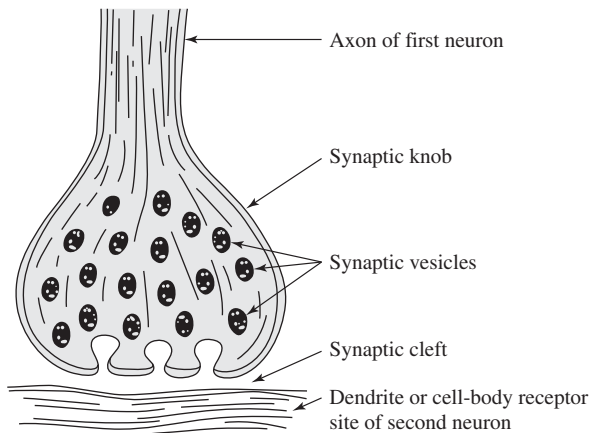


Figure 9.2 The synapse is a junction point where neuron meets neuron.

knobs elsewhere, while the second neuron has many other receptor sites that receive messages from many other neurons. Each neuron is receiving messages from hundreds or thousands of other neurons. The neuron sets off the nervous pulse; that is, it fires off its own nervous pulse only when it is properly stimulated. It was observed that, in general, the neuron will not fire in response to a single message. To fire, the neuron needs many messages arriving simultaneously or in quick succession from other neurons. In addition, the stimulating messages to fire must outweigh the messages that inhibit the neuron from firing.

In the following sections, we discuss a few mathematical models of the biological neuron. These models are referred to as the artificial neurons. Then, we analyze neural networks which are composed of the artificial neurons.

9.1 Threshold Logic Unit

A simple model of a biological neuron is the so called the *threshold logic unit* (TLU), shown in Figure 9.3. The TLU is also called the *perceptron*. A TLU is parameterized by a set of weights and a threshold (bias). The summation operation is followed by the sign activation function,

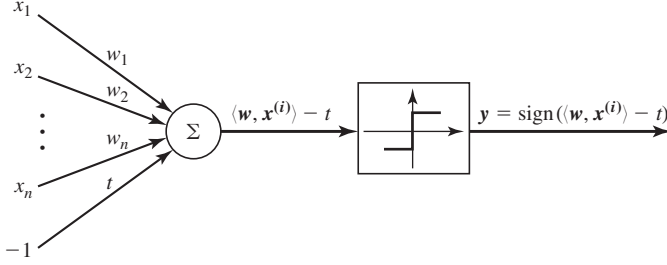


Figure 9.3 A threshold logic unit (TLU), also called the perceptron, as a model of a biological neuron.

which models the cell body of a biological neuron. The dendrites are modeled by the inputs, and their strengths are modeled by the weights. The axon corresponds to the output y . Suppose that we are given a TLU with n weights, w_1, w_2, \dots, w_n , and a threshold t . Given an input vector $\mathbf{x}^{(i)} = [x_1^{(i)} \ x_2^{(i)} \ \dots \ x_n^{(i)}]^T \in \mathbb{R}^n$, then the corresponding output of the TLU is

$$y_i = \text{sign} \left(\sum_{j=1}^n w_j x_j^{(i)} - t \right), \quad (9.1)$$

where

$$\text{sign}(z) = \begin{cases} 1 & \text{if } z \geq 0, \\ -1 & \text{if } z < 0. \end{cases}$$

We note that sometimes in the literature the sign function is undefined at 0, or it is defined as a point to a set mapping; that is, $\text{sign}(0) \in [-1, 1]$. Here, we define $\text{sign}(0) = 1$. Using the scalar product notation, we represent (9.1) as

$$y_i = \text{sign}(\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t).$$

Note that $y_i \in \{-1, 1\}$. Specifically,

$$y_i = \begin{cases} 1 & \text{if } \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t \geq 0, \\ -1 & \text{if } \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t < 0. \end{cases} \quad (9.2)$$

The goal of the TLU (perceptron) is to correctly classify the set of given input vectors $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}$ into one of the two classes, \mathcal{C}_1 or \mathcal{C}_2 . In other words, if a given vector, say $\mathbf{x}^{(i)}$, belongs to class \mathcal{C}_1 , then the corresponding output y_i should be $+1$. On the other hand, if $\mathbf{x}^{(i)}$ belongs to class \mathcal{C}_2 then the TLU's output should be -1 . Thus, we are faced with a two-class pattern classification problem. The TLU's role is to find a hyperplane that separates two decision regions. The separating hyperplane (decision boundary) has the form

$$\langle \mathbf{w}^*, \mathbf{x} \rangle - t^* = \sum_{i=1}^n w_i^* x_i - t^* = 0. \quad (9.3)$$

In the case of two variables x_1 and x_2 the decision boundary is a straight line. An illustration of such a case is given in Figure 9.4. In this case a point that lies above the boundary line is assigned to class \mathcal{C}_1 , while a point below the boundary line is assigned to class \mathcal{C}_2 . We will present an iterative algorithm for the decision boundary weights and threshold provided that a decision boundary exists. We assume that we are given p training pairs

$$(\mathbf{x}^{(1)}, d_1), (\mathbf{x}^{(2)}, d_2), \dots, (\mathbf{x}^{(p)}, d_p),$$

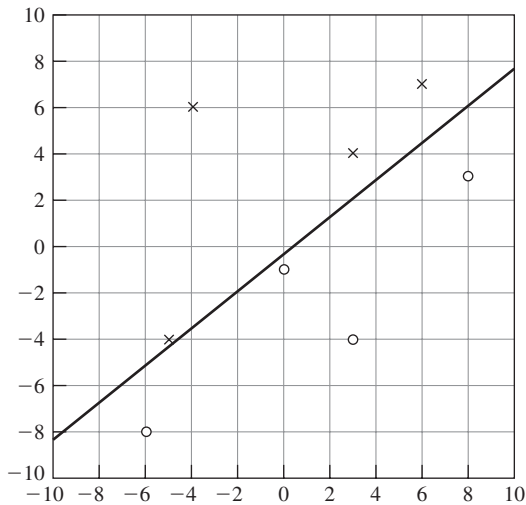


Figure 9.4 An example of separable points in two-dimensional space and a separating line found using the TLU learning algorithm.

where for $i = 1, 2, \dots, p$, the input vectors $\mathbf{x}^{(i)} \in \mathbb{R}^n$ and $d_i \in \{-1, 1\}$ represent the desired responses. We would like to determine a hyperplane (9.3) such that each input vector lies on a preassigned side of the hyperplane. If such a hyperplane exists for the given training set, then the set of input vectors is referred to as *linearly separable*, and the hyperplane is referred to as a *separating hyperplane*. If no such hyperplane exists, then the training set is referred to as *linearly nonseparable*. Finding a separating hyperplane is equivalent to finding a corresponding weight vector and a threshold. If we find such a weight vector and a threshold, then we say that the TLU correctly classifies the training set. Thus, designing a TLU that correctly classifies the given training set amounts to finding a weight vector \mathbf{w}^* and a threshold t^* such that for $i = 1, 2, \dots, p$, we have

$$\begin{aligned} \langle \mathbf{w}^*, \mathbf{x}^{(i)} \rangle - t^* &\geq 0 & \text{if } d_i = 1, \\ \langle \mathbf{w}^*, \mathbf{x}^{(i)} \rangle - t^* &< 0 & \text{if } d_i = -1. \end{aligned}$$

In our discussion, we can assume the threshold value, t , to be zero. We can achieve this by increasing the dimension of every input vector, $\mathbf{x}^{(i)}$, by augmenting it with an entry equal to -1 , and by increasing the dimension of the weight vector to $n + 1$ by augmenting it with t . Indeed,

$$y_i = \text{sign}(\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t) = \text{sign}(\langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}}^{(i)} \rangle),$$

where $\tilde{\mathbf{w}}^T = [w_1 \ w_2 \ \dots \ w_n \ t]$ and $\tilde{\mathbf{x}}^{(i)T} = [x_1^{(i)} \ x_2^{(i)} \ \dots \ x_n^{(i)} \ -1]$. In our further discussion, we assume that the given input vectors have been already augmented. In view of the above, we can formulate the training problem of a TLU as follows:

TLU Training Problem Given p training pairs $(\mathbf{x}^{(1)}, d_1), \dots, (\mathbf{x}^{(p)}, d_p)$, where $\mathbf{x}^{(i)} \in \mathbb{R}^{n+1}$ and $d_i \in \{-1, 1\}$. Determine a weight vector $\mathbf{w} \in \mathbb{R}^{n+1}$; if there exists one, such that for $i = 1, \dots, p$, we have

$$\text{sign}(\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle) = d_i.$$

Assume that the given vectors $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}$ are linearly separable—that is, there exists a weight vector \mathbf{w}^* such that

$$\langle \mathbf{w}^*, \mathbf{x}^{(i)} \rangle \geq 0 \quad \text{for every input vector } \mathbf{x}^{(i)} \text{ from class } \mathcal{C}_1$$

and

$$\langle \mathbf{w}^*, \mathbf{x}^{(i)} \rangle < 0 \quad \text{for every input vector } \mathbf{x}^{(i)} \text{ from class } \mathcal{C}_2.$$

We can now formulate the algorithm for adapting the weight vector resulting in a separating hyperplane of a given set of linearly separable points. The idea is very simple and we proceed as follows. If the given vector $\mathbf{x}^{(i)}$ is classified correctly, do not update the weights, that is,

$$\mathbf{w}(k+1) = \mathbf{w}(k).$$

Otherwise,

$$\mathbf{w}(k+1) = \mathbf{w}(k) - \mathbf{x}^{(i)} \quad \text{if } \mathbf{w}^T(k)\mathbf{x}^{(i)} \geq 0 \quad \text{and } \mathbf{x}^{(i)} \in \mathcal{C}_2 \quad (9.4)$$

and

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \mathbf{x}^{(i)} \quad \text{if } \mathbf{w}^T(k)\mathbf{x}^{(i)} < 0 \quad \text{and } \mathbf{x}^{(i)} \in \mathcal{C}_1. \quad (9.5)$$

Note that if $\mathbf{w}^T(k)\mathbf{x}^{(i)} \geq 0$ [that is, $y(k) = 1$] and $\mathbf{x}^{(i)} \in \mathcal{C}_2$ [that is, $d_i = -1$], then the error between the actual TLU output and the desired TLU output is $d_i - y(k) = -2$. Therefore, we can represent the weight update in (9.4) as

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \frac{1}{2}(d_i - y(k))\mathbf{x}^{(i)}. \quad (9.6)$$

Analogously, if $\mathbf{w}^T(k)\mathbf{x}^{(i)} < 0$ [that is, $y(k) = -1$] and $\mathbf{x}^{(i)} \in \mathcal{C}_1$, then $d_i - y(k) = 2$. Here, as in the previous case, we can represent the weight update in (9.5) as (9.6). We now summarize the above considerations in the form of the following algorithm.

TLU (PERCEPTRON) LEARNING ALGORITHM

Given a set of training pairs $\{(\mathbf{x}^{(1)}, d_1), \dots, (\mathbf{x}^{(i)}, d_i), \dots, (\mathbf{x}^{(p)}, d_p)\}$.

STEP 1 *Initialize Weight Vector and Threshold.* Set the values of the components of the initial weight vector $\mathbf{w}(0)$ and threshold $t(0)$ to zero or to small random values.

STEP 2 *Present New Input Vector and Desired Output.* Present the new input vector, $\mathbf{x}^{(i)}$, and the corresponding desired output d_i .

STEP 3 *Calculate Actual Output:*

$$y(k) = \text{sign}(\langle \mathbf{w}(k), \mathbf{x}^{(i)} \rangle - t(k)).$$

STEP 4 *Adapt Weight Vector and Threshold:*

$$\begin{aligned} \mathbf{w}(k+1) &= \mathbf{w}(k) + \frac{1}{2}(d_i - y(k))\mathbf{x}^{(i)}, \\ t(k+1) &= t(k) - \frac{1}{2}(d_i - y(k)). \end{aligned}$$

Note that the weight vector and threshold are unchanged if the correct decision is made by the TLU.

STEP 5 *Repeat by Going to Step 2.*

◆ Example 9.1

We used the above algorithm to generate a separating line for eight training pairs. For $i = 1, \dots, 4$ and $d_i = 1$,

$$[\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \mathbf{x}^{(3)} \quad \mathbf{x}^{(4)}] = \begin{bmatrix} -5 & -4 & 3 & 6 \\ -4 & 6 & 4 & 7 \end{bmatrix}.$$

The above points are labeled “ \times ” in Figure 9.4. The other four training pairs corresponding to $d_i = -1$, $i = 5, \dots, 8$, labeled “ \circ ” are

$$[\mathbf{x}^{(5)} \quad \mathbf{x}^{(6)} \quad \mathbf{x}^{(7)} \quad \mathbf{x}^{(8)}] = \begin{bmatrix} -6 & 0 & 3 & 8 \\ -8 & -1 & -4 & 3 \end{bmatrix}.$$

The weight vector and threshold were initialized randomly, where

$$\mathbf{w}(0) = [-0.4860 \quad 0.8913]^T \quad \text{and} \quad t_0 = 0.7621.$$

The TLU learning algorithm generated a separating line, shown in Figure 9.4, corresponding to the weight vector $\mathbf{w}^* = [-5.5140 \quad 6.8913]^T$ and threshold $t^* = -2.2379$.

Our goal now is to prove that the TLU learning algorithm can always find a separating hyperplane for linearly separable input vectors. In our proof, we use some technical results that we present next.

Lemma 9.1 Let $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(p)} \in \mathbb{R}^n$, be a set of p linearly separable input vectors, let $\mathbf{w} \in \mathbb{R}^n$ be a weight vector of an n -input TLU, and let $t \in \mathbb{R}$ be its threshold. Then, there exists a threshold t' such that for $i = 1, 2, \dots, p$ we have

$$\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t' \neq 0$$

and

$$\text{sign}(\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t) = \text{sign}(\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t').$$

Proof Assume, without loss of generality, that the input vectors are sorted into two sets $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(l)}\}$ and $\{\mathbf{x}^{(l+1)}, \dots, \mathbf{x}^{(p)}\}$ such that the desired response of the vectors from the first set is $d_i = -1$, while the desired response of the vectors from the second set is $d_i = 1$. Let

$$\delta = \min \{ |\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t| : 1 \leq i \leq l \}.$$

It follows from (9.2) that $\delta > 0$. Let $t' = t - \delta/2$. Then,

$$\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t' = \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t + \delta/2 \neq 0 \quad \text{for } i = 1, \dots, p,$$

and the output of the TLU with the modified threshold, t' , is the same as the output of the TLU with the old threshold, t .

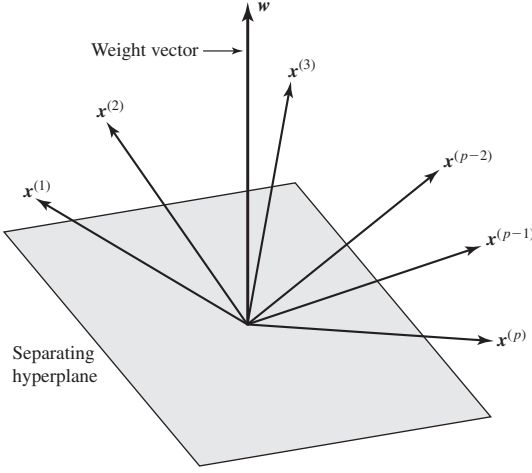


Figure 9.5 Linearly separable vectors after their preprocessing.

By the previous lemma, we can always assume that $\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t' \neq 0$ for $i = 1, \dots, p$ and that the output of a TLU can be defined as

$$y_i = \begin{cases} 1 & \text{if } \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t > 0 \\ -1 & \text{if } \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle - t < 0. \end{cases}$$

We further simplify the formulation of the TLU learning problem by preprocessing training pairs. Specifically, we replace every training pair, $(\mathbf{x}^{(i)}, d_i)$ with $d_i = -1$, with their negative, $(-\mathbf{x}^{(i)}, -d_i)$. After performing the above preprocessing, the TLU learning problem can be reformulated in a simple form that we use while analyzing the TLU learning algorithm.

Simplified TLU Learning Problem Given a set of input vectors $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}\}$. Determine a weight vector \mathbf{w} , if there exists one such that for $i = 1, \dots, p$ we have

$$\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle > 0.$$

In view of the above, for a linearly separable set of vectors, $\{\mathbf{x}^{(i)}\}$, after preprocessing there exists a hyperplane passing through the origin in \mathbb{R}^{n+1} such that all the vectors lie on one side of the hyperplane. The normal to the hyperplane is the weight vector \mathbf{w} . This is illustrated in Figure 9.5. Taking into account the Simplified TLU learning problem, we can state the TLU learning algorithm equivalently as follows:

SIMPLIFIED TLU LEARNING ALGORITHM

STEP 1 *Initialize Weight Vector.* Set the values of the components of the initial weight vector $\mathbf{w}(0)$ to zero.

STEP 2 *Test.* If $\langle \mathbf{w}(k), \mathbf{x}^{(i)} \rangle \leq 0$, go to Add, else Test next input vector.

STEP 3 *Add*

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \mathbf{x}^{(i)}, \quad k = k + 1.$$

Go to *Test*.

The algorithm is said to *converge* if for some $k < \infty$ we have

$$\langle \mathbf{w}(k), \mathbf{x}^{(i)} \rangle > 0 \quad \text{for } i = 1, \dots, p.$$

The following theorem states that if the set of input vectors is linearly separable, then the TLU learning algorithm converges.

Theorem 9.1 (The TLU Convergence Theorem) If the set of vectors $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p)}\}$ is linearly separable, then the simplified TLU learning algorithm determines a weight vector \mathbf{w} in a finite number of steps, such that $\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle > 0$ for $i = 1, 2, \dots, p$.

Proof The proof is adapted from Siu, Roychowdhury, and Kailath [260]. We begin by observing that because the set of input vectors is linearly separable, there exists a weight vector \mathbf{w}^* and a constant $\delta > 0$ such that for $i = 1, 2, \dots, p$, we have

$$\langle \mathbf{w}^*, \mathbf{x}^{(i)} \rangle > \delta.$$

To simplify further manipulations, we assume that $\|\mathbf{w}^*\| = 1$. We also assume that $\|\mathbf{x}^{(i)}\| \leq L$, for $i = 1, 2, \dots, p$; that is, the lengths of all the input vectors are bounded above by L . Let k denote the number of times the weight vector \mathbf{w} in the simplified TLU learning algorithm has been updated, and let $\mathbf{w}(k)$ be the value of the weight vector, \mathbf{w} , after the k th update. Our goal is to show that k is bounded above. To proceed, we note that if an input vector $\mathbf{x}^{(i)}$ is used for the k th update in the algorithm, then

$$\mathbf{w}(k) = \mathbf{w}(k-1) + \mathbf{x}^{(i)},$$

where $\langle \mathbf{w}(k-1), \mathbf{x}^{(i)} \rangle = \mathbf{w}^T(k-1) \mathbf{x}^{(i)} \leq 0$. We first derive an upper bound on $\|\mathbf{w}(k)\|$. We have

$$\begin{aligned} \|\mathbf{w}(k)\|^2 &= \|\mathbf{w}(k-1) + \mathbf{x}^{(i)}\|^2 \\ &= (\mathbf{w}(k-1) + \mathbf{x}^{(i)})^T (\mathbf{w}(k-1) + \mathbf{x}^{(i)}) \\ &= \mathbf{w}^T(k-1) \mathbf{w}(k-1) + \mathbf{x}^{(i)T} \mathbf{x}^{(i)} + 2\mathbf{w}^T(k-1) \mathbf{x}^{(i)} \\ &\leq \|\mathbf{w}(k-1)\|^2 + L^2, \end{aligned}$$

because $\mathbf{w}^T(k-1) \mathbf{x}^{(i)} \leq 0$ and $\|\mathbf{x}^{(i)}\|^2 \leq L^2$. Suppose that for the $(k-1)$ st update, the input vector $\mathbf{x}^{(j)}$ was used, that is,

$$\mathbf{w}(k-1) = \mathbf{w}(k-2) + \mathbf{x}^{(j)}.$$

Therefore,

$$\begin{aligned} \|\mathbf{w}(k)\|^2 &\leq \|\mathbf{w}(k-1)\|^2 + L^2 \\ &= \|\mathbf{w}(k-2) + \mathbf{x}^{(j)}\|^2 + L^2 \\ &\leq \|\mathbf{w}(k-2)\|^2 + 2L^2. \end{aligned}$$

Continuing in this manner, we obtain

$$\|\mathbf{w}(k)\|^2 \leq kL^2, \tag{9.7}$$

because by assumption, $\mathbf{w}(0) = \mathbf{0}$.

We next derive a lower bound on $\|\mathbf{w}(k)\|$. We have

$$\begin{aligned} \langle \mathbf{w}^*, \mathbf{w}(k) \rangle &= \langle \mathbf{w}^*, \mathbf{w}(k-1) \rangle + \langle \mathbf{w}^*, \mathbf{x}^{(i)} \rangle \\ &\geq \langle \mathbf{w}^*, \mathbf{w}(k-1) \rangle + \delta, \end{aligned}$$

because we assumed that $\langle \mathbf{w}^*, \mathbf{x}^{(j)} \rangle > \delta$. Observing that $\mathbf{w}(k-1) = \mathbf{w}(k-2) + \mathbf{x}^{(j)}$, we get

$$\begin{aligned}\langle \mathbf{w}^*, \mathbf{w}(k) \rangle &\geq \langle \mathbf{w}^*, \mathbf{w}(k-1) \rangle + \delta \\ &\geq \langle \mathbf{w}^*, \mathbf{w}(k-2) \rangle + 2\delta.\end{aligned}$$

Continuing in this manner, we obtain

$$\langle \mathbf{w}^*, \mathbf{w}(k) \rangle \geq k\delta. \quad (9.8)$$

Applying to the above the Cauchy–Schwarz inequality yields

$$\|\mathbf{w}^*\| \|\mathbf{w}(k)\| \geq \langle \mathbf{w}^*, \mathbf{w}(k) \rangle \geq k\delta.$$

By assumption, $\|\mathbf{w}^*\| = 1$. Hence,

$$\|\mathbf{w}(k)\| \geq k\delta.$$

Combining the obtained upper and lower bounds on $\|\mathbf{w}(k)\|^2$, given by (9.7) and (9.8), gives

$$k^2\delta^2 \leq \|\mathbf{w}(k)\|^2 \leq kL^2.$$

Therefore,

$$k \leq \frac{L^2}{\delta^2}, \quad (9.9)$$

which implies that the TLU learning algorithm converges. The proof of the TLU convergence theorem is complete.

We now present a sufficient condition for linear nonseparability.

Theorem 9.2 If there exists a positive linear combination (PLC) that equals $\mathbf{0}$ of the given set of input vectors, then the given set of input vectors is linearly nonseparable.

Proof We prove the theorem by contradiction. Assume that there exists a PLC of the given input vectors that is equal to $\mathbf{0}$, that is, there exist nonnegative constants, q_i , such that for some j , $1 \leq j \leq p$, we have $q_j > 0$, and

$$\sum_{i=1}^p q_i \mathbf{x}^{(i)} = \mathbf{0}. \quad (9.10)$$

Suppose now that the given input vectors are linearly separable; that is, there is a weight vector, \mathbf{w} , such that $\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle > 0$ for $i = 1, \dots, p$. We premultiply (9.10) by \mathbf{w}^T to get

$$\mathbf{w}^T \sum_{i=1}^p q_i \mathbf{x}^{(i)} = \sum_{i=1}^p q_i \langle \mathbf{w}, \mathbf{x}^{(i)} \rangle > 0$$

because $\langle \mathbf{w}, \mathbf{x}^{(i)} \rangle > 0$ and $q_j > 0$. However, we arrived at a contradiction because by assumption we have $\sum_{i=1}^p q_i \mathbf{x}^{(i)} = \mathbf{0}$, which completes the proof.

◆ Example 9.2

Suppose we are given the following training set:

$$\begin{aligned} \mathbf{x}^{(1)} &= \begin{bmatrix} -1 \\ -1 \end{bmatrix}, & d_1 &= -1; & \mathbf{x}^{(2)} &= \begin{bmatrix} 1 \\ 1 \end{bmatrix}, & d_2 &= -1; \\ \mathbf{x}^{(3)} &= \begin{bmatrix} -1 \\ 1 \end{bmatrix}, & d_3 &= 1; & \mathbf{x}^{(4)} &= \begin{bmatrix} 1 \\ -1 \end{bmatrix}, & d_4 &= 1. \end{aligned}$$

This training set represents the two-input Exclusive-OR (XOR) function. Let us now preprocess the training set so we can apply Theorem 9.2. By preprocessing we mean eliminating the threshold by augmenting the input vectors with -1 and then negating the input vectors for which the desired output is -1 . The preprocessed input vectors are

$$\tilde{\mathbf{x}}^{(1)} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \quad \tilde{\mathbf{x}}^{(2)} = \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix}, \quad \tilde{\mathbf{x}}^{(3)} = \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix}, \quad \tilde{\mathbf{x}}^{(4)} = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix}.$$

Note that $\sum_{i=1}^4 \tilde{\mathbf{x}}^{(i)} = \mathbf{0}$. It follows from Theorem 9.2 that the given input vectors are not linearly separable.

9.2 Identification Using Adaptive Linear Element

The adaptive linear element, also called the *adaline*, is shown in Figure 9.6. We first explain the symbols used in our discussion of its operation. Let $\mathbf{x}(t) = [x_1(t) \ \cdots \ x_n(t)]^T$ denote the adaline input. We assume that for all t

$$\|\mathbf{x}(t)\| = \sqrt{x_1^2(t) + \cdots + x_n^2(t)} \leq B_x,$$

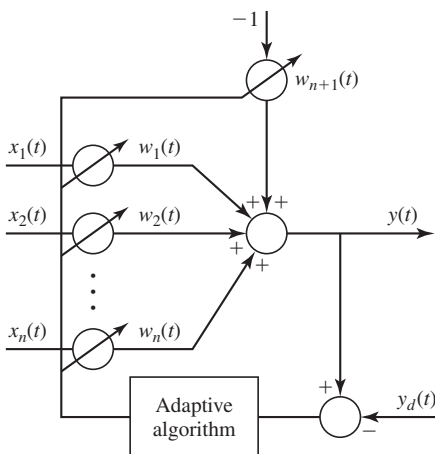


Figure 9.6 Adaline.

where $B_x > 0$ is known positive constant and $\|\cdot\| = \|\cdot\|_2$; that is, whenever we write $\|\cdot\|$ we mean the Euclidean norm $\|\cdot\|_2$. We also assume that for all t

$$\|\dot{\mathbf{x}}(t)\| = \sqrt{\dot{x}_1^2(t) + \cdots + \dot{x}_n^2(t)} \leq B_{\dot{x}}.$$

The weight vector $\mathbf{w}(t) = [w_1(t) \cdots w_n(t)]^T$. The bias weight $w_{n+1}(t)$ is connected to a constant input equal to -1 . We define $\tilde{\mathbf{x}}(t)$ to be an augmented input vector that includes the constant input to the bias weight, that is,

$$\tilde{\mathbf{x}}(t) = [x_1(t) \cdots x_n(t) \quad x_{n+1}(t)]^T.$$

Note that

$$\|\tilde{\mathbf{x}}(t)\|^2 = \tilde{\mathbf{x}}^T(t)\tilde{\mathbf{x}}(t) = \|\mathbf{x}(t)\|^2 + 1,$$

and hence $\tilde{\mathbf{x}}(t)$ is bounded away from zero at all times. The augmented vector $\tilde{\mathbf{x}}(t)$ is bounded, for all t , if and only if $\mathbf{x}(t)$ is bounded. We denote a bound of $\tilde{\mathbf{x}}(t)$ by $B_{\tilde{x}}$; that is, $\|\tilde{\mathbf{x}}(t)\| \leq B_{\tilde{x}}$ for all t . Observe that $\|\dot{\tilde{\mathbf{x}}}(t)\| = \|\dot{\mathbf{x}}(t)\|$. Thus, we can take $B_{\dot{\tilde{x}}} = B_{\dot{x}}$. Similarly, we define an augmented weight vector,

$$\tilde{\mathbf{w}}(t) = [w_1(t) \cdots w_n(t) \quad w_{n+1}(t)]^T.$$

We assume that for all t

$$\|\tilde{\mathbf{w}}(t)\| \leq B_{\tilde{w}}$$

for some constant $B_{\tilde{w}} > 0$. The output, $y(t)$, of the adaline is a linear combination of its inputs, that is,

$$y(t) = \tilde{\mathbf{x}}^T(t)\tilde{\mathbf{w}}(t).$$

Let $y_d(t)$ be the desired output of the adaline and assume that

$$|y_d(t)| \leq B_y \quad \text{and} \quad |\dot{y}_d(t)| \leq B_{\dot{y}}$$

for some positive constants B_y and $B_{\dot{y}}$. Let

$$e(t) = y(t) - y_d(t)$$

denote the error signal. The error dynamics are described by

$$\begin{aligned} \dot{e} &= \dot{y} - \dot{y}_d \\ &= \dot{\tilde{\mathbf{w}}}^T \tilde{\mathbf{x}} + \tilde{\mathbf{w}}^T \dot{\tilde{\mathbf{x}}} - \dot{y}_d. \end{aligned} \tag{9.11}$$

Our objective is to devise an adaptation algorithm for the weights $w_i(t)$, $i = 1, 2, \dots, n+1$, so that the error $e(t)$ converges to zero in finite time, say t_h , and $e(t) = 0$ for $t > t_h$. Sira-Ramírez and Colina-Morles [258] showed that we can achieve our objective if we adapt the weights according to the adaptation law,

$$\boxed{\dot{\tilde{\mathbf{w}}} = -\eta \frac{\tilde{\mathbf{x}}}{\tilde{\mathbf{x}}^T \tilde{\mathbf{x}}} \text{sign}(e(t))} \tag{9.12}$$

where the gain $\eta > 0$ is sufficiently large. We now state and prove the theorem of Sira-Ramírez and Colina-Morles [258].

Theorem 9.3 If the weight vector $\tilde{\mathbf{w}}$ is adapted according to (9.12), where

$$\eta > B_{\tilde{\mathbf{w}}} B_{\dot{\mathbf{x}}} + B_{\dot{\mathbf{y}}},$$

then for arbitrary initial error $e(0)$ the error signal $e(t)$ converges to zero in finite time, say t_h , such that

$$t_h \leq \frac{|e(0)|}{\eta - B_{\tilde{\mathbf{w}}} B_{\dot{\mathbf{x}}} - B_{\dot{\mathbf{y}}}}$$

and $e(t) = 0$ for $t > t_h$.

Proof Recall that the error dynamics are governed by (9.11), where the vector $\tilde{\mathbf{w}}(t)$ is adapted according to (9.12). Consider the following Lyapunov function candidate

$$V(e) = \frac{1}{2} e^2(t).$$

The time derivative of V evaluated on the solutions of (9.11) is

$$\begin{aligned} \dot{V}(e) &= e \dot{e} \\ &= e(\tilde{\mathbf{w}}^T \dot{\tilde{\mathbf{x}}} + \tilde{\mathbf{w}}^T \dot{\tilde{\mathbf{x}}} - \dot{y}_d). \end{aligned} \quad (9.13)$$

Substituting (9.12) into the above and taking into account that

$$\text{sign}(e) = \frac{e}{|e|} = \frac{|e|}{e},$$

we obtain

$$\begin{aligned} \dot{V} &= -\eta|e| + e(\tilde{\mathbf{w}}^T \dot{\tilde{\mathbf{x}}} - \dot{y}_d) \\ &\leq -\eta|e| + |e|(B_{\tilde{\mathbf{w}}} B_{\dot{\mathbf{x}}} + B_{\dot{\mathbf{y}}}) \\ &= (-\eta + B_{\tilde{\mathbf{w}}} B_{\dot{\mathbf{x}}} + B_{\dot{\mathbf{y}}})|e| \\ &< 0 \quad \text{for all } e \neq 0 \quad \text{and} \quad \eta > B_{\tilde{\mathbf{w}}} B_{\dot{\mathbf{x}}} + B_{\dot{\mathbf{y}}}, \end{aligned} \quad (9.14)$$

which means that $e(t) \rightarrow 0$ as t increases. It follows from our discussion above that we can take $B_{\dot{\tilde{\mathbf{x}}}} = B_{\dot{\mathbf{x}}}$. We now show that for some finite time t_h , we have $e(t_h) = 0$ and $e(t) = 0$ for $t > t_h$. Substitute (9.12) into (9.11) to get

$$\dot{e} = -\eta \text{sign}(e) + \tilde{\mathbf{w}}^T \dot{\tilde{\mathbf{x}}} - \dot{y}_d. \quad (9.15)$$

Integrating the above from $t = 0$ to $t \leq t_h$ gives

$$e(t) - e(0) = -\eta t \text{sign}(e(0)) + \int_0^t (\tilde{\mathbf{w}}^T(\tau) \dot{\tilde{\mathbf{x}}}(\tau) - \dot{y}_d(\tau)) d\tau. \quad (9.16)$$

At $t = t_h$ we have $e(t_h) = 0$, and therefore for $t = t_h$ the above equation takes the form

$$-e(0) = -\eta t_h \text{sign}(e(0)) + \int_0^{t_h} (\tilde{\mathbf{w}}^T(t) \dot{\tilde{\mathbf{x}}}(t) - \dot{y}_d(t)) dt. \quad (9.17)$$

Multiplying both sides of the above by $-\text{sign}(e(0))$ yields

$$\begin{aligned} |e(0)| &= \eta t_h - \text{sign}(e(0)) \int_0^{t_h} (\tilde{\mathbf{w}}^T(t) \dot{\tilde{\mathbf{x}}}(t) - \dot{y}_d(t)) dt \\ &\geq (\eta - (B_{\tilde{\mathbf{w}}} B_{\dot{\mathbf{x}}} + B_{\dot{\mathbf{y}}})) t_h. \end{aligned} \quad (9.18)$$

Thus, the time t_h at which $e(t)$ reaches zero value for the first time satisfies the inequality

$$t_h \leq \frac{|e(0)|}{\eta - (B_{\tilde{w}} B_{\tilde{x}} + B_{\dot{y}})}.$$

Because for $e(t) \neq 0$ we have $\frac{d}{dt} e^2(t) < 0$, we are guaranteed that we also have $e(t) = 0$ for $t > t_h$. The proof of the theorem is now complete.

If the time derivative $\dot{\tilde{x}}(t)$ is available, then we can use it in an adaptation law of \tilde{w} . Using $\dot{\tilde{x}}(t)$ in an adaptation law allows us to reduce the gain η and the bound on time t_h when the error $e(t)$ becomes zero. One such adaptation strategy that includes $\dot{\tilde{x}}(t)$ was proposed by Sira-Ramírez and Colina-Morles [258], which we present in the following theorem.

Theorem 9.4 If the adaptation law for the augmented weight vector is chosen as

$$\dot{\tilde{w}}(t) = -\frac{\tilde{x}(t) \dot{\tilde{x}}^T(t)}{\tilde{x}^T(t) \tilde{x}(t)} \tilde{w}(t) - \eta \frac{\tilde{x}(t)}{\tilde{x}^T(t) \tilde{x}(t)} \text{sign}(e(t)) \quad (9.19)$$

where $\eta > B_{\dot{y}}$, then for any given initial error $e(0)$ the error $e(t)$ converges to zero in finite time t_h such that

$$t_h \leq \frac{|e(0)|}{\eta - B_{\dot{y}}} \quad (9.20)$$

and $e(t) = 0$ for $t > t_h$.

Proof The proof of this theorem is analogous to the proof of Theorem 9.3. The only difference is that (9.15) becomes

$$\dot{e} = -\eta \text{sign}(e) - \dot{y}_d. \quad (9.21)$$

When using the adaline for control purposes, the signal to be processed by the adaline is first passed through a prefilter as shown in Figure 9.7. The prefilter is described by the equation $\dot{x} = Ax + br$, where r is the signal to be processed. The matrix A of the prefilter is chosen to be asymptotically stable. Recall that the inputs to the adaline are the components of the augmented

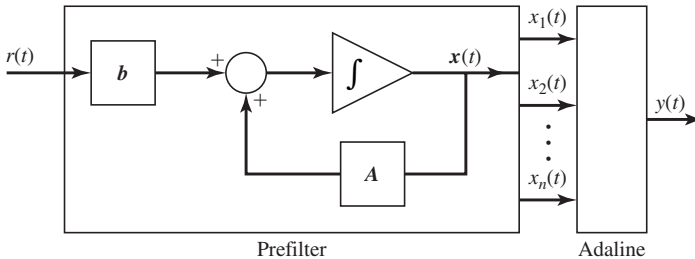


Figure 9.7 Prefilter connection to adaline.

vector $\tilde{\mathbf{x}}$, where the $(n + 1)$ st component is equal to -1 . Thus, $\dot{x}_{n+1} = 0$. Taking this into account, the prefiltered inputs to the adaline can be modeled as

$$\dot{\tilde{\mathbf{x}}} = \tilde{\mathbf{A}}\tilde{\mathbf{x}} + \tilde{\mathbf{b}}r, \quad (9.22)$$

where

$$\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0}^T & 0 \end{bmatrix} \quad \text{and} \quad \tilde{\mathbf{b}} = \begin{bmatrix} \mathbf{b} \\ 0 \end{bmatrix}.$$

Note that $\dot{\tilde{\mathbf{x}}}$ is available to us in the architecture shown in Figure 9.7. Thus, we can employ the adaptation strategy given by (9.19), which results in the following description of the structure depicted in Figure 9.7:

$$\begin{aligned} \dot{\tilde{\mathbf{x}}} &= \tilde{\mathbf{A}}\tilde{\mathbf{x}} + \tilde{\mathbf{b}}r, \\ \dot{\tilde{\mathbf{w}}} &= -\frac{\tilde{\mathbf{x}}\dot{\tilde{\mathbf{x}}}^T}{\tilde{\mathbf{x}}^T\tilde{\mathbf{x}}}\tilde{\mathbf{w}} - \eta\frac{\tilde{\mathbf{x}}}{\tilde{\mathbf{x}}^T\tilde{\mathbf{x}}}\text{sign}(e), \\ y &= \tilde{\mathbf{x}}^T\tilde{\mathbf{w}} \end{aligned} \quad (9.23)$$

We now demonstrate, with a numerical example, how the adaline system modeled by (9.23) can be applied to forward and inverse dynamics identification of dynamical systems. The nonlinear plant model is taken from Sira-Ramírez and Colina-Morles [258]. This plant model has the form

$$\begin{aligned} \dot{\alpha} &= p + \frac{\sin(\alpha)}{l}u, \\ \dot{p} &= \left(\frac{g}{l} - \frac{u^2}{l^2} \cos(\alpha) \right) - \frac{u}{l}p \cos(\alpha), \\ \dot{z} &= u, \\ y_p &= \alpha. \end{aligned} \quad (9.24)$$

The parameters g and l have the values: $g = 9.81 \text{ m/sec}^2$ and $l = 0.7 \text{ m}$. The plant input is labeled as u while the plant output is y_p . We first concern ourselves with forward dynamics identification of the above plant. A block diagram of a forward dynamics identification scheme is depicted in Figure 9.8. The input to the plant is

$$u(t) = A_1 + A_2 \cos\left(\frac{t}{2\epsilon}\right) + A_3 \sin\left(\frac{t}{\epsilon}\right),$$

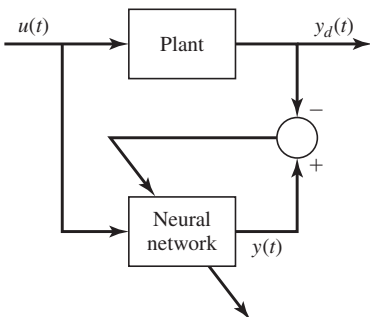


Figure 9.8 Forward dynamics identification scheme.

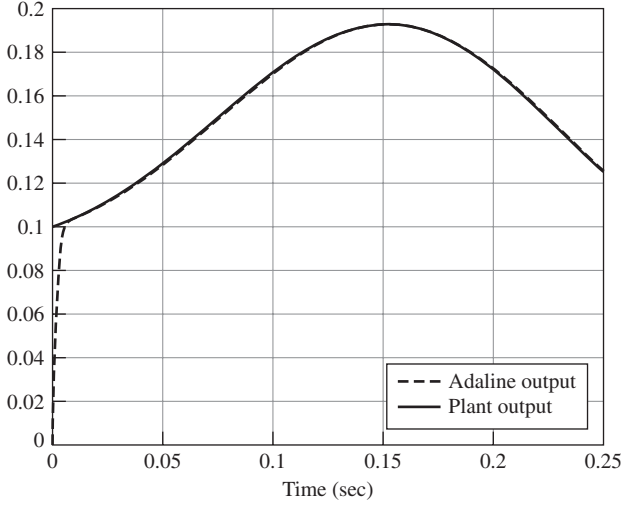


Figure 9.9 Simulation results of forward dynamics identification of (9.24).

where $A_1 = 0.4$, $A_2 = 2$, $A_3 = 3$, and $\epsilon = 0.05$. The prefilter equation (9.22) in our example has the form

$$\dot{\tilde{\mathbf{x}}} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ -1 & -3 & -3 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \tilde{\mathbf{x}} + \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} u.$$

Thus, the adaline has total four weights, including the bias weight. In our simulations, we use the following approximation of the sign function:

$$\text{sign}(e(t)) \approx \frac{e(t)}{|e(t)| + \delta}, \quad (9.25)$$

where $\delta = 0.05$. We used the above approximation to eliminate chattering and speed up the simulations. We used $\eta = 10$. The performance of the adaline as a forward dynamics identifier of the dynamical system modeled by (9.24) is shown in Figure 9.9. We used the following initial conditions:

$$\tilde{\mathbf{x}}(0) = [0.1 \quad 0 \quad 0 \quad -1]^T, \quad \tilde{\mathbf{w}}(0) = [0.1 \quad 0 \quad 0 \quad 0]^T,$$

and $[\alpha(0) \ p(0) \ z(0)] = [0.1 \ 0 \ 0]$. Increasing the gain η results in more accurate identification; that is, the error $e(t)$ goes to zero faster. Furthermore, the effects of using the approximation (9.25) are reduced for higher gain η .

In our next experiment, we employed the adaline as an inverse dynamics identifier using the scheme shown in Figure 9.10. The results of our simulation, using the same initial conditions as before, are depicted in Figure 9.11. Here we used $\eta = 90$. Increasing the gain k decreases the identification error as well as reduces the effects of the approximation of the sign function.

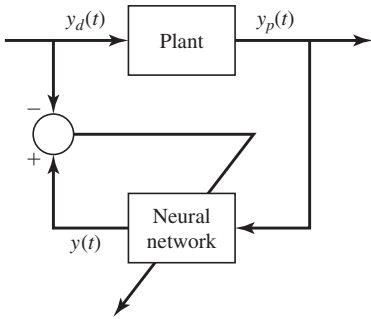


Figure 9.10 Inverse dynamics identifier scheme.

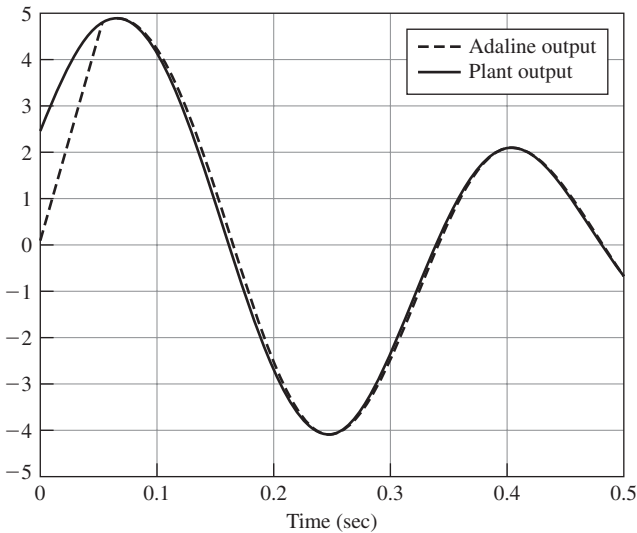


Figure 9.11 Simulation results of inverse dynamics identification of (9.24).

9.3 Backpropagation

In Section 9.1 we discussed an iterative learning algorithm for a threshold logic unit. This computational unit kept on iterating until it reached a solution, if a solution existed. In this section we present iterative learning algorithms for computational units with differentiable activation functions. This family of algorithms are called the backpropagation algorithms because they backpropagate the error between the desired and actual output of the neural network for the purpose of adjusting the network's weights in order to reduce the learning error. Specifically, a typical backpropagation algorithm consists of two passes: In the forward pass the input vectors are presented to the network, and the actual outputs of the network are calculated. Then, the error between the actual and desired outputs are calculated. In the reverse pass the error gradient with respect to the network's weights is calculated by propagating the error backwards through the network. Once the error gradient is calculated, the weights are adjusted using, for example, a descent gradient method.

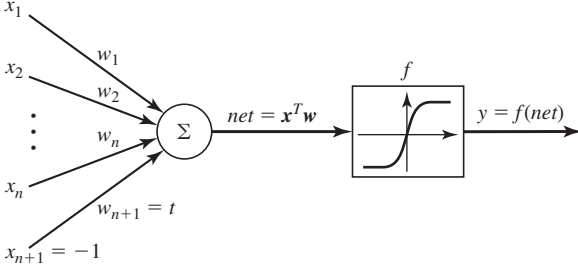


Figure 9.12 A single perceptron with a differentiable sigmoidal activation function.

We now analyze a single computational unit shown in Figure 9.12. It is also called the *perceptron*. We consider a simple case, where we train the perceptron to learn the training pair

$$(\mathbf{x}, d),$$

where $\mathbf{x} \in \mathbb{R}^{n+1}$, $x_{n+1} = (\mathbf{x})_{n+1} = -1$, and $d \in [-1, 1]$. Define the error function

$$\begin{aligned} E(\mathbf{w}) &= \frac{1}{2}(d - y)^2 \\ &= \frac{1}{2}(d - f(\text{net}))^2, \end{aligned} \quad (9.26)$$

where $\text{net} = \mathbf{x}^T \mathbf{w}$. Note that $w_{n+1} = t$, where t is a threshold, bias. The gradient of E , with respect to \mathbf{w} , is

$$\nabla E(\mathbf{w}) = -(d - f(\text{net}))f'(\text{net})\mathbf{x}. \quad (9.27)$$

Using the descent gradient algorithm with a fixed step size, $\alpha > 0$, to minimize E , we obtain

$$\begin{aligned} \mathbf{w}(k+1) &= \mathbf{w}(k) - \alpha \nabla E(\mathbf{w}(k)) \\ &= \mathbf{w}(k) + \alpha(d_k - f(\text{net}_k))f'(\text{net}_k)\mathbf{x}^{(k)}, \end{aligned} \quad (9.28)$$

where $\text{net}_k = \mathbf{x}^T(k)\mathbf{w}(k)$. Let

$$\delta^k = (d_k - f(\text{net}_k))f'(\text{net}_k). \quad (9.29)$$

Then, we can represent (9.28) as

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \alpha \delta^k \mathbf{x}(k) \quad (9.30)$$

We now generalize the above formula for a three-layer feedforward neural network shown in Figure 9.13. We wish to train the network to learn the training pairs

$$(\mathbf{x}^{(1)}, \mathbf{d}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{d}^{(2)}), \dots, (\mathbf{x}^{(p)}, \mathbf{d}^{(p)}),$$

where $\mathbf{x}^{(i)} \in \mathbb{R}^n$ and $\mathbf{d}^{(i)} \in \mathbb{R}^m$. The network has n inputs, which is the dimension of the training vectors, and m outputs. The network's inputs are connected to a layer of hidden units, which in turn is connected to a layer of output units. To construct a network to perform some task, we must set the weights of the connections appropriately. The weights specify the strength of the influence of one unit upon another. We will derive rules for adjusting the weights w_{ji}^h and w_{kj}^o such that the error function

$$E(\mathbf{w}) = \frac{1}{2} \sum_{s=1}^m (d_s - y_s)^2$$

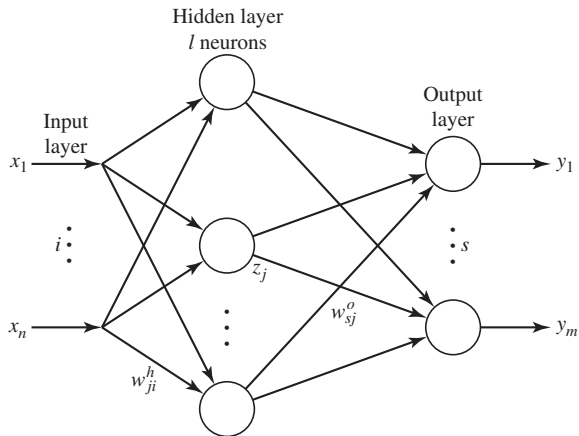


Figure 9.13 A three-layer feedforward neural network.

is minimized. The minimization is performed over the space of weights

$$\{w_{ji}^h, w_{sj}^o : i = 1, 2, \dots, n, j = 1, 2, \dots, l, s = 1, 2, \dots, m\}.$$

We first give formulas for updating the weights connecting the hidden and output layers. Let

$$net_s = \sum_{j=1}^l w_{sj}^o z_j.$$

Using the descent gradient algorithm in a similar way as in (9.30), we obtain

$$\begin{aligned} w_{sj}^o(k+1) &= w_{sj}^o(k) - \alpha_o \frac{\partial E(\mathbf{w}(k))}{\partial w_{sj}^o(k)} \\ &= w_{sj}^o(k) + \alpha_o (d_s - y_s) f'_o(net_s) z_j(k), \end{aligned} \quad (9.31)$$

where f'_o is the derivative of $f_o(net_s)$ with respect to net_s , and the values of z_j are computed by feeding into the network the input vector \mathbf{x} and then propagating it forward through the network. We have

$$z_j = f_h \left(\sum_{i=1}^n w_{ji}^h x_i \right) = f_h(net_j), \quad j = 1, 2, \dots, l.$$

To derive the update formulas for the weights in the hidden layer, we use the chain rule to compute $\partial E / \partial w_{ji}^h$, and then we apply the gradient descent algorithm:

$$\frac{\partial E}{\partial w_{ji}^h} = \frac{\partial E}{\partial z_j} \frac{\partial z_j}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}^h}, \quad (9.32)$$

where

$$\frac{\partial net_j}{\partial w_{ji}^h} = x_i, \quad (9.33)$$

$$\frac{\partial z_j}{\partial net_j} = f'_h(net_j), \quad (9.34)$$

and

$$\begin{aligned}
 \frac{\partial E}{\partial z_j} &= \frac{\partial}{\partial z_j} \left(\frac{1}{2} \sum_{s=1}^m (d_s - f_o(net_s))^2 \right) \\
 &= - \sum_{s=1}^m (d_s - f_o(net_s)) \frac{\partial f_o(net_s)}{\partial z_j} \\
 &= - \sum_{s=1}^m (d_s - y_s) f'_o(net_s) w_{sj}^o.
 \end{aligned} \tag{9.35}$$

Combining (9.32)–(9.35), we obtain rules for updating the hidden weights,

$$\begin{aligned}
 w_{ji}^h(k+1) &= w_{ji}^h(k) - \alpha_h \frac{\partial E(\mathbf{w}(k))}{\partial w_{ji}^h(k)} \\
 &= w_{ji}^h(k) + \alpha_h \left(\sum_{s=1}^m (d_s - y_s) f'_o(net_s) w_{sj}^o(k) \right) f'_h(net_j) x_i.
 \end{aligned} \tag{9.36}$$

A commonly used activation function is the sigmoid,

$$f(v) = \frac{1}{1 + e^{-v}}. \tag{9.37}$$

It is easy to check that the activation function (9.37) has the useful property,

$$\frac{df(v)}{dv} = f(v)(1 - f(v)). \tag{9.38}$$

We summarize the above derivations in the following algorithm.

BACKPROPAGATION TRAINING ALGORITHM

Given a set of training pairs $\{(\mathbf{x}^{(1)}, \mathbf{d}^{(1)}), \dots, (\mathbf{x}^{(i)}, \mathbf{d}^{(i)}), \dots, (\mathbf{x}^{(p)}, \mathbf{d}^{(p)})\}$.

STEP 1 *Initialize Weights and Thresholds.* Initialize the weights and thresholds to small random values.

STEP 2 *Present New Input Vector and Desired Output.* Present the new input vector, $\mathbf{x}^{(i)}$, and the corresponding desired output, $\mathbf{d}^{(i)}$. Calculate actual outputs.

STEP 3 *Calculate Error Gradients.*

STEP 4 *Adapt Weight Vectors and Thresholds:*

$$\begin{aligned}
 \mathbf{w}(k+1) &= \mathbf{w}(k) - \alpha \nabla_{\mathbf{w}} E(\mathbf{w}(k), \mathbf{t}(k)), \\
 \mathbf{t}(k+1) &= \mathbf{t}(k) - \alpha \nabla_{\mathbf{t}} E(\mathbf{w}(k), \mathbf{t}(k)).
 \end{aligned}$$

STEP 5 *Repeat by Returning to Step 2.*

◆ Example 9.3

We consider the two-input XOR problem. In Example 9.2 we showed that this logic function cannot be implemented using a single TLU. Here, we use a two-layer neural network with two nodes in the hidden layer and one node in the

output layer. In this example, we use the unipolar representation of the two-input XOR function; that is, the elements of the training pairs are from the set $\{0, 1\}$. Therefore, the input vectors and the corresponding desired outputs have the form

$$\begin{aligned}\mathbf{x}^{(1)} &= \begin{bmatrix} 0 \\ 0 \end{bmatrix}, & d_1 &= 0; & \mathbf{x}^{(2)} &= \begin{bmatrix} 1 \\ 1 \end{bmatrix}, & d_2 &= 0; \\ \mathbf{x}^{(3)} &= \begin{bmatrix} 0 \\ 1 \end{bmatrix}, & d_3 &= 1; & \mathbf{x}^{(4)} &= \begin{bmatrix} 1 \\ 0 \end{bmatrix}, & d_4 &= 1.\end{aligned}$$

The error function for this example takes the form

$$E = \frac{1}{2}(d_i - y_i)^2.$$

The activation function for all three neurons is the sigmoid given by (9.37). The outputs of the hidden neurons are

$$z_1 = f(w_{11}^h x_1^{(i)} + w_{12}^h x_2^{(i)} - t_1), \quad (9.39)$$

$$z_2 = f(w_{21}^h x_1^{(i)} + w_{22}^h x_2^{(i)} - t_2). \quad (9.40)$$

The network output is

$$y = f(w_{11}^o z_1 + w_{12}^o z_2 - t_o). \quad (9.41)$$

The threshold parameters are treated as learning weights. Thus, the above simple network requires updating of nine weights, where, in particular,

$$\begin{aligned}\frac{\partial E}{\partial t_1} &= (d_i - y) y' w_{11} z_1' = (d_i - y) y(1 - y) w_{11}^o z_1(1 - z_1), \\ \frac{\partial E}{\partial t_2} &= (d_i - y) y' w_{12} z_2' = (d_i - y) y(1 - y) w_{12}^o z_2(1 - z_2), \\ \frac{\partial E}{\partial t_o} &= (d_i - y) y' = (d_i - y) y(1 - y).\end{aligned}$$

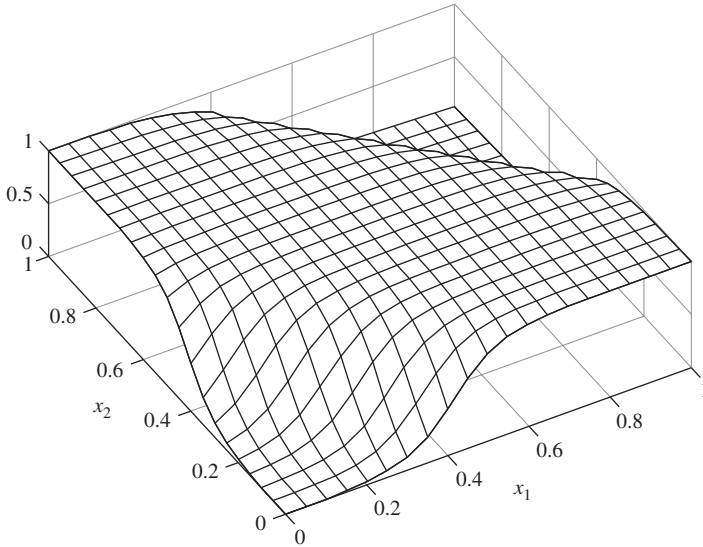
In calculating the above, we used (9.38). In our simulation the weights of the network were initialized to small random values in the interval $[-0.1, 0.1]$. The four patterns were presented in a mixed order. We used a fixed learning rate of 0.75 for all nine weights. After 50,000 iterations we obtained the following weight values:

$$\mathbf{W}^h = \begin{bmatrix} -4.9288 & -4.9334 \\ -6.7240 & -6.7526 \end{bmatrix}, \quad \mathbf{W}^o = [10.2447 \quad -10.7258],$$

and $t_1 = -7.3398$, $t_2 = -2.7853$, and $t_o = 4.8701$. In Table 9.1 we show the desired and actual responses of the trained network. In Figure 9.14 we show a plot of the function that the trained network implements.

Table 9.1 The Desired and Actual Responses of the Trained Network in Example 9.3

x_1	x_2	d_i	y_i
0	0	0	0.0087
1	1	0	0.0162
0	1	1	0.9870
1	0	1	0.9870

**Figure 9.14** The function implemented by the trained neural network to solve the two-input XOR problem in Example 9.3.

9.4 Neural Fuzzy Identifier

In this section we describe a method that can be used to construct identifiers of dynamical systems that in turn could be employed to construct neural-network-based fuzzy logic control strategies. The idea behind the method is to apply the backpropagation algorithm to a fuzzy logic system.

It follows from our discussion in Section 8.3 that a typical fuzzy logic system can be thought of as a nonlinear system described by

$$f(x_1, \dots, x_n) = f(\mathbf{x}) = \frac{\sum_{l=1}^M v_l \left(\prod_{i=1}^n \mu_{il}(x_i) \right)}{\sum_{l=1}^M \prod_{i=1}^n \mu_{il}(x_i)}, \quad (9.42)$$

where, for example,

$$\mu_{il} = \exp \left\{ -\frac{(x_i - m_{il})^2}{2\sigma_{il}^2} \right\},$$

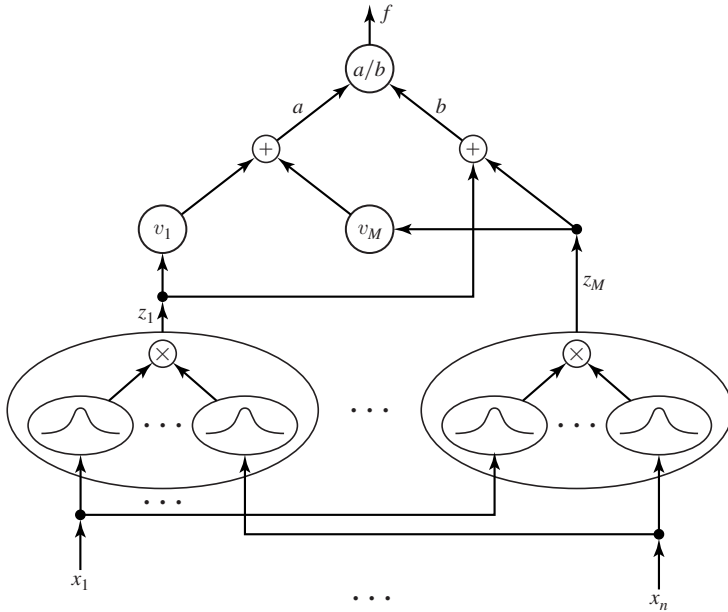


Figure 9.15 Fuzzy logic system (9.42) represented as a feedforward network.

M is the number of IF–THEN fuzzy rules, $\mathbf{x} = [x_1 \ x_2 \ \cdots \ x_n]^T$, and m_{il} and σ_{il} are the parameters of the Gaussian membership functions, while v_l is a weight parameter. These parameters will be determined using the backpropagation algorithm discussed in the previous section. To facilitate our discussion, we represent (9.42), following Wang [296, p. 169], as a feedforward network shown in Figure 9.15. The variables that characterize the feedforward network depicted in Figure 9.15 satisfy the following relations:

$$\begin{aligned}
 z_l &= \prod_{i=1}^n \exp \left\{ -\frac{(x_i - m_{il})^2}{2\sigma_{il}^2} \right\}, \quad l = 1, 2, \dots, M, \\
 b &= \sum_{l=1}^M z_l, \\
 a &= \sum_{l=1}^M v_l z_l, \\
 f &= a/b.
 \end{aligned} \tag{9.43}$$

We then use this feedforward network as a component of an identifier of dynamical systems. The identification scheme involving (9.42) is depicted in Figure 9.16. The adaptation algorithm for the identifier parameters will be derived using the backpropagation algorithm. The cost function to be minimized is

$$E = \frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}(y - f)^2,$$

where y is the output of the identified plant and f is the output of the identifier. We use a gradient algorithm to minimize E over all

$$\mathbf{w} = \{v_l, \sigma_{il}, m_{il}\}.$$

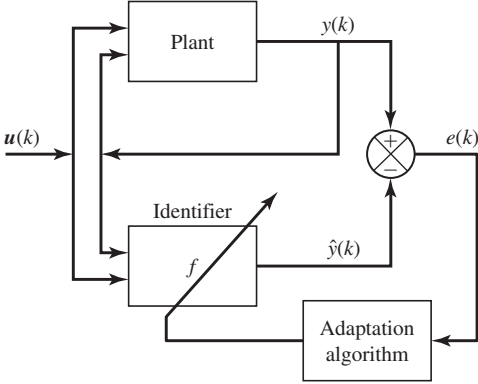


Figure 9.16 Neural fuzzy identification scheme of a nonlinear dynamical system.

The gradient algorithm that we use has the form

$$\mathbf{w}(k+1) = \mathbf{w}(k) - \alpha \nabla E(\mathbf{w}(k)), \quad (9.44)$$

where $\alpha > 0$ is a fixed step size. Note that (9.44) is just a simple descent gradient algorithm that we used in our discussion of the backpropagation algorithm.

In order to be able to use the above algorithm, we need to compute partial derivatives of E with respect to each component of \mathbf{w} . For this we first fix indices i , and l . We next compute the partial derivatives of E with respect to v_j . Then, the parameters v_j will be adapted accordingly as

$$v_j(k+1) = v_j(k) - \alpha \frac{\partial E(\mathbf{w}(k))}{\partial v_j}.$$

We note that E depends upon v_j only through a , where $f = a/b$. Hence, using the chain rule, we obtain

$$\begin{aligned} \frac{\partial E}{\partial v_j}(\mathbf{w}) &= \frac{\partial}{\partial v_j} \left(\frac{1}{2} (y - f)^2 \right) \\ &= \frac{\partial}{\partial v_j} \left(\frac{1}{2} (y - a/b)^2 \right) \\ &= -\frac{(y - f)}{b} \frac{\partial a}{\partial v_j} \\ &= \frac{f - y}{b} z_j. \end{aligned}$$

We next compute the partial derivative of E with respect to m_{st} . Observe that E depends upon m_{st} only through z_t . Using the chain rule again gives

$$\begin{aligned} \frac{\partial E}{\partial m_{st}}(\mathbf{w}) &= -(y - f) \frac{\partial f}{\partial z_t} \frac{\partial z_t}{\partial m_{st}} \\ &= -(y - f) \left(\frac{v_t b - a}{b^2} \right) z_t \frac{(x_s - m_{st})}{\sigma_{st}^2} \\ &= (f - y) \left(\frac{v_t - f}{b} \right) z_t \frac{(x_s - m_{st})}{\sigma_{st}^2}. \end{aligned}$$

Hence, the adaptation formula for m_{st} takes the form

$$m_{st}(k+1) = m_{st}(k) - \alpha(f - y) \left(\frac{v_t - f}{b} \right) z_t \frac{(x_s - m_{st}(k))}{\sigma_{st}^2(k)}.$$

To find the partial derivative of E with respect to σ_{st} , we note that E depends upon σ_{st} only through z_t . Applying the chain rule gives

$$\begin{aligned} \frac{\partial E}{\partial \sigma_{st}}(\mathbf{w}) &= -(y - f) \frac{\partial f}{\partial z_t} \frac{\partial z_t}{\partial \sigma_{st}} \\ &= (f - y) \left(\frac{v_t - f}{b} \right) z_t \frac{(x_i - m_{st})^2}{\sigma_{st}^3}. \end{aligned}$$

We are now ready to formulate an algorithm for updating on-line the weights of the neural fuzzy identifier.

ON-LINE IDENTIFICATION ALGORITHM

STEP 1 *Initialize the Identifier Parameters.* Set the values of v_l , m_{il} , and σ_{il} to small random values.

STEP 2 *Calculate the Output of the Identifier at the Time Instant k .* Compute

$$\begin{aligned} z_l &= \prod_{i=1}^n \exp\left(-\frac{(x_i - m_{il})^2}{2\sigma_{il}^2}\right), \\ b &= \sum_{l=1}^M z_l, \\ a &= \sum_{l=1}^M v_l z_l, \\ f &= a/b. \end{aligned}$$

STEP 3 *Calculate the Identification Error:*

$$e(k) = y(k) - f(k).$$

STEP 4 *Adapt the Identifier Parameters:*

$$\begin{aligned} v_l(k+1) &= v_l(k) + \alpha \frac{e(k)}{b(k)} z_l(k), \\ m_{il}(k+1) &= m_{il}(k) + \alpha \frac{e(k)}{b(k)} (v_l(k) - f(k)) z_l(k) \frac{(x_i(k) - m_{il}(k))}{\sigma_{il}^2(k)}, \\ \sigma_{il}(k+1) &= \sigma_{il}(k) + \alpha \frac{e(k)}{b(k)} (v_l(k) - f(k)) z_l(k) \frac{(x_i(k) - m_{il}(k))^2}{\sigma_{il}^3(k)}, \\ k &= k + 1. \end{aligned}$$

STEP 5 *Repeat by Returning to Step 2.*

◆ Example 9.4

The plant to be identified is the same as that considered by Wang [296, p. 175], and it is described by

$$y(k+1) = 0.3y(k) + 0.6y(k-1) + g(u(k)),$$

where the “unknown” function $g(u)$ has the form

$$g(u) = 0.6 \sin(\pi u) + 0.3 \sin(3\pi u) + 0.1 \sin(5\pi u),$$

and $u(k) = \sin(2\pi k/200)$. The identifier has the structure shown in Figure 9.16 and is described by

$$\hat{y}(k+1) = 0.3\hat{y}(k) + 0.6\hat{y}(k-1) + f(u(k)),$$

where f has the form of (9.42) with $M=10$. We chose $\alpha=0.5$, $y(0)=y(1)=0$, $f(1)=f(2)=-3$, and we started with $u(1)=0$. The identifier's parameters were initialized in a random fashion using MATLAB's command `rand`. After the identifier was trained, we tested its prediction power using the same input signal as before and the initial conditions $y(0)=y(1)=0$, $f(1)=f(2)=-3$, and $u(1)=0$. The results of this experiment are shown in Figures 9.17 and 9.18.

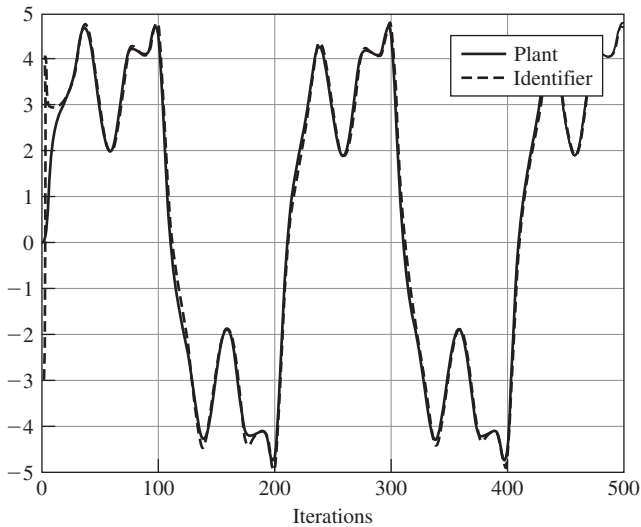


Figure 9.17 Plots of the outputs of the plant and its identifier during the learning phase in Example 9.4.

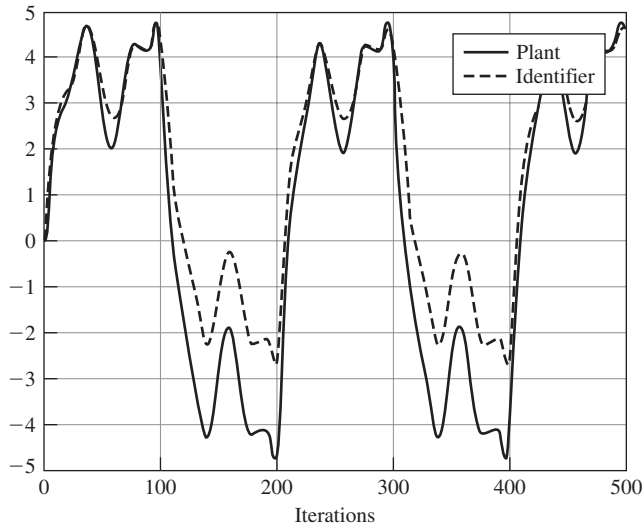


Figure 9.18 Plots of the outputs of the plant and its identifier during the testing phase in Example 9.4.

9.5 Radial-Basis Function (RBF) Networks

In this section we present radial-basis function (RBF) neural networks and their applications. We first show how an RBF network can be used for interpolation purposes. Then, we concern ourselves with the problem of identification of unknown dynamical systems using only the systems' input and output data. The RBF neural identifier is an adaptive system for which the learning algorithm is developed using Lyapunov techniques. The algorithm can add new nodes to the neural identifier to obtain the desired identification accuracy.

9.5.1 Interpolation Using RBF Networks

A diagram of a radial-basis function (RBF) network is shown in Figure 9.19. The functions $\phi(\|\mathbf{x} - \mathbf{x}^{(i)}\|)$, $i = 1, 2, \dots, N$, are called the radial-basis functions, where $\|\cdot\|$ denotes a p -norm that is usually the Euclidean 2-norm. The real parameters w_i , $i = 1, 2, \dots, N$, are the weights. A common choice for ϕ is

$$\phi(\|\mathbf{x} - \mathbf{x}^{(i)}\|) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}^{(i)}\|^2}{2\sigma_i^2}\right), \quad (9.45)$$

where $\mathbf{x}^{(i)}$ is the center of the basis function and σ_i is sometimes called its radius. The function given by (9.45) is called the Gaussian radial-basis function.

The RBF networks are used to perform interpolation. The interpolation problem can be formulated as: For a given set of N different points

$$\{\mathbf{x}^{(i)} : \mathbf{x}^{(i)} \in \mathbb{R}^n, i = 1, 2, \dots, N\}$$

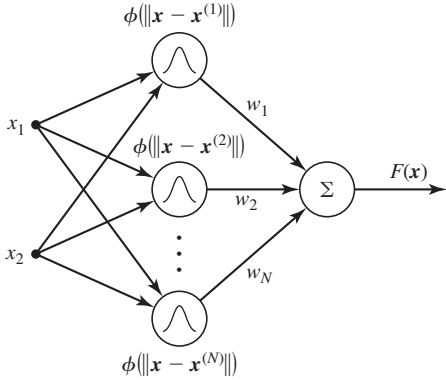


Figure 9.19 A two-input radial-basis function (RBF) network.

and a corresponding set of N real numbers

$$\{d_i : d_i \in \mathbb{R}, i = 1, 2, \dots, N\},$$

find a function $F : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$F(\mathbf{x}^{(i)}) = d_i, \quad i = 1, 2, \dots, N. \quad (9.46)$$

We attempt to solve the interpolation problem using the RBF network. The function that the RBF network implements is

$$F(\mathbf{x}) = \sum_{j=1}^N w_j \phi(\|\mathbf{x} - \mathbf{x}^{(j)}\|). \quad (9.47)$$

Substituting the interpolation conditions given by (9.46) into (9.47) gives

$$\begin{bmatrix} \phi_{11} & \phi_{12} & \cdots & \phi_{1N} \\ \phi_{21} & \phi_{22} & \cdots & \phi_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{N1} & \phi_{N2} & \cdots & \phi_{NN} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_N \end{bmatrix} = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_N \end{bmatrix}, \quad (9.48)$$

where $\phi_{ij} = \phi(\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|)$. Let $\mathbf{w} = [w_1 \ w_2 \ \cdots \ w_N]$, $\mathbf{d} = [d_1 \ d_2 \ \cdots \ d_N]$, and

$$\Phi = \begin{bmatrix} \phi_{11} & \phi_{12} & \cdots & \phi_{1N} \\ \phi_{21} & \phi_{22} & \cdots & \phi_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{N1} & \phi_{N2} & \cdots & \phi_{NN} \end{bmatrix}.$$

Using the above notation, we represent (9.48) as

$$\Phi \mathbf{w} = \mathbf{d}. \quad (9.49)$$

The matrix Φ is called the *interpolation matrix*. If the interpolation matrix is nonsingular, then we can solve (9.49) for \mathbf{w} to obtain

$$\mathbf{w} = \Phi^{-1} \mathbf{d}. \quad (9.50)$$

◆ **Example 9.5**

We construct an RBF network that interpolates data given in Table 9.2, where $\phi_{ij} = \exp(-\frac{(x^{(i)} - x^{(j)})^2}{2\sigma_i^2})$. We set $\sigma_i = 1$, for $i = 1, \dots, 5$, and evaluate (9.50) to obtain

$$\mathbf{w} = [0.8000 \quad 0.1459 \quad -0.0390 \quad 0.5737 \quad 0.2227]^T.$$

The interpolation result is shown in Figure 9.20. In Figure 9.20(a), we used a solid line to show the result of interpolating using the Gaussian radial-basis functions, while a broken line indicates the result of the straight line interpolation. The plots of the basis functions that contribute to the interpolating curve are shown in Figure 9.20(b).

Table 9.2 Interpolation Data for Example 9.5

$x^{(i)}$	d_i
0	0.8
6	0.2
7	0.4
8	0.6
10	0.3

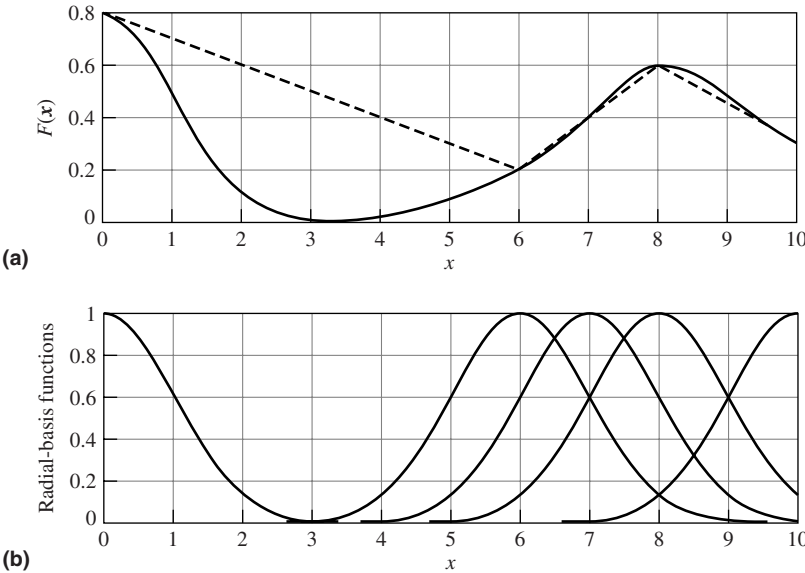


Figure 9.20 Results of interpolating data from Table 9.2 using radial-basis functions.

If the number of radial basis functions was fixed and the rows of the matrix Φ and the elements of the vector \mathbf{d} were obtained sequentially, we could use the recursive least squares algorithm, described in Section A.11, to compute the least squares estimate of \mathbf{w} .

9.5.2 Identification of a Single-Input, Single-State System

Here, we discuss the identification of a single-input, single-state dynamical system using neural network. We consider a dynamical system modeled by

$$\dot{x}(t) = f(x(t), u(t)), \quad x(0) = x_0, \quad (9.51)$$

where $f: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is an unknown, possibly nonlinear, function of x and u . We will construct a neural network to identify the function f . Adding and subtracting $ax(t)$, where $a > 0$ is a design parameter, to the right-hand side of (9.51) gives

$$\begin{aligned} \dot{x}(t) &= -ax(t) + ax(t) + f(x(t), u(t)) \\ &= -ax(t) + g(x(t), u(t)), \end{aligned} \quad (9.52)$$

where $g(x(t), u(t)) = ax(t) + f(x(t), u(t))$. The reason for the above manipulations will become apparent when we discuss learning rules of the neural network identifier of the dynamical system (9.51). This identifier is described by

$$\dot{\hat{x}}(t) = -a\hat{x}(t) + \hat{g}(x(t), u(t), \mathbf{m}_k, \sigma_k), \quad \hat{x}(0) = \hat{x}_0, \quad (9.53)$$

where \hat{x} is the state of the identification network, and \mathbf{m}_k and σ_k are adjustable parameters of the identifier. The function \hat{g} has the form

$$\hat{g}(x(t), u(t), \mathbf{m}_k, \sigma_k) = \sum_{k=1}^K w_k \phi(x, u, \mathbf{m}_k, \sigma_k), \quad (9.54)$$

where for $k = 1, 2, \dots, K$ we have $\mathbf{m}_k = [m_{xk} \ m_{uk}]$ and

$$\phi(x(t), u(t), \mathbf{m}_k, \sigma_k) = \exp\left(-\frac{1}{2\sigma_k^2}((x(t) - m_{xk})^2 + (u(t) - m_{uk})^2)\right).$$

Thus, we postulate that the function \hat{g} be generated using an RBF network. Suppose now that the basis functions are given to us. The problem of choosing basis functions is discussed later. Suppose also that we know the “optimal” weights w_k^* , $k = 1, 2, \dots, K$, that result in the least identification error between g in (9.52) and \hat{g} in (9.53), where the identification error is defined as

$$e(t) = g(x(t), u(t)) - \sum_{k=1}^K w_k^* \phi(x(t), u(t), \mathbf{m}_k, \sigma_k). \quad (9.55)$$

Using the above notation, we represent the dynamical system model (9.52) in the form

$$\dot{x}(t) = -ax(t) + \sum_{k=1}^K w_k^* \phi(x(t), u(t), \mathbf{m}_k, \sigma_k) + e(t). \quad (9.56)$$

The neural identifier, on the other hand, can be represented as

$$\dot{\hat{x}}(t) = -a\hat{x}(t) + \sum_{k=1}^K w_k \phi(x(t), u(t), \mathbf{m}_k, \sigma_k). \quad (9.57)$$

A schematic of the above identification architecture is depicted in Figure 9.21.

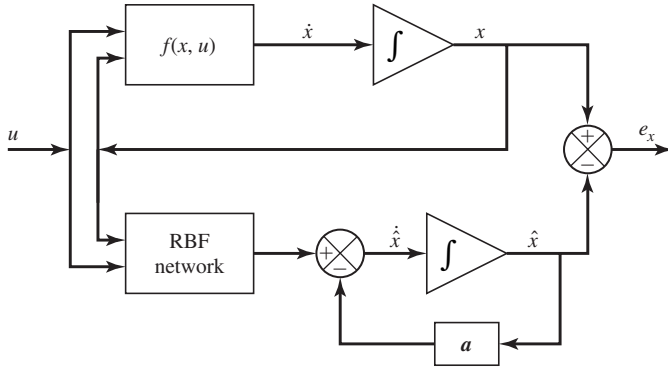


Figure 9.21 A configuration for the radial-basis function neural identifier.

9.5.3 Learning Algorithm for the RBF Identifier

We assume that both the state x and the input u of the dynamical system model (9.51), and hence that of (9.52), are bounded. We next define the state estimation error to be

$$e_x = x - \hat{x},$$

while the weights' estimation errors are

$$\xi_k = w_k^* - w_k, \quad k = 1, \dots, K.$$

We assume that the modeling error, e , defined in (9.55), is bounded by a finite constant e_K , that is, $|e| \leq e_K$. However, e_K is unknown to us. Next, we use (9.56) and (9.57) to form the differential equation describing the state error dynamics. For simplicity of notation, we omit the argument t of the time-varying functions. We have

$$\begin{aligned} \dot{e}_x &= \dot{x} - \dot{\hat{x}} \\ &= -ax + \sum_{k=1}^K w_k^* \phi(x, u, \mathbf{m}_k, \sigma_k) + e + a\hat{x} - \sum_{k=1}^K w_k \phi(x, u, \mathbf{m}_k, \sigma_k) \\ &= -ae_x + \sum_{k=1}^K \xi_k \phi(x, u, \mathbf{m}_k, \sigma_k) + e. \end{aligned} \tag{9.58}$$

We now use the second method of Lyapunov to devise a learning law for the RBF identifier that would yield a stable identification scheme in the sense that the state estimation error e_x decreases as the time increases and the weights do not diverge. For this purpose, we use the following Lyapunov function candidate:

$$V(e_x, \xi) = \frac{1}{2} \left(e_x^2 + \frac{1}{\alpha} \sum_{k=1}^K \xi_k^2 \right), \tag{9.59}$$

where $\alpha > 0$ is a design parameter and $\xi = [\xi_1 \ \xi_2 \ \dots \ \xi_K]^T$. The time derivative of V evaluated

on the trajectories of (9.58) is

$$\begin{aligned}
 \dot{V} &= e_x \dot{e}_x + \frac{1}{\alpha} \sum_{k=1}^K \xi_k \dot{\xi}_k \\
 &= e_x \left(-ae_x + \sum_{k=1}^K \xi_k \phi(x, u, \mathbf{m}_k, \sigma_k) + e \right) + \frac{1}{\alpha} \sum_{k=1}^K \xi_k \dot{\xi}_k \\
 &= -ae_x^2 + \frac{1}{\alpha} \sum_{k=1}^K \xi_k (\alpha e_x \phi(x, u, \mathbf{m}_k, \sigma_k) + \dot{\xi}_k) + e_x e.
 \end{aligned} \tag{9.60}$$

Let

$$\dot{\xi}_k = -\alpha e_x \phi(x, u, \mathbf{m}_k, \sigma_k). \tag{9.61}$$

Then (9.60) simplifies to

$$\dot{V} = -ae_x^2 + e_x e. \tag{9.62}$$

Because $\dot{w}_k = -\dot{\xi}_k$, we can represent (9.61) as

$$\boxed{\dot{w}_k = \alpha e_x \phi(x, u, \mathbf{m}_k, \sigma_k)} \tag{9.63}$$

We further evaluate (9.62) to obtain

$$\begin{aligned}
 \dot{V} &\leq -ae_x^2 + |e_x| e_K \\
 &= -a|e_x| \left(|e_x| - \frac{e_K}{a} \right)
 \end{aligned} \tag{9.64}$$

It follows from (9.62) that if the state estimation error is zero—that is, $e_K = 0$ —then

$$\dot{V} = -ae_x^2.$$

Thus in this case, \dot{V} is negative semidefinite in the $[e_x \ \xi^T]^T$ space, and hence the origin of the $[e_x \ \xi^T]^T$ space is stable, which guarantees boundedness of e_x and the weights w_k s. We will now show, using the Lyapunov-like lemma stated on page 211, that when $e = 0$ we have $e_x(t) \rightarrow 0$ as $t \rightarrow \infty$. We have to check if the assumptions of the Lyapunov-like lemma are satisfied. First, observe that V is bounded below and positive semidefinite. Next, because e_x , ξ_k as well as $\phi(x, u, \mathbf{m}_k, \sigma_k)$ are bounded, the derivative of the state estimation error, given by (9.58) for $e = 0$, is bounded as well, which implies that

$$\ddot{V} = -2ae_x \dot{e}_x$$

is bounded. This in turn implies that \dot{V} is uniformly continuous. Hence, by the Lyapunov-like lemma, $e_x(t) \rightarrow 0$ as $t \rightarrow \infty$.

In the presence of the identification error, e , if $|e_x| < e_K/a$, then it is possible that $\dot{V} > 0$, which implies that the weights w_k s may drift to infinity with time. To avoid this problem, Liu, Kadiramanathan, and Billings [187] proposed to set $\dot{w}_k = 0$ if $|e_x| < e_K/a$. Also, because e_K is unknown, a bound is imposed on the weights as well as on the error. The learning rule

becomes

$$\dot{w}_k = \begin{cases} \alpha e_x \phi(x, u, \mathbf{m}_k, \sigma_k) & \text{if } |e_x| \geq e_0 \text{ and } \|\mathbf{w}\| \leq M \\ 0 & \text{otherwise.} \end{cases} \quad (9.65)$$

where the bounds e_0 and M are design parameters, and $\mathbf{w} = [w_1 \cdots w_K]^T$. With the above modified learning law, whenever $|e_x| \geq e_0 \geq e_K/a$, the Lyapunov derivative \dot{V} is negative semidefinite. However, when $e_0 \leq |e_x| \leq e_K/a$, the weights may increase with time because it is possible that $\dot{V} > 0$. We impose an upper bound on the weights. We also impose an upper bound on the estimation error, e_x . When any of the two bounds is reached, more basis functions may be needed to improve the approximation of g . New basis functions are added to reduce the error until $\dot{V} < 0$. In what follows, we describe a procedure for adding radial-basis functions to the identification network.

9.5.4 Growing RBF Network

As we mentioned above, when $\dot{V} > 0$, the identification scheme is no longer stable, but the learning rule (9.65) guarantees that the adjustable parameters are bounded. Once the parameter adjustments can no longer reduce the estimation error and the bounds on the weights or on the estimation error are reached, additional basis functions in the RBF network should be added to reduce the modeling error. This continues until the approximation leads to the error reduction so that $\dot{V} < 0$.

In the proposed identification scheme by Liu, Kadirkamanathan, and Billings [187], the Gaussian radial-basis functions are used, whose areas of “influence” are determined by the functions’ centers and radii (widths). As the function to be approximated moves further away from a particular basis function, the contribution of that basis function decreases. This is due to the exponential form of the Gaussian basis functions. As $|z|$ increases, $\exp(-z^2)$ decreases. The maximum contribution of a particular basis function to approximating a given function occurs when the center of the basis function and the approximated function are at the same point because in this case we have $\exp(0) = 1$. This heuristic argument leads to the following rule for determining the centers of new basis functions to be added at a time instance t :

$$\mathbf{m}_{K+1} = [m_{x,K+1} \quad m_{u,K+1}]^T = [\delta_x \text{round}(x(t)/\delta_x) \quad \delta_u \text{round}(u(t)/\delta_u)]^T, \quad (9.66)$$

where the round operator rounds the elements on which it operates to the nearest integer. For example, if $\mathbf{x} = [-1.8 \ -0.3 \ 3.7]$, then $\text{round}(\mathbf{x}) = [-2 \ 0 \ 4]$. The parameters δ_x and δ_u represent the spacing of the centers of the basis functions. We can think of having a grid over the input-state space where the spacing between the vertices of the grid are determined by the parameters δ_x and δ_u . Each vertex of the grid represents the potential center of a basis function. The parameter σ_k of the new basis function is determined from the relation

$$\sigma_k = \kappa \|\mathbf{m}_k - \mathbf{m}_k^*\|, \quad (9.67)$$

where $\kappa > 0$ is a design parameter and

$$\mathbf{m}_k^* = \arg \min_{\substack{i=1, \dots, K+1 \\ i \neq k}} \{\|\mathbf{m}_k - \mathbf{m}_i\|\}. \quad (9.68)$$

Thus, \mathbf{m}_k^* is the nearest center of all basis functions to the center of the k th basis function. With the above background, we can now present the rules for adding a new basis function to the

RBF network. Initially, the network contains no basis functions. Then, once the following two conditions are satisfied, a new basis function is added. The first condition is

$$\min_{k=1,\dots,K} |x(t) - m_{xk}| > \delta_x/2 \quad \text{or} \quad \min_{k=1,\dots,K} |u(t) - m_{uk}| > \delta_u/2. \quad (9.69)$$

The above will be satisfied if the state or input, at time t , is closer to an available “unoccupied” grid vertex than to any of the existing centers of the basis functions. The basis function at the new vertex would be the closest to the existing state or input and would have the greatest influence in the approximation of g . The second condition is

$$|x(t) - \hat{x}(t)| > e_{\max}, \quad (9.70)$$

where e_{\max} is the maximum allowable state estimation error. If (9.70) is satisfied, a new basis function is added to reduce the state estimation error by reducing the modeling error.

◆ Example 9.6

The above technique for identification of dynamical systems using radial-basis function network is illustrated with the system

$$\dot{x} = 1.1(1 - x + 2xu - x^2) \exp(-0.5x^2 - u^2), \quad x(0) = 0, \quad (9.71)$$

where $u = \cos(t/100)$. This model was adapted from Liu, Kadiramanathan, and Billings [187]. There are many ways to simulate the above identification scheme. In our simulations, we discretized the above differential equation using the Euler method with the step size $dt = 0.01$ sec. The initial weights and the initial state of the identifier were set to zero. We used the following parameter values: $e_0 = 0.001$, $e_{\max} = 0.005$, $\alpha = 5$, $\kappa = 3$. In Figure 9.22, plots of the true state and its estimate

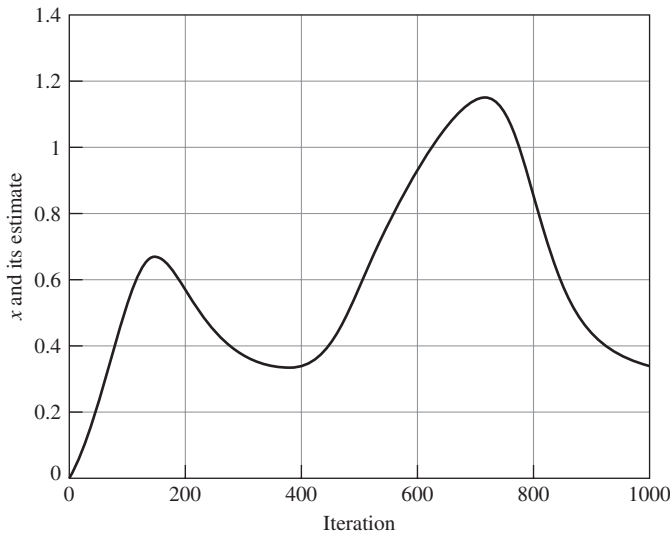


Figure 9.22 Plots of the true state and its estimate versus time. The estimated state “hugs” closely the actual state.

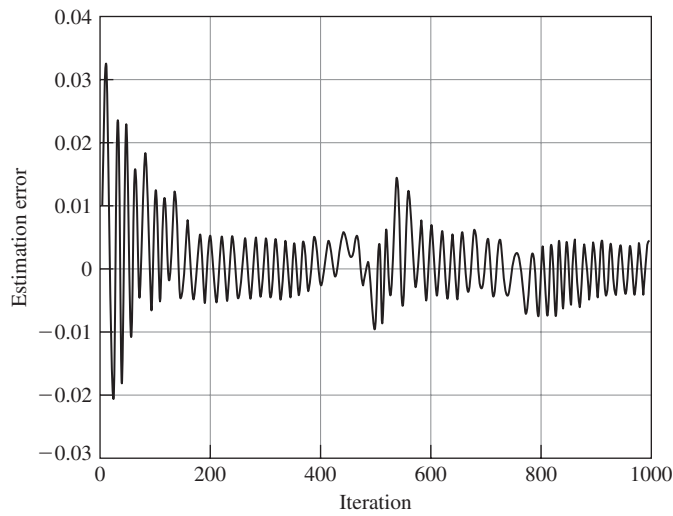


Figure 9.23 A plot of the error, e_x , between the true state and its estimate.

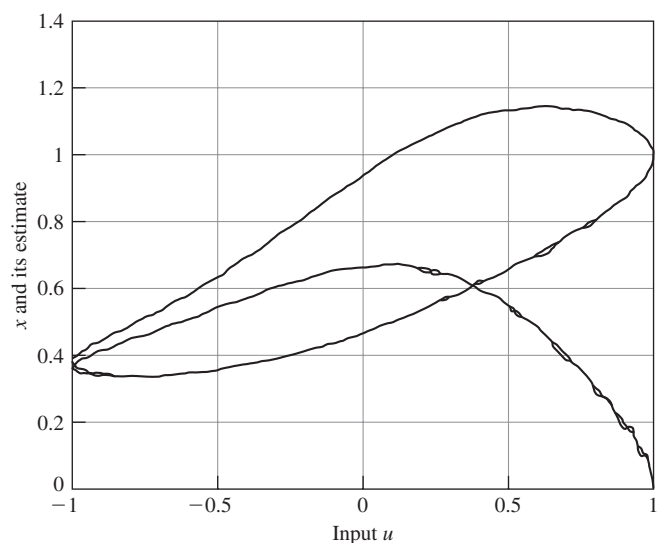


Figure 9.24 A plot of the true state and its estimate versus the input u .

versus time are shown. A plot of the state estimation error, e_x , between the true state and its estimate is shown in Figure 9.23. Initially, the error exceeds the bound e_{\max} . Adding additional basis functions drives the error back down. The plots shown in Figure 9.24 depict the true state and its estimate versus the input u . Finally, in Figure 9.25, plots of the derivative of the true state versus the true state and the derivative of the estimated state versus the state estimate are shown.

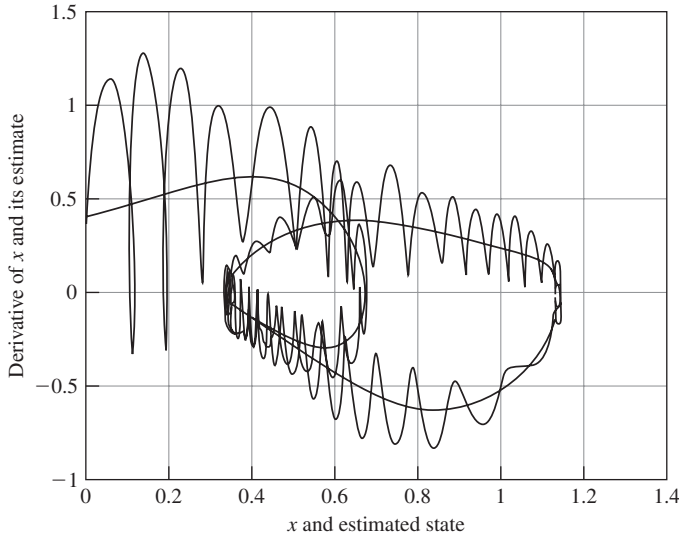


Figure 9.25 Plots of the derivative of the true state versus the true state and the derivative of the estimated state versus the state estimate of Example 9.6.

9.5.5 Identification of Multivariable Systems

We now extend the results of Subsection 9.5.2 to multi-input, multistate dynamical systems. In our derivation of the update law for the weights of the RBF network, we use a result from matrix algebra presented next. Recall the definition of the trace of a square matrix as given by (A.32). If $\mathbf{\Gamma} \in \mathbb{R}^{m \times n}$, then by (A.79) its Frobenius norm satisfies

$$\|\mathbf{\Gamma}\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n \gamma_{ij}^2 = \text{trace}(\mathbf{\Gamma}\mathbf{\Gamma}^T).$$

If the elements, γ_{ij} , of $\mathbf{\Gamma}$ are functions of time, then

$$\frac{d}{dt}(\text{trace}(\mathbf{\Gamma}\mathbf{\Gamma}^T)) = 2 \sum_{i=1}^m \sum_{j=1}^n \gamma_{ij} \dot{\gamma}_{ij} = 2 \text{trace}(\mathbf{\Gamma}\dot{\mathbf{\Gamma}}^T), \quad (9.72)$$

where $\dot{\mathbf{\Gamma}} = d\mathbf{\Gamma}/dt$. With the above background, we can proceed with the analysis of the RBF neural identifier for the multivariable systems.

The multivariable dynamical systems that we consider are modeled by

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t)), \quad \mathbf{x}(0) = \mathbf{x}_0, \quad (9.73)$$

where $\mathbf{u}(t) \in \mathbb{R}^m$, $\mathbf{x}(t) \in \mathbb{R}^n$, and $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ is the function to be identified. Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be a Hurwitz matrix—that is, the matrix whose eigenvalues are all in the open left-half plane. Because \mathbf{A} is a Hurwitz matrix, by the theorem of Lyapunov (Theorem 4.1), for any given real symmetric positive definite (r.s.p.d.) matrix \mathbf{Q} the solution \mathbf{P} to the continuous Lyapunov matrix equation

$$\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A} = -2\mathbf{Q} \quad (9.74)$$

is also (real symmetric) positive definite. We use this fact later when we develop a learning law for the RBF neural identifier.

We represent the model (9.73) in an equivalent form,

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{g}(\mathbf{x}(t), \mathbf{u}(t)), \quad \mathbf{x}(0) = \mathbf{x}_0, \quad (9.75)$$

where

$$\mathbf{g}(\mathbf{x}(t), \mathbf{u}(t)) = -\mathbf{A}\mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), \mathbf{u}(t)).$$

The k th basis function can be represented as

$$\phi_k(\mathbf{x}, \mathbf{u}) = \exp\left(-\frac{1}{2\sigma_k^2} \left\| \begin{bmatrix} \mathbf{x} \\ \mathbf{u} \end{bmatrix} - \mathbf{m}_k \right\|^2\right).$$

Let

$$\Phi_K(\mathbf{x}, \mathbf{u}) = [\phi_1(\mathbf{x}, \mathbf{u}) \quad \phi_2(\mathbf{x}, \mathbf{u}) \quad \cdots \quad \phi_K(\mathbf{x}, \mathbf{u})]^T$$

and

$$\hat{\mathbf{g}}(\mathbf{x}, \mathbf{u}, \mathbf{m}_k, \sigma_k) = \mathbf{W}_K \Phi_K,$$

where \mathbf{W}_K is the $n \times K$ weight matrix and $\mathbf{m}_k = [\mathbf{m}_{xk}^T \quad \mathbf{m}_{uk}^T]^T$ is the center of the k th radial basis function in the $[\mathbf{x}^T \quad \mathbf{u}^T]^T$ space. Note that Φ_K is a column vector with K components. The RBF neural identifier has the form

$$\dot{\hat{\mathbf{x}}}(t) = \mathbf{A}\hat{\mathbf{x}}(t) + \hat{\mathbf{g}}(\mathbf{x}, \mathbf{u}, \mathbf{m}_k, \sigma_k), \quad \hat{\mathbf{x}}(0) = \mathbf{x}_0. \quad (9.76)$$

Let \mathbf{W}_K^* be the optimal weight matrix that minimizes a norm of the identification error, that is,

$$\mathbf{W}_K^* = \arg \min_{\mathbf{W}_K} \|\mathbf{g}(\mathbf{x}(t), \mathbf{u}(t)) - \mathbf{W}_K \Phi_K(\mathbf{x}(t), \mathbf{u}(t))\|.$$

Let

$$\mathbf{e}(t) = \mathbf{g}(\mathbf{x}(t), \mathbf{u}(t)) - \mathbf{W}_K^* \Phi_K(\mathbf{x}(t), \mathbf{u}(t)).$$

Using the above definition, we represent the system model (9.75) as

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{W}_K^* \Phi_K(\mathbf{x}(t), \mathbf{u}(t)) + \mathbf{e}(t). \quad (9.77)$$

We define the state estimation error and the weight estimation error as

$$\mathbf{e}_x = \mathbf{x} - \hat{\mathbf{x}} \quad \text{and} \quad \Gamma_K = \mathbf{W}_K^* - \mathbf{W}_K.$$

Next, we form the differential equation describing the state estimation error dynamics,

$$\begin{aligned} \dot{\mathbf{e}}_x &= \dot{\mathbf{x}} - \dot{\hat{\mathbf{x}}} \\ &= \mathbf{A}\mathbf{x} + \mathbf{W}_K^* \Phi_K(\mathbf{x}, \mathbf{u}) + \mathbf{e} - \mathbf{A}\hat{\mathbf{x}} - \hat{\mathbf{g}}(\mathbf{x}, \mathbf{u}, \mathbf{m}_k, \sigma_k) \\ &= \mathbf{A}(\mathbf{x} - \hat{\mathbf{x}}) + \mathbf{W}_K^* \Phi_K(\mathbf{x}, \mathbf{u}) - \mathbf{W}_K \Phi_K(\mathbf{x}, \mathbf{u}) + \mathbf{e} \\ &= \mathbf{A}\mathbf{e}_x + \Gamma_K \Phi_K(\mathbf{x}, \mathbf{u}) + \mathbf{e}. \end{aligned} \quad (9.78)$$

We employ the following Lyapunov function candidate:

$$V = V(\mathbf{e}_x, \Gamma_K) = \frac{1}{2} \left(\mathbf{e}_x^T \mathbf{P} \mathbf{e}_x + \frac{1}{\alpha} \text{trace}(\Gamma_K \Gamma_K^T) \right). \quad (9.79)$$

Using (9.72), we evaluate the time derivative of V on the solutions of the state estimation error equation to get

$$\begin{aligned}\dot{V} &= \mathbf{e}_x^T \mathbf{P} \dot{\mathbf{e}}_x + \frac{1}{\alpha} \text{trace}(\mathbf{\Gamma}_K \dot{\mathbf{\Gamma}}_K^T) \\ &= \mathbf{e}_x^T (\mathbf{P} \mathbf{A} \mathbf{e}_x + \mathbf{P} \mathbf{\Gamma}_K \mathbf{\Phi}_K(\mathbf{x}, \mathbf{u}) + \mathbf{P} \mathbf{e}) + \frac{1}{\alpha} \text{trace}(\mathbf{\Gamma}_K \dot{\mathbf{\Gamma}}_K^T) \\ &= \frac{1}{2} \mathbf{e}_x^T (\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A}) \mathbf{e}_x + \mathbf{e}_x^T \mathbf{P} \mathbf{\Gamma}_K \mathbf{\Phi}_K(\mathbf{x}, \mathbf{u}) + \frac{1}{\alpha} \text{trace}(\mathbf{\Gamma}_K \dot{\mathbf{\Gamma}}_K^T) + \mathbf{e}_x^T \mathbf{P} \mathbf{e}. \quad (9.80)\end{aligned}$$

Using the definition of the trace operator, we obtain

$$\begin{aligned}\mathbf{e}_x^T \mathbf{P} \mathbf{\Gamma}_K \mathbf{\Phi}_K(\mathbf{x}, \mathbf{u}) &= \text{trace}(\mathbf{e}_x^T \mathbf{P} \mathbf{\Gamma}_K \mathbf{\Phi}_K(\mathbf{x}, \mathbf{u})) \\ &= \text{trace}(\mathbf{\Gamma}_K \mathbf{\Phi}_K(\mathbf{x}, \mathbf{u}) \mathbf{e}_x^T \mathbf{P}).\end{aligned} \quad (9.81)$$

Observe that

$$\dot{\mathbf{W}}_K = -\dot{\mathbf{\Gamma}}_K.$$

Using the above and (9.81) along with the Lyapunov equation (9.74), we write

$$\dot{V} = -\mathbf{e}_x^T \mathbf{Q} \mathbf{e}_x + \frac{1}{\alpha} (\text{trace}(\alpha \mathbf{\Gamma}_K \mathbf{\Phi}_K(\mathbf{x}, \mathbf{u}) \mathbf{e}_x^T \mathbf{P}) - \text{trace}(\mathbf{\Gamma}_K \dot{\mathbf{W}}_K^T)) + \mathbf{e}_x^T \mathbf{P} \mathbf{e}. \quad (9.82)$$

Consider the following weight matrix update law:

$$\boxed{\dot{\mathbf{W}}_K = \alpha \mathbf{P} \mathbf{e}_x \mathbf{\Phi}_K^T(\mathbf{x}, \mathbf{u}).} \quad (9.83)$$

With the above update law, the Lyapunov derivative (9.82) evaluates to

$$\begin{aligned}\dot{V} &= -\mathbf{e}_x^T \mathbf{Q} \mathbf{e}_x + \mathbf{e}_x^T \mathbf{P} \mathbf{e} \\ &\leq -\lambda_{\min}(\mathbf{Q}) \|\mathbf{e}_x\|^2 + \lambda_{\max}(\mathbf{P}) \|\mathbf{e}\| \|\mathbf{e}_x\| \\ &= -\lambda_{\min}(\mathbf{Q}) \left(\|\mathbf{e}_x\| - \frac{\lambda_{\max}(\mathbf{P})}{\lambda_{\min}(\mathbf{Q})} \|\mathbf{e}\| \right) \|\mathbf{e}_x\|.\end{aligned} \quad (9.84)$$

Utilizing the same type of arguments as in the scalar case (analyzed in Subsection 9.5.2) and the Lyapunov-like lemma stated on page 211, we can show that when $\mathbf{e} = \mathbf{0}$, then $\mathbf{e}_x(t) \rightarrow \mathbf{0}$ as $t \rightarrow \infty$.

If the magnitude of the state estimation error is such that $\|\mathbf{e}_x\| < \frac{\lambda_{\max}(\mathbf{P})}{\lambda_{\min}(\mathbf{Q})} \|\mathbf{e}\|$, then it may happen that $\dot{V} > 0$, which in turn may lead to the weights' drift to infinity over time. To prevent this, we modify (9.83) following the analysis of the single-input single-state case. The modified weight update law is

$$\boxed{\dot{\mathbf{W}}_K = \begin{cases} \alpha \mathbf{P} \mathbf{e}_x \mathbf{\Phi}_K^T(\mathbf{x}, \mathbf{u}) & \text{if } \|\mathbf{e}_x\| \geq e_0 \text{ and } \|\mathbf{W}_K\| \leq M \\ \mathbf{0} & \text{otherwise.} \end{cases}} \quad (9.85)$$

In the above adaptation law, e_0 and M are design parameters.

9.6 A Self-Organizing Network

Up to now we discussed neural networks that must have a teacher in order to learn presented patterns. This mode of learning is referred to as the supervised mode. In this section we analyze a neural network that can learn in the unsupervised mode and the network undergoes the self-organization process. Specifically, the network is capable of clustering training patterns. It is assumed that the number of classes, p , is known and that the network's number of weight vectors is equal to the number of classes in the training set. The weights are initialized and then normalized. After initialization, the training patterns, $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}$, are presented. In response to each presented pattern, we compute the winning neuron of the network, that is, we compute the network weight that is closest to the presented pattern. We call this computation the *similarity matching*. Let \mathbf{x} be the presented pattern, where, for simplicity, we omitted its index. Then the similarity matching can be described as

$$\|\mathbf{x} - \mathbf{w}^{(i)}\| = \min_{1 \leq j \leq p} \{\|\mathbf{x} - \mathbf{w}^{(j)}\|\},$$

where i is the index of the winning neuron. Thus, searching for the winning neuron corresponds to finding the maximum among the p scalar products,

$$\max_{j=1,2,\dots,p} \langle \mathbf{w}^{(j)}, \mathbf{x} \rangle.$$

To see that this is indeed the case, note that because the weights and the patterns are normalized—that is, $\|\mathbf{w}^{(j)}\| = \|\mathbf{x}\| = 1$ —we have

$$\begin{aligned} \|\mathbf{x} - \mathbf{w}^{(j)}\| &= (\|\mathbf{x}\|^2 - 2\langle \mathbf{w}^{(j)}, \mathbf{x} \rangle + \|\mathbf{w}^{(j)}\|^2)^{1/2} \\ &= (2 - 2\langle \mathbf{w}^{(j)}, \mathbf{x} \rangle)^{1/2}. \end{aligned}$$

Hence, minimizing $\|\mathbf{x} - \mathbf{w}^{(j)}\|$ is equivalent to maximizing the scalar product $\langle \mathbf{w}^{(j)}, \mathbf{x} \rangle$. On the other hand,

$$\langle \mathbf{w}^{(j)}, \mathbf{x} \rangle = \|\mathbf{w}^{(j)}\| \|\mathbf{x}\| \cos(\angle \mathbf{w}^{(j)}, \mathbf{x}),$$

which means that the winning neuron is “more parallel” to the input pattern \mathbf{x} than the other neurons.

After the winning neuron is declared, it is updated, while the remaining weights are not changed. This is the reason for the name of the network—the winner-take-all (WTA) network. This update algorithm is also known as the *Kohonen learning rule* after Kohonen, who first proposed this learning rule. The role of the update is to reduce the distance between the winning neuron, $\mathbf{w}^{(i)}$, and the presented pattern, \mathbf{x} . One way to accomplish this task is to apply the descent gradient method to the objective function $\frac{1}{2}\|\mathbf{x} - \mathbf{w}^{(i)}\|^2$. The gradient of $\frac{1}{2}\|\mathbf{x} - \mathbf{w}^{(i)}\|^2$ with respect to $\mathbf{w}^{(i)}$ is

$$\frac{1}{2} \nabla_{\mathbf{w}^{(i)}} \|\mathbf{x} - \mathbf{w}^{(i)}\|^2 = -(\mathbf{x} - \mathbf{w}^{(i)}).$$

Applying the descent gradient method to minimize $\frac{1}{2}\|\mathbf{x} - \mathbf{w}^{(i)}\|^2$ with a step size α_k yields

$$\mathbf{w}^{(i)}(k+1) = \mathbf{w}^{(i)}(k) + \alpha_k(\mathbf{x} - \mathbf{w}^{(i)}(k)).$$

The weights of other neurons are unchanged, that is,

$$\mathbf{w}^{(j)}(k+1) = \mathbf{w}^{(j)}(k) \quad \text{for } j \neq i.$$

We summarize our discussion in the form of the following algorithm.

WINNER-TAKE-ALL (WTA) LEARNING ALGORITHM

STEP 1 *Initialize Weights.* Generate the initial weights. Normalize each weight to have unit length: $\frac{w}{\|w\|}$.

STEP 2 *Serially Present Training Patterns.* Normalize each training pattern $\mathbf{x}^{(j)}$ to have unit length. The training patterns can be presented in a random order a fixed number of times.

STEP 3 *Select Winning Neuron.* For each WTA weight, compute the scalar product with the presented pattern. Select a weight that gives the largest scalar product. Call this the winning weight.

STEP 4 *Update Weight of Winning Neuron.* Let i be the index of the winning weight. Set

$$\begin{aligned} \mathbf{w}^{(i)}(k+1) &= \mathbf{w}^{(i)}(k) + \alpha_k (\mathbf{x}^{(j)} - \mathbf{w}^{(i)}(k)), \\ \mathbf{w}^{(j)}(k+1) &= \mathbf{w}^{(j)}(k) \quad \text{for } j \neq i. \end{aligned}$$

Normalize $\mathbf{w}^{(i)}(k+1)$ to have unit length.

STEP 5 *Repeat by Returning to Step 2.* Repeat until stopping criterion is met.

◆ Example 9.7

We take the same set of training patterns as in Zurada [320, p. 407]:

$$[\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \mathbf{x}^{(3)} \quad \mathbf{x}^{(4)} \quad \mathbf{x}^{(5)}] = \begin{bmatrix} 0.8 & 0.1736 & 0.707 & 0.342 & 0.6 \\ 0.6 & -0.9848 & 0.707 & -0.9397 & 0.8 \end{bmatrix}.$$

The above patterns are already normalized and are marked in Figure 9.26. As can be seen from Figure 9.26, we can classify the training patterns into two classes,

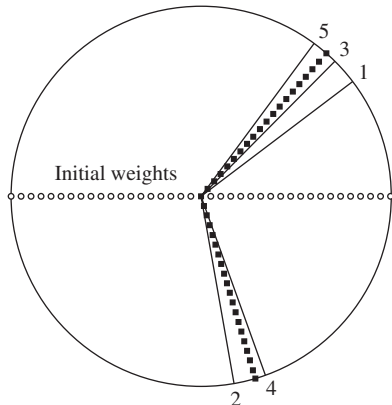


Figure 9.26 Training patterns, initial weights and the final weights of the WTA network of Example 9.7.

that is, $p = 2$. The initial normalized weights are

$$[\mathbf{w}_0^{(1)} \quad \mathbf{w}_0^{(2)}] = \begin{bmatrix} 1 & -1 \\ 0 & 0 \end{bmatrix}.$$

The initial weights are also marked in Figure 9.26. The training patterns, $\mathbf{x}^{(1)}$, $\mathbf{x}^{(2)}, \dots, \mathbf{x}^{(5)}$ are presented in ascending order and then recycled. After 200 presentations of the training set, we obtained the following final weights of the network:

$$[\mathbf{w}^{(1)} \quad \mathbf{w}^{(2)}] = \begin{bmatrix} 0.6628 & 0.2868 \\ 0.7488 & -0.9580 \end{bmatrix}.$$

The final weights are marked in Figure 9.26 with *.

9.7 Hopfield Neural Network

9.7.1 Hopfield Neural Network Modeling and Analysis

The analog Hopfield neural network is a dynamical system whose schematic is shown in Figure 9.27. It is a nonlinear interconnected system proposed by Hopfield [126] as a model of biological neural networks. Each node, or neuron, in Figure 9.27 is represented by a circle.

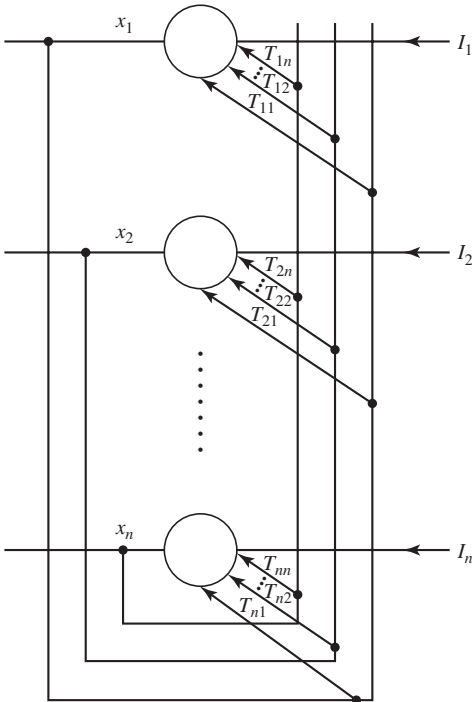


Figure 9.27 Hopfield neural network model.

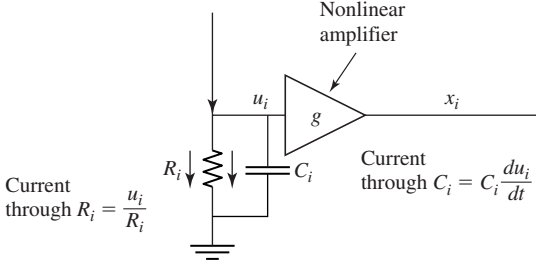


Figure 9.28 Model of a neuron in the Hopfield neural network.

structure of a neuron is depicted in Figure 9.28. The nonlinear amplifier in Figure 9.28 has an input-output characteristic described by a function $g : \mathbb{R} \rightarrow [0, 1]$. It is commonly assumed that g is continuously differentiable and strictly monotonically increasing, that is, $g(u) > g(u')$ if $u > u'$. The function g is called the *activation function*. A typical choice of g is the sigmoid function given by

$$x = g(u) = \frac{1}{1 + e^{-u/\theta}}, \quad (9.86)$$

where the positive parameter θ controls the slope of g . A circuit implementation of the Hopfield network is shown in Figure 9.29. To derive equations that model the Hopfield neural network, we use an equivalent representation of the single node of the Hopfield network shown in Figure 9.30. Applying Kirchhoff's current law at the input of each amplifier of the equivalent node representation and taking into account that the current into the amplifier is negligibly small due to the amplifier's high input resistance, we get

$$C_i \frac{du_i}{dt} + \frac{u_i}{R_i} = \frac{x_1 - u_i}{R_{i1}} + \frac{x_2 - u_i}{R_{i2}} + \cdots + \frac{x_n - u_i}{R_{in}} + I_i, \quad i = 1, 2, \dots, n, \quad (9.87)$$

where $x_i = g(u_i)$. Let

$$T_{ij} = \frac{1}{R_{ij}}, \quad \frac{1}{r_i} = \sum_{j=1}^n T_{ij} + \frac{1}{R_i}. \quad (9.88)$$

Substituting (9.88) into (9.87) yields

$$C_i \frac{du_i}{dt} + \frac{u_i}{r_i} = \sum_{j=1}^n T_{ij} x_j + I_i, \quad i = 1, 2, \dots, n \quad (9.89)$$

The system of n first-order nonlinear differential equations given by (9.89) constitutes a model of the Hopfield neural network. The Hopfield neural network can be viewed as an associative memory. This is because it can have asymptotically stable states that attract the neighboring states to develop in time. This behavior can be interpreted as an evolution of an imperfect, corrupted pattern toward the correct, stored pattern—hence the term associative memory. More on neural associative memory can be found in Subsection 9.9.1.

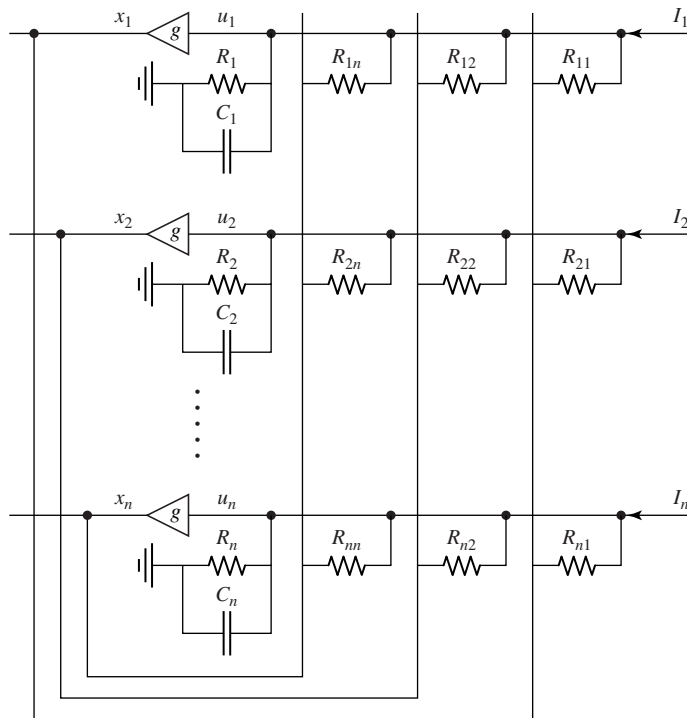


Figure 9.29 A circuit realization of the Hopfield neural network.

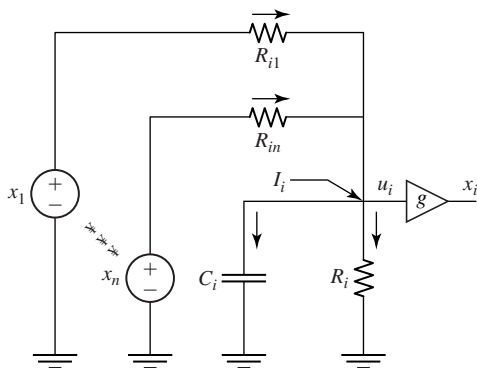


Figure 9.30 Equivalent representation of a single node in the circuit realization of the Hopfield neural network.

In further analysis, we use an equivalent representation of (9.89). To obtain this equivalent representation of the Hopfield model, we first note that the activation function g is, by assumption, strictly increasing and differentiable. Therefore g^{-1} exists and is also differentiable—see the inverse function theorem on page 390. In Figure 9.31, we show plots of $x = g(u)$ and $u = g^{-1}(x)$ for $\theta = 0.2$. We rewrite equations (9.89) as

$$C_i \frac{dg^{-1}(x_i)}{dx_i} \frac{dx_i}{dt} + \frac{g^{-1}(x_i)}{r_i} = \sum_{j=1}^n T_{ij} x_j + I_i, \quad i = 1, 2, \dots, n. \quad (9.90)$$

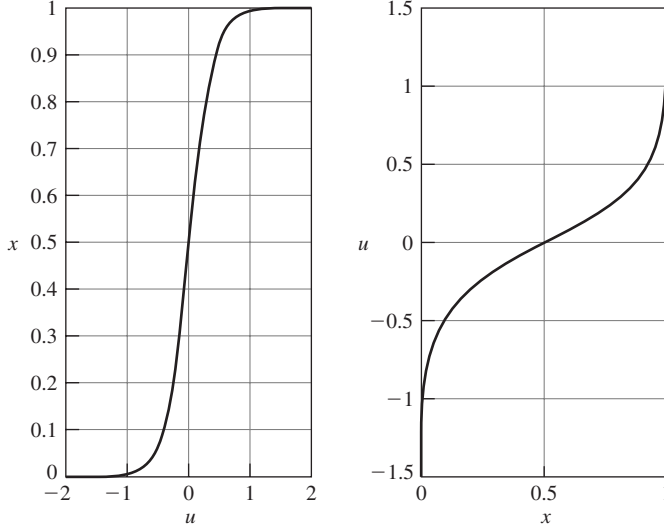


Figure 9.31 Plots of $x = g(u)$ and $u = g^{-1}(x)$.

Let $\mathbf{T} \in \mathbb{R}^{n \times n}$ be the matrix whose elements are T_{ij} and $\mathbf{v} = [I_1 \ I_2 \ \dots \ I_n]^T \in \mathbb{R}^n$. Using this notation, we represent (9.89) in matrix form:

$$\begin{bmatrix} C_1 \frac{dg^{-1}(x_1)}{dx_1} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & C_n \frac{dg^{-1}(x_n)}{dx_n} \end{bmatrix} \dot{\mathbf{x}}(t) = \mathbf{T}\mathbf{x} + \mathbf{v} - \begin{bmatrix} \frac{g^{-1}(x_1)}{r_1} \\ \vdots \\ \frac{g^{-1}(x_n)}{r_n} \end{bmatrix}. \quad (9.91)$$

Associated with the Hopfield neural network is its computational energy function, $E : \mathbb{R}^n \rightarrow \mathbb{R}$, given by

$$\begin{aligned} E(\mathbf{x}) &= -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n T_{ij} x_i x_j - \sum_{i=1}^n I_i x_i + \sum_{i=1}^n \int_0^{x_i} \frac{g^{-1}(s)}{r_i} ds \\ &= -\frac{1}{2} \mathbf{x}^T \mathbf{T} \mathbf{x} - \mathbf{v}^T \mathbf{x} + \sum_{i=1}^n \int_0^{x_i} \frac{g^{-1}(s)}{r_i} ds. \end{aligned} \quad (9.92)$$

Suppose now that $\mathbf{T} = \mathbf{T}^T$. Then, the gradient of E is

$$\nabla E(\mathbf{x}) = -\mathbf{T}\mathbf{x} - \mathbf{v} + \begin{bmatrix} \frac{g^{-1}(x_1)}{r_1} & \dots & \frac{g^{-1}(x_n)}{r_n} \end{bmatrix}^T. \quad (9.93)$$

Let

$$\mathbf{S}(\mathbf{x}) = \begin{bmatrix} C_1 \frac{dg^{-1}(x_1)}{dx_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & C_n \frac{dg^{-1}(x_n)}{dx_n} \end{bmatrix}. \quad (9.94)$$

We now combine (9.93) and (9.94) to obtain an equivalent form of the Hopfield network modeling equations:

$$\boxed{\mathbf{S}(\mathbf{x})\dot{\mathbf{x}}(t) = -\nabla E(\mathbf{x}(t))} \quad (9.95)$$

We analyze the dynamics of the Hopfield network modeled by (9.95) using the computational energy function. We first show that E is bounded. Taking norms of (9.92) yields

$$|E(\mathbf{x})| \leq \frac{1}{2} \|\mathbf{T}\| \|\mathbf{x}\|^2 + \|\mathbf{v}\| \|\mathbf{x}\| + \sum_{i=1}^n \frac{1}{r_i} \left| \int_0^{x_i} g^{-1}(s) ds \right|. \quad (9.96)$$

The components of the vector \mathbf{x} are the outputs of the nonlinear amplifiers, and they take values from the interval $[0, 1]$. Hence, $\|\mathbf{x}\| = \|\mathbf{x}\|_2 \leq \sqrt{n}$. Thus, to show that $E(\mathbf{x})$ is bounded, it is enough to show that the integral $\int_0^{x_i} g^{-1}(s) ds$ is bounded. We first find the expression for g^{-1} . Recall that $g(u)$ is given by (9.86). Therefore,

$$\frac{1}{x} = 1 + e^{-u/\theta},$$

or

$$\frac{1}{x} - 1 = e^{-u/\theta}, \quad (9.97)$$

Taking logarithms of both sides of the above yields

$$\ln\left(\frac{1-x}{x}\right) = -\frac{u}{\theta}.$$

Hence,

$$u = g^{-1}(x) = \theta(\ln(x) - \ln(1-x)). \quad (9.98)$$

Using the symmetry property of g^{-1} , we obtain

$$\begin{aligned} \left| \int_0^{x_i} g^{-1}(s) ds \right| &\leq \int_0^{x_i} |g^{-1}(s)| ds \\ &\leq -2 \int_0^{1/2} g^{-1}(s) ds. \end{aligned} \quad (9.99)$$

Recall that

$$\int \ln(s) ds = s \ln(s) - s.$$

Hence,

$$\int \ln(1-s) ds = -(1-s) \ln(1-s) - s.$$

Using the above, we obtain

$$\begin{aligned}
 \int_0^{1/2} g^{-1}(s) ds &= \int_0^{1/2} \theta (\ln(s) - \ln(1-s)) ds \\
 &= \theta (s \ln(s) - s + (1-s) \ln(1-s) + s) \Big|_0^{1/2} \\
 &= \theta (s \ln(s) + (1-s) \ln(1-s)) \Big|_0^{1/2}.
 \end{aligned} \tag{9.100}$$

Note that

$$s \ln(s) = \frac{\ln(s)}{1/s}.$$

Applying L'Hôpital's rule to the above yields

$$\lim_{s \rightarrow 0} \frac{\ln(s)}{1/s} = \lim_{s \rightarrow 0} \frac{1/s}{-1/s^2} = 0.$$

Applying the above to (9.100), we get

$$\int_0^{1/2} g^{-1}(s) ds = \theta \ln(1/2). \tag{9.101}$$

Combining (9.99) and (9.101) yields

$$\left| \int_0^{x_i} g^{-1}(s) ds \right| \leq -2\theta \ln(1/2), \tag{9.102}$$

which completes the proof of the statement that E is bounded.

Let $\mathbf{x}(t)$ denote the trajectory of the Hopfield network model (9.95). We will show that if $\mathbf{T} = \mathbf{T}^T$, then

$$\frac{d}{dt} E(\mathbf{x}(t)) \leq 0$$

on the Hopfield network trajectories. Using the chain rule gives

$$\frac{d}{dt} E(\mathbf{x}(t)) = (\nabla E(\mathbf{x}(t)))^T \dot{\mathbf{x}}(t). \tag{9.103}$$

It follows from (9.95) that

$$\nabla E(\mathbf{x}(t)) = -\mathbf{S}(\mathbf{x}(t))\dot{\mathbf{x}}(t). \tag{9.104}$$

Substituting (9.104) into (9.103), we obtain

$$\begin{aligned}
 \frac{d}{dt} E(\mathbf{x}(t)) &= -\dot{\mathbf{x}}^T(t) \mathbf{S}(\mathbf{x}(t)) \dot{\mathbf{x}}(t) \\
 &= -\sum_{i=1}^n C_i \frac{dg^{-1}(x_i(t))}{dx_i} \dot{x}_i^2,
 \end{aligned} \tag{9.105}$$

where $C_i > 0$ and

$$\frac{dg^{-1}(x_i)}{dx_i} = \theta(1/x_i + 1/(1-x_i)) > 0 \tag{9.106}$$

because $x_i \in (0, 1)$. Hence,

$$\frac{d}{dt} E(\mathbf{x}(t)) \leq 0$$

on the trajectories of the Hopfield model (9.95).

9.7.2 Analog-to-Digital Converter

We will now show how the Hopfield neural network can be used to solve some engineering problems. Specifically, we will discuss the construction of an analog-to-digital (A/D) converter using the Hopfield neural network. The idea was put forward by Tank and Hopfield [278]. In A/D conversion, an analog value is transformed into its corresponding binary representation. The conversion that we will discuss consists of two steps. The first step is to round the given analog value to the nearest integer. The second step is to convert the integer into its binary equivalent. For example, 3.4 is first rounded to 3, then 3 is converted into its binary equivalent “11” if we have a 2-bit converter, “011” if we have a 3-bit converter, or “0011” if we have a 4-bit converter. Therefore, an N -bit A/D converter can be viewed as a device that assigns to a given analog input a corresponding vertex of the hypercube $[0, 1]^N$.

The first step in A/D conversion, rounding of a given number to its nearest integer, can be organized as solving an optimization problem by a neural network. This entails formulation of a suitable cost, or penalty, function to be minimized. In fact, finding a suitable cost function will be the first step in constructing a neural A/D converter.

To better understand the idea behind solving A/D conversion problems using neural networks, we first consider constructing a 2-bit A/D converter. Let V_1 be the most significant bit (MSB), let V_0 be the least significant bit (LSB), and let x be an analog value to be converted. In the actual neural network, V_1 and V_0 are the outputs of the two neurons in the neural A/D converter. We first consider a very simple cost function,

$$E = \frac{1}{2}(V_1 2^1 + V_0 2^0 - x)^2. \quad (9.107)$$

A plot of the above function, for an analog input of $x = 1.5$, is shown in Figure 9.32(a). Note that E is defined over the two-dimensional hypercube $[0, 1]^2$. However, constructing a circuit minimizing the cost function (9.107) would not yield good results because V_0 and V_1 would take values from the interval $[0, 1]$ rather than the set $\{0, 1\}$. To overcome the problem, we add to (9.107) two terms whose purpose is to penalize for noninteger values of V_0 and V_1 . The modified cost function is

$$E = \frac{1}{2}(V_1 2^1 + V_0 2^0 - x)^2 + a V_1(V_1 - 1) + b V_0(V_0 - 1), \quad (9.108)$$

where a and b are weighting coefficients. Tank and Hopfield [278] suggest selecting a and b so that the quadratic terms V_0^2 and V_1^2 are canceled out. In particular, if we set $a = -2$ and $b = -0.5$, then (9.108) becomes

$$E = \frac{1}{2}x^2 + 2V_0V_1 - 2V_1x - V_0x + 2V_1 + \frac{1}{2}V_0. \quad (9.109)$$

The term used to cancel the quadratic terms is plotted in Figure 9.32(b). The final cost function given by (9.109) is plotted in Figure 9.32(c). To convert an analog input x into its binary equivalent, we minimize the cost function (9.109). One can interpret this process as follows. A given analog input x defines a corresponding cost function E . We can view E as a landscape. Now, a ball placed somewhere on this landscape will move toward the lowest point, which should correspond to one of the vertices of the hypercube $[0, 1]^2$.

Our next task is to construct differential equations that model a neural network implementing an A/D converter. The parameters of the state equations of the network will be chosen so that the time derivative of E evaluated on the solutions of the equations describing the network is

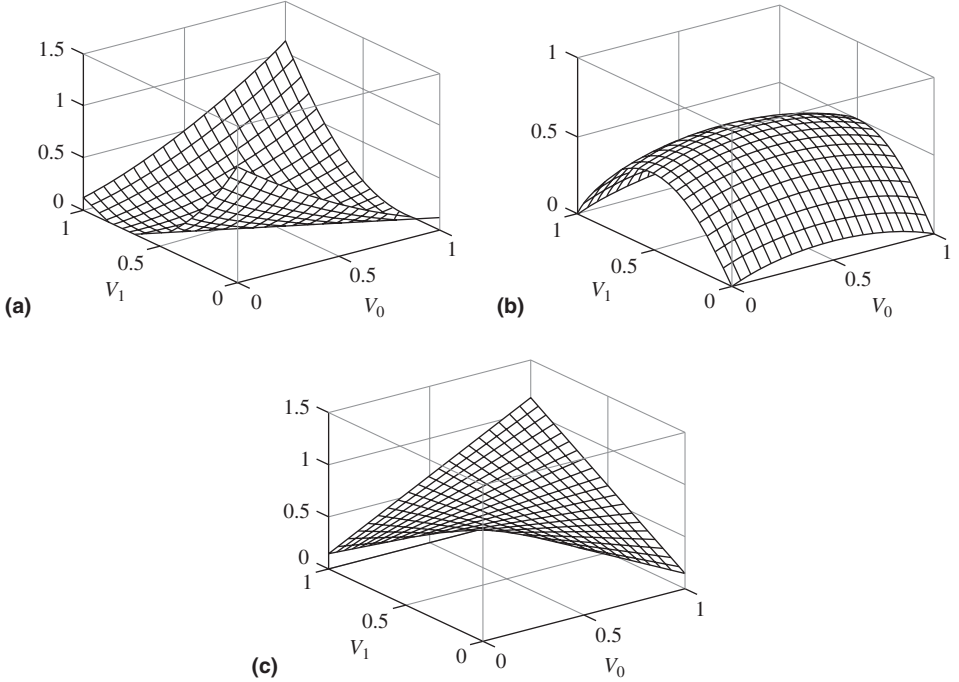


Figure 9.32 Constructing the cost function for a 2-bit A/D converter. Plots of the cost function and its components. (a) Plot of (9.107) for an analog input of $x = 1.5$. (b) Plot of the term used to cancel the quadratic terms in (9.107). (c) Plot of the cost function (9.109).

negative. We first consider a general case of an N -bit A/D converter. After this, we analyze in detail state equations of a 2-bit A/D converter. Note that

$$\frac{dE}{dt} = \frac{\partial E}{\partial V_0} \dot{V}_0 + \frac{\partial E}{\partial V_1} \dot{V}_1 + \cdots + \frac{\partial E}{\partial V_{N-1}} \dot{V}_{N-1},$$

where $V_i = g_i(u_i)$. Hence, similarly as in (9.105), it is enough to set

$$\frac{\partial E}{\partial V_i} = -C_i \dot{u}_i = -C_i \frac{dg_i^{-1}(V_i)}{dV_i} \dot{V}_i, \quad i = 0, 1, \dots, N-1 \quad (9.110)$$

to obtain

$$\frac{dE}{dt} = - \sum_{i=0}^{N-1} C_i \frac{dg_i^{-1}(V_i)}{dV_i} \dot{V}_i^2 \leq 0, \quad (9.111)$$

whereas in (9.106) we have $\frac{dg_i^{-1}(V_i)}{dV_i} > 0$. Thus, if the state equations of an N -bit A/D converter satisfy (9.110), then $dE/dt \leq 0$ on the trajectories of the converter. Using (9.110), we can write the state equations that model the 2-bit A/D converter consisting of two neurons. Let $V_0 = g_0(u_0)$ and $V_1 = g_1(u_1)$. Then, substituting (9.109) into (9.110) and performing required

differentiation, we obtain

$$\begin{aligned} C_0 \dot{u}_0 &= -\frac{\partial E}{\partial V_0} = -(2V_1 - x + 0.5) = -(2g_1(u_1) - x + 0.5), \\ C_1 \dot{u}_1 &= -\frac{\partial E}{\partial V_1} = -2(V_0 - x + 1) = -2(g_0(u_0) - x + 1). \end{aligned} \quad (9.112)$$

For simplicity, we assume that $C_0 = C_1 = 1$ and that $g_0 = g_1 = g$. Note that because C is present in both state equations, assigning $C_0 = C_1 = 1$ does not distort the dynamical behavior of the network; it only changes the time scale of both equations by the same factor. Observing that for $i = 1, 2$,

$$\dot{u}_i = \frac{dg_i^{-1}(V_i)}{dV_i} \dot{V}_i,$$

we represent (9.112) as

$$\begin{aligned} \dot{V}_0 &= -(2V_1 - x + 0.5) \left(\frac{dg^{-1}(V_0)}{dV_0} \right)^{-1}, \\ \dot{V}_1 &= -2(V_0 - x + 1) \left(\frac{dg^{-1}(V_1)}{dV_1} \right)^{-1}. \end{aligned} \quad (9.113)$$

We display a plot of the cost function E for the analog value $x = 1.5$ in Figure 9.33(a) and a phase portrait of the 2-bit A/D converter in Figure 9.33(b). We chose the activation function

$$V_i = g(u_i) = \frac{1}{1 + e^{-4u_i}}.$$

Using (9.98), we find

$$u_i = g^{-1}(V_i) = \frac{1}{4}(\ln(V_i) - \ln(1 - V_i)).$$

Thus,

$$\frac{dg^{-1}(V_i)}{dV_i} = \frac{1}{4V_i(1 - V_i)}.$$

Note that for $x = 1.5$, the cost function has two minimizers in the hypercube $[0, 1]^2$ that are separated by a “ridge.” Initial conditions on each side of this ridge yield a different final value of the state vector $[V_0, V_1]$. Depending on the initial conditions of the state vector of the A/D converter, the analog value $x = 1.5$ will be converted to either the binary “01” or the binary “10.” Thus, initial conditions influence the dynamical behavior of the A/D converter. Also, the activation nonlinearity influences the dynamical behavior of the A/D converter via the multiplicative term $(\frac{dg^{-1}(V_i)}{dV_i})^{-1}$ present in the state equations (9.113). One can see that the multiplicative factor influences the direction in which the converter state trajectory evolves.

Employing the above technique, we can construct higher-order A/D converters. For example, to construct a 3-bit A/D converter, we begin with the cost function

$$E = \frac{1}{2}(V_2^2 + V_1^2 + V_0^2 - x)^2 + aV_2(V_2 - 1) + bV_1(V_1 - 1) + cV_0(V_0 - 1).$$

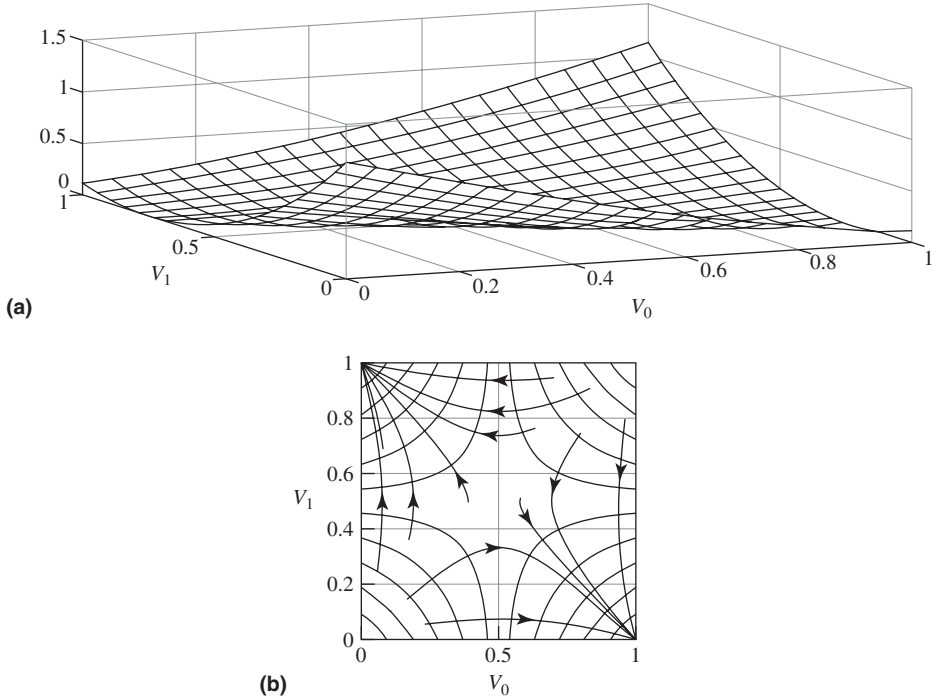


Figure 9.33 A plot of the cost function and a phase portrait of the 2-bit A/D converter for the analog value $x = 1.5$.

To cancel the quadratic terms, select $a = -8$, $b = -2$, $c = -1/2$. The state equations of the converter are obtained from the equations

$$\dot{V}_i = -\frac{1}{C_i} \frac{\partial E}{\partial V_i} \left(\frac{dg_i^{-1}(V_i)}{dV_i} \right)^{-1}, \quad i = 0, 1, 2.$$

In our analysis of A/D neural converters, we observed that while the 2-bit converter yielded correct results, the 3-bit and 4-bit converters showed some deficiency in their performance. The errors observed for these two converters were quite significant and increased with the number of bits and consequently with network size.

9.8 Hopfield Network Stability Analysis

In this section we apply Lyapunov's second method to analyze stability of the Hopfield neural network, in the sense of Lyapunov. The neural network is viewed as a large interconnected dynamical system—that is, as an interconnection of many neurons. First we analyze stability properties of a single neuron, and then we proceed with the stability analysis of the whole network in terms of the properties of the single neurons. Interconnection information is used in conjunction with the properties of the individual neurons to obtain sufficient conditions for asymptotic stability of the whole network.

9.8.1 Hopfield Network Model Analysis

Recall the equations describing the Hopfield neural net. They are given by (9.89), that is,

$$C_i \frac{du_i}{dt} + \frac{u_i}{r_i} = \sum_{j=1}^n T_{ij} x_j + I_i, \quad i = 1, 2, \dots, n,$$

where $x_i = g_i(u_i)$. In our further analysis, we assume that the activation functions $g_i : \mathbb{R} \rightarrow (-1, 1)$, $i = 1, 2, \dots, n$, are continuously differentiable and strictly monotonically increasing. We also assume that

$$u_i g_i(u_i) > 0 \quad \text{for } u_i \neq 0$$

and

$$g_i(0) = 0.$$

To proceed further, we divide both sides of the above equations describing the Hopfield network by $C_i > 0$ to get

$$\dot{u}_i = -b_i u_i + \sum_{j=1}^n A_{ij} g_j(u_j) + U_i, \quad i = 1, 2, \dots, n, \quad (9.114)$$

where $b_i = 1/r_i C_i$, $A_{ij} = T_{ij}/C_i$, and $U_i = I_i/C_i$. We are interested in the behavior of solutions of the above model near its equilibrium points. For $U_i = 0$ the equilibrium points are obtained by solving the set of algebraic equations

$$-b_i u_i + \sum_{j=1}^n A_{ij} g_j(u_j) = 0, \quad i = 1, 2, \dots, n.$$

In some applications the external inputs U_i are held constant, that is, $U_i = c_i = \text{constant}$. In such cases the equilibrium points are obtained by solving for u_i the set of the following algebraic equations:

$$-b_i u_i + \sum_{j=1}^n A_{ij} g_j(u_j) + c_i = 0, \quad i = 1, 2, \dots, n.$$

However, we can always transform the algebraic equations with nonzero constant external inputs into the algebraic equations where the external inputs are set to zero. Indeed, for $i = 1, 2, \dots, n$ let

$$u_i = y_i + \frac{c_i}{b_i}$$

and

$$G_i(y_i) = g_i\left(y_i + \frac{c_i}{b_i}\right).$$

Substituting the above into (9.114) and assuming that $U_i = c_i$ yields

$$\dot{y}_i = -b_i y_i + \sum_{j=1}^n A_{ij} G_j(y_j), \quad i = 1, 2, \dots, n. \quad (9.115)$$

The equilibrium points of the above system of differential equations are obtained by solving the following algebraic equations:

$$\begin{aligned} -b_i y_i + \sum_{j=1}^n A_{ij} G_j(y_j) &= -b_i y_i + \sum_{j=1}^n A_{ij} g_j\left(y_j + \frac{c_j}{b_j}\right) \\ &= 0, \quad i = 1, 2, \dots, n. \end{aligned}$$

Thus, after the transformation, the algebraic equations with external constant inputs used for computing the equilibrium points assume the same form as the algebraic equations corresponding to zero external inputs. Hence, without loss of generality, we can assume in our further analysis that the external inputs are zero and that the network is modeled by the equations

$$\boxed{\dot{u}_i = -b_i u_i + \sum_{j=1}^n A_{ij} g_j(u_j), \quad i = 1, 2, \dots, n} \quad (9.116)$$

Throughout, we assume that an isolated equilibrium point of interest is located at the origin of the state space \mathbb{R}^n . If this is not the case, then by an appropriate change of variables we can arrange for the equilibrium point of interest to be translated to the origin of \mathbb{R}^n . Let $\mathbf{u} = [u_1 \cdots u_n]^T \in \mathbb{R}^n$ and suppose $\mathbf{u}^* \neq \mathbf{0}$ is an isolated equilibrium point of interest of (9.116). For $i = 1, 2, \dots, n$ let

$$x_i = u_i - u_i^* \quad \text{and} \quad h_i(u_i) = h_i(x_i + u_i^*) = g_i(u_i) - g_i(u_i^*).$$

Substituting the above into (9.116) yields

$$\begin{aligned} \dot{x}_i &= -b_i x_i - b_i u_i^* + \sum_{j=1}^n A_{ij} (h_j(x_j + u_j^*) + g_j(u_j^*)) \\ &= -b_i x_i + \sum_{j=1}^n A_{ij} h_j(x_j + u_j^*) + \left\{ -b_i u_i^* + \sum_{j=1}^n A_{ij} g_j(u_j^*) \right\}. \end{aligned}$$

Note that, by the definition of the equilibrium point, the expression in the curly brackets vanishes. Thus we have

$$\dot{x}_i = -b_i x_i + \sum_{j=1}^n A_{ij} h_j(x_j + u_j^*), \quad i = 1, 2, \dots, n. \quad (9.117)$$

Let

$$G_j(x_j) = h_j(x_j + u_j^*).$$

Then, differential equations (9.117) can be written as

$$\dot{x}_i = -b_i x_i + \sum_{j=1}^n A_{ij} G_j(x_j), \quad i = 1, 2, \dots, n,$$

where $G_j : \mathbb{R} \rightarrow (c_1, c_2) \subseteq (-2, 2)$, $c_1 < 0 < c_2$, $G_j(0) = 0$, G_j is continuously differentiable and strictly monotonically increasing in x_j . Furthermore $x_j G_j(x_j) > 0$ for all $x_j \neq 0$. We now

represent the above model of the Hopfield neural network in the form

$$\dot{x}_i = -b_i x_i + A_{ii} G_i(x_i) + \sum_{\substack{j=1 \\ j \neq i}}^n A_{ij} G_j(x_j), \quad i = 1, 2, \dots, n.$$

We shall view the terms

$$\sum_{\substack{j=1 \\ j \neq i}}^n A_{ij} G_j(x_j)$$

as the interconnections between the i th neuron and the remaining neurons in the network. We shall refer to the equations

$$\dot{x}_i = -b_i x_i + A_{ii} G_i(x_i), \quad i = 1, 2, \dots, n \quad (9.118)$$

as models of free isolated neurons. In what follows we analyze stability properties of an isolated single neuron.

9.8.2 Single-Neuron Stability Analysis

A single neuron in the Hopfield network, when the external input is set to zero, is modeled by a first-order differential equation:

$$\boxed{\dot{x}_i = -b_i x_i + A_{ii} G_i(x_i)} \quad (9.119)$$

Consider now, as in Michel et al. [203], the following positive definite Lyapunov function candidate

$$V_i(x_i) = \frac{1}{2} x_i^2.$$

Its time derivative evaluated along the trajectories of the single neuron model is

$$\begin{aligned} \dot{V}_i(x_i(t)) &= x_i \dot{x}_i \\ &= x_i (-b_i x_i + A_{ii} G_i(x_i)). \end{aligned}$$

If there exists $r_i > 0$ such that for all $x_i \in (-r_i, r_i)$

$$\begin{aligned} -b_i x_i + A_{ii} G_i(x_i) &> 0 & \text{if } x_i < 0, \\ -b_i x_i + A_{ii} G_i(x_i) &= 0 & \text{if } x_i = 0, \\ -b_i x_i + A_{ii} G_i(x_i) &< 0 & \text{if } x_i > 0, \end{aligned}$$

then for $x_i \neq 0$ and $x_i \in (-r_i, r_i)$ we have

$$\dot{V}_i(x_i(t)) < 0$$

and $\dot{V}_i(0) = 0$. Hence the null solution $x_i = 0$ is uniformly asymptotically stable. Figure 9.34 illustrates this case. Now, if there exists $r_i > 0$ such that for $x_i \in (-r_i, r_i)$ we have

$$\begin{aligned} -b_i x_i + A_{ii} G_i(x_i) &< 0 & \text{if } x_i < 0, \\ -b_i x_i + A_{ii} G_i(x_i) &= 0 & \text{if } x_i = 0, \\ -b_i x_i + A_{ii} G_i(x_i) &> 0 & \text{if } x_i > 0, \end{aligned}$$

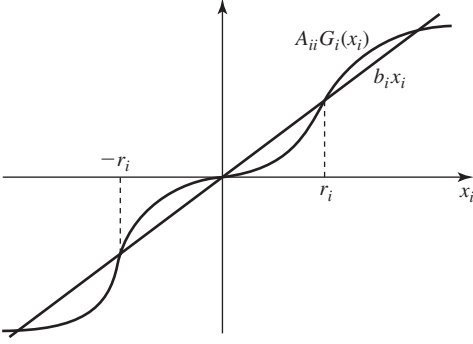


Figure 9.34 Characteristics of a single neuron model whose null solution is uniformly asymptotically stable.

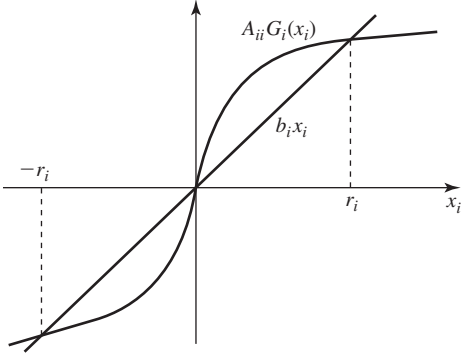


Figure 9.35 Characteristics of single-neuron model whose null solution is unstable.

then in addition to V_i being positive definite we have

$$\dot{V}_i(x_i(t)) > 0$$

for $x_i \neq 0$ and $x_i \in (-r_i, r_i)$, and $\dot{V}_i(0) = 0$. Thus the null solution $x_i = 0$ is unstable. Figure 9.35 illustrates the case when the null solution of a single neuron is unstable. Assume now that for some $r_i > 0$ there exist constants $\sigma_{i1} > 0$ and $\sigma_{i2} > 0$ such that for all $x_i \in (-r_i, r_i)$, $x_i \neq 0$ we have

$$\sigma_{i1} x_i^2 < x_i G_i(x_i) < \sigma_{i2} x_i^2, \quad (9.120)$$

or, equivalently,

$$\sigma_{i1} < \frac{G_i(x_i)}{x_i} < \sigma_{i2}.$$

If in addition to the sufficient conditions for asymptotic stability we also have

$$-b_i + A_{ii} \delta_i < 0, \quad (9.121)$$

where

$$\delta_i = \begin{cases} \sigma_{i1} & \text{if } A_{ii} < 0, \\ \sigma_{i2} & \text{if } A_{ii} > 0, \end{cases} \quad (9.122)$$

then the null solution of the single-neuron model is exponentially stable. Indeed, for $x_i \neq 0$ we have

$$\begin{aligned}
 \dot{V}_i(x_i(t)) &= x_i \dot{x}_i \\
 &= x_i(-b_i x_i + A_{ii} G_i(x_i)) \\
 &= \left(-b_i + A_{ii} \frac{G_i(x_i)}{x_i}\right) x_i^2 \\
 &\leq (-b_i + A_{ii} \delta_i) x_i^2 \\
 &= 2(-b_i + A_{ii} \delta_i) V_i(x_i(t)).
 \end{aligned}$$

Applying Theorem A.20 yields

$$\begin{aligned}
 V_i(x_i(t)) &= \frac{x_i^2(t)}{2} \\
 &\leq \exp\{2(-b_i + A_{ii} \delta_i)(t - t_0)\} V_i(x_i(t_0)) \\
 &= \exp\{2(-b_i + A_{ii} \delta_i)(t - t_0)\} \frac{x_i^2(t_0)}{2}.
 \end{aligned}$$

Hence

$$|x_i(t)| \leq \exp\{(-b_i + A_{ii} \delta_i)(t - t_0)\} |x_i(t_0)|.$$

From the above we conclude that the null solution is not only asymptotically stable but even more, it is exponentially stable.

9.8.3 Stability Analysis of the Network

Lyapunov's second method in conjunction with interconnection information will now be used in the stability analysis of the Hopfield neural model viewed as an interconnected dynamical system. Applications of Lyapunov's second method to a general class of interconnected systems were investigated by Bailey [16]. Here, we use Lyapunov's second method to investigate stability properties of the null solution of the neural network described by

$$\boxed{\dot{x}_i = -b_i x_i + \sum_{j=1}^n A_{ij} G_j(x_j), \quad i = 1, 2, \dots, n} \quad (9.123)$$

We impose the following assumptions on the interconnections:

A1: There exist positive constants r_i , r_j , and a_{ij} such that for $|x_i| < r_i$ and $|x_j| < r_j$, $i, j = 1, 2, \dots, n$, we have

$$x_i A_{ij} G_j(x_j) \leq |x_i| a_{ij} |x_j|.$$

A2: Let S be an $n \times n$ real matrix whose elements are

$$s_{ij} = \begin{cases} \alpha_i(-b_i + a_{ii}) & \text{if } i = j, \\ \frac{1}{2}(\alpha_i a_{ij} + \alpha_j a_{ji}) & \text{if } i \neq j. \end{cases}$$

We assume that there exists a vector $\alpha = [\alpha_1 \dots \alpha_n]^T$, where $\alpha_i > 0$ for $i = 1, 2, \dots, n$, such that the matrix S is negative definite.

We now state and prove a theorem concerned with the exponential stability of the neural network model. This theorem can be found in Michel et al. [203].

Theorem 9.5 The null solution of the neural network model (9.123) is uniformly exponentially stable if assumptions A1 and A2 are satisfied.

Proof Consider the following Lyapunov function candidate:

$$\begin{aligned} V(\mathbf{x}) &= \frac{1}{2} \sum_{i=1}^n \alpha_i x_i^2 \\ &= \frac{1}{2} \mathbf{x}^T \begin{bmatrix} \alpha_1 & \dots & 0 \\ & \alpha_2 & \\ \vdots & & \ddots & \vdots \\ 0 & \dots & & \alpha_n \end{bmatrix} \mathbf{x}, \end{aligned}$$

where $\alpha_i > 0$, $i = 1, 2, \dots, n$, are defined in assumption A2. Note that V is real symmetric positive definite. The time derivative of V evaluated along the trajectories of the neural model (9.123) is

$$\begin{aligned} \dot{V}(\mathbf{x}(t)) &= \sum_{i=1}^n \alpha_i x_i \dot{x}_i \\ &= \sum_{i=1}^n \alpha_i x_i \left(-b_i x_i + \sum_{j=1}^n A_{ij} G_j(x_j) \right) \\ &\leq \sum_{i=1}^n \alpha_i \left(-b_i x_i^2 + |x_i| \sum_{j=1}^n a_{ij} |x_j| \right) \\ &= |\mathbf{x}|^T \mathbf{R} |\mathbf{x}|, \end{aligned}$$

where $|\mathbf{x}| = [|x_1| \ |x_2| \ \dots \ |x_n|]^T$, and the matrix $\mathbf{R} \in \mathbb{R}^{n \times n}$ has its elements defined as

$$r_{ij} = \begin{cases} \alpha_i(-b_i + a_{ii}) & \text{if } i = j, \\ \alpha_i a_{ij} & \text{if } i \neq j. \end{cases}$$

Note that

$$\begin{aligned} |\mathbf{x}|^T \mathbf{R} |\mathbf{x}| &= |\mathbf{x}|^T \left(\frac{\mathbf{R} + \mathbf{R}^T}{2} \right) |\mathbf{x}| \\ &= |\mathbf{x}|^T \mathbf{S} |\mathbf{x}| \\ &\leq \lambda_{\max}(\mathbf{S}) \|\mathbf{x}\|^2, \end{aligned}$$

where the matrix \mathbf{S} is defined in assumption A2. By assumption A2 the matrix \mathbf{S} is real symmetric negative definite. Hence $\lambda_{\max}(\mathbf{S})$ is real and negative. Thus

$$\dot{V}(\mathbf{x}(t)) \leq \lambda_{\max}(\mathbf{S}) \|\mathbf{x}(t)\|^2 < 0$$

in a neighborhood of $\mathbf{x} = \mathbf{0}$ of radius

$$r = \min_{1 \leq i \leq n} r_i.$$

By Proposition 4.1, the null solution is uniformly exponentially stable.

In our stability analysis of the Hopfield neural network we do not require that the parameters A_{ij} , and hence the interconnections T_{ij} , be symmetric. The symmetry requirement on the interconnections was usually imposed in the early stability analyses of this network.

Note also that a necessary condition for the matrix S to be negative definite is that its diagonal elements be negative, that is,

$$-b_i + a_{ii} < 0, \quad i = 1, 2, \dots, n. \quad (9.124)$$

However, this condition is not sufficient for the matrix S to be negative definite. The above necessary condition, though, can shed some light on the influence of the stability properties of the individual neurons on the stability of the whole network. To show how the stability of the individual neurons may affect the stability of the whole network, we set $i = j$ in A1 to obtain

$$x_i A_{ii} G_i(x_i) \leq a_{ii} |x_i|^2. \quad (9.125)$$

Applying now the sector condition, $\sigma_{i1} x_i^2 < x_i G_i(x_i) < \sigma_{i2} x_i^2$, to $x_i A_{ii} G_i(x_i)$ gives

$$x_i A_{ii} G_i(x_i) \leq A_{ii} \delta_i |x_i|^2, \quad (9.126)$$

where δ_i is defined in (9.122). Thus the sector condition (9.120) is sufficient for (9.125) to hold and we can write

$$a_{ii} = A_{ii} \delta_i.$$

Substituting the above into (9.124) yields

$$-b_i + A_{ii} \delta_i < 0, \quad i = 1, 2, \dots, n,$$

which are sufficient conditions for the exponential stability of the individual neurons. Thus, if we first find constants a_{ii} , $i = 1, 2, \dots, n$, and then use them to construct the matrix S , then for so-formed matrix S to be negative definite it is necessary that the above sufficient conditions for the exponential stability of the individual neurons hold.

We now present a corollary to Theorem 9.5. Before stating the corollary, we need to perform some algebraic manipulations. We first combine assumption A2 with the sector condition $\sigma_{i1} x_i^2 < x_i G_i(x_i) < \sigma_{i2} x_i^2$. It follows from the sector condition that for $x_j \neq 0$,

$$\left| \frac{G_j(x_j)}{x_j} \right| < \sigma_{j2}.$$

Applying the above when $i \neq j$ yields

$$\begin{aligned} x_i A_{ij} G_j(x_j) &\leq |A_{ij}| |x_i G_j(x_j)| \\ &= |A_{ij}| \left| x_i \frac{G_j(x_j)}{x_j} x_j \right| \\ &\leq |A_{ij}| |x_i| |x_j| \sigma_{j2}. \end{aligned}$$

When $i = j$ we obtain

$$x_i A_{ii} G_i(x_i) \leq \delta_i A_{ii} x_i^2,$$

where δ_i is defined as above. We use the results of the above manipulations in the following assumption:

A3: There exist positive constants r_i and r_j such that for $|x_i| < r_i$ and $|x_j| < r_j$, the interconnections of the neural network model (9.123) satisfy the conditions

$$\begin{aligned} x_i A_{ii} G_i(x_i) &\leq \delta_i A_{ii} x_i^2 & \text{for } i = j, \\ x_i A_{ij} G_j(x_j) &\leq |x_i| |A_{ij}| |x_j| \sigma_{j2} & \text{for } i \neq j. \end{aligned}$$

We shall use another assumption.

A4: There exists a vector $\alpha = [\alpha_1 \cdots \alpha_n]^T \in \mathbb{R}^n$, $\alpha_i > 0$, such that the matrix $S \in \mathbb{R}^{n \times n}$ is negative definite, where the elements of S are defined as

$$s_{ij} = \begin{cases} \alpha_i(-b_i + \delta_i A_{ii}) & \text{if } i = j, \\ \frac{1}{2}(\alpha_i |A_{ij}| \sigma_{j2} + \alpha_j |A_{ji}| \sigma_{i2}) & \text{if } i \neq j. \end{cases}$$

We are now ready to state and prove a corollary to Theorem 9.5. This result also comes from Michel et al. [203].

Corollary 9.1 The null solution of the neural network model (9.123) is uniformly exponentially stable if the assumptions A3 and A4 are satisfied.

Proof Consider the following quadratic positive definite Lyapunov function candidate

$$V(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^n \alpha_i x_i^2,$$

where $\alpha_i > 0$, $i = 1, 2, \dots, n$, are defined in assumption A4. Taking into account the assumptions A3 and A4, we evaluate the time derivative of V along the trajectories of the neural network (9.123) to get

$$\begin{aligned} \dot{V}(\mathbf{x}(t)) &= \sum_{i=1}^n \alpha_i x_i \left(-b_i x_i + \sum_{j=1}^n A_{ij} G_j(x_j) \right) \\ &= \sum_{i=1}^n \alpha_i x_i \left(-b_i x_i + A_{ii} G_i(x_i) + \sum_{\substack{j=1 \\ j \neq i}}^n A_{ij} G_j(x_j) \right) \\ &\leq \sum_{i=1}^n \alpha_i \left((-b_i + \delta_i A_{ii}) x_i^2 + \sum_{\substack{j=1 \\ j \neq i}}^n |A_{ij}| |x_i| |x_j| \sigma_{j2} \right) \\ &= \frac{1}{2} |\mathbf{x}|^T \mathbf{S} |\mathbf{x}| \\ &< 0 \quad \text{for } \mathbf{x} \neq \mathbf{0}. \end{aligned}$$

Using the above, the fact that V is a quadratic function, and Proposition 4.1, we conclude that the null solution is uniformly exponentially stable.

We illustrate the above corollary with a numerical example involving a neural network consisting of two neurons.

◆ Example 9.8

The neural net that we analyze in this example is described by the equations

$$\begin{bmatrix} \dot{u}_1 \\ \dot{u}_2 \end{bmatrix} = \begin{bmatrix} -1.4u_1 + 1.2g_2(u_2) \\ -1.1u_2 + g_1(u_1) \end{bmatrix}, \quad (9.127)$$

where

$$g_i(u_i) = \frac{2}{\pi} \arctan\left(\lambda \frac{\pi}{2} u_i\right), \quad i = 1, 2. \quad (9.128)$$

The value of the parameter λ was set to be equal to 1.4. In Figure 9.36, we show a plot of the above activation nonlinearity. The network has three equilibrium states:

$$\mathbf{u}^{(1)} = \begin{bmatrix} -0.4098 \\ -0.4245 \end{bmatrix}, \quad \mathbf{u}^{(2)} = \begin{bmatrix} 0.4098 \\ 0.4245 \end{bmatrix}, \quad \text{and} \quad \mathbf{u}^{(3)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

When computing the equilibrium points, we began by finding the first coordinates of them. This can be done, for example, by combining the two algebraic equations

$$-1.4u_1 + 1.2g_2(u_2) = 0 \quad \text{and} \quad -1.1u_2 + g_1(u_1) = 0$$

and then solving them for u_1 . A plot of the function obtained by substituting the second of the above equations into the first is shown in Figure 9.37. We can use Lyapunov's first method, as we did in Example 4.8, to show that the equilibrium states $\mathbf{u}^{(1)}$ and $\mathbf{u}^{(2)}$ are asymptotically stable, while the equilibrium $\mathbf{u}^{(3)}$ is unstable.

We now apply Corollary 9.1 to show that the equilibrium states $\mathbf{u}^{(1)}$ and $\mathbf{u}^{(2)}$ are in fact uniformly exponentially stable. We analyze the equilibrium state $\mathbf{u}^{(1)}$; the analysis of $\mathbf{u}^{(2)}$ is identical. We first translate the equilibrium state of interest

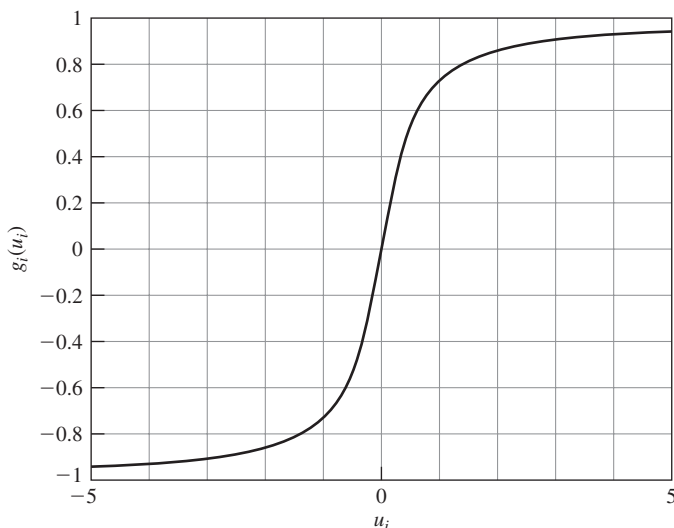


Figure 9.36 A plot of the activation nonlinearity (9.128).

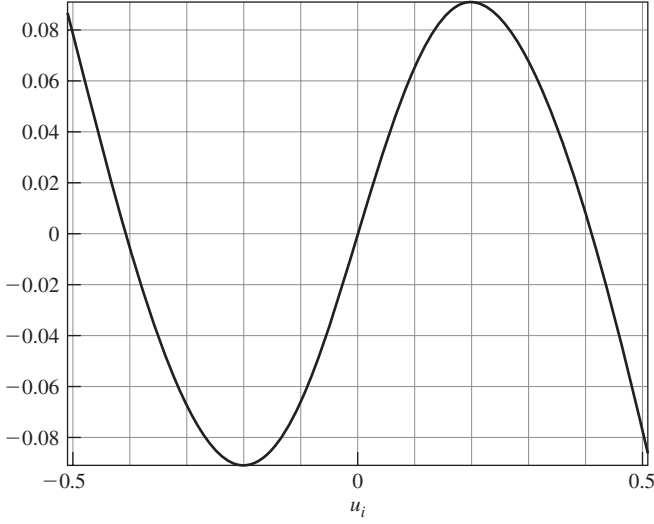


Figure 9.37 A plot of the function used to find the first coordinates of the equilibrium points of the neural net of Example 9.8.

to the origin of \mathbb{R}^2 using the translation

$$\mathbf{x} = \mathbf{u} - \mathbf{u}^{(1)}$$

Let

$$G_i(x_i) = g_i(x_i + u_i^{(1)}) - g_i(u_i^{(1)}), \quad i = 1, 2.$$

Taking into account the above relation, the neural network model in the new coordinates takes the form

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} -1.4x_1 + 1.2G_2(x_2) \\ -1.1x_2 + G_1(x_1) \end{bmatrix}. \quad (9.129)$$

The equilibrium state $\mathbf{x}_e = \mathbf{0}$ of the above net corresponds to the equilibrium state $\mathbf{u}_e = \mathbf{u}^{(1)}$ of the net in the u coordinates. To be able to apply Corollary 9.1, we need to find σ_{i2} for $i = 1, 2$ such that

$$0 < \frac{G_i(x_i)}{x_i} < \sigma_{i2} \quad \text{for } |x_i| < r,$$

where $r > 0$. In this example, we took $r = \min\{|u_1^{(1)}|, |u_2^{(1)}|\} = 0.4098$. Graphs of $G_i(x_i)/x_i$ versus x_i , for $i = 1, 2$, are shown in Figure 9.38. We calculated $\sigma_{12} = 1.1394$ and $\sigma_{22} = 1.1165$. We next formed the matrix $\mathbf{S} \in \mathbb{R}^{2 \times 2}$ defined in A4. Selecting $\alpha_1 = \alpha_2 = 1$ and taking into account that $A_{11} = 0$, $A_{12} = 1.4$, and $A_{21} = 1$, we obtain

$$\begin{aligned} \mathbf{S} &= \begin{bmatrix} -1.4\alpha_1 & (\alpha_1|A_{12}|\sigma_{22} + \alpha_2|A_{21}|\sigma_{12})/2 \\ (\alpha_2|A_{21}|\sigma_{12} + \alpha_1|A_{12}|\sigma_{22})/2 & -1.1\alpha_2 \end{bmatrix} \\ &= \begin{bmatrix} -1.4000 & 1.2396 \\ 1.2396 & -1.1000 \end{bmatrix}. \end{aligned}$$

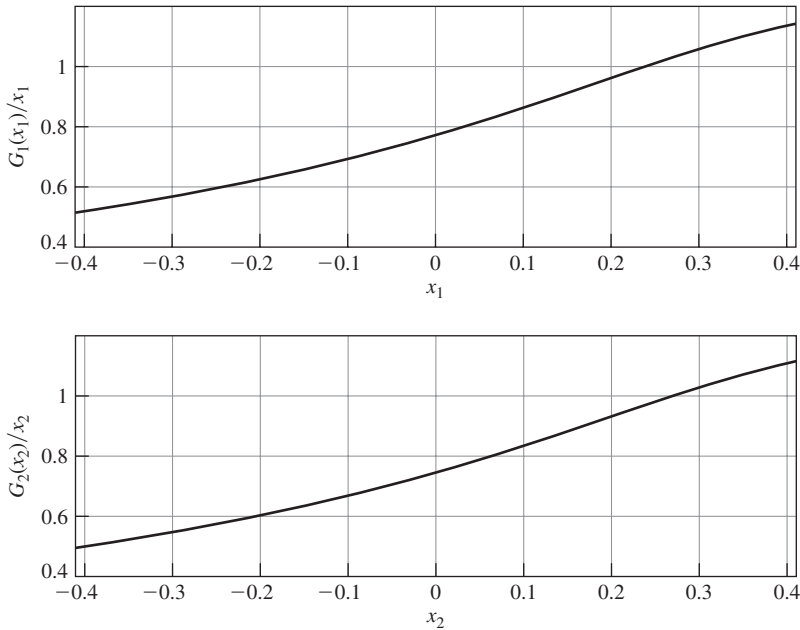


Figure 9.38 Plots of $G_i(x_i)/x_i$, $i = 1, 2$, in Example 9.8.

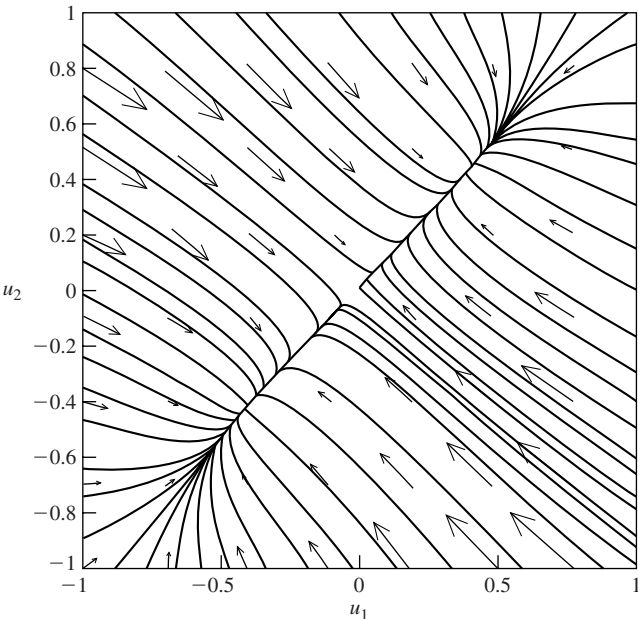


Figure 9.39 Phase plane portrait of the neural net of Example 9.8.

The eigenvalues of S are located at -2.4987 and -0.0013 . Hence, the matrix S is negative definite, and therefore the equilibrium state $\mathbf{x}_e = \mathbf{0}$ of the net (9.129), and hence $\mathbf{u}_e = \mathbf{u}^{(1)}$ of the net (9.127), is uniformly exponentially stable. In Figure 9.39, we show a state plane portrait of the neural network modeled by (9.127).

9.9 Brain-State-in-a-Box (BSB) Models

9.9.1 Associative Memories

In this section we analyze a class of neural network models that can be used to construct the so-called associative memories. We distinguish between two types of association: autoassociation and heteroassociation. In autoassociation, after storing a set of patterns by a neural network, where patterns are usually represented as vectors, a distorted (noisy) pattern of a stored pattern is subsequently presented to the network. The task of the neural network is to retrieve (recall) the original stored pattern from the noisy pattern. In heteroassociation an arbitrary set of input patterns is paired with another arbitrary set of output patterns. Operation of an associative memory is characterized by two stages: storage phase, where patterns are being stored by the neural network, and recall phase, where memorized patterns are being retrieved in response to a noisy pattern presented to the network. Figure 9.40 illustrates the recall phase, where the input pattern represents a noisy version of stored pattern \mathbf{x}_j . The pattern associator produces the output pattern \mathbf{y} in response to the input pattern \mathbf{x} . For a perfect recall in an autoassociative memory, we should have $\mathbf{y} = \mathbf{x}_j$. When $\mathbf{y} \neq \mathbf{x}_j$, the associative memory is said to have made a recall error (or error in recall). We now analyze a static linear autoassociator. For the sake of simplicity, we assume that the given patterns, $\mathbf{v}^{(j)}$, $j = 1, \dots, r$, are orthogonal, that is, $\langle \mathbf{v}^{(j)}, \mathbf{v}^{(k)} \rangle = 0$ for $k \neq j$. Suppose that we form the weight matrix as the sum of the outer products of the given patterns, that is,

$$\mathbf{W} = \sum_{j=1}^r \mathbf{v}^{(j)} \mathbf{v}^{(j)T} = \mathbf{V} \mathbf{V}^T,$$

where $\mathbf{V} = [\mathbf{v}^{(1)} \ \mathbf{v}^{(2)} \ \dots \ \mathbf{v}^{(r)}]$. Let \mathbf{v} be one of the given patterns. We present \mathbf{v} to the linear associator. Then, because the patterns $\mathbf{v}^{(j)}$, $j = 1, \dots, r$, are mutually orthogonal, the output of the associator is

$$\mathbf{W} \mathbf{v} = \sum_{j=1}^r \mathbf{v}^{(j)} \mathbf{v}^{(j)T} \mathbf{v} = \mathbf{v} \mathbf{v}^T \mathbf{v} = c \mathbf{v},$$

where $c = \mathbf{v}^T \mathbf{v}$. Thus, we obtained back the presented pattern.

We will now show that our linear associator can reconstruct the missing part of a stored pattern. Assume that

$$\mathbf{v} = \mathbf{v}' + \mathbf{v}'';$$



Figure 9.40 Operation of pattern associator.

that is, the given pattern \mathbf{v} is composed of two parts. Suppose that \mathbf{v}' and \mathbf{v}'' are orthogonal to each other and that \mathbf{v}' is also orthogonal to the remaining patterns in the set. Taking this into account while presenting \mathbf{v}' to our linear associator gives

$$\begin{aligned} \mathbf{W}\mathbf{v}' &= \sum_{j=1}^r \mathbf{v}^{(j)} \mathbf{v}^{(j)T} \mathbf{v}' \\ &= \mathbf{v} \mathbf{v}^T \mathbf{v}' \\ &= (\mathbf{v}' + \mathbf{v}'')(\mathbf{v}' + \mathbf{v}'')^T \mathbf{v}' \\ &= (\mathbf{v}' + \mathbf{v}'') \mathbf{v}'^T \mathbf{v}' \\ &= d \mathbf{v}, \end{aligned}$$

where $d = \mathbf{v}'^T \mathbf{v}'$. Thus the linear associator was able to reconstruct the missing part of the presented pattern. In the above discussion, we assumed that the patterns are mutually orthogonal. If this is not the case, the output of the associator will contain an extra component called “cross-talk.” This linear static associator has very limited capability. However, we shall use the idea of the outer product to construct neural associative memories.

A neural network that can be used as an associative memory is the Brain-State-in-a-Box (BSB) neural model, which was proposed by Anderson et al. [9] in 1977. It can recognize a pattern from its given noisy version. For this reason the BSB neural model is often referred to as an associative memory. The BSB neural model is a discrete dynamical system whose dynamics are described by the difference equation

$$\mathbf{x}(k+1) = \mathbf{g}(\mathbf{x}(k) + \alpha \mathbf{W} \mathbf{x}(k)) \quad (9.130)$$

subject to the initial condition

$$\mathbf{x}(0) = \mathbf{x}_0,$$

where $\mathbf{x}(k) \in \mathbb{R}^n$ is the state of the BSB network at time k , $\alpha > 0$ is the step size, and $\mathbf{W} \in \mathbb{R}^{n \times n}$ is a symmetric weight matrix determined during the training procedure. The function $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is a vector-valued function whose i th component $(\mathbf{g})_i$ is defined as

$$(\mathbf{g}(\mathbf{x}))_i = \begin{cases} 1 & \text{if } x_i \geq 1 \\ x_i & \text{if } -1 < x_i < 1 \\ -1 & \text{if } x_i \leq -1. \end{cases}$$

The above nonlinear activation function is sometimes referred to as the linear saturating nonlinearity. A plot of such a function is shown in Figure 9.41. The BSB model is named from the fact that the network trajectories are constrained to be in the hypercube $H_n = [-1 \ 1]^n$. In Figure 9.42 we show a state plane portrait of the BSB model (9.130) with the weight matrix

$$\mathbf{W} = \begin{bmatrix} 7 & -3 \\ -3 & 10 \end{bmatrix} \quad (9.131)$$

and the step size is $\alpha = 0.1$. The state space of this two-dimensional neural net is the square $H_2 = [-1 \ 1]^2$. Each vertex of the square has its basin of attraction. In this example all four vertices are asymptotically stable equilibrium states of the BSB neural net. In Figure 9.43 we

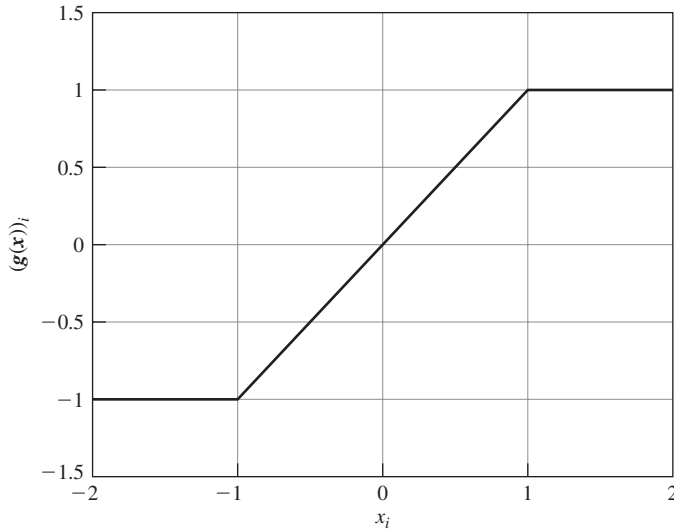


Figure 9.41 A plot of linear saturating nonlinearity used in the BSB model as an activation function.

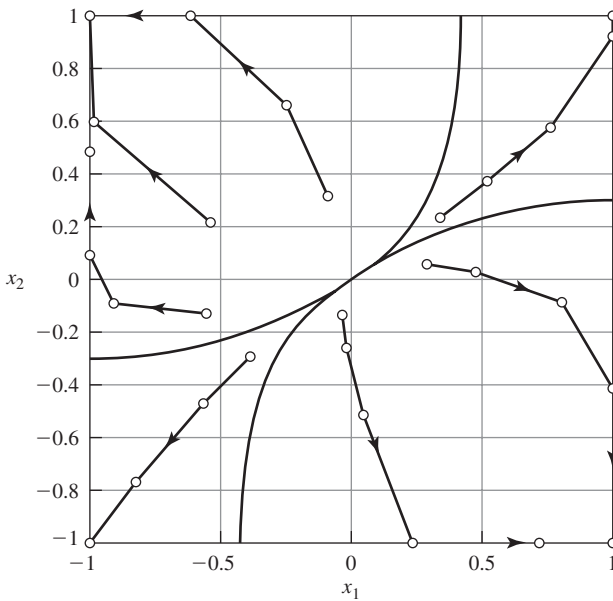


Figure 9.42 Trajectories of the BSB model with the weight matrix W given by (9.131).

show a state plane portrait of the BSB model (9.130) with the weight matrix

$$W = \begin{bmatrix} 5 & -3 \\ -3 & 2 \end{bmatrix} \quad (9.132)$$

and the step size $\alpha = 0.1$. In this example, only two vertices are asymptotically stable equilibrium states of the BSB neural net. When the BSB model is used as an associative memory,

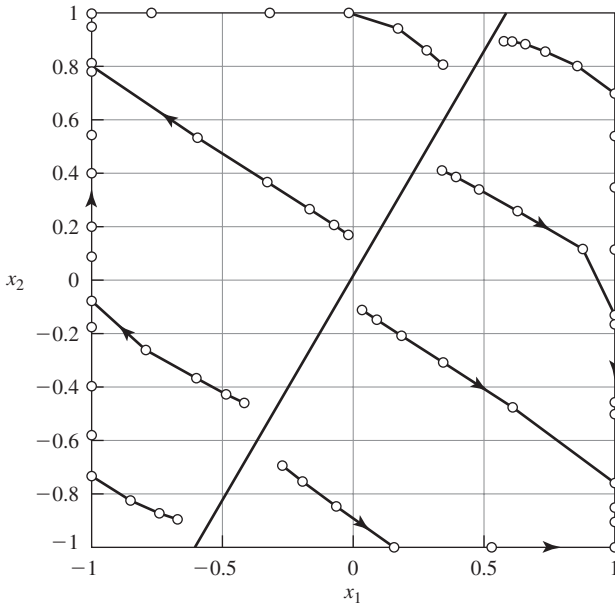


Figure 9.43 Trajectories of the BSB model with the weight matrix W given by (9.132).

asymptotically stable equilibrium states of the model represent stored patterns. Noisy patterns are represented by the states in the basins of attraction of the stable equilibrium states. Thus, suppose we are given a noisy, or imperfect, pattern represented by an initial condition \mathbf{x}_0 that is located in a basin of attraction of a stable equilibrium state corresponding to the correct pattern. Then, the corresponding trajectory of the BSB model evolves in time to this stable equilibrium state representing the correct stored pattern. An associative memory is able to recall all the information stored in the memory when it is excited with a sufficient large portion of that memory's information. This is in contrast to the address addressable memory (AAM) used in digital computers. In this type of memory, a desired set of information is recalled only when the entire correct address of the information is given. For this reason, we can view associative memories as content addressable memories (CAM).

In summary, realization of associative memories with neural networks requires storage of a set of given patterns as asymptotically stable equilibrium states of the network. The states in basins of attraction about each stable equilibrium state correspond to noisy versions of the stored patterns. When a noisy version of a stored pattern is sufficiently close to it, the network trajectory should converge to the equilibrium state corresponding to the correct pattern. Neural network associative memories should possess the following useful properties:

1. Each given pattern to be stored should be stored as an asymptotically stable equilibrium state of the network.
2. The number of asymptotically stable equilibrium states of the network that do not correspond to the given stored patterns—that is, the number of spurious states—should be minimal.
3. A nonsymmetric interconnection structure, which eases difficulties in implementing associative memory using neural networks.

4. The ability to control the extent of the basin of attraction about the equilibrium state corresponding to each stored pattern.
5. Learning capability, that is, the ability to add patterns to be stored as asymptotically stable equilibrium states to an existing set of stored patterns without affecting the existing equilibria.
6. Forgetting capability, that is, the ability to delete specified equilibria from a given set of stored equilibria without affecting the rest of the equilibria.
7. A high storage and retrieval efficiency, that is, the ability to efficiently store and retrieve a large number (compared to the order n of the network) of patterns.

9.9.2 Analysis of BSB Models

The BSB neural net has a symmetric weight matrix, which is not a desirable property of a neural associative memory. Hui and Žak [130] modified the BSB model to allow for nonsymmetric weight matrix as well as to better control the extent of basins of attraction. A modified form of the BSB neural model, which we refer to as the generalized Brain-State-in-a-Box (gBSB) neural model, is well-suited for implementing associative memories. Its dynamics are described by the equation

$$\mathbf{x}(k+1) = \mathbf{g}((\mathbf{I}_n + \alpha \mathbf{W})\mathbf{x}(k) + \alpha \mathbf{b}), \quad (9.133)$$

where \mathbf{I}_n is the $n \times n$ identity matrix, $\mathbf{b} \in \mathbb{R}^n$, and the weight matrix $\mathbf{W} \in \mathbb{R}^{n \times n}$ need not be symmetric as in the BSB model. Thus, the original BSB model can be viewed as a special case of the gBSB model when the interconnection matrix \mathbf{W} is symmetric and $\mathbf{b} = \mathbf{0}$. We add that the analysis of the BSB model behavior on the boundary regions of the hypercube H_n can be reduced to the study of the reduced-order gBSB-type models rather than that of the BSB type.

We now describe an algorithm for computing the weight matrix \mathbf{W} and the vector \mathbf{b} that result in the given patterns being stored as asymptotically stable vertices of the hypercube H_n . In further discussion, we use the following notation and definitions. Let

$$\mathbf{L}(\mathbf{x}) = (\mathbf{I}_n + \alpha \mathbf{W})\mathbf{x} + \alpha \mathbf{b}$$

and

$$\mathbf{T}(\mathbf{x}) = \mathbf{g}(\mathbf{L}(\mathbf{x})).$$

Definition 9.1 A point $\mathbf{x} \in H_n$ is an equilibrium state of $\mathbf{x}(k+1) = \mathbf{T}(\mathbf{x}(k))$ if $\mathbf{x} = \mathbf{T}(\mathbf{x})$.

Note that an equilibrium state is a constant solution to the system $\mathbf{x}(k+1) = \mathbf{T}(\mathbf{x}(k))$.

Let $(\mathbf{L}(\mathbf{x}))_i$ be the i th component of $\mathbf{L}(\mathbf{x})$ and let $v_i^{(j)}$ be the i th component of the vector $\mathbf{v}^{(j)}$. Let $\mathbf{v} = [v_1 \cdots v_n]^T \in \{-1, 1\}^n$; that is, \mathbf{v} is a vertex of H_n because for $i = 1, 2, \dots, n$ we have $v_i = \pm 1$. For a vertex \mathbf{v} of H_n to be an equilibrium state of the gBSB net, we must have

$$v_i = (\mathbf{g}(\mathbf{L}(\mathbf{v})))_i, \quad i = 1, 2, \dots, n;$$

that is, if $v_i = 1$, then we should have

$$(\mathbf{L}(\mathbf{v}))_i \geq 1, \quad (9.134)$$

and for $v_i = -1$ we should have

$$(\mathbf{L}(\mathbf{v}))_i \leq -1. \quad (9.135)$$

Combining (9.134) and (9.135), we obtain the following necessary and sufficient conditions for the vertex \mathbf{v} to be an equilibrium state.

Theorem 9.6 A vertex \mathbf{v} of the hypercube H_n is an equilibrium state of the gBSB neural net if and only if

$$(\mathbf{L}(\mathbf{v}))_i v_i \geq 1, \quad i = 1, 2, \dots, n.$$

Equivalently, a vertex \mathbf{v} of the hypercube H_n is not an equilibrium state of the gBSB neural net if and only if

$$(\mathbf{L}(\mathbf{v}))_i v_i < 1 \quad \text{for some } i = 1, 2, \dots, n.$$

◆ Example 9.9

We verify that the vertex $\mathbf{v}^{(1)} = [1 \ 1]^T$ of H_2 is not an equilibrium state of the BSB net with the weight matrix given by (9.132) and the step size $\alpha = 0.1$. Note that for the BSB net we have

$$\mathbf{L}(\mathbf{x}) = (\mathbf{I}_n + \alpha \mathbf{W}) \mathbf{x}.$$

We first compute $\mathbf{L}(\mathbf{v}^{(1)})$ to get

$$\mathbf{L}(\mathbf{v}^{(1)}) = (\mathbf{I}_2 + \alpha \mathbf{W}) \mathbf{v}^{(1)} = \begin{bmatrix} 1.2 \\ 0.9 \end{bmatrix}.$$

Note that because

$$(\mathbf{L}(\mathbf{v}^{(1)}))_2 v_2^{(1)} = 0.9 < 1,$$

the vertex $\mathbf{v}^{(1)}$ is not an equilibrium state of the above BSB net.

On the other hand, the vertex $\mathbf{v}^{(2)} = [1 \ -1]^T$ is an equilibrium state of this BSB model because

$$(\mathbf{L}(\mathbf{v}^{(2)}))_1 v_1^{(2)} = 1.8 \geq 1 \quad \text{and} \quad (\mathbf{L}(\mathbf{v}^{(2)}))_2 v_2^{(2)} = 1.5 \geq 1;$$

that is, the necessary and sufficient conditions of Theorem 9.6 are satisfied.

Using the fact that $v_i^2 = 1$ and $\alpha > 0$, we reformulate Theorem 9.6 as follows.

Theorem 9.7 A vertex \mathbf{v} of the hypercube H_n is an equilibrium state of the gBSB neural model if and only if

$$(\mathbf{W}\mathbf{v} + \mathbf{b})_i v_i \geq 0, \quad i = 1, 2, \dots, n.$$

In our analysis we use stability definitions in the sense of Lyapunov, which we now recall.

Definition 9.2 An equilibrium state \mathbf{x}^* of $\mathbf{x}(k+1) = \mathbf{T}(\mathbf{x}(k))$ is (uniformly) stable if for every $\varepsilon > 0$ there is a $\delta = \delta(\varepsilon) > 0$ such that if $\|\mathbf{x}(0) - \mathbf{x}^*\| < \delta$, then $\|\mathbf{x}(k) - \mathbf{x}^*\| < \varepsilon$ for all $k \geq 0$, where $\|\cdot\|$ may be any p -norm of a vector—for example, the Euclidean norm $\|\cdot\|_2$.

The above definition says that if an equilibrium is stable, then a trajectory starting sufficiently close to the equilibrium will stay close to the equilibrium for all future times. If the equilibrium is not stable, then it is called unstable.

Definition 9.3 The equilibrium solution \mathbf{x}^* is asymptotically stable if it is stable and if there exists an $\eta > 0$ such that if $\|\mathbf{x}(0) - \mathbf{x}^*\| < \eta$, then $\lim_{k \rightarrow \infty} \mathbf{x}(k) = \mathbf{x}^*$.

The following theorem gives us a sufficient condition for a vertex of H_n to be an asymptotically stable equilibrium state.

Theorem 9.8 Let \mathbf{v} be a vertex of the hypercube H_n . If

$$(\mathbf{L}(\mathbf{v}))_i v_i > 1, \quad i = 1, 2, \dots, n,$$

then \mathbf{v} is an asymptotically stable equilibrium state of the gBSB neural model.

Proof Because $(\mathbf{L}(\mathbf{x}))_i x_i$ is a continuous function, there exists a neighborhood $N(\mathbf{v})$ about \mathbf{v} such that

$$(\mathbf{L}(\mathbf{x}))_i x_i > 1 \quad \text{for all } i = 1, 2, \dots, n \text{ and } \mathbf{x} \in N(\mathbf{v}).$$

Therefore, for any $\mathbf{x} \in N(\mathbf{v})$ we have

$$\mathbf{g}(\mathbf{L}(\mathbf{x})) = \mathbf{v},$$

which means that \mathbf{v} is indeed asymptotically stable equilibrium state.

Taking into account the fact that $v_i^2 = 1$ and $\alpha > 0$, we reformulate the above theorem as follows.

Theorem 9.9 Let \mathbf{v} be a vertex of the hypercube H_n . If

$$(\mathbf{W}\mathbf{v} + \mathbf{b})_i v_i > 0, \quad i = 1, 2, \dots, n,$$

then \mathbf{v} is an asymptotically stable equilibrium state of the gBSB neural model.

9.9.3 Synthesis of Neural Associative Memory

In our further discussion, we use a special class of matrices that are characterized in the following definition.

Definition 9.4 A matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$ is said to be row diagonal dominant if

$$d_{ii} \geq \sum_{\substack{j=1 \\ j \neq i}}^n |d_{ij}|, \quad i = 1, 2, \dots, n.$$

It is strongly row diagonal dominant if

$$d_{ii} > \sum_{\substack{j=1 \\ j \neq i}}^n |d_{ij}|, \quad i = 1, 2, \dots, n.$$

Now, suppose that we are given a matrix $\mathbf{V} \in \mathbb{R}^{n \times r}$. Its pseudo-inverse, or the Moore–Penrose inverse, is defined as follows.

Definition 9.5 Given $V \in \mathbb{R}^{n \times r}$, a matrix $V^\dagger \in \mathbb{R}^{r \times n}$ is called a pseudo-inverse of V if

1. $VV^\dagger V = V$
2. $V^\dagger V V^\dagger = V^\dagger$
3. $(VV^\dagger)^T = VV^\dagger$
4. $(V^\dagger V)^T = V^\dagger V$

The pseudo-inverse has the following useful properties:

- $(V^T)^\dagger = (V^\dagger)^T$
- $(V^\dagger)^\dagger = V$

Let $v^{(1)}, \dots, v^{(r)}$ be given vertices of H_n that we wish to store as asymptotically stable equilibrium states of the gBSB neural net; that is, we are interested in building a neural autoassociative memory. We now present a number of technical results that will lead to formulation of a design algorithm for the weight matrix of the gBSB neural associative memory. Let $D \in \mathbb{R}^{n \times n}$ satisfy the condition

$$d_{ii} > \sum_{\substack{k=1 \\ k \neq i}}^n |d_{ik}| + |b_i|, \quad i = 1, 2, \dots, n. \quad (9.136)$$

We now present an expression for constructing the weight matrix W ensuring that the given patterns are stored as asymptotically stable equilibrium states of the gBSB neural net.

Theorem 9.10 Let $V = [v^{(1)} \dots v^{(r)}] \in \{-1, 1\}^{n \times r}$. If

$$W = DVV^\dagger,$$

where the matrix D satisfies (9.136), then $v^{(1)}, \dots, v^{(r)}$ are asymptotically stable states of the gBSB net.

Proof Using the first property of the Moore–Penrose inverse of V , we obtain

$$WV = DVV^\dagger V = DV.$$

Hence, for $j = 1, \dots, r$ we have

$$Wv^{(j)} = Dv^{(j)}. \quad (9.137)$$

Substituting the above into the condition of Theorem 9.9 and using (9.136) gives

$$\begin{aligned} (Wv^{(j)} + b)_i v_i^{(j)} &= (Dv^{(j)} + b)_i v_i^{(j)} \\ &= \left(d_{ii} v_i^{(j)} + \sum_{\substack{k=1 \\ k \neq i}}^n d_{ik} v_k^{(j)} + b_i \right) v_i^{(j)} \\ &\geq d_{ii} - \sum_{\substack{k=1 \\ k \neq i}}^n |d_{ik}| - |b_i| \\ &> 0. \end{aligned}$$

By Theorem 9.9, $\mathbf{v}^{(j)}$ is asymptotically stable. The above argument holds for all $j = 1, \dots, r$. The proof is complete.

Note that the weight matrix $\mathbf{W} = \mathbf{D}\mathbf{V}\mathbf{V}^\dagger$ does not have to be symmetric thanks to the presence of the matrix \mathbf{D} in the formula for \mathbf{W} . The matrix \mathbf{D} is a design parameter, and therefore it can be chosen to be nonsymmetric so that \mathbf{W} becomes nonsymmetric.

Corollary 9.2 Every vertex that is contained in $\text{Range}(\mathbf{V})$ is also stored as an asymptotically stable equilibrium state.

Proof Let $\mathbf{v} \in \{-1, 1\}^n$ be such that $\mathbf{v} \in \text{Range}(\mathbf{V})$. Hence, there is a vector $\boldsymbol{\gamma} \in \mathbb{R}^r$ such that $\mathbf{v} = \mathbf{V}\boldsymbol{\gamma}$. Using the above and the first property of the pseudo-inverse matrix of \mathbf{V} yields

$$\begin{aligned}\mathbf{W}\mathbf{v} &= \mathbf{D}\mathbf{V}\mathbf{V}^\dagger \mathbf{v} \\ &= \mathbf{D}\mathbf{V}\mathbf{V}^\dagger \mathbf{V}\boldsymbol{\gamma} \\ &= \mathbf{D}\mathbf{V}\boldsymbol{\gamma} \\ &= \mathbf{D}\mathbf{v}.\end{aligned}\tag{9.138}$$

Substituting the above into the condition of Theorem 9.9 and performing manipulation similar to the ones in the proof of Theorem 9.10, we arrive at the result.

Theorem 9.11 Under the conditions of Theorem 9.10, every vertex \mathbf{v} of H_n that is contained in $\text{Range}(\mathbf{V})^\perp$ such that for some $i = 1, \dots, n$,

$$b_i v_i < 0$$

is not an equilibrium state of the gBSB net.

Proof It follows from the definition of the pseudo-inverse that if \mathbf{V}^\dagger is a pseudo-inverse matrix of $\mathbf{V} \in \mathbb{R}^{n \times r}$, then there exists a square matrix \mathbf{M} such that

$$\mathbf{V}^\dagger = \mathbf{M}\mathbf{V}^T.\tag{9.139}$$

Suppose that $\mathbf{v} \in \text{Range}(\mathbf{V})^\perp$. Then, $\mathbf{V}^T \mathbf{v} = \mathbf{0}$, and therefore, by (9.139),

$$\mathbf{V}^\dagger \mathbf{v} = \mathbf{0}.\tag{9.140}$$

Hence, $\mathbf{W}\mathbf{v} = \mathbf{D}\mathbf{V}\mathbf{V}^\dagger \mathbf{v} = \mathbf{0}$, and

$$(\mathbf{W}\mathbf{v} + \mathbf{b})_i v_i = b_i v_i, \quad i = 1, \dots, n.\tag{9.141}$$

If for some i we have $b_i v_i < 0$, then

$$(\mathbf{W}\mathbf{v} + \mathbf{b})_i v_i < 0,$$

which implies that \mathbf{v} is not an equilibrium point of the gBSB neural net.

The above results are clear and easy to apply to construct the weight matrix \mathbf{W} . However, if the diagonal elements of the matrix \mathbf{D} are too large, we may end up storing all of the vertices of the hypercube H_n . This follows from the following result concerning the weight matrix.

Theorem 9.12 If the weight matrix \mathbf{W} of the gBSB neural net is such that for all $i = 1, \dots, n$,

$$w_{ii} > \sum_{\substack{k=1 \\ k \neq i}}^n |w_{ik}| + |b_i|,$$

then all the vertices of H_n are asymptotically stable equilibrium states of the gBSB net.

Proof Let \mathbf{v} be an arbitrary vertex of H_n , that is, $\mathbf{v} \in \{-1, 1\}^n$. Then, using the assumption on the weight matrix, we obtain

$$\begin{aligned} (\mathbf{W}\mathbf{v} + \mathbf{b})_i v_i &= w_{ii} + \sum_{\substack{k=1 \\ k \neq i}}^n w_{ik} v_k v_i + b_i v_i \\ &\geq w_{ii} - \sum_{\substack{k=1 \\ k \neq i}}^n |w_{ik}| - |b_i| \\ &> 0. \end{aligned}$$

It follows from Theorem 9.9 that all vertices of H_n are asymptotically stable states of the gBSB net.

We now provide an “improved” method for constructing the weight matrix. The method is based on the following theorem.

Theorem 9.13 Let $\mathbf{V} = [\mathbf{v}^{(1)} \ \dots \ \mathbf{v}^{(r)}] \in \{-1, 1\}^{n \times r}$. Let $\mathbf{B} = [\mathbf{b} \ \dots \ \mathbf{b}] \in \mathbb{R}^{n \times r}$ be the matrix consisting of the column vector \mathbf{b} repeated r times, $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\beta_j = \sum_{i=1}^r (\mathbf{V}^\dagger \mathbf{v}^{(j)})_i$, and $\beta = -\max_j |\beta_j|$, where $j = 1, 2, \dots, r$. Suppose

$$d_{ii} > \sum_{\substack{k=1 \\ k \neq i}}^n |d_{ik}| + |b_i(1 - \beta)|, \quad i = 1, 2, \dots, n.$$

If

$$\mathbf{W} = (\mathbf{D}\mathbf{V} - \mathbf{B})\mathbf{V}^\dagger + \mathbf{A}(\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger), \quad (9.142)$$

then all the patterns $\mathbf{v}^{(j)}$, $j = 1, 2, \dots, r$, are stored as asymptotically stable equilibrium states of the gBSB neural net.

Proof Using the first property of the Moore–Penrose inverse of \mathbf{V} , we obtain

$$\begin{aligned} \mathbf{W}\mathbf{V} &= (\mathbf{D}\mathbf{V} - \mathbf{B})\mathbf{V}^\dagger \mathbf{V} + \mathbf{A}(\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger)\mathbf{V} \\ &= \mathbf{D}\mathbf{V} - \mathbf{B}\mathbf{V}^\dagger \mathbf{V} \\ &= \mathbf{D}\mathbf{V} - [\mathbf{B}\mathbf{V}^\dagger \mathbf{v}^{(1)} \ \mathbf{B}\mathbf{V}^\dagger \mathbf{v}^{(2)} \ \dots \ \mathbf{B}\mathbf{V}^\dagger \mathbf{v}^{(r)}] \\ &= \mathbf{D}\mathbf{V} - [\mathbf{b}\beta_1 \ \mathbf{b}\beta_2 \ \dots \ \mathbf{b}\beta_r]. \end{aligned}$$

The j th column of the above is

$$\mathbf{W}\mathbf{v}^{(j)} = \mathbf{D}\mathbf{v}^{(j)} - \mathbf{b}\beta_j, \quad j = 1, 2, \dots, r.$$

We now show that a sufficient condition for a pattern $\mathbf{v}^{(j)}$ to be an asymptotically stable equilibrium of $\mathbf{x}(k+1) = \mathbf{T}(\mathbf{x}(k))$ is satisfied. Performing simple manipulations and taking into account the assumptions, for $j = 1, 2, \dots, r$ we obtain

$$\begin{aligned}
 (\mathbf{W}\mathbf{v}^{(j)} + \mathbf{b})_i v_i^{(j)} &= (\mathbf{D}\mathbf{v}^{(j)} - \mathbf{b}\beta_j + \mathbf{b})_i v_i^{(j)} \\
 &= d_{ii} + \sum_{\substack{k=1 \\ k \neq i}}^n d_{ik} v_k^{(j)} v_i^{(j)} + b_i(1 - \beta_j) v_i^{(j)} \\
 &\geq d_{ii} - \sum_{\substack{k=1 \\ k \neq i}}^n |d_{ik}| - \left| b_i \left(1 - \left(-\max_j |\beta_j| \right) \right) \right| \\
 &= d_{ii} - \sum_{\substack{k=1 \\ k \neq i}}^n |d_{ik}| - |b_i(1 - \beta)| \\
 &> 0,
 \end{aligned}$$

where the last inequality follows from the assumption on the matrix \mathbf{D} . Hence, by Theorem 9.9, the given patterns are asymptotically stable equilibrium states of the gBSB model.

If the patterns are linearly independent and $r < n$, then $\text{rank } \mathbf{V} = r$ and therefore $\mathbf{V}^\dagger \mathbf{V} = \mathbf{I}_r$. In this case $\beta_j = 1$ for all j and

$$\mathbf{W}\mathbf{v}^{(j)} = \mathbf{D}\mathbf{v}^{(j)} - \mathbf{b}, \quad j = 1, 2, \dots, r,$$

Thus in this case, a sufficient condition for the patterns $\mathbf{v}^{(j)}$, $j = 1, 2, \dots, r$, to be stored as asymptotically stable equilibrium states of the gBSB model is that for $i = 1, 2, \dots, n$,

$$d_{ii} > \sum_{\substack{k=1 \\ k \neq i}}^n |d_{ik}|.$$

The above condition indeed guarantees that all the desired patterns are stored as asymptotically stable equilibria. As we mentioned before, if d_{ii} 's are too large, then we may end up storing many other vertices of H_n in addition to the given patterns that we want to be stored. To minimize the number of spurious patterns stored, we impose the constraints of the form

$$d_{ii} < \sum_{\substack{k=1 \\ k \neq i}}^n |d_{ik}| + |b_i|, \quad i = 1, 2, \dots, n,$$

on the matrix \mathbf{D} while synthesizing the weight matrix \mathbf{W} .

We now discuss the role of different terms in the expression (9.142) on the weight matrix in storing the desired patterns and minimizing the number of spurious states. First, the term $(\mathbf{D}\mathbf{V} - \mathbf{B})\mathbf{V}^\dagger$ is responsible for the given patterns to be stored as asymptotically stable equilibrium states of the gBSB model. This is because for a given pattern, say $\mathbf{v}^{(j)}$, we have $\mathbf{\Lambda}(\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger)\mathbf{v}^{(j)} = \mathbf{0}$, where $\mathbf{0}$ is a vector of zeros. The role of the term $\mathbf{\Lambda}(\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger)$ is to destabilize patterns that are not in $\text{Span}[\mathbf{v}^{(1)} \dots \mathbf{v}^{(r)}]$. Specifically, if a pattern \mathbf{y} is

orthogonal to $\text{Span}[\mathbf{v}^{(1)} \ \dots \ \mathbf{v}^{(r)}]$, then the term $\mathbf{\Lambda}(\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger)$ will not only destabilize the vertex of H_n corresponding to the pattern \mathbf{y} , if it is an equilibrium state of the gBSB model, it will make this state a nonequilibrium of the gBSB model. To see this, consider the expression

$$\mathbf{W}\mathbf{y} = (\mathbf{D}\mathbf{V} - \mathbf{B})\mathbf{V}^\dagger\mathbf{y} + \mathbf{\Lambda}(\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger)\mathbf{y}.$$

Because

$$\mathbf{y} \perp \text{Span}[\mathbf{v}^{(1)} \ \dots \ \mathbf{v}^{(r)}],$$

we have

$$\mathbf{W}\mathbf{y} = \mathbf{\Lambda}(\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger)\mathbf{y} = \mathbf{\Lambda}\mathbf{y}.$$

Therefore,

$$(\mathbf{W}\mathbf{y} + \mathbf{b})_i y_i = (\mathbf{\Lambda}\mathbf{y} + \mathbf{b})_i y_i = \lambda_{ii} + \sum_{\substack{j=1 \\ j \neq i}}^n \lambda_{ij} y_j y_i + b_i y_i.$$

Hence, if we choose $\mathbf{\Lambda}$ such that for some $i = 1, 2, \dots, n$,

$$\lambda_{ii} < - \sum_{\substack{j=1 \\ j \neq i}}^n |\lambda_{ij}| - |b_i|,$$

then

$$(\mathbf{W}\mathbf{y} + \mathbf{b})_i y_i < 0 \quad \text{for some } i = 1, 2, \dots, n.$$

Therefore, the vertex \mathbf{y} cannot be an equilibrium state. The inclusion of the matrix \mathbf{B} in the synthesis formula for the weight matrix \mathbf{W} was to prevent an automatic storage of the negatives of the given patterns as well as to somehow help to direct the gBSB model trajectories toward the given patterns. For this reason, we select \mathbf{b} to be in $\text{Span}\{\mathbf{V}\}$. We recommend making \mathbf{b} a weighted sum of $\mathbf{v}^{(i)}$, $i = 1, 2, \dots, r$, with positive weight coefficients. The above results can be used as guides for selecting the components of the matrices \mathbf{D} , $\mathbf{\Lambda}$, the vector \mathbf{b} , and consequently the matrix \mathbf{B} . We summarize the above considerations in the form of an algorithm for synthesizing the weight matrix \mathbf{W} .

WEIGHT MATRIX CONSTRUCTION ALGORITHM

Given: $\mathbf{v}^{(j)}$, $j = 1, 2, \dots, r$.

STEP 1 Form $\mathbf{B} = [\mathbf{b} \ \mathbf{b} \ \dots \ \mathbf{b}]$, where

$$\mathbf{b} = \sum_{j=1}^r \epsilon_j \mathbf{v}^{(j)}, \quad \epsilon_j > 0, \quad j = 1, 2, \dots, r.$$

STEP 2 Select \mathbf{D} such that for $i = 1, 2, \dots, n$ we have

$$d_{ii} > \sum_{\substack{k=1 \\ k \neq i}}^n |d_{ik}| \quad \text{and} \quad d_{ii} < \sum_{\substack{k=1 \\ k \neq i}}^n |d_{ik}| + |b_i|.$$

STEP 3 Choose \mathbf{A} such that for $i = 1, 2, \dots, n$,

$$\lambda_{ii} < - \sum_{\substack{j=1 \\ j \neq i}}^n |\lambda_{ij}| - |b_i|.$$

STEP 4 Compute

$$\mathbf{W} = (\mathbf{D}\mathbf{V} - \mathbf{B})\mathbf{V}^\dagger + \mathbf{A}(\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger).$$

STEP 5 Identify spurious equilibrium patterns. You can use the condition of Theorem 9.7 for this purpose.

When the set of patterns to be stored is not linearly independent, we can still use the preceding algorithm to compute the weight matrix \mathbf{W} . However, the second constraint on the matrix \mathbf{D} may lead to a violation of the condition on \mathbf{D} in Theorem 9.13. In such a case we cannot guarantee that all the given patterns will be stored, nor we will have any control on the number of spurious states. We may then relax this condition. If attempts in synthesizing associative memory for a set of linearly dependent patterns result in unacceptable designs, we may try a different approach. We can start by picking a maximal set of patterns that are linearly independent from the given set of patterns. Let us assume that there are s patterns to be stored and that r of them form a maximal linearly independent subset. Suppose, for notational simplicity, that the patterns $\mathbf{v}^{(q)}$, $q = 1, \dots, r$, form a linearly independent set. We then apply the above algorithm to generate the weight matrix \mathbf{W} based on the linearly independent patterns $\mathbf{v}^{(q)}$, $q = 1, \dots, r$. After that, we modify the obtained weight matrix and the vector \mathbf{b} to ensure that the remaining patterns $\mathbf{v}^{(j)}$, $j = r + 1, r + 2, \dots, s$, will also be stored. We will focus our attention on ensuring that the pattern $\mathbf{v}^{(r+1)}$ is stored as an asymptotically stable state. The remaining patterns $\mathbf{v}^{(j)}$, $j = r + 2, r + 3, \dots, s$, can be stored using the same procedure. Because

$$\mathbf{v}^{(r+1)} \in \text{Span}\{\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(r)}\},$$

that is, $\mathbf{v}^{(r+1)}$ is a linear combination of $\mathbf{v}^{(i)}$, $i = 1, 2, \dots, r$, it can be represented as

$$\mathbf{v}^{(r+1)} = \sum_{q=1}^r \gamma_q \mathbf{v}^{(q)} = \mathbf{V}\boldsymbol{\gamma},$$

where $\boldsymbol{\gamma} = [\gamma_1 \ \dots \ \gamma_r]^T \in \mathbb{R}^r$. Then, using (9.142), we obtain

$$\begin{aligned} \mathbf{W}\mathbf{v}^{(r+1)} &= \mathbf{D}\mathbf{V}\boldsymbol{\gamma} - \mathbf{B}\boldsymbol{\gamma} \\ &= \mathbf{D}\mathbf{v}^{(r+1)} - \left(\sum_{q=1}^r \gamma_q \right) \mathbf{b}. \end{aligned}$$

Hence, if

$$\begin{aligned} (\mathbf{W}\mathbf{v}^{(r+1)} + \mathbf{b})v_i^{(r+1)} &= \left(\mathbf{D}\mathbf{v}^{(r+1)} - \left[\sum_{q=1}^r \gamma_q - 1 \right] \mathbf{b} \right)_i v_i^{(r+1)} \\ &> 0, \quad i = 1, 2, \dots, n, \end{aligned}$$

then the vertex $\mathbf{v}^{(r+1)}$ is asymptotically stable. Observe that because of the assumptions on the matrix \mathbf{D} , we have

$$\text{sign}(\mathbf{D}\mathbf{v}^{(r+1)})_i = \text{sign}(v_i^{(r+1)}).$$

With the above in mind, we obtain stronger sufficient conditions for $\mathbf{v}^{(r+1)}$ to be asymptotically stable. They are

$$\text{sign} \left[\left(1 - \sum_{q=1}^r \gamma_q \right) \mathbf{b} \right]_i = \text{sign}(v_i^{(r+1)}), \quad i = 1, 2, \dots, n.$$

We can represent the above conditions as

$$\text{sign } b_i = \text{sign}(v_i^{(r+1)}) \text{sign} \left(1 - \sum_{q=1}^r \gamma_q \right), \quad i = 1, 2, \dots, n.$$

Thus, if we adjust the components of \mathbf{b} so that the above conditions are satisfied, then we are guaranteed that $\mathbf{v}^{(r+1)}$ is stored as an asymptotically stable equilibrium state. We may try to further adjust the components of \mathbf{b} , if possible, to store the remaining patterns. Then, using the same matrices \mathbf{D} , $\mathbf{\Lambda}$, and \mathbf{V} used to compute the original weight matrix \mathbf{W} , we can choose an appropriate \mathbf{b} to generate a new weight matrix \mathbf{W} .

If \mathbf{b} cannot be adjusted, employing the above technique, to accommodate all the patterns, we may then have to use the following alternative procedure. First, recall the stability conditions that the patterns $\mathbf{v}^{(r+1)}, \dots, \mathbf{v}^{(s)}$ should satisfy:

$$\left(\mathbf{D}\mathbf{v}^{(k)} - \left[\sum_{q=1}^r \gamma_q - 1 \right] \mathbf{b} \right)_i v_i^{(k)} > 0, \quad i = 1, 2, \dots, n, \quad k = r+1, r+2, \dots, s.$$

Using the requirement for the matrix \mathbf{D} to be diagonally dominant, we arrive at stronger conditions of the form

$$d_{ii} - \sum_{\substack{j=1 \\ j \neq i}}^n |d_{ij}| + \left(1 - \sum_{q=1}^r \gamma_q \right) b_i v_i^{(k)} > 0, \quad i = 1, \dots, n, \quad k = r+1, \dots, s.$$

Thus, we may use the above conditions and try to decrease the magnitude of the b_i 's until the conditions are met for all $k = r+1, \dots, s$. However, we are not guaranteed to always satisfy the above conditions. If the conditions are satisfied, then employing the same matrices \mathbf{D} , $\mathbf{\Lambda}$, and \mathbf{V} used to compute the original weight matrix \mathbf{W} , we choose an appropriate \mathbf{b} to generate a new weight matrix \mathbf{W} .

◆ Example 9.10

We synthesize a gBSB-based associative memory to store the following three patterns:

$$\mathbf{v}^{(1)} = \begin{bmatrix} -1 \\ -1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad \mathbf{v}^{(2)} = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \\ -1 \end{bmatrix}, \quad \text{and} \quad \mathbf{v}^{(3)} = \begin{bmatrix} 1 \\ -1 \\ 1 \\ -1 \\ 1 \end{bmatrix}.$$

We first choose matrices \mathbf{D} and $\mathbf{\Lambda}$. Let

$$\mathbf{D} = \begin{bmatrix} 1 & 0 & 0 & 0 & -0.5 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & -0.2 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}, \quad \mathbf{\Lambda} = -2\mathbf{I}_5.$$

Set $\mathbf{V} = [\mathbf{v}^{(1)} \ \mathbf{v}^{(2)} \ \mathbf{v}^{(3)}]$ and $\mathbf{B} = [\mathbf{b} \ \mathbf{b} \ \mathbf{b}]$, where

$$\mathbf{b} = [0.6 \quad -0.5 \quad 0.7 \quad -0.4 \quad 0.8]^T.$$

Note that $\text{rank } \mathbf{V} = 2$. Let the step size $\alpha = 0.9$. Using the above, we compute the weight matrix \mathbf{W} , where

$$\mathbf{W} = \begin{bmatrix} -0.6500 & 0.2667 & -0.2667 & -1.3500 & -0.2667 \\ 0.1250 & -1.0833 & -0.9167 & -0.1250 & -0.9167 \\ -0.1750 & -0.9500 & -1.0500 & 0.1750 & 0.9500 \\ -1.4000 & -0.0667 & 0.0667 & -0.6000 & 0.0667 \\ -0.2000 & -0.8667 & 0.8667 & 0.2000 & -1.1333 \end{bmatrix}.$$

We then check which vertices of H_5 satisfy the sufficiency condition of Theorem 9.8 for asymptotic stability. In other words, we wish to identify the spurious states. In this example, only the given patterns satisfy the condition of Theorem 9.8.

◆ Example 9.11

In this example, we use a gBSB neural net with the state in the hypercube $\{-1, 1\}^{25}$. We construct the weight matrix to store the four letters of alphabet whose bit maps are shown in Figure 9.44. The bit maps are 5×5 with elements from the set $\{-1, 1\}$, with -1 corresponding to black and 1 to white. The pattern vectors are obtained from the bit maps by stacking their columns. For example, if $A = [\mathbf{a}_1 \ \mathbf{a}_2 \ \mathbf{a}_3 \ \mathbf{a}_4 \ \mathbf{a}_5]$, where \mathbf{a}_i , $i = 1, 2, \dots, 5$, is the i th column of A , then the corresponding pattern vector is obtained from A by applying the stacking

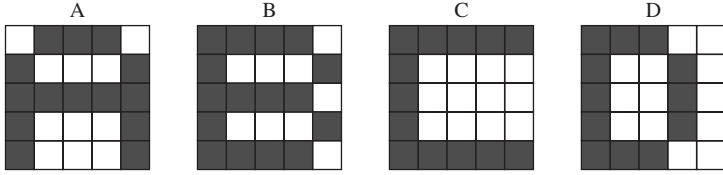


Figure 9.44 Bit maps of the characters *A*, *B*, *C*, and *D* stored in the gBSB neural net of Example 9.11.

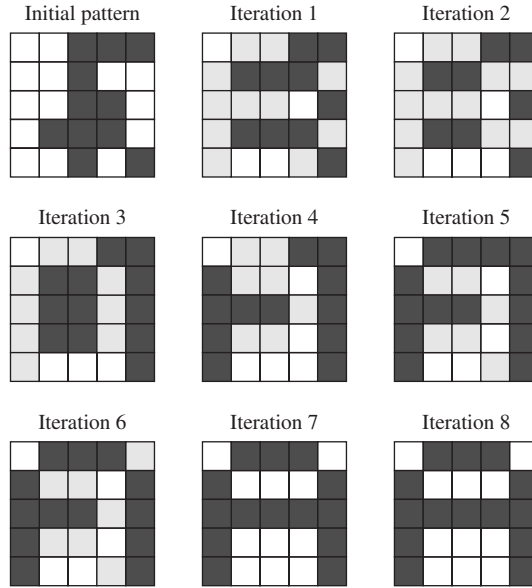


Figure 9.45 Snapshots of the iterations of the gBSB net of Example 9.11 converging toward the letter *A*.

operator $\mathbf{s} : \mathbb{R}^{5 \times 5} \rightarrow \mathbb{R}^{25}$ to get

$$\mathbf{s}(A) = A(:) = \begin{bmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_5 \end{bmatrix},$$

where $A(:)$ is the MATLAB's command for the stacking operation on the matrix A . The weight matrix, \mathbf{W} , of the gBSB net was constructed using the weight matrix construction algorithm, where

$$\mathbf{V} = [A(:) \quad B(:) \quad C(:) \quad D(:)].$$

The values of other parameters chosen in this example were: $\mathbf{D} = 1.1 \mathbf{I}_{25}$, $\mathbf{\Lambda} = -2 \mathbf{I}_{25}$, $\mathbf{b} = 0.3 A(:) + 0.1 B(:)$, and $\alpha = 0.5$. Snapshots of the network iterates for different initial patterns are shown in Figures 9.45–9.48.

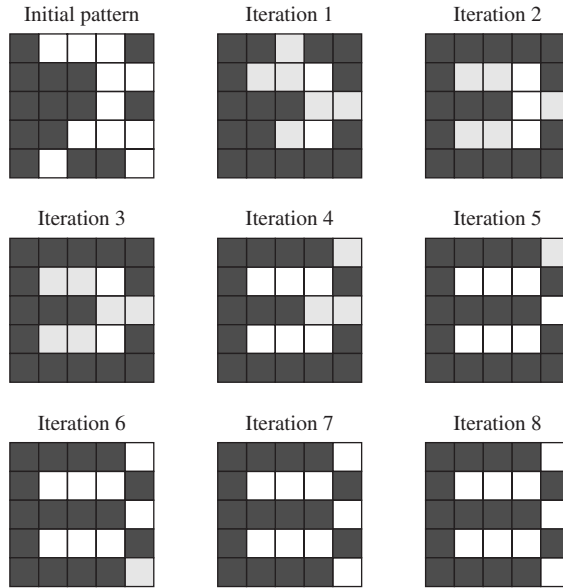


Figure 9.46 Snapshots of the iterations of the gBSB net of Example 9.11 converging toward the letter *B*.

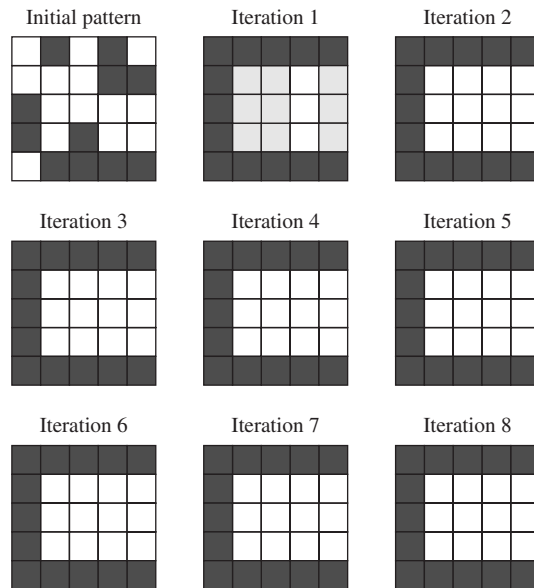


Figure 9.47 Snapshots of the iterations of the gBSB net of Example 9.11 converging toward the letter *C*.

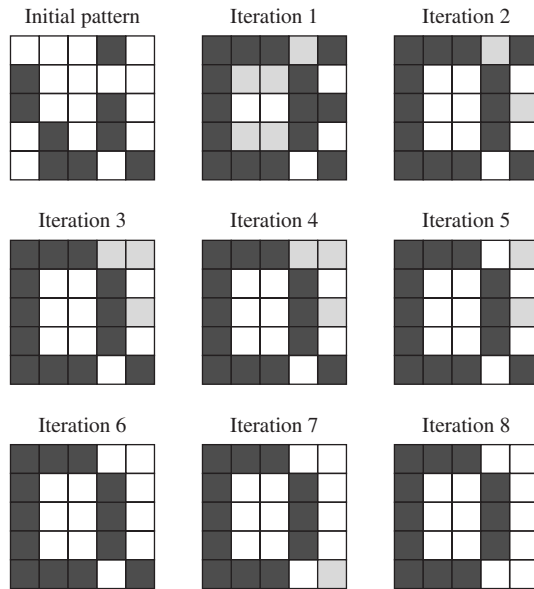


Figure 9.48 Snapshots of the iterations of the gBSB net of Example 9.11 converging toward the letter *D*.

9.9.4 Learning

A pattern is learned by a memory if its noisy or incomplete version presented to the memory is mapped back to this pattern. In our discussion of the proposed iterative learning, we use the following theorem, which allows one to iteratively compute the pseudo-inverse of a given matrix. The theorem can be found in Campbell and Meyer [40, p. 54].

Theorem 9.14 Let $U \in \mathbb{R}^{n \times (r+1)}$ be partitioned by slicing off its last column, so that

$$U = \begin{bmatrix} V & y \end{bmatrix}.$$

Let

$$\begin{aligned} k &= V^\dagger y, \\ u &= (I_n - VV^\dagger)y, \\ z &= \begin{cases} u^\dagger & \text{if } u \neq 0, \\ \frac{k^T V^\dagger}{1+k^T k} & \text{if } u = 0. \end{cases} \end{aligned}$$

Then,

$$U^\dagger = \begin{bmatrix} V^\dagger - kz \\ z \end{bmatrix}.$$

Learning a new pattern requires updating the weight matrix W . This, in turn, involves augmenting the matrices V and B and recomputing V^\dagger . Augmenting V and B is straightforward,

and updating V^\dagger can be performed recursively using the above theorem. The above procedure guarantees that the new pattern will be stored as an asymptotically stable equilibrium state if it is linearly independent of the learned patterns. If, however, the pattern to be learned is a linear combination of the stored patterns, then we do not update the matrix V but rather modify the vector \mathbf{b} as described at the end of the previous subsection. After that we update the weight matrix W . We now illustrate the iterative learning process on a numerical example.

◆ Example 9.12

Consider Example 9.10. Suppose that we wish to add

$$\mathbf{v}^{(4)} = [-1 \quad -1 \quad -1 \quad -1 \quad 1]^T$$

to the previously designed memory. We augment the pattern matrix to include $\mathbf{v}^{(4)}$ in it. Thus, we get $V = [\mathbf{v}^{(1)} \quad \mathbf{v}^{(2)} \quad \mathbf{v}^{(3)} \quad \mathbf{v}^{(4)}]$. We also augment the matrix B by adding extra \mathbf{b} column to it. We then update V^\dagger using Theorem 9.14 and recompute the weight matrix W to obtain

$$W = \begin{bmatrix} 0.1714 & 0.8143 & 0.8286 & -0.5286 & -0.8143 \\ 0.4643 & -0.8571 & -0.4643 & 0.2143 & -1.1429 \\ 0.7786 & -0.3143 & 0.2214 & 1.1286 & 0.3143 \\ -0.8286 & 0.3143 & 0.8286 & -0.0286 & -0.3143 \\ -0.4857 & -1.0571 & 0.4857 & -0.0857 & -0.9429 \end{bmatrix}.$$

Next, we check which vertices of H_5 satisfy the condition of Theorem 9.8 for asymptotic stability. We find that only the given four patterns satisfy the condition of Theorem 9.8.

9.9.5 Forgetting

A previously stored pattern is forgotten or deleted from the memory if a stimulus that is a perturbed version of that pattern is presented to the memory and is not mapped back to the stored pattern. We now discuss the case when we wish to update the weight matrix in such a way that the specified patterns are removed from memory. We refer to the process of removing patterns from the associative memory as forgetting. Our discussion uses the following theorem, which can be considered as the inverse of Theorem 9.14. A general version of the theorem, appears in Boullion and Odell [31, p. 14].

Theorem 9.15 Let $U \in \mathbb{R}^{n \times (r+1)}$ be partitioned by slicing off its last column, so that

$$U = [\mathbf{V} \quad \mathbf{y}],$$

and let its pseudo-inverse $U^\dagger \in \mathbb{R}^{(r+1) \times n}$ be partitioned by slicing off its last row, so that

$$U^\dagger = \begin{bmatrix} \mathbf{S} \\ \mathbf{z} \end{bmatrix}.$$

Then,

$$\mathbf{V}^\dagger = \begin{cases} \mathbf{S} + \mathbf{S} \frac{\mathbf{y}\mathbf{z}}{1-\mathbf{z}\mathbf{y}} & \text{if } \mathbf{z}\mathbf{y} \neq 1 \\ \mathbf{S} - \mathbf{S}\mathbf{z}^\dagger\mathbf{z} & \text{if } \mathbf{z}\mathbf{y} = 1. \end{cases}$$

Proof Referring to Theorem 9.14, we first show that $\mathbf{u} = \mathbf{0}$ if and only if $\mathbf{z}\mathbf{y} \neq 1$, or, equivalently, $\mathbf{u} \neq \mathbf{0}$ if and only if $\mathbf{z}\mathbf{y} = 1$. Indeed, if $\mathbf{u} = \mathbf{0}$, then

$$\mathbf{z} = \frac{\mathbf{k}^T \mathbf{V}^\dagger}{1 + \mathbf{k}^T \mathbf{k}}.$$

Hence,

$$\mathbf{z}\mathbf{y} = \frac{\mathbf{k}^T \mathbf{V}^\dagger \mathbf{y}}{1 + \mathbf{k}^T \mathbf{k}} = \frac{\mathbf{k}^T \mathbf{k}}{1 + \mathbf{k}^T \mathbf{k}} < 1.$$

On the other hand, if $\mathbf{u} \neq \mathbf{0}$, then

$$\mathbf{z} = \mathbf{u}^\dagger = \{(\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger) \mathbf{y}\}^\dagger.$$

Taking into account the fact that $\mathbf{u} \neq \mathbf{0}$ and that $\mathbf{u} \in \mathbb{R}^n$, its pseudo-inverse is

$$\mathbf{u}^\dagger = (\mathbf{u}^T \mathbf{u})^{-1} \mathbf{u}^T = (\mathbf{y}^T (\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger) \mathbf{y})^{-1} \mathbf{y}^T (\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger),$$

because $(\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger)$ is an orthogonal projector and hence

$$(\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger)^2 = (\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger).$$

Therefore, $\mathbf{z}\mathbf{y} = \mathbf{u}^\dagger \mathbf{y} = 1$.

We now prove the expressions for \mathbf{V}^\dagger . We start with the case when $\mathbf{z}\mathbf{y} \neq 1$. It follows from Theorem 9.14 that

$$\mathbf{S} = \mathbf{V}^\dagger - \mathbf{V}^\dagger \mathbf{y}\mathbf{z} = \mathbf{V}^\dagger (\mathbf{I}_n - \mathbf{y}\mathbf{z}). \quad (9.143)$$

Because $\mathbf{z}\mathbf{y} \neq 1$, the matrix $(\mathbf{I}_n - \mathbf{y}\mathbf{z})$ is invertible. Using the matrix inversion formula (A.30) from Section A.3 gives

$$(\mathbf{I}_n - \mathbf{y}\mathbf{z})^{-1} = \mathbf{I}_n + \frac{\mathbf{y}\mathbf{z}}{1 - \mathbf{z}\mathbf{y}}. \quad (9.144)$$

Applying formula (9.144) to (9.143) yields

$$\mathbf{V}^\dagger = \mathbf{S} + \mathbf{S} \frac{\mathbf{y}\mathbf{z}}{1 - \mathbf{z}\mathbf{y}}.$$

We now prove the expression for \mathbf{V}^\dagger for the case when $\mathbf{z}\mathbf{y} = 1$. It follows from Theorem 9.14 that

$$\mathbf{S} = \mathbf{V}^\dagger - \mathbf{k}\mathbf{z}, \quad (9.145)$$

and hence, because $\mathbf{z} \neq \mathbf{0}$,

$$\mathbf{S}\mathbf{z}^\dagger = \mathbf{V}^\dagger \mathbf{z}^\dagger - \mathbf{k}. \quad (9.146)$$

But in this case, $\mathbf{u} \neq \mathbf{0}$. Therefore,

$$\mathbf{z}^\dagger = (\mathbf{u}^\dagger)^\dagger = \mathbf{u} = (\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger) \mathbf{y}.$$

Substituting the above expression for \mathbf{z}^\dagger into (9.146), we get

$$\begin{aligned} \mathbf{S}\mathbf{z}^\dagger &= \mathbf{V}^\dagger (\mathbf{I}_n - \mathbf{V}\mathbf{V}^\dagger) \mathbf{y} - \mathbf{k} \\ &= (\mathbf{V}^\dagger - \mathbf{V}^\dagger \mathbf{V} \mathbf{V}^\dagger) \mathbf{y} - \mathbf{k} \\ &= -\mathbf{k}. \end{aligned}$$

We now substitute the above into (9.145) to obtain $\mathbf{V}^\dagger = \mathbf{S} - \mathbf{S}\mathbf{z}^\dagger\mathbf{z}$ which completes the proof of the theorem.

We can delete a pattern from the memory by first removing it from \mathbf{V} , then updating \mathbf{V}^\dagger using Theorem 9.15, removing one column from \mathbf{B} , and recomputing the weight matrix \mathbf{W} . We illustrate the iterative forgetting process with the following numerical example.

◆ Example 9.13

We consider Example 9.12. Suppose that we wish to delete

$$\mathbf{v}^{(3)} = [1 \quad -1 \quad 1 \quad -1 \quad 1]^T$$

from the previously designed memory. We can use Theorem 9.15 to update \mathbf{V}^\dagger —that is, to find the pseudo-inverse of $\mathbf{V} = [\mathbf{v}^{(1)} \quad \mathbf{v}^{(2)} \quad \mathbf{v}^{(4)}]$. We also update the matrix \mathbf{B} by deleting from it one column. We leave the remaining data unchanged. Recomputing the weight matrix \mathbf{W} gives

$$\mathbf{W} = \begin{bmatrix} -0.7333 & 1.2667 & 0.1500 & 0.1500 & -1.2667 \\ 0.9167 & -1.0833 & -0.1250 & -0.1250 & -0.9167 \\ 0.0500 & 0.0500 & -0.3250 & 1.6750 & -0.0500 \\ -0.0667 & -0.0667 & 1.4000 & -0.6000 & 0.0667 \\ -0.8667 & -0.8667 & 0.2000 & 0.2000 & -1.1333 \end{bmatrix}.$$

We then check which vertices of H_5 satisfy the condition of Theorem 9.8 for asymptotic stability. We find that only $\mathbf{v}^{(1)}$, $\mathbf{v}^{(2)}$, and $\mathbf{v}^{(4)}$ satisfy the condition of Theorem 9.8, which is what we desire.

Notes

This chapter gives an introduction to computing using artificial neural networks from a dynamical system perspective. The disclaimer that we found in the book by Hertz, Krogh, and Palmer [120, p. xix] also applies to this chapter as well as to the whole field of artificial neural networks: “The adjective ‘neural’ is used because much of the inspiration for such networks comes from neuroscience, *not* because we are concerned with networks of real neurons. Brain modeling is a different field and, though we sometimes describe biological analogies, our prime concern is with what *artificial* networks can do, and why.” The readers interested in principles underlying the organization of biological neurons into circuits and networks within different regions of the brain may consult the book edited by Shepherd [257]. For textbooks with an engineering approach to neural networks, we recommend Hassoun [116], Hagan, Demuth, and Beale [110], Zurada [320], Haykin [118], Fausett [79], Khanna [158], or Harvey [115]. Applications of artificial neural networks to signal processing and solving optimization problems are presented in the book by Cichocki and Unbehauen [52]. A psychological slant to the subject is presented in the book by Rumelhart, McClelland, and the PDP Research Group [247], and Anderson [8]. Impatient readers wishing to get a feel of artificial neural networks may want to consult a tutorial by Lippmann [186] and by Hush and Horne [137].

According to Haykin [118, p. 38], “The modern area of neural networks began with the pioneering work of McCulloch and Pitts (1943). McCulloch was a psychiatrist and neuroanatomist by training; he spent some 20 years thinking about the representation of an event in the nervous system. Pitts was a mathematical prodigy, who joined McCulloch in 1942.” In their paper [200], McCulloch and Pitts, assuming that the activity of the neuron is an “all-or-none” process, developed a logical calculus of neural events. They argued that, with addition of memory, a neural network consisting of a sufficiently large number of their neurons can compute anything that can be computed using the principles of logic. McCulloch and Pitts’ paper influenced John von Neumann, who used idealized switch-delay elements derived from the McCulloch–Pitts neuron model in the construction of the EDVAC (Electronic Discrete Variable Automatic Computer), which was developed out of the ENIAC (Electronic Numerical Integrator and Computer). The ENIAC was the first general-purpose computer, and it was invented by J. Presper Eckert and John Mauchly.

More recent investigators of the computational processes in the nervous system compare the brain to a hybrid, analog–digital computer [163]. According to Koch [163, p. 210], individual nerve cells communicate between each other via streams of digital pulses. Information is, on the other hand, processed in the analog domain. The resulting signal is then converted back into digital pulses and conveyed to the neurons.

For other analog learning algorithms for adaptive linear elements (adalines) in the spirit of Section 9.2, see Sira-Ramírez, Colina-Morles, and Rivas-Echeverria [259].

As we mentioned before, the Brain-State-in-a-Box (BSB) neural model was proposed by Anderson et al. in 1977 [9]. The dynamics of the BSB model and its variants were further analyzed by Greenberg [107], Golden [102, 103], Grossberg [108], Hui and Žak [130], Hakl [113], Hui, Lillo, and Žak [129], Schultz [254], Lillo et al. [184], Anderson [8], Perfetti [229], Hassoun [116], Žak, Lillo, and Hui [312], Chan and Žak [44, 45], and Park, Cho, and Park [220]. A continuous counterpart of the BSB model referred to as the linear system operating in a saturated mode was analyzed in Li, Michel, and Porod [183].

One of the first techniques used to design neural network associative memories is the *outer product method*. This method was first used by Hopfield [125, 126]. The outer product method results in a network with a symmetric interconnection structure that does not necessarily store the desired patterns as equilibrium states of the network. The network designed using this method can efficiently store up to $0.15n$ arbitrary patterns and has learning capabilities. The *projection learning rule* proposed by Personnaz, Guyon, and Dreyfus [230, 231] produces a network with a symmetric interconnection structure that always stores the desired patterns as equilibrium states, although not necessarily asymptotically stable. The network designed with this method can effectively store up to $0.5n$ desired patterns and has learning capabilities. The *eigenstructure method* used by Li, Michel, and Porod [183] yields a neural network associative memory capable of storing any given pattern as an asymptotically stable equilibrium state. The network designed with this method is capable of effectively storing a number of patterns that may exceed the order of the network. However, the resulting network will have a symmetric interconnection structure. Furthermore, there are no provisions for learning for networks designed by the eigenstructure method. In a series of papers, Michel, Farrell, and Porod [203], Michel and Farrell [202], Michel, Farrell, and Sun [204], and Yen and Michel [307] modified the eigenstructure method. The modified method can be used to design networks capable of storing on the order of $0.5n$ patterns as asymptotically stable equilibrium states. Networks designed by this method need not have symmetric interconnecting structures. Furthermore, learning is incorporated into the

method. Networks designed by the previous methods use energy function techniques. A benefit of using energy functions is that the resulting networks are globally stable; that is, the networks' trajectories always tend to some equilibrium states. On the other hand, if the modified method yields a nonsymmetric interconnecting structure, the resulting network is not guaranteed to be globally stable.

EXERCISES

- 9.1** Write a MATLAB program that implements the TLU (Perceptron) learning algorithm for generating a separating line for points on a plane. Test your program on the following data. For $i = 1, \dots, 4$ and $d_i = 1$,

$$[\mathbf{x}^{(1)} \quad \mathbf{x}^{(2)} \quad \mathbf{x}^{(3)} \quad \mathbf{x}^{(4)}] = \begin{bmatrix} 5 & -7 & 8 & -1 \\ 5 & 6 & -3 & 2 \end{bmatrix}.$$

The above points constitute class \mathcal{C}_1 . You can label the above points using, for example, “×”. Class \mathcal{C}_2 consists of the points corresponding to $d_i = -1$. They are

$$[\mathbf{x}^{(5)} \quad \mathbf{x}^{(6)} \quad \mathbf{x}^{(7)}] = \begin{bmatrix} -6 & -4 & 7 \\ 1 & -4 & -5 \end{bmatrix}$$

You can label these points using “○”. The initial weight vector and threshold are

$$\mathbf{w}(0) = [-5 \quad -1]^T \quad \text{and} \quad t_0 = 5.$$

Plot the initial separating line as well as the final separating line.

- 9.2** Write a MATLAB program that implements the WTA algorithm. Test your program on the following data. The initial weights of the WTA network are

$$[\mathbf{w}_0^{(1)} \quad \mathbf{w}_0^{(2)} \quad \mathbf{w}_0^{(3)}] = \begin{bmatrix} 0.4403 & 0.7411 & -0.8866 \\ 0.8979 & -0.6714 & 0.4625 \end{bmatrix}, \quad \alpha = 0.5.$$

The training patterns are

$$\mathbf{X} = \begin{bmatrix} -0.2 & 0.2 & -0.1 & 0.98 & 0.98 & .8 & -0.6 & -0.8 & -0.7 \\ 0.98 & 0.98 & .9 & 0.2 & -0.2 & 0.1 & -0.8 & -0.6 & -0.7 \end{bmatrix}$$

Present the training set 200 times. Find the weights of the network after training.

- 9.3** Use the identification scheme depicted in Figure 9.16 and the neural fuzzy on-line identification algorithm from Section 9.4 to identify the plant described by

$$y(k+1) = \frac{y(k)(y(k-1)+2)(y(k)+2.5)}{8.5 + y^2(k) + y^2(k-1)} + u(k),$$

where

$$u(k) = 2 \cos(2\pi k/100) \quad \text{for } k \leq 200$$

and

$$u(k) = 1.2 \sin(2\pi/20) \quad \text{for } k > 200.$$

The identifier is to be described by

$$\hat{y}(k+1) = f(k) + u(k),$$

where f has the form of (9.42) with $M = 100$. Choose $\alpha = 0.5$, $y(0) = y(1) = 0$, and $f(1) = f(2) = 0$, and start with $u(1) = 0$. (The above plant was also used by Lin and Lee [185, p. 731] to test their identification scheme.)

- 9.4** Write a MATLAB program that implements the learning algorithm for the RBF identifier discussed in Section 9.5.
- 9.5** Prove the following lemma that can be found in Park, Cho, and Park [220].

Lemma 9.2 Let $\mathbf{v} \in \{-1, 1\}^n$ be an asymptotically stable equilibrium state of the gBSB neural net. Let \mathbf{x} be a vertex of H_n that differs from \mathbf{v} in the component x_i for some $i = 1, \dots, n$. If

$$(\mathbf{W}\mathbf{x} + \mathbf{b})_i v_i > 0,$$

then \mathbf{x} is not an equilibrium state of the gBSB neural net.

- 9.6** Use Lemma 9.2 from Exercise 9.5 to prove the following proposition that can be found in Park, Cho, and Park [220].

Proposition 9.1 Let $\mathbf{v} \in \{-1, 1\}^n$ be an asymptotically stable equilibrium state of the gBSB neural net. If $w_{ij} = 0$ for $i = 1, \dots, n$, then none of the vertices of H_n that differ from \mathbf{v} in just one component is an equilibrium state of the gBSB neural net.

- 9.7** Use Lemma 9.2 from Exercise 9.5 to prove the following proposition that can be found in Park, Cho, and Park [220]. Before stating the proposition, we define the *Hamming distance* between two vectors whose components are -1 or 1 as the number of components in which the two vectors differ. Given two vectors \mathbf{x} and \mathbf{y} , the Hamming distance between them is denoted $HD(\mathbf{x}, \mathbf{y})$.

Proposition 9.2 Let $\mathbf{v} \in \{-1, 1\}^n$ be an asymptotically stable equilibrium state of the gBSB neural net and let k be an integer from the set $\{1, 2, \dots, n\}$. If

$$\left(\sum_{j=1}^n w_{ij} x_j + b_i \right) v_i > 2k \max_j |w_{ij}|,$$

then a vertex \mathbf{x} of H_n , such that $0 < HD(\mathbf{x}, \mathbf{v}) \leq k$, is not an equilibrium state of the gBSB neural net.



CHAPTER 10

Genetic and Evolutionary Algorithms

The word “algorithm” itself is quite interesting; at first glance it may look as though someone intended to write “logarithm” but jumbled up the first four letters. The word did not appear in *Webster’s New World Dictionary* as late as 1957; we find only the older form “algorism” with its ancient meaning, i.e., the process of doing arithmetic using Arabic numerals.” . . . “Finally, historians of mathematics found the true origin of the word algorism: it comes from the name of a famous Persian textbook author, Abu Ja’far Mohammed ibn Mūsā al-Khowārizmī (c. 825) . . . Al-Khowārizmī wrote the celebrated book *Kitab al jabr w’al-muqabala* (“Rules of restoration and reduction”) . . . The modern meaning for algorithm is quite similar to that of *recipe, process, method, technique, procedure, routine*, except that the word “algorithm” connotes something just a little different.

—*The Art of Computer Programming, Volume 1* [162, pp. 1–3]

Genetic algorithms, introduced by Holland [124] in the 1960s, have been originally proposed to model adaptive processes, both natural and man-made. In this chapter we discuss applications of genetic algorithms to solving optimization problems. Genetic algorithms were inspired by biological genetics and evolutionary processes. This is the reason for biological terms being used in the context of genetic algorithms. One of the variants of genetic algorithms, called the evolutionary algorithm, is also presented.

10.1 Genetics as an Inspiration for an Optimization Approach

The *gene* is the basic unit of inheritance. Genes are pieces of DNA, which is the material inside the cell nucleus, that carry the genetic instructions for making living organisms. Many traits are determined by genes inherited from the parents of the individual. Consider gene \mathcal{A} that determines a certain trait, such as seed coat texture—round or wrinkled. Thus gene \mathcal{A} has two alternative forms. These different forms of the same gene are called *alleles*, or *allelic genes*. The allelic forms of each gene are symbolized by the uppercase and lowercase of the same letter. In

our example, gene A has two alleles: allele A and allele a . One of the alleles controls the round seed coat texture, while the other allele controls the wrinkled texture. Genes are organized into *chromosomes*, which are small packages of genes in the nucleus of a cell. Different kinds of organisms have different numbers of chromosomes. Humans have 23 pairs of chromosomes, including two sex chromosomes. In mating, each parent contributes one chromosome to each pair of the chromosomes of an offspring, so the offspring gets half of his or her chromosomes from his or her mother and half from his or her father. The assemblage of genes, encrypted as a string of DNA molecules, is called the *genome*. The human genome consists of about 35,000 genes. The particular gene constitution is called the *genotype*. The expression of the genotype is the *phenotype*, which is the way we look and behave.

Genetic algorithms are a result of an effort to model adaptation phenomena in natural and artificial systems. The term *adaptation* can be viewed as the composition of two words: *ad* + *aptare*, meaning “to fit to” [124, p. 1]. “Just what are adaptation’s salient features? We can see at once that adaptation, whatever its context, involves a progressive modification of some structure or structures. These structures constitute the grist of the adaptive process, being largely determined by the field of study. Careful observation of successive structural modifications generally reveals a basic set of structural modifiers or operators; repeated action of these operators yields the observed modification sequences” [124, p. 3].

In this chapter we describe two genetic algorithms and apply them to solving optimization problems. An optimization problem of a real-valued function, $f : \Omega \rightarrow \mathbb{R}$, requires finding an element $x \in \Omega$ such that f is maximized or minimized. The function f is called the *fitness function*. Suppose that we wish to maximize f over the parameter space Ω . The parameter space, Ω , of an underlying optimization problem is also called the *phenotype space*. A classical genetic algorithm works with binary representations of the points from the parameter space, rather than directly with the elements of Ω . We will refer to the classical genetic algorithms as the canonical algorithms. The strings of binary digits (1s and 0s) corresponding to these points are called chromosomes. The space of chromosomes is called the *genotype space*. In general, the chromosomes in the genotype space are abstract mathematical objects on which the so-called genetic operators work. After a stopping criterion for the algorithm is satisfied, a decoding function transforms a solution chromosome from the genotype space into the phenotype space. Genetic algorithms are derivative-free optimization algorithms that use only the fitness function. They do not use the fitness function derivatives. The higher the fitness of a chromosome, the better the candidate solution encoded by the chromosome. Genetic algorithms work with populations of chromosomes rather than individual chromosomes. A genetic algorithm is an iterative process that at each iteration uses the information obtained from testing a current population of chromosomes to direct its search into promising areas of the search space. Thus, generating a suitable population of chromosomes at each iteration is the key to the algorithm performance. In the simplest form, creating a new population of chromosomes to be tested involves three genetic operators: selection, crossover, and mutation. These genetic operators work on genotypes that encode phenotypes.

The role of the selection operator is to pick a population of chromosomes for further processing, called reproduction. Chromosomes are selected into a mating pool for reproduction with a probability proportional to their fitness.

The crossover operator randomly exchanges substrings between two parent chromosomes to create two offspring. The simplest crossover operator is the one-point crossover, where a crossing site is randomly picked and the corresponding substrings are exchanged, as illustrated

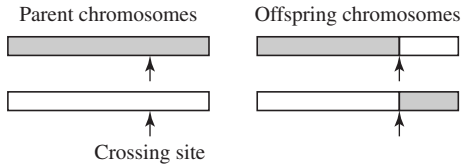


Figure 10.1 The one-point crossover.

in Figure 10.1. For example, the two parent chromosomes

1010101 and 1111010

could have the crossing site after the fourth bit in each to produce the two offspring

1010010 and 1111101.

The mutation operator randomly flips an arbitrary bit in a chromosome. For example, the chromosome 010101 might undergo a mutation in its third position to yield a new chromosome 011101.

We now provide a description of the canonical genetic algorithm using the above components.

CANONICAL GENETIC ALGORITHM (GA)

For a given optimization problem characterized by a fitness function f :

STEP 1 Choose a method of encoding candidate solutions into chromosomes; that is, decide upon a map between the phenotype and genotype spaces. (A commonly used encoding method is to concatenate the binary strings representing each real parameter value x_i .)

STEP 2 Start with a randomly generated population of chromosomes.

STEP 3 Evaluate the fitness of each chromosome in the population.

STEP 4 Replace the current population with a new population by applying selection, crossover, and mutation operations on the current population.

STEP 5 If a stopping criterion is satisfied, stop and choose the fittest chromosome from the population as a solution, else go to Step 3.

Another type of a genetic algorithm that we analyze in this chapter attempts to avoid the need for a decoding function by using the chromosome representation as close as possible to their representation in the phenotype space.

10.2 Implementing a Canonical Genetic Algorithm

We now discuss implementation issues of the canonical genetic algorithm (GA) when applied to solving optimization problems. We consider optimization problems of the form

$$\text{maximize } g(\mathbf{x}), \quad \mathbf{x} \in \Omega \subset \mathbb{R}^n. \quad (10.1)$$

A maximizing point or maximizer is denoted \mathbf{x}^* . We note that there is no loss of generality by considering only maximization problems, for if we are asked to minimize g , we can transform

the minimization problem into the maximization problem by means of a simple transformation,

$$\min_x g(\mathbf{x}) = -\max_x (-g(\mathbf{x})). \quad (10.2)$$

Also note that if \mathbf{x}^* is a minimizer of $g(\mathbf{x})$, then \mathbf{x}^* is a maximizer of $-g(\mathbf{x})$, and vice versa. Once an objective function of the maximization problem is given, the next step is to transform the objective function into a fitness function. For example, we can use the transformation

$$f(\mathbf{x}) = g(\mathbf{x}) - g_{\min}, \quad (10.3)$$

where $f(\mathbf{x})$ denotes a fitness function and g_{\min} is the minimum value of $g(\mathbf{x})$ for the current population of chromosomes. Note that the constructed fitness function is indeed nonnegative on the chromosomes of the current population.

We now discuss the implementation of Step 1 of the canonical GA presented in the previous section. We first analyze the case when the fitness function, f , is a function of only one parameter x . Suppose that

$$x \in [x_{\min}, x_{\max}] = \Omega.$$

A common encoding method, which transforms points from the phenotype space into a genotype space, is the binary encoding. In our example the phenotype space is the interval $[x_{\min}, x_{\max}]$. If we decide to use l bits to represent a point x , then this corresponds to mapping the interval $[x_{\min}, x_{\max}]$ into the discrete genotype space $[0, 2^l - 1]$ consisting of 2^l binary strings. Thus, the number of bits l determines the size of the genotype space and hence the resolution, or accuracy, of the subsequent search. We can control the resolution by appropriately choosing the parameter l . Let r denote the desired resolution. Then, the parameter l that results in the required resolution can be found by solving the equation

$$r = \frac{x_{\max} - x_{\min}}{2^l - 1}. \quad (10.4)$$

If the fitness function depends on many parameters, then we can encode the elements of the phenotype space by concatenating the binary strings representing each component of $\mathbf{x} = [x_1 \ x_2 \ \cdots \ x_n]^T$. For example, let $f = f(x_1, x_2)$, where $x_1 \in [x_{1,\min}, x_{1,\max}]$ and $x_2 \in [x_{2,\min}, x_{2,\max}]$. If we use 10 bits to encode each variable, then the genotype space will consist of 20-bit binary strings and the size of the genotype space is $2^{10} \times 2^{10} = 1,048,576$.

After initializing randomly a population of chromosomes, we evaluate the fitness $f(\mathbf{x})$ of each chromosome \mathbf{x} in the population. This requires decoding from the genotype space into the phenotype space. Suppose, as before, that the fitness function depends on two parameters x_1 and x_2 and that we used the standard encoding method described above. We now have to convert a given binary string into two real numbers $x_1 \in [x_{1,\min}, x_{1,\max}]$ and $x_2 \in [x_{2,\min}, x_{2,\max}]$. We can perform the required decoding process in three steps. First note that the first l bits of the chromosome \mathbf{x} correspond to x_1 while the remaining bits correspond to x_2 . We first convert the strings of l bits into integers. Thus, for example, if $l = 4$ and the chromosome is 10010111, which is an 8-bit string, then the first 4 bits yield

$$x_{1,int} = 1 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 = 9,$$

while the remaining 4 bits yield

$$x_{2,int} = 0 \times 2^3 + 1 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 = 7.$$

Then, we convert $x_{1,int}$ from the range $[0, 2^l - 1]$ into the range $[0, (x_{1,max} - x_{1,min})]$ by multiplying $x_{1,int}$ by

$$\frac{x_{1,max} - x_{1,min}}{2^l - 1}. \quad (10.5)$$

To obtain the corresponding value of x_1 , we add $x_{1,min}$ to (10.5) to get

$$x_1 = x_{1,min} + x_{1,int} \frac{x_{1,max} - x_{1,min}}{2^l - 1}. \quad (10.6)$$

We obtain x_2 proceeding in an analogous manner. In general, if \mathbf{x} has n components, the mapping from the phenotype space into the genotype space has the form

$$x_i = x_{i,min} + x_{i,int} \frac{x_{i,max} - x_{i,min}}{2^l - 1}, \quad i = 1, 2, \dots, n \quad (10.7)$$

The next step of the GA algorithm is the application of the selection operator to pick chromosomes from the current population for reproduction to create the new population. In the canonical GA, a pair of chromosomes is selected from the current population with the probability of selection being proportional to the fitness values of the chromosomes. Selection is done with replacement, which means that the same chromosome can be selected more than once for reproduction. A common method of selecting parent chromosomes is the so-called *roulette-wheel method*. In this method each chromosome is given a slice of a wheel. The area of the wheel slice is proportional to the chromosome's fitness. The wheel is spun, and the chromosome is selected on whose wheel slice the ball comes to rest. The roulette-wheel method can be implemented in the following steps:

1. Calculate the fitness of each chromosome in the current population.
2. Calculate the total fitness of the current population by summing the fitnesses of each chromosome in the current population.
3. Divide the fitness of each chromosome by the total fitness.
4. Calculate the cumulative fitness of the current population using the scaled fitness of each chromosome found in the previous step. You can use MATLAB's command `cumsum` for this purpose.
5. Return the first population member whose corresponding cumulative fitness is greater than the random number from the interval $(0, 1)$. The following MATLAB command can be used here: `find(cum_fitness - rand > 0)`, where `cum_fitness` is the vector of scaled cumulative fitnesses found in step 4.
6. Repeat the above step to get the second parent chromosome.

Using the parent chromosomes picked by the selection operator, we form two offspring by crossing over, with probability p_c , at a randomly chosen crossing site. The uniform probability p_c is called the crossover probability or crossover rate. If no crossover takes place, two offspring chromosomes are formed that are exact copies of their parent chromosomes. One can also use a multipoint crossover operator. The two-point crossover operation is illustrated in Figure 10.2.

The offspring chromosomes obtained from the crossover are then acted upon by the mutation operator. Mutation is performed at each bit of the offspring chromosomes with probability p_m . The probability p_m is called the mutation probability or mutation rate. The mutation operator

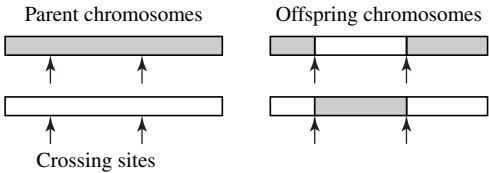


Figure 10.2 The two-point crossover.

randomly flips some of the bits in a chromosome. For example, the chromosome 101010 might be mutated in its third bit to yield 100010. The mutation probability is usually very small—for example, $p_m = 0.01$.

Each iteration, also called a *generation*, of the GA replaces the current population with a new one. A typical number of iterations of a GA algorithm is from 50 to 500 or more. The entire set of generations is called a *run*.

◆ **Example 10.1**

We apply the canonical GA described above to maximize the function

$$f(x) = -15 \sin^2(2x) - (x - 2)^2 + 160, \quad |x| \leq 10. \quad (10.8)$$

This function has many local maximizers and minimizers, as can be seen from its plot shown in Figure 10.3. The canonical GA was applied to maximize the function given by (10.8) with the chromosomes composed of 12 bits, $p_c = 0.75$, and $p_m = 0.01$. In Figure 10.4 a summary of a typical run is shown. In this run the best candidate solution was $x = 1.565324$ whose corresponding objective value is $f(x) = 159.809259$. The best solution that was observed in other runs resulted in $f(1.575092) = 159.818346$.

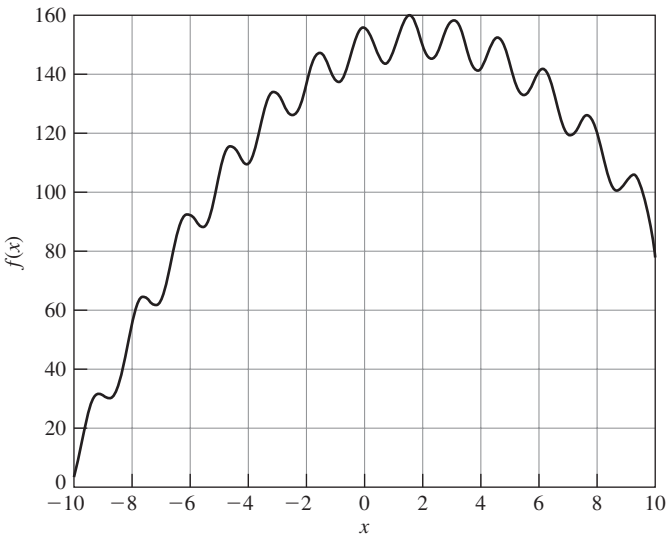


Figure 10.3 A plot of the function given by (10.8).

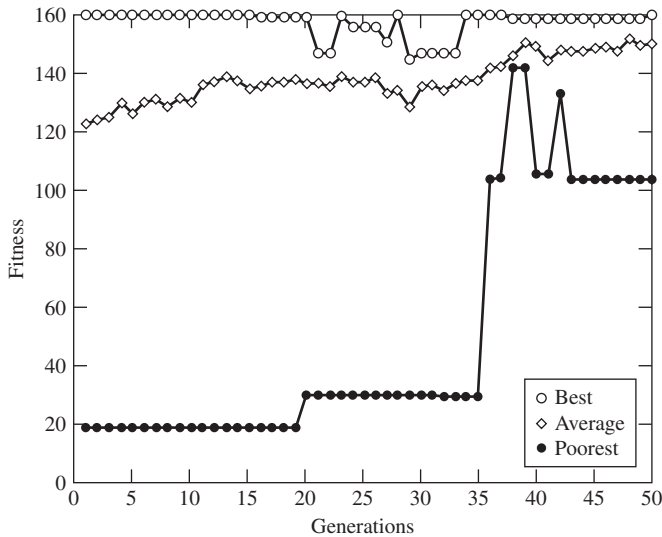


Figure 10.4 Summary of a sample run of the canonical GA when searching for a maximizer of the function given by (10.8).

In the next section we will try to get answers to some whys and hows of genetic algorithms.

10.3 Analysis of the Canonical Genetic Algorithm

In the analysis of genetic algorithms, the notion of a *schema* will be useful. A schema is a set of bit strings. It can be described using a template made up of ones, zeros, and asterisks, which represent “don’t cares.” Thus, the asterisk $*$ matches either 0 or 1. For example, the schema $H = *0*$ contains the binary strings 000, 001, 100, and 101. Note that a bit string of length l belongs to 2^l different schemas. For example, the string 10 of length $l = 2$ belongs to the schemas $**$, $*0$, $1*$, and 10 . Similarly, when $l = 3$, then any bit string of length $l = 3$ belongs to $2^3 = 8$ schemas. Elements of a set of N bit strings of length l are members of between 2^l and $N \times 2^l$ schemas. If all the strings are identical, then they belong to 2^l schemas; otherwise there may be up to $N \times 2^l$ schemas to which the N strings of length l can belong.

The parameters that characterize a given schema are its *order* and *defining length*. The order of a schema H , denoted $o(H)$, is the number of its fixed bits. For example, the order of $H = *01*1$ is $o(H) = 3$ because H has three bits fixed. The defining length of the schema H , denoted $d(H)$, is the distance between its first and the last fixed bits, which is the difference between the positions of the rightmost fixed symbol and the leftmost fixed symbol. The defining length of $H = *0*1$ is $d(H) = 4 - 2 = 2$. In further discussion, the term “schema” will be used to denote a set of binary strings represented by a template, as well as the template itself.

The capability of a chromosome to solve a given optimization problem—that is, the chromosome performance—is measured by its fitness. Similarly, we can characterize the performance of a schema by its average fitness. We now derive an expression for the average fitness of a given schema. Let x be a chromosome that matches the template H , that is, $x \in H$. Then, we say

that “ \mathbf{x} is an instance of H .” Let H be a schema with at least one instance in the k th generation. Denote by $e(H, k)$ the number of instances of H in the k th generation, and let $\bar{u}(H, k)$ be the average fitness of H in the k th generation, that is,

$$\bar{u}(H, k) = \frac{\sum_{\mathbf{x} \in H} f(\mathbf{x})}{e(H, k)}. \quad (10.9)$$

When the genetic algorithm is explicitly evaluating the fitness of N chromosomes from the current generation, it is implicitly evaluating the average fitness of up to $N \times 2^l$ schemas to which the N chromosomes belong.

We now find $\mathcal{E}(H, k + 1)$, the expected number of instances of the schema H in the $(k + 1)$ st generation, given the k th generation. We first analyze the effect of the selection operator only. Assume that the roulette-wheel method is used in selecting parent chromosomes to the mating pool for reproduction. The probability that the chromosome \mathbf{x} is going to be selected to the mating pool, $M(k)$, is

$$\frac{f(\mathbf{x})}{\sum_{i=1}^N f(\mathbf{x}_i^{(k)})},$$

where $\mathbf{x}_i^{(k)}$ refers to the i th chromosome in the k th generation and $f(\mathbf{x})$ is the fitness of the chromosome \mathbf{x} . Therefore, the expected number of chromosomes equal to \mathbf{x} that are selected to the mating pool is

$$N \frac{f(\mathbf{x})}{\sum_{i=1}^N f(\mathbf{x}_i^{(k)})} = \frac{f(\mathbf{x})}{\bar{f}(k)},$$

where

$$\bar{f}(k) = \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}_i^{(k)})$$

is the average fitness of the k th generation. Let $m(k)$ be the number of chromosomes in $M(k)$ that are instances of H . The expected value, $\mathcal{M}(H, k)$, of the number of chromosomes in the mating pool that are instances of H is

$$\frac{\sum_{\mathbf{x} \in H} f(\mathbf{x})}{\bar{f}(k)}. \quad (10.10)$$

Substituting (10.9) into (10.10) yields an alternative expression for the expected value of the number of chromosomes in the mating pool that are instances of H :

$$\boxed{\mathcal{M}(H, k) = \frac{\bar{u}(H, k)}{\bar{f}(k)} m(H, k)} \quad (10.11)$$

It follows from the above equation that a schema whose average fitness is greater than the population average will be found in an increasing number of chromosomes.

We will now quantify the effects of the one-point crossover on a given schema. Given two chromosomes, where at least one of them is an instance of the schema H , we say that the schema H survives crossover operation if at least one of the offspring is an instance of H . We will give a lower bound on the probability $S_c(H)$ that the schema H will survive the one-point crossover. A schema survives if the crossing site occurs outside the defining length. Crossovers that take place within the defining length can destroy the schema. For example, let $H = *00 * *$, and let the

parent chromosomes be $p_1 = 10011$ and $p_2 = 11100$. Note that $p_1 \in H$, while $p_2 \notin H$. For the crossing site after the third position, there are two offspring chromosomes: $s_1 = 10000$, which is an instance of H , and $s_2 = 11111$, which is not an instance of H . Thus, the schema H survived this crossover. Suppose now that the crossing site is after the second position, which is within the defining length of the schema H . The offspring chromosomes are $s_1 = 10100$ and $s_2 = 11011$. None of the offspring chromosome is an instance of H , which means that H was destroyed by the crossover. If both parent chromosomes are instances of H , then the offspring chromosomes are also instances of H irrespective of the location of the crossing site. An upper bound on the probability $D_c(H)$ that the schema H will be destroyed by the one-point crossover is

$$D_c(H) \leq p_c \frac{d(H)}{l-1}. \quad (10.12)$$

Therefore, a lower bound on the probability $S_c(H)$ that H survives is

$$S_c(H) = 1 - D_c(H) \geq 1 - p_c \frac{d(H)}{l-1} \quad (10.13)$$

Thus, the probability of survival under one-point crossover is higher for short schemas.

A schema H survives under mutation of its element, if all defined bits of this element of H remain unchanged. Let $S_m(H)$ denote the probability that the schema H survives under mutation of an instance of H . For each defined bit, the probability that it will not be mutated is $1 - p_m$. The probability that no defined bit of H will be mutated is obtained by multiplying $1 - p_m$ by itself $o(H)$ times. Therefore,

$$S_m(H) = (1 - p_m)^{o(H)} \quad (10.14)$$

We now combine (10.11), (10.13), and (10.14) to obtain a lower bound on the expected number of instances of H in the $(k+1)$ st generation, given the k -generation:

$$\mathcal{E}(H, k+1) \geq \frac{\bar{u}(H, k)}{\bar{f}(k)} m(H, k) \left(1 - p_c \frac{d(H)}{l-1}\right) (1 - p_m)^{o(H)} \quad (10.15)$$

The above is known as the *schema theorem*. It follows from our discussion that the canonical genetic algorithm implicitly processes, in parallel, schemas. This implicit computation is performed with no operation being explicitly performed on schemas. To illustrate the last statement, consider the string 101. As the fitness of this string is evaluated, partial information is obtained about the expected fitness of all possible schemas to which the evaluated string belongs. In our example, computing the fitness of the string 101 gives partial information about the fitness of strings in the schemas $***$, $**1$, $*0*$, $1**$, $*01$, $1*1$, \dots .

10.4 Simple Evolutionary Algorithm (EA)

Canonical genetic algorithms use binary representations of the elements of the parameter space. The reason behind this is the schema theory, discussed in the previous section, which allows one to analyze canonical genetic algorithms. When using a binary representation of the elements of the parameter space, canonical algorithms use encoding and decoding functions. A binary

representation may have some disadvantages though. Indeed, let $h : \{0, 1\}^l \rightarrow \Omega$ denote a decoding function used in the canonical genetic algorithm implementation. Then the original problem of optimizing $f : \Omega \rightarrow \mathbb{R}$ is transformed into the problem of optimizing the combined objective function

$$f(h) : \{0, 1\}^l \rightarrow \mathbb{R}.$$

The combined objective function to be optimized may be more complex than the original objective function. For example, it may introduce extra nonlinearities or add extra optimizers, which can make the search process for a global optimizer more difficult.

In contrast to the canonical genetic algorithms, in the optimization procedure discussed next, the parameter space is no longer coded using the binary representation. The genetic operators work directly on real-valued vectors. A term used to describe an optimization algorithm using genetic-type operators working on the real-valued vectors is an *evolutionary algorithm*. In what follows, we present a very simple evolutionary algorithm for solving optimization problems of the form given by (10.1).

SIMPLE EVOLUTIONARY ALGORITHM (EA)

For a given optimization problem of the form (10.1), perform the following steps:

STEP 1 Initialize a population of candidate solutions (chromosomes) to the problem. Each chromosome is a real-valued vector of the same length as the vector \mathbf{x} of decision variables.

STEP 2 Evaluate the fitness of each candidate solution in the population.

STEP 3 Replace the current population with a new population by applying selection, crossover, and mutation operations on the current population. The roulette-wheel selection method can be used for selection. A possible implementation of the crossover operation is to use the average crossover operator. This operator takes two parents, say \mathbf{x}_j and \mathbf{x}_k , and returns one offspring

$$\mathbf{x} = \frac{1}{2}(\mathbf{x}_j + \mathbf{x}_k).$$

Another implementation of the crossover operator is a convex combination of two parents. Specifically, let $r \in (0, 1)$. Then, the offspring is

$$\mathbf{x} = r\mathbf{x}_j + (1 - r)\mathbf{x}_k.$$

As for the mutation operation, let \mathbf{x} be a chromosome to be mutated. Then, an implementation of the mutation operation could have the form

$$\mathbf{x}' = \mathbf{x} + r\mathbf{d},$$

where \mathbf{d} is a randomly generated feasible vector and $r \in (0, 1)$. Note that if the feasible set is convex, then the above genetic operators produce chromosomes that are elements of the feasible set. Otherwise, conditions should be imposed to ensure that candidate solutions are the elements of the feasible set.

STEP 4 If a stopping criterion is satisfied, terminate the algorithm and choose the fittest chromosome from the population as a solution, else go to Step 2.

In addition, both in the canonical genetic algorithm and in the evolutionary algorithm, the elitist strategy, introduced by K. DeJong in 1975, can be incorporated. This strategy copies the best chromosomes from the old population into the next population to prevent losing the best solutions in the succeeding iterations. As observed by Mitchell [206], among others, the elitist strategy may lead to the domination of a population by a “super chromosome.” However, it appears that the elitist strategy improves genetic and evolutionary algorithms performance.

10.5 Evolutionary Fuzzy Logic Controllers

In this section we employ evolutionary algorithms to design tracking controllers. We describe three design techniques for evolutionary fuzzy logic controllers for the step-lane-change maneuver of a ground vehicle. The presentation of this section is based on the paper by Lee and Žak [179]. We mention that the step-lane-change maneuver is an important aspect in the design of the Intelligent Vehicle/Highway System (IVHS), the highway of the future.

An essential element in combining an evolutionary algorithm (EA) and a fuzzy logic controller (FLC) is the encoding method used to represent a fuzzy rule base as a chromosome on which the EA operates. Selecting an appropriate encoding procedure is a critical factor that affects the efficiency of the EA and the final solution. Another essential element is the fitness function. Each chromosome defines a fuzzy rule base. After conducting a simulation of the structure depicted in Figure 10.5 and measuring the tracking error, the fitness of the chromosome representing the fuzzy rule base of the simulated FLC is calculated. With all the chromosomes assigned a fitness value, genetic operators yield a new generation and the entire procedure is repeated until a stopping criterion, usually the number of iterations or the fitness value, is satisfied.

10.5.1 Vehicle Model and Control Objective

We use a “bicycle model” of a front-wheel-steering vehicle. Modeling of a ground vehicle is covered in detail in Subsection 1.7.2. The standard coordinate system of the Society of Automotive Engineers (SAE) is used. A schematic of the model used is shown in Figure 10.6. The modeling equations are

$$\begin{aligned}\dot{x}_1 &= -v_x x_3 + \frac{1}{m}(F_{yf} + F_{yr}), \\ \dot{x}_2 &= x_3, \\ \dot{x}_3 &= \frac{2}{I}(l_1 F_{yf} - l_2 F_{yr}), \\ \dot{x}_4 &= -v_x x_2 - x_1,\end{aligned}$$

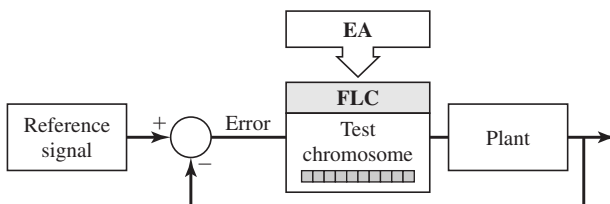


Figure 10.5 A configuration used to design a fuzzy logic controller using an evolutionary algorithm.

Table 10.1 Vehicle Data and Symbols Used in Modeling Equations

δ	Steering angle of the front wheel	
F_{yf}	Lateral force of the front tire	
F_{yr}	Lateral force of the rear tire	
m	Total mass of the vehicle	1280 kg
I	Moment of the inertia of the vehicle	2500 kg·m ²
l_1	Distance from center of gravity to front axle	1.2 m
l_2	Distance from center of gravity to rear axle	1.22 m
v_x	Longitudinal velocity of the vehicle	18.3 m/sec

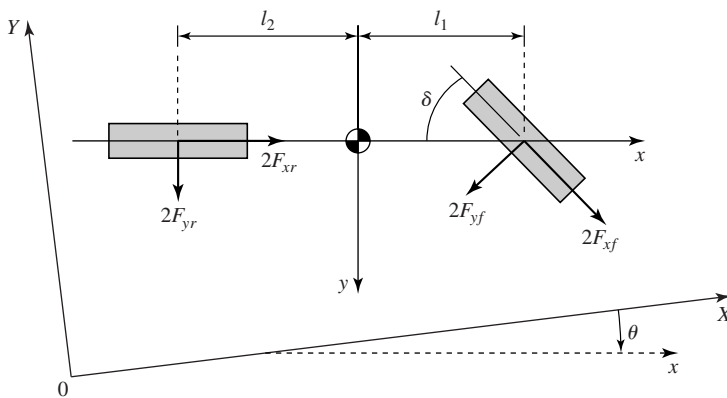


Figure 10.6 Schematic of a bicycle model of a front-wheel-steering vehicle.

where $x_1 = \dot{y}$ is the lateral velocity of the vehicle, $x_2 = \theta$ is the yaw angle, $x_3 = \dot{\theta}$ is the yaw velocity, and $x_4 = Y$ is the lateral position of the vehicle. Other parameters and their values are given in Table 10.1. The lateral forces $F_{yf} = F_{yf}(\alpha_f)$ and $F_{yr} = F_{yr}(\alpha_r)$ are nonlinear functions of respective slip angles with a typical plot of a lateral force versus tire slip angle shown in Figure 1.15. Recall that the slip angles for the front and rear tires are

$$\alpha_f = \delta - \frac{l_1 x_3 + x_1}{v_x} \quad \text{and} \quad \alpha_r = \frac{l_2 x_3 - x_1}{v_x}.$$

The input is the steering angle δ of the front wheel and the output is the lateral position of the vehicle, $x_4 = Y$.

The step-lane-change maneuver, along with cruising, traveling path selection, and highway exiting and entering, is one of the essential elements of automated vehicle control. However, automation of these tasks requires a roadway reference system that involves sensing for lane detection. In our simulation, we assume that the road reference system is already installed. Our objective is to design a fuzzy logic controller that forces the vehicle to track a given reference lateral position signal r . The reference signal is generated by a reference signal generator whose transfer function is $1/(s^3 + 1.75s^2 + 2.15s + 1)$. The tracking error e is the difference between the reference lateral position signal r and the actual lateral position Y of the vehicle. The role of

the fuzzy controller is to produce the steering wheel angle δ using the information of the tracking error e and the change in the tracking error \dot{e} so that the vehicle tracks the reference signal. The SIMULINK block diagram of the closed-loop vehicle system is depicted in Figure 10.7. In the following, we describe three different methods for an evolutionary algorithm (EA)-aided fuzzy logic controller (FLC) design.

10.5.2 Case 1: EA Tunes Fuzzy Rules

It is assumed in this case that the membership functions are fixed. In our simulation experiment, we use the membership functions shown in Figure 10.8. The labels of fuzzy sets LN, N, Z, P, and LP denote large negative, negative, zero, positive, and large positive. There are five fuzzy

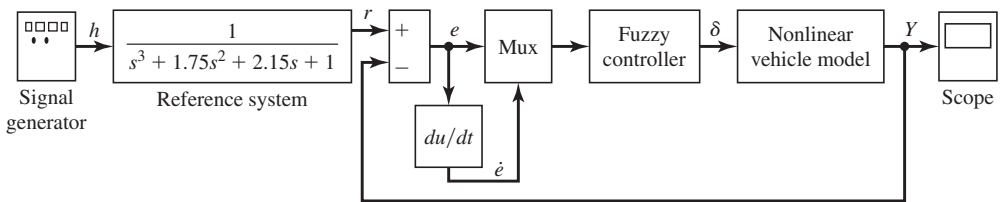


Figure 10.7 SIMULINK block diagram of the closed-loop vehicle system driven by a fuzzy logic controller.

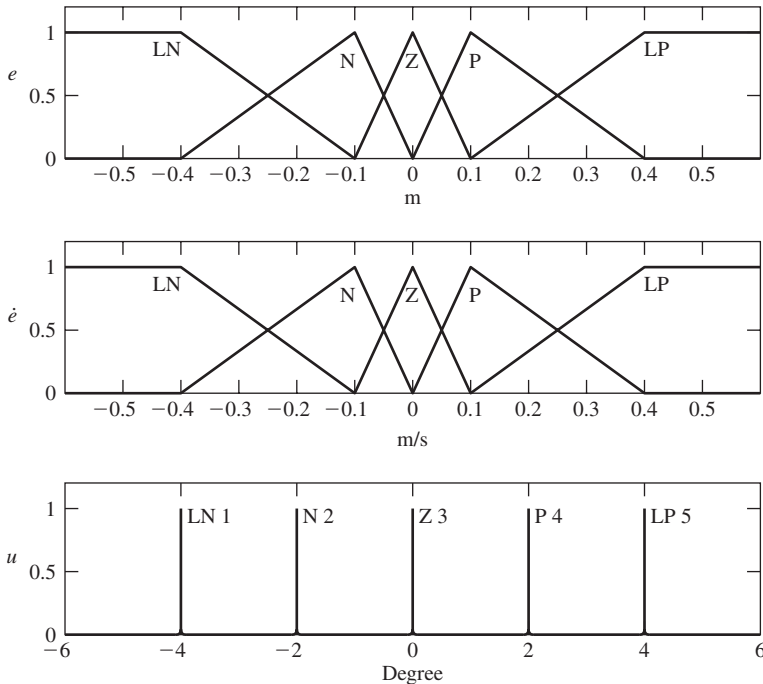


Figure 10.8 Fuzzy membership functions for e , \dot{e} , and u .

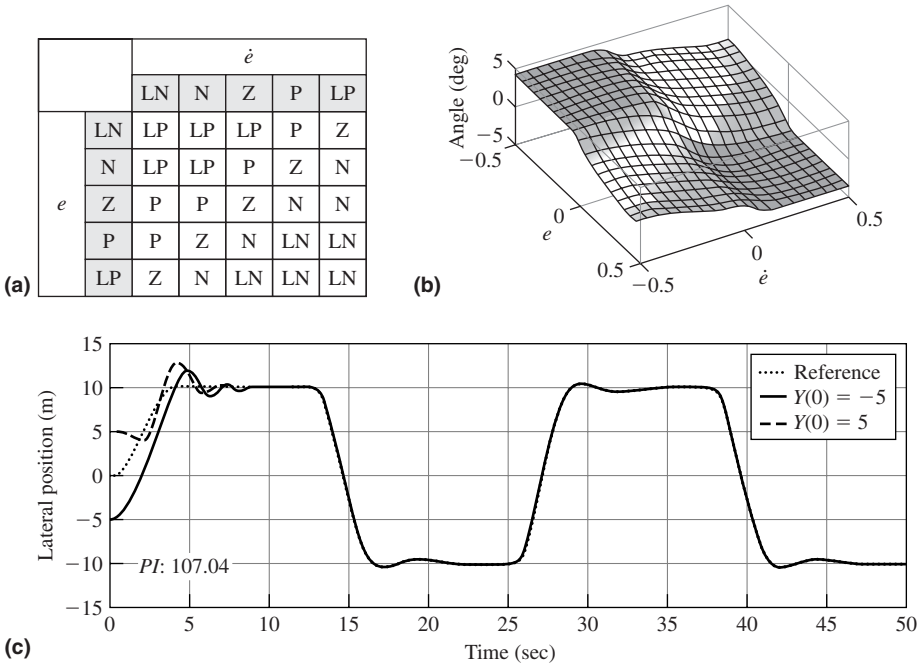


Figure 10.9 FLC constructed using a fuzzy rule base obtained heuristically. (a) Fuzzy rules table. (b) Control surface. (c) Tracking performance.

sets for e and five fuzzy sets for \dot{e} ; hence we can have 25 different rules of the form

IF e is LN AND \dot{e} is LN then δ is SP.
:
IF e is LP AND \dot{e} is LP then δ is SN.

These rules can be arranged into a rules table as shown in Figure 10.9(a). We mention here that the rules shown in Figure 10.9(a) were arrived at heuristically. Using the rules table, we construct an FLC using the center average defuzzifier given by (8.2). The control surface derived using the rules in Figure 10.9(a) is shown in Figure 10.9(b). Tracking performance for two different initial conditions $Y(0) = 5$ m and $Y(0) = -5$ m is illustrated in Figure 10.9(c). Tracking errors for each initial condition were calculated. We labeled the tracking error for the first initial condition as e_1 while the tracking error of the second initial condition was labeled as e_2 . We then constructed the performance index (PI) as the sum of the integrals of the square of the error (ISE) over the simulation time:

$$PI = \int_0^{50} (e_1^2(t) + e_2^2(t)) dt.$$

The value of the PI in the above experiment is 107.04. The performance index, PI , was evaluated using the circuit shown in Figure 10.10.

Of course, we could include more error trajectories in the expression for the PI . However, for the sake of simplicity, we decided to use only two tracking error trajectories in our simulation experiments.

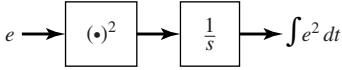


Figure 10.10 Circuit used to evaluate the performance index, PI .

We now employ an EA to generate a fuzzy rules table automatically. In order to represent the whole fuzzy rules table as a chromosome, we encode the labels of the fuzzy singletons corresponding to the controller outputs as integers, where LN is represented by 1, N by 2, Z by 3, P by 4, and LP by 5. This encoding of the fuzzy singletons of the control signal u is depicted in Figure 10.8. The actual control signal, however, is produced as the weighted sum of the fuzzy singletons using the center average defuzzifier. Because there are 25 rules, each chromosome is a string composed of 25 integers. For example, the rules table of Figure 10.9(a) can be represented as a chromosome of the form

5	5	5	4	3	5	5	4	3	2	4	4	3	2	2	4	3	2	1	1	3	2	1	1	1	1
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

In our experiment, we decided upon the population size of 30. The roulette-wheel method was used to select chromosomes to the mating pool. We employed the single-point crossover operator with the crossover rate $p_c = 0.9$. The role of the mutation operator was to replace each element of a chromosome randomly with an integer from the set $\{1, \dots, 5\}$. We used the mutation rate $p_m = 0.05$.

We also employed an elitist strategy where two best individuals of the current generation were passed unchanged to the next generation. The fitness function employed was the inverse of the PI discussed above, that is,

$$\text{Fitness} = \frac{1}{PI} = \frac{1}{\int_0^{50} (e_1^2(t) + e_2^2(t)) dt}. \quad (10.16)$$

Recall that $e_1(t)$ is the tracking error with the initial lateral position at 5 m, and $e_2(t)$ is the tracking error with the initial lateral position at -5 m. The input to the reference system was a square wave with the amplitude of 10 m and the frequency of 0.04 Hz. The rules table of the best FLC, after 50 iterations of the EA, is shown in Figure 10.11(a). Figure 10.11(b) shows the fuzzy control surface based on the rules table of Figure 10.11(a). Two different initial conditions, $Y(0) = 5$ m and $Y(0) = -5$ m, were used to test the tracking performance of the fuzzy controller; Figure 10.11(c) depicts the results. The vehicle lateral position becomes almost indistinguishable from the reference position after about 7 seconds. The value of the PI is 97.10, which is an improvement over the performance of the FLC obtained using heuristically generated fuzzy rules.

10.5.3 Case 2: EA Tunes Fuzzy Membership Functions

In this case, fuzzy rules are assumed to be known and fixed. The role of an EA is to tune the input and output membership functions. We use the heuristically obtained rules table of Figure 10.9(a) as the given fuzzy rule base. The fuzzy sets for each input can be completely described using two real numbers if the corresponding fuzzy membership functions are symmetric with respect to the origin and, at each point of their universe of discourse, the sum of the values of the membership functions equals 1. The controller output fuzzy sets are also assumed to be symmetric with respect to the origin. Under the above assumptions, the controller output fuzzy singletons can be completely characterized using just two parameters. Hence, six real parameters can completely characterize all the input and output membership functions. So the chromosomes are strings of

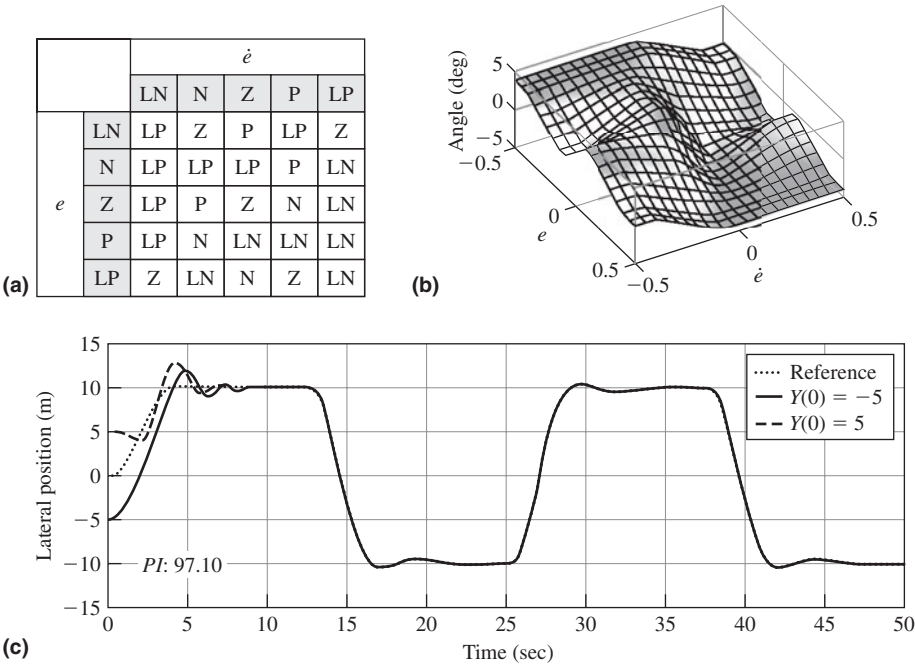


Figure 10.11 The best FLC of the 50th generation using a fuzzy rule base automatically generated by the EA. (a) Fuzzy rules table. (b) Control surface. (c) Tracking performance.

just six real parameters. For example, the fuzzy sets of Figure 10.8 can be represented as the string:

$$\begin{bmatrix} 0.1 & 0.4 & 0.1 & 0.4 & 2 & 4 \end{bmatrix}$$

We now proceed to describe the EA that operates on the above-described fuzzy rule base. The population size, the selection method, the crossover operator and its rate, the mutation operator and its rate, and the fitness function, as well as other parameters, are the same as in the previous case. The resulting fuzzy sets, after 50 iterations of the EA, are shown in Figure 10.12. The control surface and the tracking performance of the closed-loop system driven by the FLC constructed using fuzzy sets of Figure 10.12 are depicted in Figure 10.13. The *PI* value is 43.51. The tracking error of the lateral position almost disappears in about 4 seconds.

10.5.4 Case 3: EA Tunes Fuzzy Rules and Membership Functions

We observed in the previous cases that changing fuzzy rules or modifying the shape of membership functions can improve the performance of the resulting FLC. We now analyze an EA that simultaneously operates on the fuzzy rules and fuzzy membership functions. This approach is a modification of the technique described in references 138 and 149.

In general, an i th fuzzy rule can have the form

$$\text{IF } x_1 \text{ IS } F_1^i \text{ AND } \cdots \text{ AND } x_n \text{ IS } F_n^i, \text{ THEN } u_1 \text{ IS } \theta_1^i \text{ AND } \cdots \text{ AND } u_m \text{ IS } \theta_m^i,$$

where $1 \leq i \leq r$, the x_j 's are the state variables where $1 \leq j \leq n$, and the u_k 's are the controller's

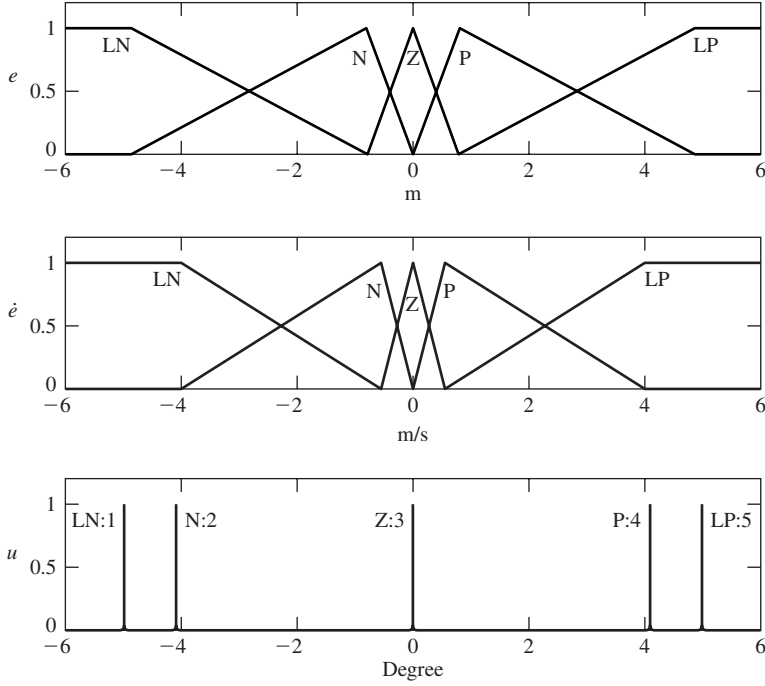


Figure 10.12 Input and output fuzzy membership functions of the best FLC of the 50th generation obtained by tuning fuzzy membership functions.

outputs, where $1 \leq k \leq m$. The F_j^i 's are fuzzy sets with trapezoidal membership functions, and the θ_k^i 's are the controller's output variables whose membership functions are singletons. A membership function, $\mu_{F_j^i}(x_j)$, of the fuzzy set, F_j^i , can be completely characterized using four parameters c_{ij} , d_{ij} , l_{ij} , and r_{ij} . It can be described as

$$\mu_{F_j^i}(x_j) = \begin{cases} \frac{x_j - l_{ij}}{c_{ij} - l_{ij} - d_{ij}} & \text{if } l_{ij} < x_j < c_{ij} - d_{ij}, \\ 1 & \text{if } c_{ij} - d_{ij} \leq x_j < c_{ij} + d_{ij}, \\ -\frac{x_j - r_{ij}}{r_{ij} - c_{ij} - d_{ij}} & \text{if } c_{ij} + d_{ij} < x_j < r_{ij}, \\ 0 & \text{otherwise,} \end{cases} \quad (10.17)$$

or, equivalently,

$$\mu_{F_j^i}(x_j) = \max \left(0, \min \left(1, \frac{x_j - l_{ij}}{c_{ij} - l_{ij} - d_{ij}}, -\frac{x_j - r_{ij}}{r_{ij} - c_{ij} - d_{ij}} \right) \right). \quad (10.18)$$

Figure 10.14 shows a generic fuzzy membership function of the form given by (10.17) or, equivalently, by (10.18). The fuzzy singletons, on the other hand, are described as

$$\mu_{\theta_k^i}(z) = \begin{cases} 1 & \text{if } z = \theta_k^i, \\ 0 & \text{otherwise.} \end{cases} \quad (10.19)$$

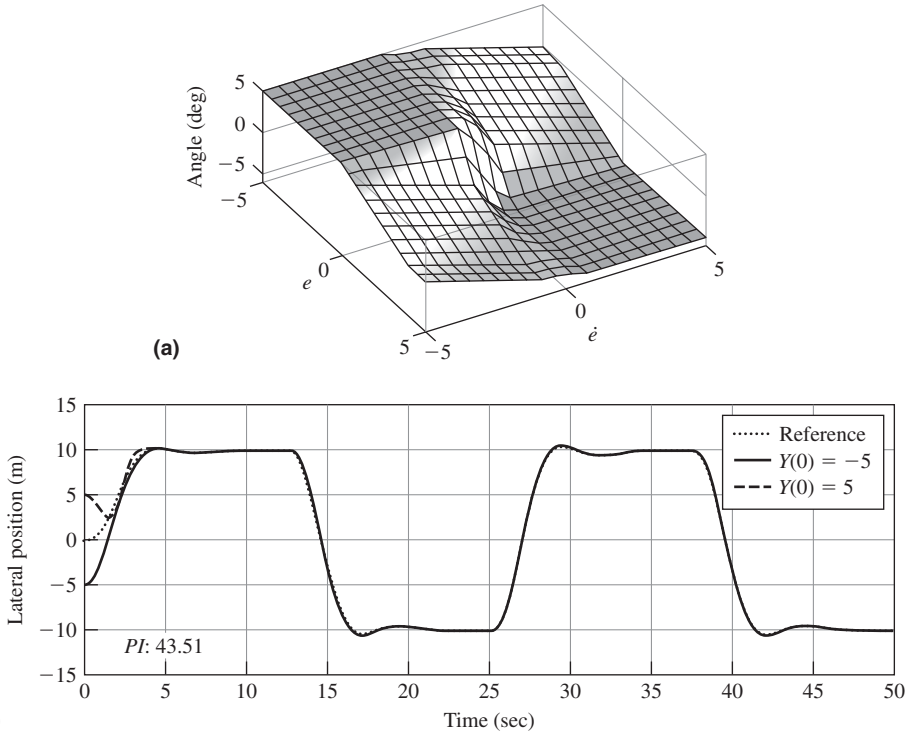


Figure 10.13 The best FLC of the 50th generation obtained by tuning fuzzy membership functions. (a) Control surface. (b) Tracking performance.

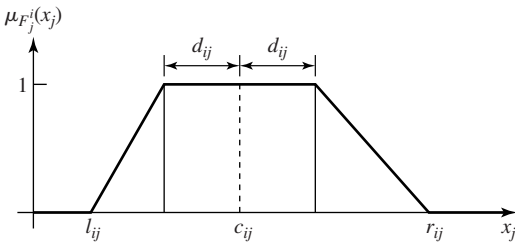


Figure 10.14 An example of a trapezoidal membership function.

The controller's signal u_k is calculated using the center average defuzzification formula,

$$u_k = \frac{\sum_{i=1}^r w_i \theta_k^i}{\sum_{i=1}^r w_i}, \quad (10.20)$$

where r is the number of fuzzy rules and $w_i = \mu_{F_1^i} \mu_{F_2^i} \times \cdots \times \mu_{F_n^i}$.

We now present a method of representing a given fuzzy rule base as a chromosome. In order to tune fuzzy rules and to alter the shape of membership functions, we require that a chromosome contains all the information pertaining to the structure of the fuzzy rules and the shape of membership functions explicitly. A chromosome is composed of two substructures: One substructure contains the fuzzy rules' structure matrix, whereas the other substructure contains the information concerning the parameters of the membership functions of the fuzzy rule base.

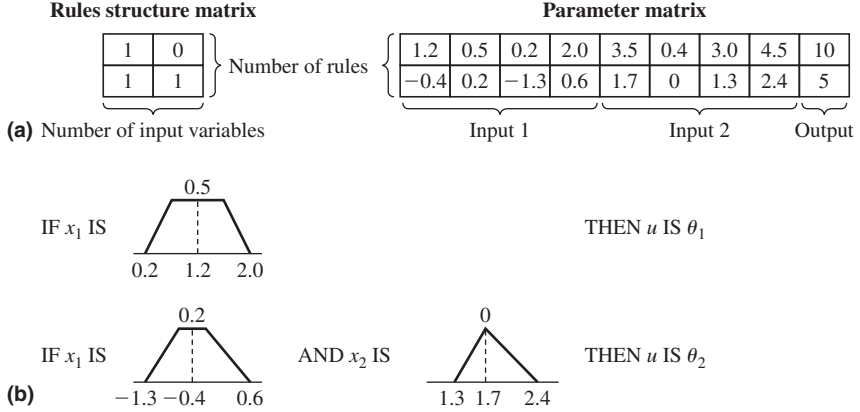


Figure 10.15 An example of a chromosome structure and its corresponding fuzzy rule base where the controller has two inputs, x_1 and x_2 , and one output u .

The fuzzy rules' structure matrix contains the information related to the absence or presence of the premise part of the fuzzy rule corresponding to the particular controller input. Hence the elements of this matrix are from the set $\{0, 1\}$.

The parameter matrix contains the parameters that describe fuzzy membership functions. As previously mentioned, there are four parameters for each fuzzy trapezoidal membership function and one for each fuzzy singleton. The parameter matrix has the same number of rows as the number of rules, and the number of columns is equal to four times the number of the controller inputs plus the number of the controller's outputs. Figure 10.15(a) shows an example of the rules' structure matrix and the parameter matrix for the controller with two inputs and one output whose operation is described by two IF-THEN rules. The fuzzy rules and membership functions corresponding to the rules' structure matrix and the parameter matrix in Figure 10.15(a) are depicted in Figure 10.15(b). Note that a 0 in the second column of the first row of the rules' structure matrix means that the premise part of the first fuzzy rule corresponding to the second controller input is absent; that is, the first rule does not involve the controller's second input.

The fitness function has the form

$$\text{Fitness} = \frac{1}{PI + \alpha \times (\text{number of rules})}$$

$$= \frac{1}{\int_0^{50} (e_1^2(t) + e_2^2(t)) dt + \alpha \times (\text{number of rules})}, \quad (10.21)$$

where $e_1(t)$ is the tracking error with the initial lateral position at 5 m, and $e_2(t)$ is the tracking error with the initial lateral position at -5 m. The reference signal used has the amplitude of 10 m and the frequency of 0.04 Hz over 50 seconds. The design parameter α is employed to penalize for the excessive number of fuzzy rules. In our simulations, $\alpha = 1$.

We now describe genetic operators utilized by the EA. Three genetic operators are used with the probability p_g . We first discuss a genetic operator that modifies the number of fuzzy rules. Let r_{new} be the number of fuzzy rules after the genetic operator is applied, and let r_{old} be the number of fuzzy rules before the genetic operator is applied. Then,

$$r_{\text{new}} = r_{\text{old}} + \Delta r,$$

where Δr is the correction term that we analyze next. The number of fuzzy rules cannot exceed

a prespecified bound r_{\max} . It is reasonable to assume that, on the average, r_{old} is about $r_{\max}/2$. Therefore, we wish to construct the correction term so that

$$\Delta r \in [-r_{\max}/2, r_{\max}/2].$$

Let Ψ be a random variable with $N(0, 1)$. The probability that Ψ will take its values in the interval $(-3, 3)$ is 0.99. Let $\text{integer}(\cdot)$ be the operator that produces the nearest integer value of its argument. Then, the random variable, $\text{integer}(\Psi r_{\max}/6)$, will take its value in the interval $(-\text{integer}(r_{\max}/2), \text{integer}(r_{\max}/2))$ with the probability 0.99. In view of the above, we use the following expression for updating the number of fuzzy rules:

$$r_{\text{new}} = r_{\text{old}} + \text{integer}(\Psi r_{\max}/6). \quad (10.22)$$

As follows from the discussion above, this particular form of the genetic operator will ensure that the resultant variation range of the rules' number is reasonable with respect to the maximum number r_{\max} of the rules allowed. If the number of new rules, r_{new} , falls outside the allowed range, then we set $r_{\text{new}} = r_{\text{old}}$.

The second genetic operator acts upon the elements of the rules' structure matrix by taking the binary complement of each element—that is, $1 \rightarrow 0$ and $0 \rightarrow 1$ with the probability p_g .

The third genetic operator acts upon the parameter matrix as per the following expression:

$$p_{ij}^{\text{new}} = p_{ij}^{\text{old}} + \eta_{ij} \Psi e^{1-F_{\text{old}}/F_{\max}}, \quad (10.23)$$

where p_{ij}^{old} is the current value of the (ij) th element of the parameter matrix, p_{ij}^{new} is the updated value of the (ij) th element of the parameter matrix, Ψ is a random variable with $N(0, 1)$, F_{old} is the fitness value of the chromosome being operated upon, and F_{\max} is the maximum fitness value of the chromosomes of the current population. Each design parameter η_{ij} was chosen according to the following recipe:

$$\eta_{ij} = (\text{maximum of the range of } p_{ij} - \text{minimum of the range of } p_{ij})/20.$$

The term $e^{1-F_{\text{old}}/F_{\max}}$ in equation (10.23) allows a chromosome with low fitness to vary the elements of the parameter matrix in a larger range than for the chromosomes with high fitness values.

In our implementation of the EA, we employ the elitist strategy to ensure that 20% of the fittest chromosomes in each generation are passed directly to the next generation without being acted upon by the genetic operators. The genetic operators, described above, are applied to the top 20% of the chromosomes in the current generation to produce the remaining 80% of the chromosomes of the new population.

After applying genetic operators to the chromosomes of the current generation, postprocessing may be required to ensure that the resulting chromosomes are feasible. If a zero row in the rules' structure matrix is created, then that row is deleted and the corresponding row in the parameter matrix is also deleted. The following conditions must be satisfied by the parameter matrix for the chromosome containing this matrix to be feasible:

$$\begin{aligned} d_{ij} &\geq 0, \\ l_{ij} &\leq c_{ij} - d_{ij}, \\ r_{ij} &\geq c_{ij} + d_{ij}. \end{aligned}$$

Whenever any one of the above conditions is violated by a parameter, then that parameter is

assigned a boundary value. For example, if $d_{ij} < 0$, then d_{ij} is set to 0, and if $l_{ij} > c_{ij} - d_{ij}$, then l_{ij} is set to $c_{ij} - d_{ij}$, etc.

The above-described EA is used to construct a fuzzy logic controller for the front-wheel-steering ground vehicle model, whose dynamics are characterized by a certain type of symmetry that allows us to simplify computations. Specifically, suppose that in the rule base we have a fuzzy rule of the form

IF x_1 is F_1 and x_2 is F_2 , THEN u is θ ,

where the fuzzy sets F_1 and F_2 are specified by the membership functions $\mu_{F_1}(c_1, d_1, r_1, l_1)$ and $\mu_{F_2}(c_2, d_2, r_2, l_2)$, respectively. The parameters c_i, d_i, r_i , and $l_i, i = 1, 2$, characterize membership functions of the form shown in Figure 10.14. Because of the symmetry in the vehicle dynamics, the rule base also contains the rule

IF x_1 is \tilde{F}_1 and x_2 is \tilde{F}_2 , THEN u is $-\theta$,

where \tilde{F}_1 and \tilde{F}_2 have membership functions

$$\mu_{\tilde{F}_1} = \mu_{F_1}(-c_1, d_1, -l_1, -r_1) \quad \text{and} \quad \mu_{\tilde{F}_2} = \mu_{F_2}(-c_2, d_2, -l_2, -r_2),$$

respectively.

The size of the population in our simulations was 30, the maximum number of rules was $r_{\max} = 10$, and the probability p_g was equal to 0.1. A fuzzy rule base obtained after 100 iterations is shown in Figure 10.16. Only four rules were generated. The fuzzy rule base shown

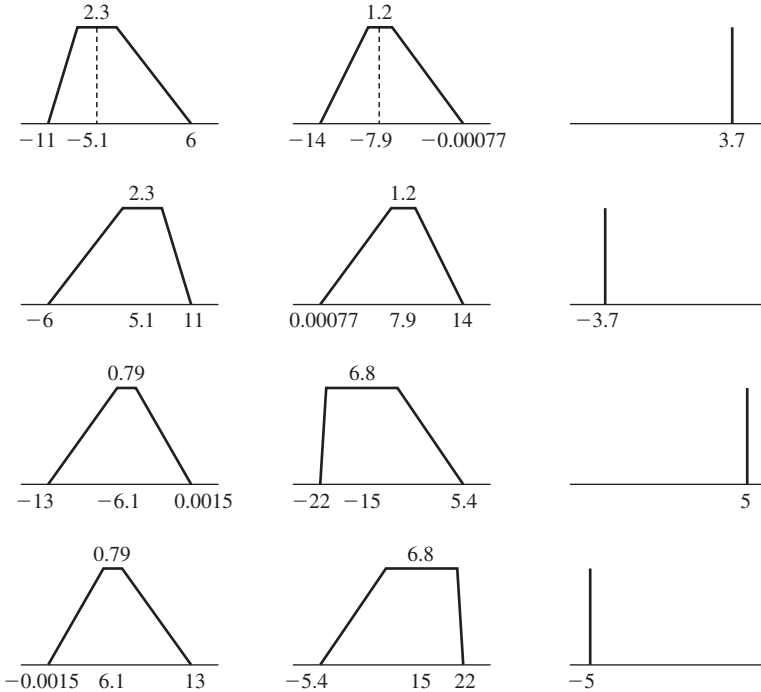


Figure 10.16 The best fuzzy rule base after 100 iterations of the EA operating on fuzzy membership functions and fuzzy rule base.

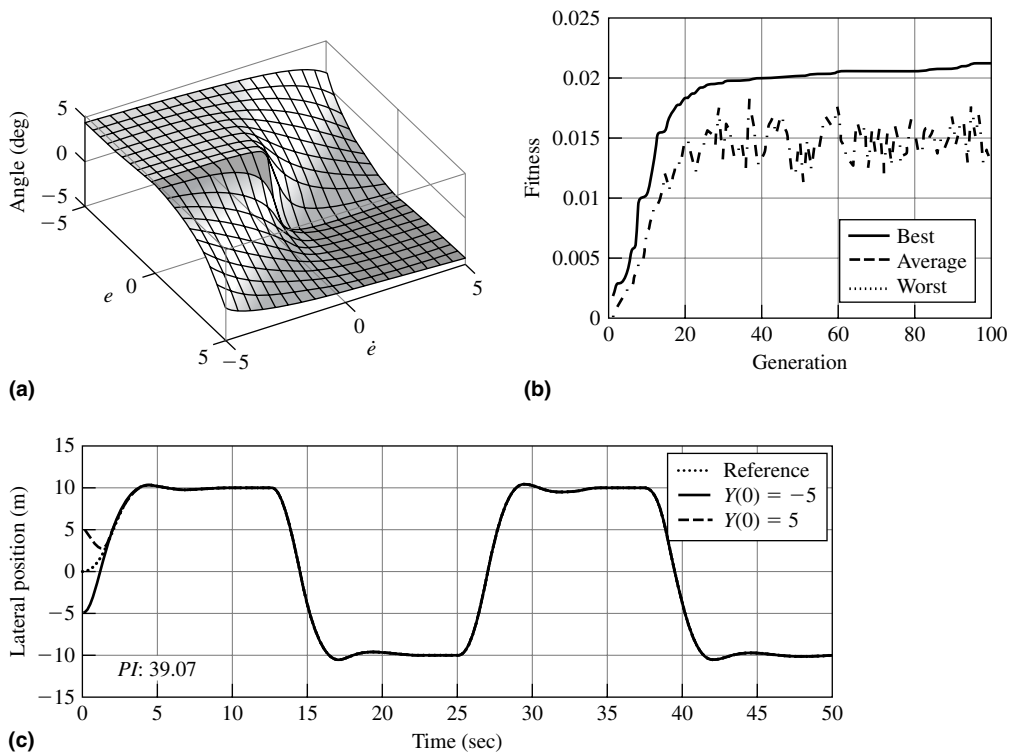


Figure 10.17 Summary of the simulation run of the closed-loop system driven by the fuzzy logic controller obtained from the fuzzy rule base shown in Figure 10.16. (a) Control surface. (b) Fitness values of the best, average, and worst chromosomes. (c) Tracking performance.

in Figure 10.16 was then used to construct a fuzzy logic controller that was connected to the vehicle model in the closed-loop configuration shown in Figure 10.7. The fuzzy control surface constructed using the obtained fuzzy rule base is shown in Figure 10.17(a). The results of a typical simulation run, along with the fitness values of the best, average, and worst chromosomes in each generation, are shown in Figure 10.17(b), while the tracking performance is illustrated in Figure 10.17(c). As in the previous cases, the simulations were performed for two different initial conditions, and excellent tracking was achieved in both cases. The value of the *PI* was reduced to 38.06, which is lower than in the two previous designs.

Notes

Basic principles of genetics have their roots in the experiments performed by a nineteenth-century monk, Gregor Mendel, in his monastery garden in Brno, now in the Czech Republic. Mendel worked with the garden pea (*Pisum sativum*) because it possesses traits that are easily recognizable and appear in widely contrasting forms. In addition, the garden pea lends itself well to controlled crosses. Mendel published his findings in 1866 in a paper entitled, “Experiments in Plant Hybridization.” According to Elseth and Baumgardner [74, p. 3], “The publication of

this paper marked the birth of genetics as a science. The principles stated in the paper still stand today as the cornerstone of modern genetics.” The importance of Mendel’s work went completely unrecognized for 34 years. “In 1900, three botanists, Hugo de Vries (Dutch), Carl Correns (German), and Erich von Tschermak-Seysenegg (Austrian), simultaneously discovered Mendel’s paper on heredity, recognized its importance, and cited it in their own publications” [74, p. 7]. We add that the term *gene* was coined by W. L. Johannsen in 1903.

A seminal text on genetic algorithms by Holland [124] is definitely worth perusing. However, Holland’s monograph is not intended as an introduction to genetic algorithms. A very well written introductory text to genetic algorithms is Mitchell’s book [206]. A number of engineering design examples using genetic algorithms can be found in Gen and Cheng [96]. Bäck, Hammel, and Schwefel [15] survey the history as well as the current state of evolutionary computation, which is the term that encompasses evolution strategies, evolutionary programming, and genetic algorithms. Applications of genetic algorithms to automatic programming using evolving computer programs are described by Koza [166]. Automatic programming with genetic algorithms is called by Koza “genetic programming” (GP). Robust controller design with genetic algorithms are reported by Patton and Liu [224] and by Ge, Lee, and Zhu [94]. Kozek, Roska, and Chua [168], Pal and Bhandari [219], and Zamparelli [315] use genetic algorithms to train cellular neural networks. Kim and Myung [159] apply genetic algorithms to solve constrained optimization problems, while Koza et al. [167] employ genetic algorithms for the automatic design of electrical circuits. Yang, Hachino, and Tsuji [306] combine least squares techniques with a genetic algorithm to on-line identification of time-delay systems.

EXERCISES

- 10.1** Recall that the “elitism” is a genetic operator that forces the GA, or EA, to retain some number of the best individuals from the current population in the new population. Incorporate the elitism operator into your GA program and test it on the function (10.8). You may try, for example, to keep the best two individuals from the current population in the new population.
- 10.2** Implement the simple EA from Section 10.4 and test it on
- (a) the function (10.8);
 - (b) the function $f(x) = x + |\sin(32x)|$, where $0 \leq x < \pi$.
- The function in (b) comes from Mitchell [206, p. 9].
- 10.3** Implement the simple EA from Section 10.4 that maximizes the MATLAB’s peaks function

$$f(x, y) = 3(1 - x)^2 \exp(-x^2 - (y + 1)^2) - 10\left(\frac{x}{5} - x^3 - y^5\right) \exp(-x^2 - y^2) - \frac{1}{3} \exp(-(x + 1)^2 - y^2),$$

where $-3 \leq x, y \leq 3$.

10.4 Implement the simple EA from Section 10.4 that maximizes the function

$$f(x, y) = 3 + y \sin(0.7x) + x \cos(1.5y),$$

where $-3.0 \leq x \leq 12.1$ and $4.1 \leq y \leq 5.8$.

10.5 Consider the classic traveling salesman problem (TSP): Suppose a salesman must visit clients in different cities and then return home. What is the shortest tour through those cities, visiting each one once and only once? Implement in MATLAB a simple variant of the genetic algorithm that optimizes the salesman route for 10 cities. How many different possible paths are there for 10 cities? Following Fogel [88, p. 29], use the following representation of any candidate solution: For a given map of the cities to be visited (including the salesman’s home base), assign to each city a positive integer. Thus, for example, if there were six cities, one possible solution might be [1 2 3 4 5 6]. The above represents an order of progression of the salesman’s trip. Note that because the problem requires a round trip, the first coordinate of the vector representing the first city of the trip is also the last city visited by the salesman. Because the trip is a closed loop, it does not matter from which city the salesman starts his trip. As far as a crossover operator implementation is concerned, we recommend using Fogel’s idea [88, p. 29], where a single-parent operator produces an offspring by inverting the visiting order between two randomly chosen coordinates of the parent vector. As an example of this implementation of the single-parent crossover operator, suppose that the parent is as above and that the randomly chosen coordinates—that is, the inversion points—are the second and fifth ones. The offspring would then be [1 5 3 4 2 6].

Plot a map with the cities and mark the obtained salesman’s route. The coordinates of the cities are given in Table 10.2.

Table 10.2 Locations of the Cities in Exercise 10.5

<i>x</i>	0.4306	3.7094	6.9330	9.3582	4.7758	1.2910	4.83831	9.4560	3.6774	3.2849
<i>y</i>	7.7288	2.9727	1.7785	6.9080	2.6394	4.5774	8.43692	8.8150	7.0002	7.5569

10.6 Consider a closed-loop system with a fuzzy logic controller (FLC) in the loop shown in Figure 10.18. Suppose that we use five membership functions for the FLC inputs and its output. A possible set of rules is given in Table 10.3, where LN denotes large negative, N means negative, ZE stands for zero, P means positive, and LP refers to large positive. Implement a variant of the simple EA from Section 10.4 that generates the set of rules to be used to synthesize the FLC. Assign to each fuzzy set a positive integer; for example, let LN be encoded as 1, N as 2, ZE as 3, P as 4, and LP as 5. Then,

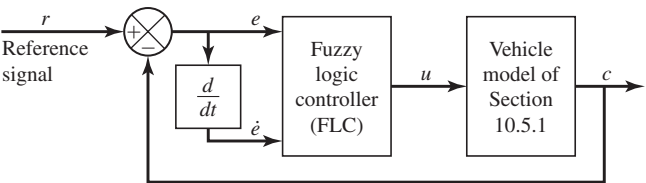


Figure 10.18 Closed-loop system with a fuzzy logic controller in the loop.

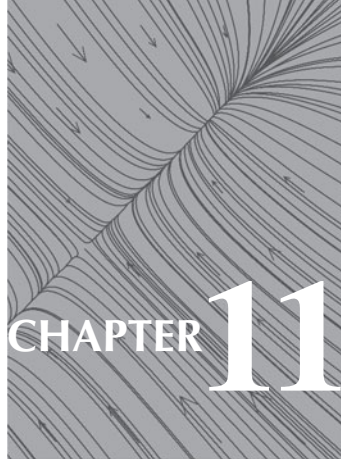
Table 10.3 Rules Matrix to Be Used to Construct a Fuzzy Logic Controller in Exercise 10.6

$e \backslash \dot{e}$	LN	N	ZE	P	LP
LN	LN	LN	LN	N	ZE
N	LN	LN	N	ZE	P
ZE	N	N	ZE	P	P
P	N	ZE	P	LP	LP
LP	ZE	P	LP	LP	LP

a candidate solution, which is shown in Table 10.3, can be represented in a vector form as

[1 1 2 2 3 1 1 2 3 4 1 2 3 4 5 2 3 4 5 5 3 4 4 5 5];

that is, a vector representation of the candidate solution given by the rules matrix is obtained using the stacking operation of the columns of the rules matrix. Use the one-point crossover operator similar to the one in the canonical GA. The mutation operation can be implemented by altering a randomly chosen coordinate in the given chromosome—that is, the vector representing a candidate solution. You can use the roulette-wheel method to determine the replacement. For example, if the seventh component was selected to be altered, then you can use the roulette-wheel method to determine an integer from the set {1, 2, 3, 4, 5} that will replace the “1” that occupies the seventh coordinate of the above chromosome. Chan, Lee, and Leung [47] used a similar version of the EA to generate fuzzy rules to synthesize a fuzzy logic target tracking algorithm. You can use the fitness function given by (10.16).



Chaotic Systems and Fractals

An apparent paradox is that chaos is deterministic, generated by fixed rules which do not themselves involve any elements of change. We even speak of deterministic chaos. In principle, the future is completely determined by the past; but in practice small uncertainties, much like minute errors of measurement which enter into calculations, are amplified, with the effect that even though the behavior is predictable in the short term, it is unpredictable over the long run.

—*Chaos and Fractals: New Frontiers of Science* [227, p. 11]

11.1 Chaotic Systems Are Dynamical Systems with Wild Behavior

The objective of this chapter is to show the reader “the wild things that simple nonlinear equations can do” [198, p. 459]. These simple nonlinear models with very complicated dynamics are examples of *chaotic systems*, which are deterministic nonlinear dynamical systems that exhibit a random-like behavior. A discrete-time model of population dynamics discussed in the next section, Newton’s method for finding roots of certain polynomials, or a jetting balloon after its release are examples of dynamical chaotic systems. Trajectories of chaotic dynamical systems are sensitive to initial conditions in the sense that starting from slightly different initial conditions the trajectories diverge exponentially. Another property of a chaotic dynamical system is *loss of information* about initial conditions. We follow Moon [209, p. 4] to explain this property. Suppose that we are given a state-plane portrait of a chaotic dynamical system in the (x, y) -plane. Assume that we specify the x coordinate with accuracy Δx and the y coordinate with accuracy Δy . Then, we could form a grid in the (x, y) -plane, shown in Figure 11.1, where each area would be of size $\Delta x \Delta y$. Thus, if we select initial conditions to the stated accuracy, then we know that these initial conditions are located somewhere in the shaded box in the (x, y) -plane corresponding to the initial time t_0 . For a chaotic system, this uncertainty grows in time to, say, N boxes at a later time t_1 , as illustrated in Figure 11.1. The uncertainty growth, as mentioned above, is exponential and for continuous-time systems can be described by

$$N(t) \approx e^{\mu(t-t_0)} N(t_0), \quad (11.1)$$

where the constant μ is a measure of sensitivity with respect to initial conditions. Lyapunov exponents, which we discuss in a subsequent section, can be used to measure the initial condition

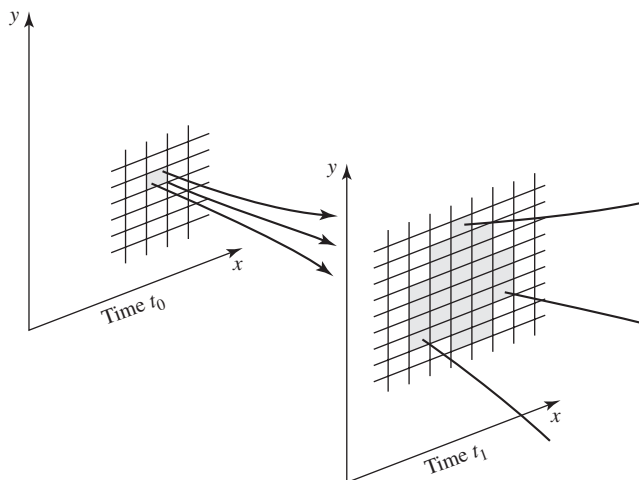


Figure 11.1 An illustration of the growth of uncertainty in a chaotic dynamical system.

sensitivity. Chaotic behavior can arise in low-order nonlinear dynamical systems. In the case of continuous-time systems, chaotic behavior can only occur for systems of dimension three or higher. However, in discrete-time systems, chaos can occur even in one-dimensional systems, as we will see in the following section.

11.2 Chaotic Behavior of the Logistic Equation

11.2.1 The Logistic Equation—An Example from Ecology

The so-called logistic equation, or logistic map, whose behavior we analyze in this section, arises in the analysis of the growth of a population, for example. Before proceeding, we define a few terms we will use in our discussion.

Definition 11.1 A map is a function whose domain (input) space and range (output) space are the same. Let $\mathbf{x}(0)$ be a point in \mathbb{R}^n and f be a map on \mathbb{R}^n . The orbit, or trajectory, of $\mathbf{x}(0)$ under f is the set of points

$$\{\mathbf{x}(0), f(\mathbf{x}(0)), f^2(\mathbf{x}(0)), \dots, f^k(\mathbf{x}(0)), \dots\},$$

where f^k is the k th iterate of f . The starting point, $\mathbf{x}(0)$, of the orbit is called the initial value of the orbit.

In what follows, we derive a discrete-time version of the logistic equation. The continuous-time logistic equation was analyzed in Subsection 1.7.5. We first analyze a situation where there are no predators and there is an unlimited supply of food. For the sake of argument, we consider a population of rabbits. Assume that, on the average, each rabbit gives birth to σ new rabbits in one unit of time and that no rabbit dies. Let $N(k)$ be the number of rabbits at time $t = k$; then the change in the population over a time period between times k and $k + 1$ is

$$N(k + 1) - N(k) = \sigma N(k). \quad (11.2)$$

Thus,

$$N(k+1) = (1 + \sigma)N(k). \quad (11.3)$$

Equation (11.3) is a very simple linear model of population growth. Given the number $N(0)$ of rabbits at time $k = 0$, the number of rabbits at time k is

$$N(k) = (1 + \sigma)^k N(0). \quad (11.4)$$

One can see from (11.4) that the linear model given by (11.3) predicts the exponential growth of population. We now try to make the above model more realistic. We assume, as before, that the number of births is proportional to the size of the population. Let b be the birth rate; then the number of births at time k is $bN(k)$. We also assume that the number of deaths is proportional to the population size. Let d be the death rate; then the number of deaths at time k is $dN(k)$. The change in population in a time period between times k and $k+1$ is

$$N(k+1) - N(k) = bN(k) - dN(k). \quad (11.5)$$

We rewrite equation (11.5) as

$$N(k+1) = (1 + b - d)N(k). \quad (11.6)$$

Let $\sigma = b - d$. Then, equation (11.6) has the same form as equation (11.3). In both cases, the population grows exponentially. We refer to the constant σ as the *unrestricted growth rate*. Over short periods of time, the above equation may be sufficient for purposes of an analysis of a population. Over long periods of time, the growth rates may change as the size of a population changes.

We now alter the above linear models by introducing the concept of saturation. We first assume that the environment of the population can only support a certain number of species. We denote this number by K , and we call it the *carrying capacity*. Thus, if $N(k) > K$, then there will not be enough food or space available, and more species will die than will be born. In such a situation the growth rate should be negative. If the population at time k is $N(k) < K$, then the growth rate should be positive. Finally, if $N(k)$ is much smaller than the carrying capacity, then the growth rate should be close to the unrestricted growth rate σ . A possible function that satisfies the above requirements has the form

$$\sigma \left(1 - \frac{N(k)}{K} \right). \quad (11.7)$$

The change in the size of population then can be described as

$$N(k+1) - N(k) = \sigma \left(1 - \frac{N(k)}{K} \right) N(k), \quad (11.8)$$

or, equivalently

$$N(k+1) = (1 + \sigma)N(k) - \frac{\sigma}{K}N^2(k). \quad (11.9)$$

Equation (11.9) is called the *logistic equation*. The term $\frac{\sigma}{K}N^2(k)$ is referred to as the *damping term*, because it dampens the growth of the population. To simplify (11.9), we scale the variable $N(k)$ by introducing the new variable

$$x(k) = \frac{\sigma}{\rho K}N(k),$$

where

$$\rho = 1 + \sigma.$$

Then, equation (11.9) becomes

$$x(k+1) = \rho x(k)(1 - x(k)) \quad (11.10)$$

We will now show that if $0 < \rho \leq 4$ and $0 \leq x(0) \leq 1$, then $0 \leq x(k) \leq 1$ for $k = 1, 2, \dots$. Let

$$f(x) = \rho x(1 - x).$$

Then,

$$\frac{df(x)}{dx} = \rho - 2\rho x = 0$$

for $x = 1/2$ and $f_{\max}(x) = \rho/4$. Hence, for $k = 1, 2, \dots$,

$$0 \leq x(k) \leq 1 \quad \text{if} \quad 0 < \rho \leq 4 \quad \text{and} \quad 0 \leq x(0) \leq 1.$$

11.2.2 Stability Analysis of the Logistic Map

We begin our investigation of the logistic equation (11.10) by seeking its equilibrium states (fixed points). These can be found by setting $x(k+1)$ equal to $x(k)$ in (11.10) to obtain

$$x(k)(\rho x(k) + (1 - \rho)) = 0.$$

It follows that the equilibrium solutions are

$$x^* = 0 \quad \text{and} \quad x^* = \frac{\rho - 1}{\rho}.$$

We will now study attracting properties of the fixed points as a function of the parameter ρ . Roughly speaking, a fixed point $x^* = f(x^*)$ of the system $x(k+1) = f(x(k))$ is attracting if, for a starting point close to x^* , the next iterate is still closer. Formally, a fixed point x^* is attracting if there is an $\varepsilon > 0$ such that when $|x(0) - x^*| < \varepsilon$, then $\lim_{k \rightarrow \infty} x(k) = x^*$.

A fixed point x^* is repelling if for any $\varepsilon > 0$ no matter how small, if $0 < |x(0) - x^*| < \varepsilon$, then $|x(k) - x^*| > \varepsilon$ for some k .

Theorem 11.1 Suppose $x^* = f(x^*)$ is a fixed point of $x(k+1) = f(x(k))$, where f is continuously differentiable. Then, x^* is attracting if

$$|f'(x^*)| < 1$$

and is repelling if

$$|f'(x^*)| > 1.$$

The case when $|f'(x^*)| = 1$ is inconclusive.

Proof We first show that if $x^* = f(x^*)$ is such that $|f'(x^*)| < 1$, then x^* is attracting. Because by assumption $|f'(x^*)| < 1$ and f' is continuous, there exists an interval $(x^* - \varepsilon, x^* + \varepsilon)$ about x^* such that for some $L > 0$,

$$|f'(x^*)| < L < 1 \quad \text{and} \quad |f'(x)| < L \quad \text{for every } x \in (x^* - \varepsilon, x^* + \varepsilon).$$

Let $x(k) \in (x^* - \varepsilon, x^* + \varepsilon)$. We will show that $x(k+1)$ is closer to x^* than $x(k)$. We have

$$|x(k+1) - x^*| = |f(x(k)) - f(x^*)|$$

because $x^* = f(x^*)$ is a fixed point. Using the mean value theorem yields

$$|x(k+1) - x^*| = |f(x(k)) - f(x^*)| = |f'(\tilde{x})||x(k) - x^*|$$

for some $\tilde{x} \in (x(k), x^*)$ or $\tilde{x} \in (x^*, x(k))$ if $x^* < x(k)$. Since $x(k) \in (x^* - \varepsilon, x^* + \varepsilon)$, then $\tilde{x} \in (x^* - \varepsilon, x^* + \varepsilon)$ also. Thus, $|f'(\tilde{x})| < L < 1$, and therefore

$$|x(k+1) - x^*| < L|x(k) - x^*|.$$

From the above, we conclude that if $x(0) \in (x^* - \varepsilon, x^* + \varepsilon)$, then

$$|x(k) - x^*| < L^k|x(0) - x^*|.$$

Because $L < 1$, we have $\lim_{k \rightarrow \infty} |x(k) - x^*| = 0$, that is,

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

which means that x^* is an attracting fixed point.

Similarly, if $x^* = f(x^*)$ is such that $|f'(x^*)| > 1$, then it is a repelling fixed point. Indeed, if $|f'(x^*)| > 1$, then because f' is continuous, there is an interval $(x^* - \varepsilon, x^* + \varepsilon)$ about x^* and a constant $L > 1$ such that $|f'(x)| > L$ for all $x \in (x^* - \varepsilon, x^* + \varepsilon)$. We will now show that if $x(k) \in (x^* - \varepsilon, x^* + \varepsilon)$, then $x(k+1)$ will be further away from x^* than $x(k)$. Using the mean value theorem, we obtain

$$\begin{aligned} |x(k+1) - x^*| &= |f(x(k)) - f(x^*)| \\ &= |f'(\tilde{x})||x(k) - x^*| \\ &> L|x(k) - x^*|. \end{aligned} \tag{11.11}$$

If $x(k+1)$ is no longer in $(x^* - \varepsilon, x^* + \varepsilon)$, we are done. Otherwise, we repeat the process to get

$$\begin{aligned} |x(k+2) - x^*| &= |f(x(k+1)) - f(x^*)| \\ &> L|x(k+1) - x^*| \\ &> L^2|x(k) - x^*|. \end{aligned}$$

If $x(k+2)$ is no longer in the interval $(x^* - \varepsilon, x^* + \varepsilon)$, we are done. If not, we repeat the process to obtain

$$|x(k+m) - x^*| > L^m|x(k) - x^*|.$$

Since $(x^* - \varepsilon, x^* + \varepsilon)$ is a finite interval and $L^m \rightarrow \infty$ as $m \rightarrow \infty$, at some time, $x(k+m)$ must leave the interval, which means that x^* is a repelling fixed point. The proof of the theorem is complete.

Equation (11.11) tells us that the trajectory $\{x(k)\}$ diverges exponentially from x^* as $k \rightarrow \infty$. Thus, a trajectory with an initial condition close to x^* for the case when $|f'(x^*)| = L > 1$ will separate from x^* at a multiplicative rate of approximately $L > 1$ per iteration until it moves significantly away from x^* . In other words, the distance $|x(k) - x^*|$ for $x(0)$ close to x^* will grow exponentially; that is, this distance will be magnified approximately by L for each iteration. To quantify this multiplicative rate of separation of points close to x^* , the Lyapunov exponent term is introduced that we will discuss in Section 11.4.

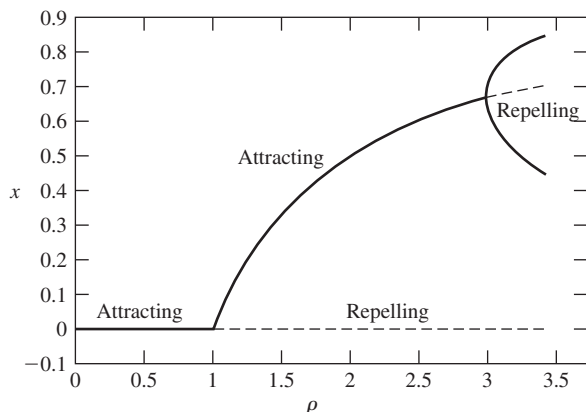


Figure 11.2 Exchange of attracting properties of the fixed points of the logistic equation for $0 < \rho \leq 3.44$.

We now consider the logistic equation (11.10). Recall that $f = \rho x(1 - x)$ denotes the right-hand side of the logistic equation. The derivative of f with respect to x is

$$f'(x) = \rho(1 - 2x).$$

Therefore, $f'(0) = \rho$. Since $|f'(0)| < 1$ when $\rho < 1$, by Theorem 11.1, the fixed point $x^* = 0$ is attracting for $\rho < 1$ and repelling for $\rho > 1$. The fixed point $x^* = (\rho - 1)/\rho$ is outside of the range of x for $\rho < 1$. Thus, if $0 < \rho < 1$, we have just one attractive fixed point. For $\rho > 1$, the fixed point $x^* = 0$ is repelling and $x^* = (\rho - 1)/\rho$ comes in within the range of the variable x . The derivative of f evaluated at $x^* = (\rho - 1)/\rho$ is

$$f'\left(\frac{\rho - 1}{\rho}\right) = 2 - \rho.$$

By Theorem 11.1, the fixed point $x^* = (\rho - 1)/\rho$ is attracting if $|2 - \rho| < 1$ —that is, if

$$1 < \rho < 3.$$

Thus, for $\rho = 1$ there is an exchange of attracting properties from one fixed point to the other. This change in the qualitative behavior at a certain parameter value is called a *bifurcation*. In Figure 11.2, the fixed points $x^* = 0$ and $x^* = (\rho - 1)/\rho$ are shown and the ranges of ρ for which they are attracting. For $\rho > 3$, there are no attracting fixed points, and the solutions of the logistic equation become very complex. We discuss the logistic equation dynamics for $\rho > 3$ after displaying its solutions for $\rho = 0.7$, $\rho = 1.4$, $\rho = 2.8$, and $\rho = 3.6$ in Figure 11.3. The initial condition in all four cases was $x(0) = 0.9$.

Another way of displaying a solution of a first-order difference equation is a *cobweb*. To construct a cobweb for the given first-order discrete-time equation, $x(k + 1) = f(x(k))$, we first draw the function $y = f(x)$ representing the right-hand side of the difference equation. Then, we draw the straight line $y = x$. The fixed points correspond to the points of intersection of these two graphs. The solution sequence, represented by the cobweb, starting at $x(0)$ is obtained as follows. Starting from $x(0)$, draw upward vertical line to intersect the graph of $f(x)$. This corresponds to the calculation of $x(1)$. The value of $x(1)$ is transferred to the straight line $y = x$. The x coordinate of the intersection is $x(1)$, which is then used to generate $x(2)$ using the above procedure. The process is repeated over and over again. In Figure 11.4 the cobwebs corresponding to the solutions displayed in Figure 11.3 are presented.

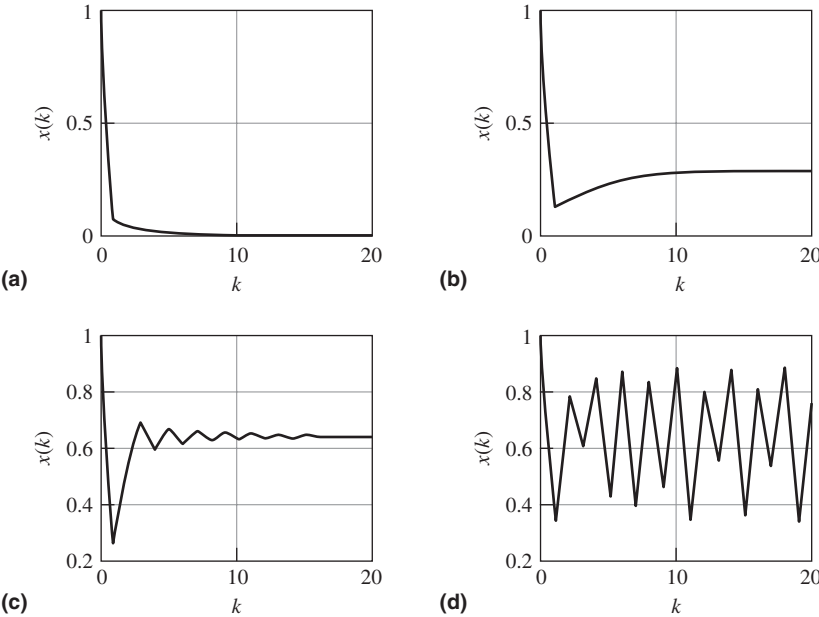


Figure 11.3 Solutions of the logistic equation for different values of the parameter ρ : (a) $\rho = 0.7$, (b) $\rho = 1.4$, (c) $\rho = 2.8$, and (d) $\rho = 3.6$.

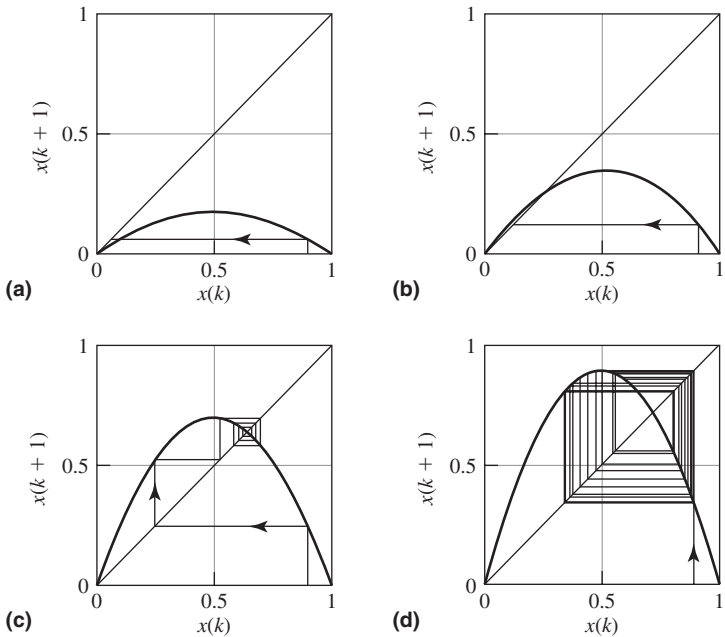


Figure 11.4 Cobwebs for the logistic equation: (a) $\rho = 0.7$, (b) $\rho = 1.4$, (c) $\rho = 2.8$, and (d) $\rho = 3.6$.

As we mentioned before, for $\rho > 3$ there are no attracting fixed points. As ρ increases, the solutions of the logistic equation exhibit increasing complexity. For ρ just slightly larger than 3, the solution sequence settles down into a steady oscillation of period 2 as shown in Figure 11.5. Then, as ρ is further increased, periodic solutions of period 4, 8, 16, \dots appear. In Figure 11.6 a solution of period 4 is shown corresponding to $\rho = 3.5$. For $\rho = 3.9$ the solution is shown in Figure 11.7.

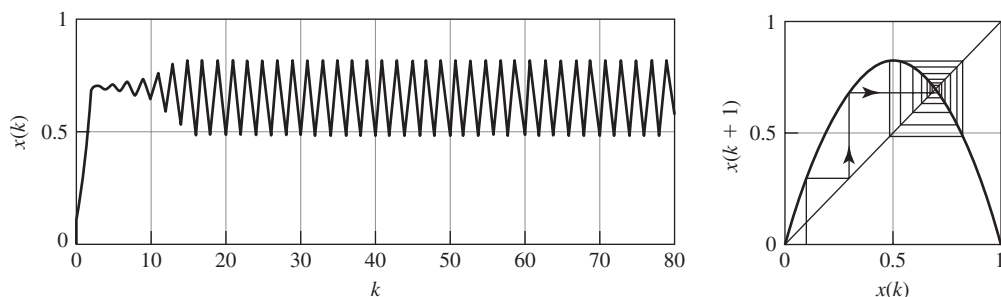


Figure 11.5 A two-cycle solution and its cobweb for $\rho = 3.3$ and $x(0) = 0.1$.

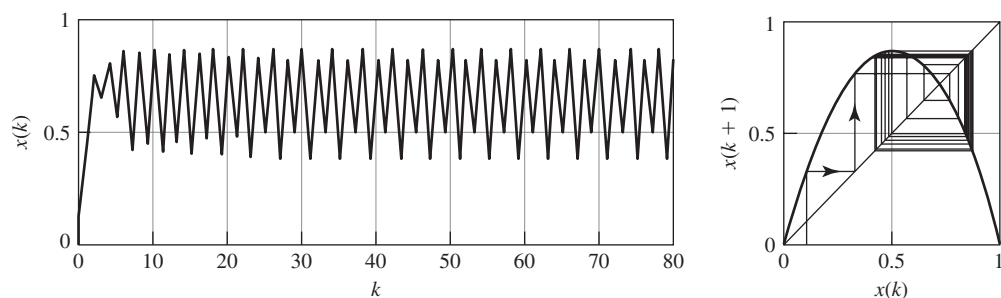


Figure 11.6 A period-four solution and its cobweb for $\rho = 3.5$ and $x(0) = 0.1$.

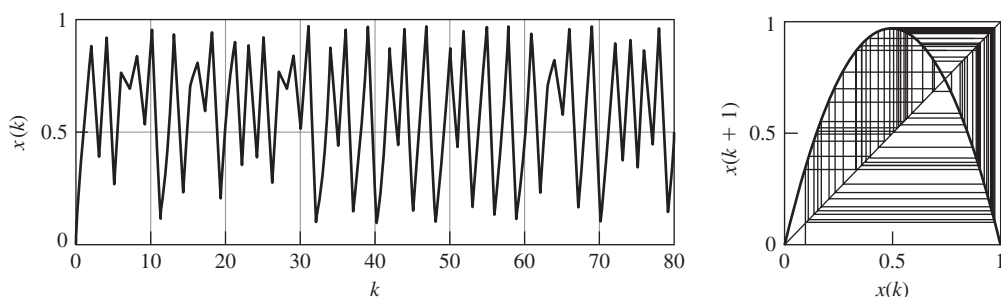


Figure 11.7 Aperiodic solution and its cobweb for $\rho = 3.9$ and $x(0) = 0.1$.

We now analyze the solutions for $\rho > 3$. To simplify the analysis, we first make the substitution

$$x(k) = z(k) + \frac{\rho - 1}{\rho}, \quad (11.12)$$

which is equivalent to shifting the fixed point $x^* = (\rho - 1)/\rho$ to the origin in the (ρ, z) coordinates. Note that

$$x(k + 1) = z(k + 1) + \frac{\rho - 1}{\rho}. \quad (11.13)$$

Substituting (11.12) and (11.13) into the logistic equation (11.10) yields

$$z(k + 1) + \frac{\rho - 1}{\rho} = \rho \left(z(k) + \frac{\rho - 1}{\rho} \right) \left(1 - z(k) - \frac{\rho - 1}{\rho} \right). \quad (11.14)$$

After performing some manipulations, we represent (11.14) as

$$\boxed{z(k + 1) = (2 - \rho)z(k) - \rho z^2(k)} \quad (11.15)$$

The fixed points of (11.15) are $z^* = 0$ corresponding to $x^* = (\rho - 1)/\rho$, and $z^* = (1 - \rho)/\rho$ corresponding to $x^* = 0$. We investigate the attracting properties of the fixed point $z^* = 0$ for $\rho > 3$. We see from Figure 11.5 that for ρ just greater than 3, the solution of the logistic equation goes to an attracting period-two cycle. Two numbers p and q form a period-2 cycle for a first-order discrete system if, when $x(k) = p$, then $x(k + 1) = q$, $x(k + 2) = p$, and so on. A period-2 cycle is attracting if there are intervals (p_1, p_2) and (q_1, q_2) about p and q , respectively, such that:

1. If $x(0) \in (p_1, p_2)$, then $\lim_{k \rightarrow \infty} x(2k) = p$ and $\lim_{k \rightarrow \infty} x(2k + 1) = q$.
2. If $x(0) \in (q_1, q_2)$, then $\lim_{k \rightarrow \infty} x(2k) = q$ and $\lim_{k \rightarrow \infty} x(2k + 1) = p$.

The fixed points p and q of the period-2 cycle of equation (11.15) can be found by solving the equation

$$z = g(g(z)), \quad (11.16)$$

where $g = (2 - \rho)z - \rho z^2$. We can use MATLAB's function `compose` to form (11.16), and then use the function `solve` to solve (11.16). We obtain four roots. Two of them are the fixed points that we analyzed above. The remaining two roots are the fixed points of the period-2 cycle:

$$p = -\frac{1}{2\rho}(\rho - 3 + \sqrt{\rho^2 - 2\rho - 3}) \quad (11.17)$$

and

$$q = -\frac{1}{2\rho}(\rho - 3 - \sqrt{\rho^2 - 2\rho - 3}). \quad (11.18)$$

Once we have found the 2-cycle, we next determine the range of ρ for which it is attracting. A 2-cycle p and q of the system $z(k + 1) = g(z(k))$ is attracting if the fixed points p and q of the map $y(k + 1) = h(y(k))$ are attracting, where $h(y) = g(g(y))$. By Theorem 11.1, the fixed points $p = h(p)$ and $q = h(q)$ are attracting if $|h'(p)| < 1$ and $|h'(q)| < 1$. Using the chain rule gives

$$h'(y) = \frac{d}{dy}g(g(y)) = g'(g(y))g'(y). \quad (11.19)$$

Because $g(p) = q$ and $g(q) = p$, we obtain

$$h'(p) = g'(g(p))g'(p) = g'(q)g'(p). \quad (11.20)$$

Similarly,

$$h'(q) = g'(p)g'(q). \quad (11.21)$$

Therefore, $|h'(p)| < 1$ and $|h'(q)| < 1$ when

$$|g'(p)g'(q)| < 1. \quad (11.22)$$

Using Theorem 11.1 and the arguments similar to the ones above, we conclude that the 2-cycle represented by p and q is repelling if

$$|g'(p)g'(q)| > 1. \quad (11.23)$$

We next substitute (11.17) and (11.18) into $g'(p)g'(q)$ and evaluate the resulting expression using MATLAB's function `subs`, to obtain

$$g'(p)g'(q) = (2 - \rho - 2\rho p)(2 - \rho - 2\rho q) = 4 - \rho^2 + 2\rho. \quad (11.24)$$

Thus, the 2-cycle represented by the pair (p, q) is attracting if

$$-1 < 4 - \rho^2 + 2\rho < 1, \quad (11.25)$$

where $\rho > 3$. Solving (11.25) yields

$$3 < \rho < 1 + \sqrt{6}. \quad (11.26)$$

When $\rho > 1 + \sqrt{6}$, an attracting 4-cycle appears. The bifurcation diagram in Figure 11.8 shows that as the value of ρ increases, the period doubling process takes place. Each of the branches of the bifurcation diagram produces a period-2 cycle of its own. The process continues as the parameter ρ approaches the value $\rho = 3.5699$, until no periodic cycles are created. For $\rho = 3.5699$ the behavior is not cyclic anymore. The behavior of the solutions of the logistic equation for $\rho \geq 3.5699$ is complex indeed, as can be seen, for example, by inspecting the bifurcation diagram shown in Figure 11.8. Various windows of ranges of $\rho > 3.5699$ correspond

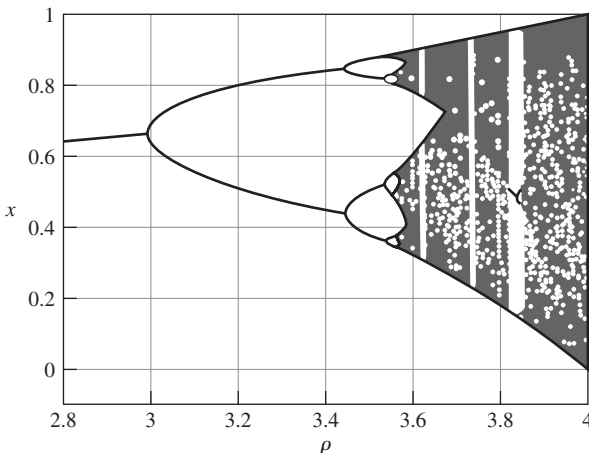


Figure 11.8 Bifurcation diagram of the logistic equation.

to different periodicities of solutions as well as to aperiodic solutions. These aperiodic solutions are called *chaotic solutions* or *chaotic motions*. Here, we quote a descriptive definition of chaotic motion due to Thompson and Stewart [281, p. 197]: “... one way to define chaos is in the negative—recurrent behavior that is not an equilibrium, a cycle or even a quasi-periodic motion—there is more to be said about chaos. Chaotic motion has some aspect that is provably as random as a coin toss. The randomness arises from sensitive dependence on imperfectly known initial conditions, resulting for example in broadband noise in the power spectra of observable time series. This seems remarkable because the dynamical system needs no stochastic input to achieve this.” A precise definition of chaotic motion is given in Section 11.4.

11.2.3 Period Doubling to Chaos

We will now describe one of the ways in which chaotic motion can arise—by repeated flip bifurcations, that is, through period doubling. Such a behavior arises for $3 \leq \rho \leq 3.5699$. We will show that there is a sequence of numbers

$$c_1 < c_2 < \cdots < c_n < \cdots \quad (11.27)$$

such that when

$$c_n < \rho < c_{n+1},$$

then the logistic equation has an attracting 2^n -cycle. Thus, for example, if $c_8 < \rho < c_9$, then the solution will settle into an attracting 2^8 -cycle. We already have $c_1 = 3$ and $c_2 = 1 + \sqrt{6} = 3.4495$. We are interested in finding the remaining elements of the sequence (11.27). The sequence (11.27) is convergent, where

$$\lim_{n \rightarrow \infty} c_n = c_\infty = 3.5699.$$

To find the elements of the sequence (11.27), we could, in principle, form composite maps $f^{2^n}(x)$ and use stability results from the previous subsection. We have observed that if the parameter ρ is slightly greater than $1 + \sqrt{6}$, then the 2-cycle is repelling and the solution goes to an attracting 4-cycle. A cobweb of an attracting 4-cycle is shown in Figure 11.9. We could find the four points

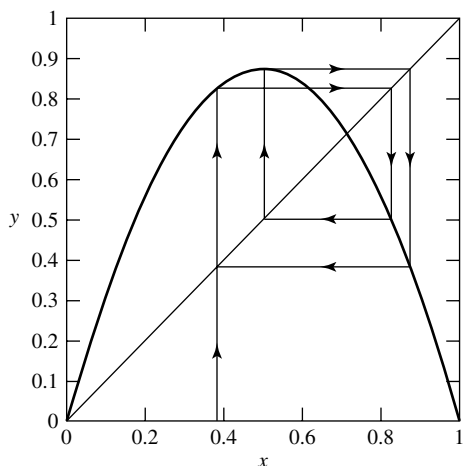


Figure 11.9 Cobweb of an attracting 4-cycle for $\rho = 3.5$.

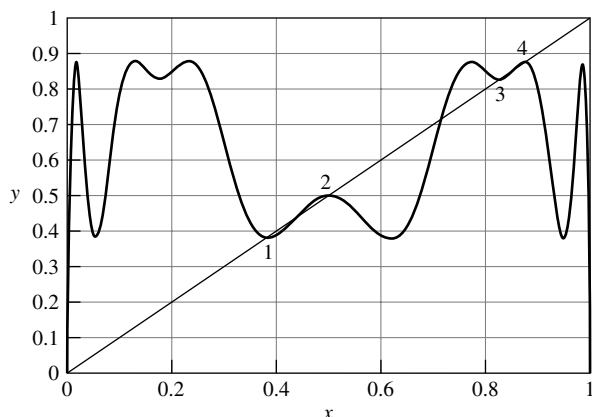


Figure 11.10 A graphical solution of $x = f^4(x)$, where $f(x) = 3.5x(1 - x)$. Marked points, corresponding to slopes less than one, form an attracting 4-cycle.

corresponding to an attracting 4-cycle by solving the algebraic equation

$$x = f(f(f(f(x)))), \quad (11.28)$$

where $f(x) = \rho x(1 - x)$. A graph of the function $y = f^4(x)$ is shown in Figure 11.10. The solution to (11.28) gives eight points. Four of these points, corresponding to $y = f^4(x)$ having negative slopes when intersecting the line $y = x$, form an attracting 4-cycle. Using results from the previous subsection, we could find the range of ρ resulting in attracting 4-cycle solutions. We could then proceed to find the range of ρ corresponding to solutions going to attracting 8-cycles, then 16-cycles, and so forth. However, this approach is only feasible for $n = 1$ and $n = 2$. We will have to take a different approach to compute the elements of the sequence (11.27). We will use an iteration technique consisting of four steps:

1. Finding the second iterate, $z(k + 2) = g(g(z(k)))$, of the logistic map (11.15).
2. Shifting one of the points of the 2-cycle of the second iterate to the origin.
3. Neglecting cubic and quartic terms in the shifted map.
4. Scaling the map obtained in step 3.

Performing iterations consisting of the above four steps will allow us to find the elements of the sequence (11.27).

The second iterate of (11.15) is

$$\begin{aligned} z(k + 2) &= (2 - \rho)z(k + 1) - \rho z^2(k + 1) \\ &= (2 - \rho)((2 - \rho)z(k) - \rho z^2(k)) - \rho((2 - \rho)z(k) - \rho z^2(k))^2 \\ &= (2 - \rho)^2 z(k) - ((2 - \rho)\rho + \rho(2 - \rho)^2)z^2(k) \\ &\quad + 2\rho^2(2 - \rho)z^3(k) - \rho^3 z^4(k). \end{aligned} \quad (11.29)$$

The points p and q , given by (11.17) and (11.18), are fixed points of the map (11.29). We select one of the points, p or q , to work with. Suppose that we choose the point p . Consider the transformation

$$v(k) = z(k) - p. \quad (11.30)$$

Note that

$$v(k+2) = z(k+2) - p. \quad (11.31)$$

Substituting (11.30) and (11.31) into (11.29), performing manipulations, and neglecting terms of order three and four, we obtain

$$v(k+2) = (4+2\rho-\rho^2)v(k) - \rho(\rho^2-2\rho-3-3\sqrt{\rho^2-2\rho-3})v^2(k). \quad (11.32)$$

One can easily verify the above result using MATLAB's function `subs` applied to (11.29)–(11.31). We will be studying fixed points of (11.32). We replace $k+2$, the argument of v on the left-hand side of (11.32), with $k+1$. Let

$$C_1(p) = \rho^2 - 2\rho - 3 - 3\sqrt{\rho^2 - 2\rho - 3}. \quad (11.33)$$

Then, (11.32) can be represented as

$$v(k+1) = (4+2\rho-\rho^2)v(k) - \rho C_1(p)v^2(k). \quad (11.34)$$

Note that if we used q in (11.30) instead of p , we would get

$$v(k+1) = (4+2\rho-\rho^2)v(k) - \rho C_1(q)v^2(k), \quad (11.35)$$

where now $C_1(q) = \rho^2 - 2\rho - 3 + 3\sqrt{\rho^2 - 2\rho - 3}$. In the next step, we transform the v map given by (11.34) into the same format as the z -map given by (11.15). This will allow us to apply stability properties of the z -map to the v -map. We introduce an iterated parameter ρ_1 such that

$$2 - \rho_1 = 4 + 2\rho_0 - \rho_0^2, \quad (11.36)$$

and represent (11.34) as

$$v(k+1) = (2 - \rho_1)v(k) - \rho_0 C_1(p)v^2(k). \quad (11.37)$$

Finally, we scale (11.37) using the transformation

$$s(k) = \frac{\rho_0 C_1(p)}{\rho_1} v(k). \quad (11.38)$$

Substituting (11.38) into (11.37) gives

$$\frac{\rho_1}{\rho_0 C_1(p)} s(k+1) = \frac{\rho_1(2 - \rho_1)}{\rho_0 C_1(p)} s(k) - \frac{\rho_1^2 \rho_0 C_1(p)}{\rho_0^2 C_1^2(p)} s^2(k). \quad (11.39)$$

Simplifying (11.39) yields

$$s(k+1) = (2 - \rho_1)s(k) - \rho_1 s^2(k) \quad (11.40)$$

The above s -map has the same format as the z -map given by (11.15) and analyzed in the previous subsection. It follows from these results that the s -map given by (11.40) has an attracting 2-cycle for $3 < \rho_1 < 1 + \sqrt{6}$, which corresponds to an attracting 4-cycle of the z -map given by (11.15). Using (11.36), we conclude that $3 < \rho_1 < 1 + \sqrt{6}$ corresponds to

$$1 + \sqrt{6} < \rho_0 < 1 + \sqrt{4 + \sqrt{6}} = 3.5396.$$

Indeed, representing (11.36) as

$$\rho_0^2 - 2\rho_0 - (2 + \rho_1) = 0$$

and using the quadratic formula, we get

$$\rho_0 = 1 + \sqrt{3 + \rho_1},$$

where we took the solution that gives $\rho_0 > 0$. The above concludes the first iteration. The result of this iteration is the element $c_3 = 1 + \sqrt{4 + \sqrt{6}} = 3.5396$ of the sequence (11.27) and the scaling parameter $C_1(p)$ whose properties will be investigated in the following subsection.

We then proceed with the next iteration. In particular, executing the steps of the iteration method on the s -map yields the new map that we call the t -map. This map has the form

$$t(k+1) = (4 + 2\rho_1 - \rho_1^2)t(k) - \rho_1 C_2 t^2(k). \quad (11.41)$$

Introducing the iterated parameter, ρ_2 , such that

$$\rho_2 = \rho_1^2 - 2\rho_1 - 2, \quad (11.42)$$

along with scaling the t -map using the transformation

$$u(k) = \frac{\rho_1 C_2}{\rho_2} t(k),$$

yields the u -map of the form

$$u(k+1) = (2 - \rho_2)u(k) - \rho_2 u^2(k). \quad (11.43)$$

The above map has an attracting 2-cycle for $3 < \rho_2 < 1 + \sqrt{6}$, which corresponds to an attracting 4-cycle of the s -map, which in turn corresponds to an attracting 8-cycle of the z -map. We calculate the range of the parameter ρ corresponding to attracting 8-cycles of the z -map by working backwards the recursive relation (11.42) to get

$$c_4 = 1 + \sqrt{3 + \rho_1} = 1 + \sqrt{3 + 1 + \sqrt{3 + \rho_2}} = 1 + \sqrt{4 + \sqrt{4 + \sqrt{6}}} = 3.5573.$$

In general, after $i + 1$ iterations we obtain

$$\rho_{i+1} = \rho_i^2 - 2\rho_i - 2 \quad (11.44)$$

and work backwards to find the term c_{i+3} of the sequence (11.27). We can also compute the elements of the sequence (11.27) using the recursive relation

$$\rho_i = \rho_{i+1}^2 - 2\rho_{i+1} - 2, \quad i = 1, 2, \dots \quad (11.45)$$

Solving the above, using the quadratic formula, yields

$$\rho_{i+1} = 1 + \sqrt{3 + \rho_i}.$$

This means that the elements of the sequence (11.27) satisfy the difference equation

$$c_{i+1} = 1 + \sqrt{3 + c_i}, \quad c(1) = 3. \quad (11.46)$$

The equation (11.46) has two equilibrium points that are obtained by solving the algebraic equation

$$c = 1 + \sqrt{3 + c}.$$

The equilibrium point

$$c_{eq} = \frac{3 + \sqrt{17}}{2}$$

is attracting, and we have

$$\lim_{i \rightarrow \infty} c_i = c_\infty = c_{eq} = 3.5616.$$

An analysis without truncation when computing second iterates of the logistic equation yields $c_\infty = 3.5699$. Such an analysis can be found in the paper by Feigenbaum [81]. After the parameter ρ reaches the limiting value $\rho = 3.5699$, no stable solution patterns of the logistic equation emerge. We obtain a chaotic behavior.

11.2.4 The Feigenbaum Numbers

The sequence (11.27) and the scaling factors

$$\frac{\rho_i C_{i+1}}{\rho_{i+1}} \quad (11.47)$$

have certain properties that are universal for all quadratic maps. Here, we perform an approximate analysis of these properties following Kahn [145, Chapter 16].

We begin by rewriting equation (11.45) as

$$(\rho_{i+1} - 1)^2 = 3 + \rho_i. \quad (11.48)$$

Hence,

$$\rho_{i+1} - 1 = \sqrt{3 + \rho_i} \quad (11.49)$$

and

$$\rho_i - 1 = \sqrt{3 + \rho_{i-1}}. \quad (11.50)$$

Subtracting (11.50) from (11.49) gives

$$\rho_{i+1} - \rho_i = \sqrt{3 + \rho_i} - \sqrt{3 + \rho_{i-1}}. \quad (11.51)$$

Summing (11.50) and (11.49) yields

$$\rho_{i+1} + \rho_i - 2 = \sqrt{3 + \rho_i} + \sqrt{3 + \rho_{i-1}}. \quad (11.52)$$

Multiplying (11.51) and (11.52), we get

$$(\rho_{i+1} - \rho_i)(\rho_{i+1} + \rho_i - 2) = \rho_i - \rho_{i-1}. \quad (11.53)$$

We rewrite (11.53) as

$$\rho_{i+1} + \rho_i - 2 = \frac{\rho_i - \rho_{i-1}}{\rho_{i+1} - \rho_i}. \quad (11.54)$$

Let

$$\boxed{\delta_i = \frac{\rho_i - \rho_{i-1}}{\rho_{i+1} - \rho_i}} \quad (11.55)$$

The limiting value of δ_i as i tends to infinity is called the *Feigenbaum number* δ . To find its value, we use (11.55) to represent equation (11.54) as

$$\delta_i = \rho_{i+1} + \rho_i - 2. \quad (11.56)$$

We obtain an approximate value for the Feigenbaum number δ because of the truncation procedure that we used to obtain (11.47). We have

$$\delta_{\text{approx}} = \lim_{i \rightarrow \infty} (\rho_{i+1} + \rho_i - 2) = 2c_{\infty} - 2 = 5.1231 \quad (11.57)$$

The exact value of the Feigenbaum number δ is $\delta_{\text{exact}} = 4.669201 \dots$

The other parameter that characterizes our iteration method is the scaling factor given by (11.47). The limiting value of (11.47) as i tends to infinity is called the *Feigenbaum number* α . Its approximate value is

$$\begin{aligned} \alpha_{\text{approx}} &= \lim_{i \rightarrow \infty} \frac{\rho_i C_{i+1}(p)}{\rho_{i+1}} \\ &= \frac{c_{\infty} C_{\infty}(p)}{c_{\infty}} \\ &= c_{\infty}^2 - 2c_{\infty} - 3 - 3\sqrt{c_{\infty}^2 - 2c_{\infty} - 3} \\ &= -2.24 \end{aligned} \quad (11.58)$$

The exact value of the Feigenbaum number α is $\alpha_{\text{exact}} = -2.50 \dots$

11.3 Fractals

This section is devoted to modeling of irregular shapes. These non-Euclidean objects are called *fractals*. The term *fractal* was coined by Mandelbrot [195], and it comes from the Latin adjective *fractus*, meaning “irregular.” The corresponding Latin verb is *frangere*, which means “to break to create irregular fragments.” An example of a simple fractal is the von Koch snowflake curve shown in Figure 11.11 along with an iterative procedure for constructing the curve. We start with dividing a line segment into thirds and then replacing the middle segment with two equal

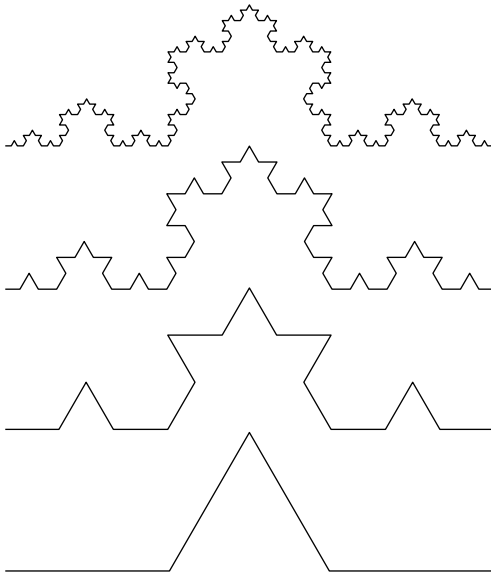


Figure 11.11 The first four stages in constructing the von Koch snowflake curve.

segments that form two sides of an equilateral triangle. The outcome of this stage is a shape consisting of four segments. In the next stage, each of the four segments is divided into thirds, and the middle segments are replaced with two equal segments that form two sides of equilateral triangle. The process is repeated over and over again. The resulting curve is characterized by *self-similarity* or *scaling*.

One of the characteristics of a fractal is its dimension. A line segment is a one-dimensional Euclidean object, and we can divide this line segment into N identical parts. In such a case, the individual parts can be viewed as scaled down, by the factor $r = 1/N$, versions of the original object. Similarly, a square, a two-dimensional Euclidean object, can be divided into N self-similar parts. These parts can be viewed as scaled down copies by the factor $r = 1/\sqrt{N}$ of the original square. In general, a D -dimensional object can be divided into N copies of itself. Each copy can be viewed as a scaled down version by the factor

$$r = \frac{1}{N^{1/D}} \quad (11.59)$$

of the parent object. Equation (11.59) can be represented as

$$r^D N = 1. \quad (11.60)$$

We now generalize the above arguments to the objects that are not necessarily Euclidean. Taking logarithms of both sides of (11.60) yields the expression for the *fractal self-similarity dimension*:

$$D = \frac{\log N}{\log (1/r)} \quad (11.61)$$

We now use (11.61) to compute the dimension of the von Koch curve. Any segment of the curve is composed of four subsegments, that is, $N = 4$. Each subsegment is a scaled-down copy of its parent by the factor $r = 1/3$. Hence, the self-similarity dimension of the von Koch curve is

$$D = \frac{\log N}{\log (1/r)} = \frac{\ln 4}{\ln 3} = 1.26. \quad (11.62)$$

We can see from the above that the fractal dimension, unlike the dimension of a common Euclidean pattern, need not be an integer.

Another example of a fractal is the *Cantor set*, sometimes called *Cantor's ternary set*. The Cantor set consists of the points of the unit interval $[0, 1]$ that remain after “middle third” intervals have been successively removed. This set is illustrated in Figure 11.12. It is an example of so-called *disconnected fractal* or *fractal dust*. For Cantor's set, we have $N = 2$ and $1/r = 3$.

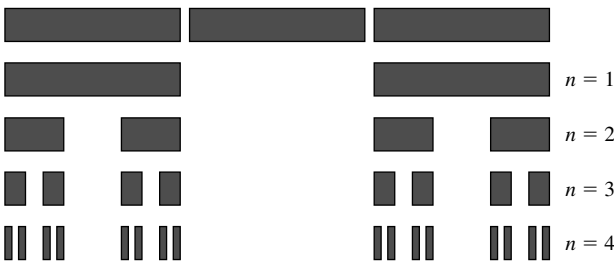


Figure 11.12 The first four iterations in constructing Cantor's ternary set.

Hence, the self-similarity dimension of the Cantor set is

$$D = \frac{\log N}{\log(1/r)} = \frac{\ln 2}{\ln 3} = 0.631. \quad (11.63)$$

Other types of fractal dimension are nicely explained in Peitgen, Jürgens, and Saupe [227].

11.3.1 The Mandelbrot Set

The Mandelbrot set, denoted \mathcal{M} , is the set of all complex $c \in \mathbb{C}$ for which the sequence

$$0, c, c^2 + c, (c^2 + c)^2 + c, \dots \quad (11.64)$$

does not tend to ∞ as n tends to ∞ . Thus, the Mandelbrot set is generated by the map in the complex plane,

$$z \mapsto z^2 + c, \quad c \in \mathbb{C}, \quad (11.65)$$

or, equivalently,

$$z(k+1) = z^2(k) + c, \quad c \in \mathbb{C}. \quad (11.66)$$

Let

$$z(k) = x(k) + jy(k) \quad \text{and} \quad c = a + jb.$$

Then,

$$\begin{aligned} z(k+1) &= x(k+1) + jy(k+1) \\ &= x^2(k) - y^2(k) + a + j(2x(k)y(k) + b) \\ &= x^2(k) + 2jx(k)y(k) + (jy(k))^2 + a + jb \\ &= z^2(k) + c. \end{aligned}$$

Thus, we can express the first-order quadratic complex dynamical system (11.66) as a second-order real dynamical system of the form

$$\begin{aligned} x(k+1) &= x^2(k) - y^2(k) + a, \\ y(k+1) &= 2x(k)y(k) + b. \end{aligned} \quad (11.67)$$

The representations (11.66) and (11.67) are indeed equivalent.

In terms of the dynamical system (11.67), the Mandelbrot set is the set of parameter values, a and b , for which the solutions of (11.67) starting at $x(0) = y(0) = 0$ remain bounded. The Mandelbrot set is shown in Figure 11.13.

11.3.2 Julia Sets

A filled Julia set corresponding to the given values of the parameters a and b of the dynamical system (11.67) is the set of initial conditions $(x(0), y(0))$ resulting in bounded system trajectories. For each value of $c = a + jb$ there is a different Julia set. In order to construct a Julia set, pick first a point (a, b) , say in the Mandelbrot set. Then, pick an initial condition $(x(0), y(0))$ and compute the resulting trajectory of (11.67) starting from the chosen initial condition $(x(0), y(0))$. Plot only the initial conditions $(x(0), y(0))$ that give rise to bounded trajectories. The set of plotted points is the filled Julia set corresponding to the particular parameter values (a, b) . The boundary of the filled Julia set is the Julia set. Two examples of filled Julia sets are shown in

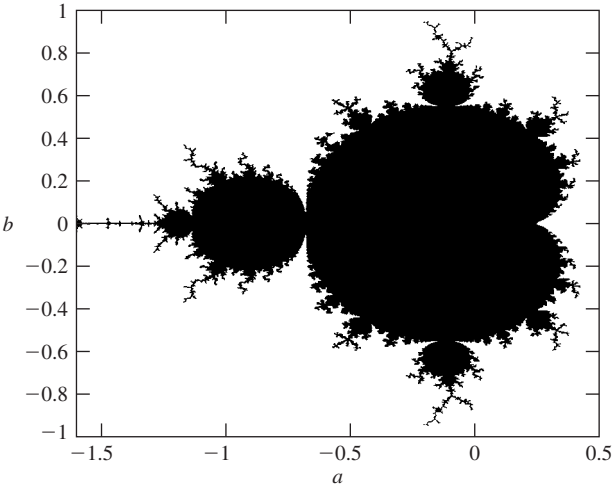


Figure 11.13 The Mandelbrot set for $x(0) = y(0) = 0$.

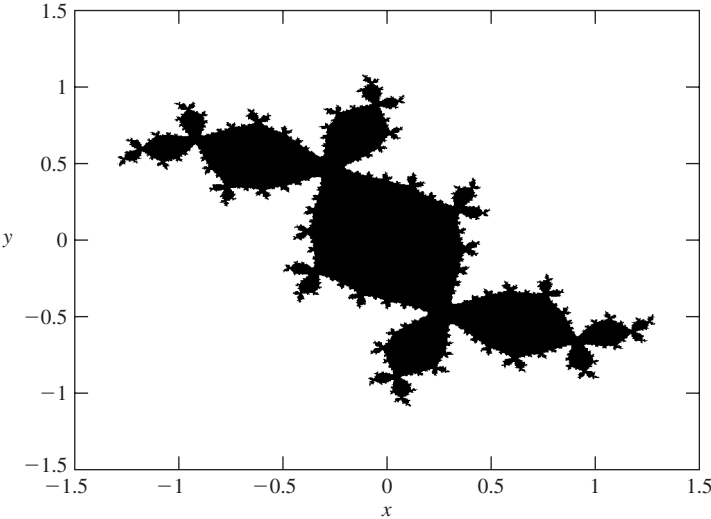


Figure 11.14 Filled Julia set for $c = -0.12256117 + j0.74486177$.

Figures 11.14 and 11.15. The filled Julia set shown in Figure 11.14 is called “Douady’s Rabbit” by Devaney and Keen [63, p. 61]. Peitgen and Saupe [228, p. 199] observed that in a sense the Mandelbrot set can be considered as a one-page dictionary of all Julia sets. This means that the Mandelbrot set stores infinitely many different Julia sets.

11.3.3 Iterated Function Systems

An iterated function system (IFS) can be used to model natural objects such as clouds, smoke, leaves, or ferns. An IFS uses affine transformations that we discuss next. A two-dimensional affine transformation $w : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ from \mathbb{R}^2 into itself is defined by

$$w(x) = Ax + b, \tag{11.68}$$

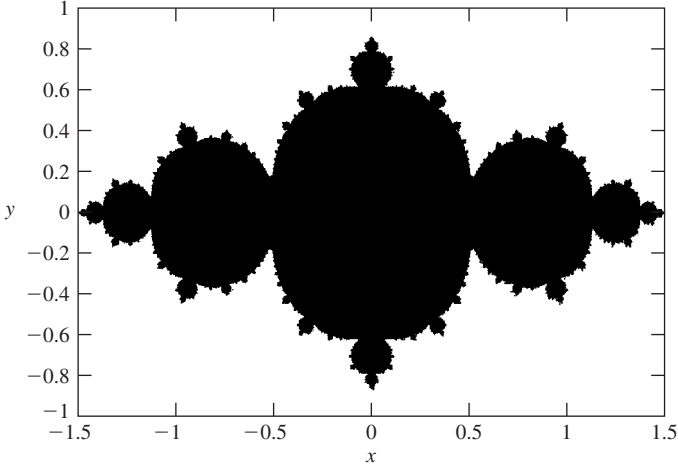


Figure 11.15 Filled Julia set for $c = -0.75$.

where $A \in \mathbb{R}^{2 \times 2}$ and $b \in \mathbb{R}^2$. The affine transformation (11.68) satisfies the condition

$$\|w(x) - w(y)\| \leq L\|x - y\|, \quad (11.69)$$

where the Lipschitz constant L is

$$L = \|A\|.$$

The affine transformation w is contractive if $L < 1$. It is a symmetry if

$$\|w(x) - w(y)\| = \|x - y\|.$$

The affine transformation w is expansive if $L > 1$.

A two-dimensional IFS consists of a set of N affine transformations denoted

$$\{w_1, w_2, \dots, w_N\},$$

each transformation mapping \mathbb{R}^2 into itself, together with a set of probabilities

$$\{p_1, p_2, \dots, p_N\},$$

where each $p_i \geq 0$ and

$$p_1 + p_2 + \dots + p_N = 1.$$

That is, an IFS is the set

$$\{w_n, p_n : n = 1, 2, \dots, N\}.$$

Let $L_n, n = 1, 2, \dots, N$, denote the Lipschitz constant for w_n . Then, we say that the IFS obeys the average contractivity condition if

$$L_1^{p_1} L_2^{p_2} \dots L_N^{p_N} < 1. \quad (11.70)$$

Interesting mathematical properties of the IFS that satisfy the condition (11.70) are reported by Barnsley [21].

An IFS operates as follows. A set of N affine transformations is chosen and an initial point is picked. Then, the affine transformations are iteratively applied over and over again in a random

order with probabilities p_n , for example, proportional to the determinants of A_n . The result is a sequence of points from \mathbb{R}^2 :

$$\mathbf{x}(0), \mathbf{x}(1), \dots$$

As an example, suppose that $N = 4$, where

$$\begin{aligned} A_1 &= \begin{bmatrix} 0.02 & 0 \\ 0 & 0.21 \end{bmatrix}, & b_1 &= \begin{bmatrix} -0.15 \\ -0.8 \end{bmatrix}, \\ A_2 &= \begin{bmatrix} 0.8496 & 0.0255 \\ -0.0255 & 0.8496 \end{bmatrix}, & b_2 &= \begin{bmatrix} 0 \\ 0.15 \end{bmatrix}, \\ A_3 &= \begin{bmatrix} -0.1554 & 0.235 \\ 0.1958 & 0.1865 \end{bmatrix}, & b_3 &= \begin{bmatrix} 0.05 \\ -0.55 \end{bmatrix}, \\ A_4 &= \begin{bmatrix} 0.1554 & -0.235 \\ 0.1958 & 0.1865 \end{bmatrix}, & b_4 &= \begin{bmatrix} -0.35 \\ -0.4 \end{bmatrix}. \end{aligned}$$

The above data comes from Bunde and Havlin [38]. The elements of the set of probabilities

$$\{p_1, p_2, p_3, p_4\}$$

were constructed as

$$p_n = \frac{\det(A_n)}{\sum_{n=1}^N \det(A_n)}, \quad n = 1, 2, 3, 4.$$

The initial point was chosen to be $\mathbf{x}(0) = [-0.1 \ -0.2]^T$. Then, the roulette-wheel method was used to select the affine transformation to operate on $\mathbf{x}(0)$. The roulette-wheel method operates as follows. Each affine transformation is assigned a slice of a wheel, where the slice is proportional to the probability p_n of the particular transformation w_n . A random number from the interval $[0, 1]$ is generated, which corresponds to spinning the wheel. The transformation under the wheel's marker is selected. Suppose that w_k was selected. Then w_k operates on $\mathbf{x}(0)$ and the result is

$$\mathbf{x}(1) = A_k \mathbf{x}(0) + b_k.$$

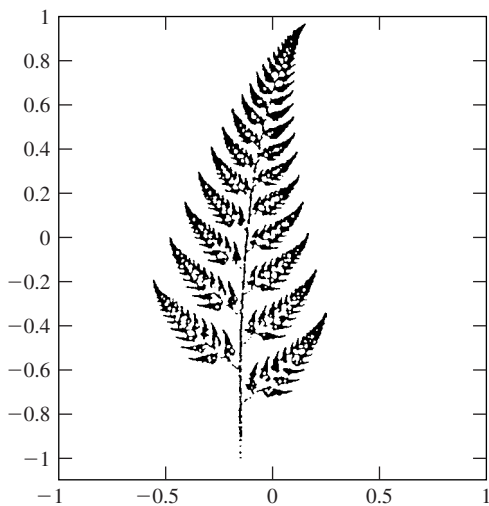


Figure 11.16 The Barnsley fern.

The roulette wheel is spun again to select the affine transformation to act upon $\mathbf{x}(1)$ in order to obtain $\mathbf{x}(2)$, and so on.

The result of 30,000 iterations is shown in Figure 11.16. The object is called the Barnsley fern because it was M. Barnsley [21] who first generated the fern fractal using an IFS.

11.4 Lyapunov Exponents

Trajectories of a chaotic dynamical system are characterized by a sensitive dependence upon initial conditions, which means that the distance between the adjacent trajectories grows exponentially with time as illustrated in Figure 11.17. For a discrete dynamical system, if x^* is a fixed point of a one-dimensional map f and if $|f'(x^*)| = a > 1$, then by Theorem 11.1, trajectories starting near x^* will separate from x^* at a multiplicative rate of approximately a per iteration until they move significantly far away from x^* . Let x_1 be a point that constitutes the k -cycle. Then, for a k -cycle, the trajectory starting from a point close to x_1 , will separate from x_1 at a rate approximately equal to the product of the derivatives of the k points of the k -cycle. Let this product be equal to A and suppose that $A > 1$. Then, it takes k iterations to separate a trajectory starting close to x_1 by a distance A from x_1 . We can say that in the case of a k -cycle the average rate of separation is $A^{1/k}$ per iteration. We now are ready to introduce the term *Lyapunov number* to describe the average multiplicative rate of separation of points close to x_1 . For a k -cycle, the Lyapunov number $L(x_1)$ is

$$\begin{aligned} L(x_1) &= (|(f^k)'(x_1)|)^{1/k} \\ &= (|f'(x_1)| |f'(x_2)| \cdots |f'(x_k)|)^{1/k}. \end{aligned} \quad (11.71)$$

The concept of the Lyapunov number is now extended to nonperiodic trajectories.

Definition 11.2 Let f be a differentiable map on the real line \mathbb{R} and let $\{x(0), x(1), \dots\}$ be a trajectory starting from $x(0)$. Then, the Lyapunov number $L(x(0))$ of the trajectory starting from $x(0)$ is

$$L(x(0)) = \lim_{N \rightarrow \infty} (|f'(x(0))| \cdots |f'(x(N-1))|)^{1/N}$$

if the limit exists.

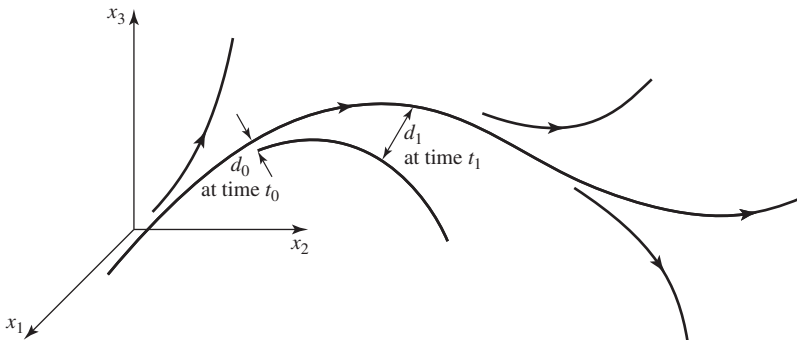


Figure 11.17 The change in distance between adjacent chaotic trajectories that is used to compute the largest Lyapunov exponent.

The Lyapunov exponent, $\mu(x(0))$, is then defined as the natural logarithm of the Lyapunov number; that is,

$$\mu(x(0)) = \ln L(x(0)) = \lim_{N \rightarrow \infty} \frac{\ln |f'(x(0))| + \cdots + \ln |f'(x(N-1))|}{N} \quad (11.72)$$

if the limit exists.

We now look at the Lyapunov exponent from another perspective. Let d_0 denote the distance between the two initial states of a continuous dynamical system. Then, for a chaotic motion, at a later time,

$$d(t) = e^{\mu t} d_0. \quad (11.73)$$

The measure of the trajectories' divergence is obtained by averaging the exponential growth at many points along a trajectory. This is because the chaotic trajectories' divergence can only be measured locally since if the trajectories are bounded, $d(t)$ in (11.73) cannot go to infinity. To define the measure of the divergence of chaotic trajectories, first a reference trajectory is obtained. Next a point on an adjacent trajectory is selected and $d(t)/d_0$ is measured. This is illustrated in Figure 11.17. Then, a new $d_0(t)$ is selected on new adjacent trajectory. Using the collected data, the Lyapunov exponent is computed as

$$\mu(x(0)) = \lim_{N \rightarrow \infty} \frac{1}{t_N - t_0} \sum_{k=1}^N \ln \frac{d(t_k)}{d_0(t_{k-1})}. \quad (11.74)$$

There are n Lyapunov exponents for an n -dimensional nonlinear dynamical system, just as there are n eigenvalues that govern behavior of an n -dimensional linear system. To define these Lyapunov exponents for an n -dimensional nonlinear dynamical system model, first choose an n -dimensional sphere centered at a point on a reference trajectory. Let the radius of this sphere be δ_0 . We can view this sphere as an n -dimensional ellipsoid with n semiaxes of equal length $\delta_i(0) = \delta_0$. Next, at a later time t construct an n -dimensional ellipsoid with the property that all the trajectories emanating from the previously chosen sphere pass through this ellipsoid. Let the n semiaxes of the ellipsoid be denoted $\delta_i(t)$. Then, we can compute the Lyapunov exponents as

$$\mu_i(x(0)) = \lim_{t \rightarrow \infty} \left(\lim_{\delta_i(0) \rightarrow 0} \left\{ \frac{1}{t} \ln \frac{\delta_i(t)}{\delta_i(0)} \right\} \right), \quad i = 1, 2, \dots, n \quad (11.75)$$

if the limits exist. It is common to order the Lyapunov exponents so that

$$\mu_1 \geq \mu_2 \geq \dots \geq \mu_n.$$

Thus, the Lyapunov exponents are numbers that are used to measure the exponential attraction, or separation, in time of two neighboring trajectories of a dynamical system.

◆ Example 11.1

This example was adapted from Moon [209, p. 193]. Suppose that we are given a one-dimensional map

$$x(k+1) = f(x(k)). \quad (11.76)$$

We will demonstrate again that in the regions where f is differentiable, the Lyapunov exponent is determined by $|f'|$. Indeed, let $x(k)$ and $x(k) + \varepsilon$ be two

initial conditions. Then, we have

$$d_0(k-1) = \varepsilon \quad \text{and} \quad d(k) = |f(x(k) + \varepsilon) - f(x(k))| \approx |f'|_{x=x(k)} |\varepsilon|.$$

Using (11.74), we obtain

$$\begin{aligned} \mu(x(0)) &= \lim_{N \rightarrow \infty} \frac{1}{t_N - t_0} \sum_{k=1}^N \ln \frac{d(t_k)}{d_0(t_{k-1})} \\ &= \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \ln |f'(x(k))|, \end{aligned}$$

which is the same as (11.72).

The notion of the Lyapunov exponent will be used to formally define chaotic motion. The other term that we will use in this definition is an asymptotically periodic trajectory that we define next. This definition is adapted from Alligood, Sauer, and Yorke [7, p. 108].

Definition 11.3 A trajectory $\{x(0), \dots, x(N), \dots\}$ of a smooth map f is called *asymptotically periodic* if it converges to a periodic trajectory as $N \rightarrow \infty$; that is, there exists a periodic trajectory $\{y_1, y_2, \dots, y_k, y_1, y_2, \dots\}$ such that

$$\lim_{N \rightarrow \infty} \|x(N) - y(N)\| = 0.$$

Using the definition of Lyapunov exponent and the above definition, we can define a chaotic trajectory. This definition is adapted from Alligood et al. [7, p. 110].

Definition 11.4 Let f be a map of the real line \mathbb{R} , and let $\{x(0), x(1), \dots\}$ be a bounded trajectory of the map f . This trajectory is chaotic if:

1. The trajectory is not asymptotically periodic.
2. The Lyapunov exponent $\mu(x(0))$ is greater than zero.

For a map on \mathbb{R}^n , each trajectory has n Lyapunov numbers. The Lyapunov numbers, in this case, measure the rates of separation from the current point of the trajectory along n orthogonal directions. The first direction is the one where the separation between the nearby points is the greatest. The second direction is the one in which the separation is the greatest from all directions perpendicular to the first direction. The third direction will be the one where the separation is the greatest among directions perpendicular to the first two directions, and so on. We illustrate the above in two dimensions in Figure 11.18, where a unit disk is centered at $x(0)$. After N iterates

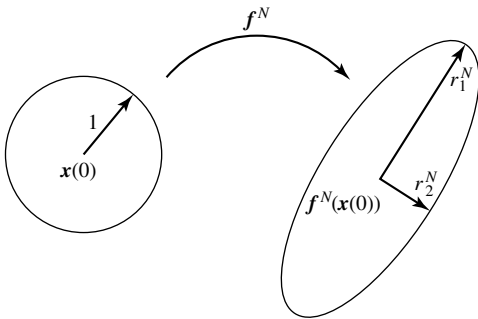


Figure 11.18 The unit disk is mapped into an ellipse by a two-dimensional map.

of the two-dimensional map f , the disk evolves into an ellipse. The per-iteration changes of the ellipse axes are the Lyapunov numbers of $\mathbf{x}(0)$. Thus, the Lyapunov numbers measure the contraction or expansion near the trajectory of $\mathbf{x}(0)$ during the first N iterations. The natural logarithm of each Lyapunov number is a Lyapunov exponent.

Let $Df(\mathbf{x}(0))$ be the first derivative of f evaluated at $\mathbf{x}(0)$. Let $J_N = Df^N(\mathbf{x}(0))$ be the first derivative of the N th iterate of f . We denote by B the unit sphere centered at $\mathbf{x}(0)$. With the above notation, we are ready to give the definition of Lyapunov exponents for a map on \mathbb{R}^n . This definition is adapted from Alligood et al. [7, p. 195].

Definition 11.5 Let f be a smooth map on \mathbb{R}^n and let r_k^N be the length of the k th longest orthogonal axis of the ellipsoid $J_N B$ for a trajectory with the initial state $\mathbf{x}(0)$. The k th Lyapunov number is

$$L_k(\mathbf{x}(0)) = \lim_{N \rightarrow \infty} (r_k^N)^{1/N}$$

if the limit exists. The k th Lyapunov exponent of $\mathbf{x}(0)$ is

$$\mu_k(\mathbf{x}(0)) = \ln L_k(\mathbf{x}(0)).$$

Using the above definition, we can extend the definition of chaotic trajectories for maps on \mathbb{R}^n . This definition is adapted from Alligood et al. [7, p. 196].

Definition 11.6 Let f be a map on \mathbb{R}^n , $n \geq 1$, and let $\{\mathbf{x}(0), \mathbf{x}(1), \mathbf{x}(2), \dots\}$ be a bounded trajectory of the map f . The trajectory is chaotic if:

1. It is not asymptotically periodic.
2. No Lyapunov exponent is equal to zero.
3. $\mu_1(\mathbf{x}(0)) > 0$.

For continuous systems, the above definition of a chaotic trajectory has to be modified—see Alligood et al. [7, pp. 385–386].

To numerically compute Lyapunov exponents, a reference trajectory of a given nonlinear system is constructed and the ellipsoid semiaxes lengths r_k 's are found using linearized models. We now analyze this process in some detail. For this, suppose that we are given an n th-order continuous dynamical system model

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}), \quad (11.77)$$

where $\mathbf{x}(t) \in \mathbb{R}^n$ is the system's state at time t , and $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^n$ is called the vector field. The solution to (11.77) subject to initial condition $\mathbf{x}(t_0) = \mathbf{x}_0$ is called a trajectory and is denoted $\phi_t(\mathbf{x}_0)$. For every t , the mapping

$$\phi_t: \mathbb{R}^n \rightarrow \mathbb{R}^n,$$

maps the state space \mathbb{R}^n into itself. The mapping ϕ_t is called the *flow* of the system (11.77). For example, if (11.77) is linear—that is, $\mathbf{f}(\mathbf{x}) = \mathbf{A}\mathbf{x}$, where $\mathbf{A} \in \mathbb{R}^{n \times n}$ —then

$$\phi_t(\mathbf{x}_0) = e^{\mathbf{A}(t-t_0)} \mathbf{x}_0.$$

Consider now an equilibrium state \mathbf{x}_e of (11.77). The local behavior of (11.77) in a neighborhood of \mathbf{x}_e is determined by the linearized model of (11.77) about the equilibrium state \mathbf{x}_e . The linearized model is

$$\frac{d}{dt} \Delta \mathbf{x} = D\mathbf{f}(\mathbf{x}_e) \Delta \mathbf{x}, \quad (11.78)$$

where $\Delta \mathbf{x} = \mathbf{x} - \mathbf{x}_e$. We can also linearize the system (11.77) about its trajectory $\phi_t(\mathbf{x}_0)$. The linearized model about the trajectory $\phi_t(\mathbf{x}_0)$ is

$$\frac{d}{dt} \Delta \mathbf{x} = D\mathbf{f}(\phi_t(\mathbf{x}_0)) \Delta \mathbf{x}, \quad \Delta \mathbf{x}(t_0) = \Delta \mathbf{x}_0. \quad (11.79)$$

The solution to (11.79) is

$$\Delta \mathbf{x}(t) = \Phi_t(\mathbf{x}_0) \Delta \mathbf{x}_0.$$

We call the $n \times n$ matrix $\Phi_t(\mathbf{x}_0)$ the *transition matrix* of the system (11.79). The transition matrix satisfies the equation

$$\frac{d}{dt} \Phi = D\mathbf{f}(\phi_t(\mathbf{x}_0)) \Phi, \quad \Phi_0 = I_n.$$

Let

$$\{m_i(t)\}_{i=1}^n$$

be the eigenvalues of the transition matrix $\Phi_t(\mathbf{x}_0)$. Then, the Lyapunov exponents of \mathbf{x}_0 are

$$\mu_i(\mathbf{x}_0) = \lim_{t \rightarrow \infty} \frac{1}{t} \ln |m_i(t)|, \quad i = 1, 2, \dots, n \quad (11.80)$$

if the limits exist. The above is an alternative to (11.75) representation of the Lyapunov exponents.

◆ Example 11.2

The Lyapunov exponents of an equilibrium state \mathbf{x}_e of the system (11.77), with $t_0 = 0$, can be found by first linearizing (11.77) about the equilibrium state to obtain a linearized model of the form (11.78). The transition matrix of (11.78) is

$$\Phi_t(\mathbf{x}_e) = e^{D\mathbf{f}(\mathbf{x}_e)t}.$$

The i th eigenvalue of the transition matrix $\Phi_t(\mathbf{x}_e)$ is $m_i(t) = e^{\lambda_i t}$, where λ_i is the i th eigenvalue of the matrix $D\mathbf{f}(\mathbf{x}_e)$. Therefore, the Lyapunov exponents of \mathbf{x}_e are

$$\begin{aligned} \mu_i(\mathbf{x}_e) &= \lim_{t \rightarrow \infty} \frac{1}{t} \ln |m_i(t)| \\ &= \lim_{t \rightarrow \infty} \frac{1}{t} \ln |e^{\Re \lambda_i t} e^{j \Im \lambda_i t}| \\ &= \lim_{t \rightarrow \infty} \frac{1}{t} \ln (|e^{\Re \lambda_i t}| |\cos(\Im \lambda_i t) + j \sin(\Im \lambda_i t)|) \\ &= \lim_{t \rightarrow \infty} \frac{1}{t} \Re \lambda_i \\ &= \Re \lambda_i. \end{aligned} \quad (11.81)$$

Thus, the Lyapunov exponents in this case are the real parts of the eigenvalues of the linearized system about the equilibrium state. Note that if \mathbf{x}_0 is in the basin of attraction of \mathbf{x}_e , that is, if

$$\lim_{t \rightarrow \infty} \phi_t(\mathbf{x}_0) = \mathbf{x}_e,$$

then the Lyapunov exponents of \mathbf{x}_0 are the same as the Lyapunov exponents of \mathbf{x}_e .

◆ Example 11.3

Suppose that we are given a discrete-time dynamical system model,

$$\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k)), \quad \mathbf{x}(0) = \mathbf{x}_0. \quad (11.82)$$

Let \mathbf{x}_e be the equilibrium state of the model (11.82), that is,

$$\mathbf{x}_e = \mathbf{f}(\mathbf{x}_e).$$

The linearized model is

$$\Delta \mathbf{x}(k+1) = D\mathbf{f}(\mathbf{x}_e)\Delta \mathbf{x}(k).$$

The state transition matrix of the linearized model is

$$\Phi_k(\mathbf{x}_e) = (D\mathbf{f}(\mathbf{x}_e))^k \quad (11.83)$$

If λ_i is the i th eigenvalue of the matrix $D\mathbf{f}(\mathbf{x}_e)$, then λ_i^k is the i th eigenvalue of the transition matrix (11.83). The Lyapunov exponents of \mathbf{x}_e , in this case, are

$$\begin{aligned} \mu_i(\mathbf{x}_e) &= \lim_{k \rightarrow \infty} \frac{1}{k} \ln |m_i(k)| \\ &= \lim_{k \rightarrow \infty} \frac{1}{k} \ln |\lambda_i|^k \\ &= \ln |\lambda_i|. \end{aligned} \quad (11.84)$$

In summary, Lyapunov exponents can be used to characterize a chaotic motion. Specifically, a necessary condition for a chaotic motion to occur in a dynamical system with bounded trajectories is that the largest Lyapunov exponent be positive. An efficient numerical method for evaluating the largest Lyapunov exponent can be found, for example, in the paper by Parker and Chua [221]; see also their book [222] on practical numerical algorithms for chaotic systems.

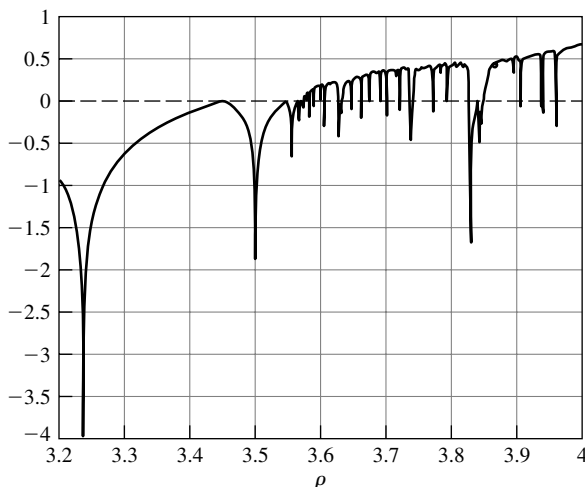


Figure 11.19 Lyapunov exponent versus parameter ρ for the logistic equation (11.10).

In Figure 11.19 we show a plot of the Lyapunov exponent versus the parameter ρ for the logistic equation (11.10).

11.5 Discretization Chaos

In this section we consider examples of chaotic systems that are discrete-time analogs of continuous-time systems, where the latter systems do not exhibit any chaotic behavior. We begin by considering a method of obtaining discrete models from continuous models. Suppose that we are given a scalar differential equation

$$\dot{x} = ax, \quad x(t_0) = x_0. \quad (11.85)$$

The solution to the above differential equation is

$$x(t) = e^{a(t-t_0)}x(t_0). \quad (11.86)$$

Let $\Delta t = t - t_0$. Then, we can express (11.86) as

$$x(1) = e^{a\Delta t}x(0), \quad (11.87)$$

where $x(1) = x(t_0 + \Delta t)$ and $x(0) = x(t_0)$. More generally,

$$x(k+1) = e^{a\Delta t}x(k), \quad (11.88)$$

where $x(k+1) = x((k+1)\Delta t)$ and $x(k) = x(k\Delta t)$. The discrete model (11.88) of the continuous differential equation (11.85) is exact in the sense that the values of (11.86) and (11.88) are the same at the sampling times. We generalize the above method to a class of nonlinear differential equations. Consider a scalar differential equation

$$\dot{x}(t) = x(t)h(x(t)). \quad (11.89)$$

Dividing both sides of the above by x and then integrating yields

$$\int \frac{dx}{x} = \int h(x(t)) dt. \quad (11.90)$$

Suppose now that we hold $h(x(t))$ constant on the interval $[k\Delta t, (k+1)\Delta t]$. Then, evaluating (11.90) over the above interval gives

$$\ln \left(\frac{x((k+1)\Delta t)}{x(k\Delta t)} \right) = h(x(k\Delta t))\Delta t.$$

Using the notation introduced after (11.88), we represent the above as

$$x(k+1) = e^{h(x(k))\Delta t}x(k). \quad (11.91)$$

◆ Example 11.4

Consider the system of differential equations

$$\begin{aligned} \dot{x}_1 &= \gamma x_1(1 - x_1 - 0.5x_2), \\ \dot{x}_2 &= 2\gamma x_2(1 - 0.5x_1 - x_2), \end{aligned} \quad (11.92)$$

where $\gamma = 2.5$. The above system models a pair of competing species. For a discussion of modeling of competing species population dynamics, the reader is

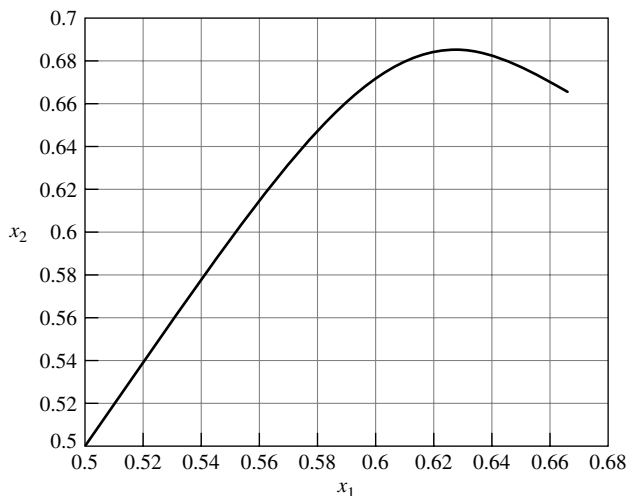


Figure 11.20 A trajectory of the system (11.92).

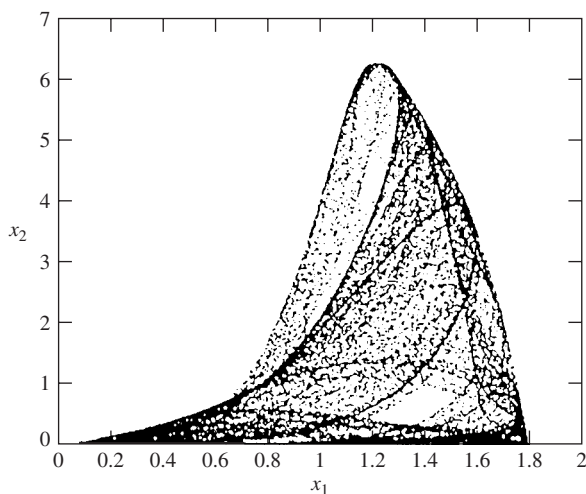


Figure 11.21 A trajectory of the system (11.93).

referred to Subsection 1.7.5. The system (11.92) was analyzed by Grantham and Athalye [106]. A trajectory of the system (11.92) for $x_1(0) = x_2(0) = 0.5$ is shown in Figure 11.20. We next discretize (11.92) using the formula (11.91). For $\Delta t = 1$, we obtain the following discrete-time analog:

$$\begin{aligned} x_1(k+1) &= x_1(k) \exp(\gamma(1 - x_1(k) - 0.5x_2(k))), \\ x_2(k+1) &= x_2(k) \exp(2\gamma(1 - 0.5x_1(k) - x_2(k))). \end{aligned} \quad (11.93)$$

The trajectory of the discrete model (11.93) starting from the same initial condition as its continuous-time analog is shown in Figure 11.21. The continuous-time model trajectory asymptotically converges to the equilibrium state $[2/3 \ 2/3]^T$, while the discrete-time trajectory wanders around chaotically.

11.6 Controlling Chaotic Systems

A trajectory of a chaotic system behaves in a random-like fashion. The set on which the chaotic motion is taking place attracts trajectories starting from a set of initial conditions. Roughly speaking, this attracting set is called a *chaotic attractor*. It is not possible to predict when, or if, a trajectory will pass through or near a given point in the chaotic attractor. However, with the probability equal to one, arbitrarily small neighborhood of any given point in the chaotic attractor will be visited by the chaotic system's trajectory in a finite time. This property of a chaotic system is referred to as *ergodicity*. Using the ergodicity property, one can propose a control algorithm whose task it is to bring the chaotic system's trajectory to a prespecified state in the system's state space. The idea behind the chaotic control algorithm, which we discuss next, is to set up "fishing nets" in the chaotic attractor in a neighborhood of an operating point to "catch" the system trajectory. Once the trajectory is "caught," a control law is used to bring the trajectory to the desired operating point. The controller is constructed using linear quadratic regulator (LQR) technique. The chaotic algorithm is illustrated with two examples: the Hénon map and the bouncing ball system.

11.6.1 Ingredients of Chaotic Control Algorithm

Chaotic motion is characterized by extreme sensitivity to initial conditions. This might suggest that it would be impossible, for example, to stabilize a chaotic system with respect to a desired operating point, because any perturbation resulting from applying a control action would lead to a chaotic motion's exponential growth in time. However, it turns out that precisely this property of a chaotic motion can be used to convert a chaotic attractor to a periodic motion. A chaotic attractor, as observed by Ott, Grebogi, and Yorke [217], typically has embedded within it an infinite number of unstable periodic orbits. These periodic orbits are dense in the attractor in the sense that they pass through any neighborhood, however small, of any point in the attractor. Ott et al. [217] noticed that these properties could be advantageous in achieving control objectives. They proposed a method, now called the OGY method, which involves converting, using small time-dependent perturbations of a system parameter, a chaotic attractor into a periodic orbit. The idea behind this method is to exploit the ergodicity property of a chaotic system. Specifically, all of the unstable periodic orbits of the chaotic attractor will eventually be visited by the system trajectory. Once the trajectory gets to a neighborhood of the desired unstable periodic orbit, a small perturbation is applied to stabilize the orbit of interest.

The chaotic control algorithm, which we discuss next, can be viewed as a variant of the OGY method. This algorithm was presented by Vincent [290] and is characterized by two ingredients: a chaotic attractor and a controllable target. These components of the chaotic algorithm are illustrated in Figure 11.22. We assume that either the system under consideration is chaotic or

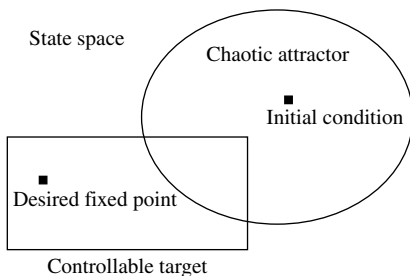


Figure 11.22 A pictorial representation of the chaotic algorithm.

that it can be made chaotic using an open-loop control. Next, we select a fixed point of interest and assume that a feedback law can be constructed such that the closed-loop system is nonchaotic and the chosen fixed point is asymptotically stable. The region of attraction of the fixed point should have a nonempty intersection with the chaotic attractor. A controllable target can be any subset of the region of attraction of the fixed point that has a nonempty intersection with the chaotic attractor. Note that the fixed point does not have to be located in the chaotic attractor. If we start the system from any initial state in the chaotic attractor, the resulting chaotic trajectory will at some time arrive at the intersection region of the chaotic attractor and the controllable target. The chaotic algorithm keeps track when the trajectory reaches the intersection region. When this happens, the open-loop control, which was used to create chaotic motion, is turned off and the feedback control is turned on to drive the system toward the fixed point. The region of attraction intersecting the chaotic attractor plays the role of a “fishing net” for the chaotic trajectory.

11.6.2 Chaotic Control Algorithm

We consider two classes of single-input dynamical systems. The first class is modeled by a nonlinear difference equation,

$$\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k), u(k)), \quad (11.94)$$

and the other class is modeled by a nonlinear differential equation,

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), u(t)). \quad (11.95)$$

The control u , in both cases, is bounded. The control u is either a specified function of time or a specified function of the system state. In the first case we have an open-loop control, whereas in the latter case we have a closed-loop control. We assume that for a given system there exists an open-loop control such that the system has a chaotic attractor. In addition, we assume that for a specified constant control $u = u_e$ there is a corresponding equilibrium state located in a vicinity of the chaotic attractor. The equilibrium state satisfies the equation

$$\mathbf{x}_e = \mathbf{f}(\mathbf{x}_e, u_e)$$

for the discrete-time system, and it satisfies the relation

$$\mathbf{0} = \mathbf{f}(\mathbf{x}_e, u_e)$$

for the continuous-time system.

We now briefly describe the way we obtain a controllable target for the system of interest. We first linearize the system about the selected equilibrium state. The linearized discrete-time system has the form

$$\Delta \mathbf{x}(k+1) = \mathbf{A} \Delta \mathbf{x}(k) + \mathbf{b} \Delta u(k),$$

where $\Delta \mathbf{x}(k) = \mathbf{x}(k) - \mathbf{x}_e$ and $\Delta u(k) = u(k) - u_e$. The linearized continuous-time system is

$$\frac{d}{dt} \Delta \mathbf{x}(t) = \mathbf{A} \Delta \mathbf{x}(t) + \mathbf{b} \Delta u(t),$$

where $\Delta \mathbf{x}(t) = \mathbf{x}(t) - \mathbf{x}_e$ and $\Delta u(t) = u(t) - u_e$. Then, the LQR method, described in Chapter 5, is used to construct a state variable feedback control law, $\Delta u(\Delta \mathbf{x}) = -\mathbf{k} \Delta \mathbf{x}$, that guarantees the origin of the linearized system,

$$\Delta \mathbf{x}(k+1) = (\mathbf{A} - \mathbf{b}\mathbf{k}) \Delta \mathbf{x}(k), \quad (11.96)$$

or

$$\frac{d}{dt} \Delta \mathbf{x}(t) = (\mathbf{A} - \mathbf{b}k) \Delta \mathbf{x}(t), \quad (11.97)$$

to be asymptotically stable. This implies that for the corresponding nonlinear system driven by the state feedback controller, the equilibrium state about which the system was linearized, is asymptotically stable. We use a Lyapunov function obtained from the linearized system to estimate the region of asymptotic stability (RAS) of the closed-loop nonlinear system. The Lyapunov function has the form $V(\Delta \mathbf{x}) = \Delta \mathbf{x}^T \mathbf{P}_l \mathbf{x}$, where \mathbf{P}_l is obtained by solving the corresponding Lyapunov equation for the optimal closed-loop linear system given by (11.96) or by (11.97). This Lyapunov function is then used to determine a controllable target for the nonlinear system driven by a bounded controller of the form

$$u = -\text{sat}(\mathbf{k} \Delta \mathbf{x}) + u_e, \quad (11.98)$$

where

$$\text{sat}(v) = \begin{cases} u_{\max} - u_e & \text{if } v > u_{\max} - u_e, \\ v & \text{if } u_{\min} - u_e \leq v \leq u_{\max} - u_e, \\ u_{\min} - u_e & \text{if } v < u_{\min} - u_e, \end{cases}$$

and the bounds u_{\min} and u_{\max} are specified by a designer. Using the above control law, we then seek the largest level curve $V(\mathbf{x}) = V_{\max}$ such that for the points $\{\mathbf{x} : V(\mathbf{x}) \leq V_{\max}\}$, for the discrete-time system we have

$$\begin{aligned} \Delta V(k) &= V(\mathbf{x}(k+1)) - V(\mathbf{x}(k)) \\ &= (\mathbf{x}(k+1) - \mathbf{x}_e)^T \mathbf{P}_l (\mathbf{x}(k+1) - \mathbf{x}_e) - (\mathbf{x}(k) - \mathbf{x}_e)^T (k) \mathbf{P}_l (\mathbf{x}(k) - \mathbf{x}_e) \\ &= (\mathbf{f}(\mathbf{x}(k), u(k)) - \mathbf{x}_e)^T \mathbf{P}_l (\mathbf{f}(\mathbf{x}(k), u(k)) - \mathbf{x}_e) - (\mathbf{x}(k) - \mathbf{x}_e)^T \mathbf{P}_l (\mathbf{x}(k) - \mathbf{x}_e) \\ &< 0, \end{aligned}$$

or for the continuous-time system we obtain

$$\begin{aligned} \dot{V}(\mathbf{x}) &= 2(\mathbf{x} - \mathbf{x}_e)^T \mathbf{P}_l \dot{\mathbf{x}} \\ &= 2(\mathbf{x} - \mathbf{x}_e)^T \mathbf{P}_l \mathbf{f}(\mathbf{x}, u) \\ &< 0. \end{aligned}$$

We next describe the chaotic algorithm. When the given nonlinear system trajectory is in the region $\{\mathbf{x} : V(\mathbf{x}) > V_{\max}\}$, an open-loop control is applied to keep the system in the chaotic attractor until the system reaches a point in the set $\{\mathbf{x} : V(\mathbf{x}) \leq V_{\max}\}$, where the open-loop control is turned off and the feedback saturating control law is applied to force the system trajectory to asymptotically approach the equilibrium state. More formally, we state the chaotic algorithm in the form of the following algorithm:

```

IF  $V(\mathbf{x}) > V_{\max}$ 
  THEN  $u =$  open-loop control to keep the system in the chaotic attractor
ELSE  $u = -\text{sat}(\mathbf{k} \Delta \mathbf{x}) + u_e$ 
END

```

We illustrate the above algorithm when applied to the Hénon map.

◆ Example 11.5

The dynamical system known as the Hénon map is described by the difference equations

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} -1.4x_1^2(k) + x_2(k) + 1 \\ 0.3x_1(k) \end{bmatrix}. \quad (11.99)$$

It was named after the French astronomer M. Hénon, who discovered chaotic behavior of the above system in 1974. The orbit of $\mathbf{x}(0) = [1 \ 0]^T$ of the Hénon map is shown in Figure 11.23. We modify the Hénon map by adding the control input into the first equation so that the map becomes

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} -1.4x_1^2(k) + x_2(k) + 1 + u(k) \\ 0.3x_1(k) \end{bmatrix}. \quad (11.100)$$

There are two fixed states corresponding to $u(k) = u_e = 0$. They are obtained by solving two algebraic equations:

$$\begin{aligned} x_{1,eq} &= -1.4x_{1,eq}^2 + x_{2,eq} + 1, \\ x_{2,eq} &= 0.3x_{1,eq}. \end{aligned}$$

Solving the above gives

$$[0.6314 \ 0.1894]^T \quad \text{and} \quad [-1.1314 \ -0.3394]^T.$$

We select to work with the first of the above fixed states. We linearize the

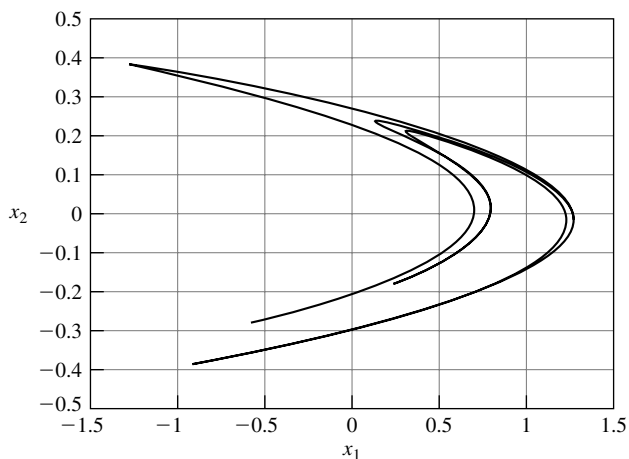


Figure 11.23 An orbit of the Hénon map.

system (11.100) about the first fixed state to obtain

$$\mathbf{A} = \left[\begin{array}{cc} -2.8x_1 & 1 \\ 0.3 & 0 \end{array} \right] \bigg|_{x_1=0.6314} = \left[\begin{array}{cc} -1.7678 & 1 \\ 0.3 & 0 \end{array} \right], \quad \mathbf{b} = \left[\begin{array}{c} 1 \\ 0 \end{array} \right].$$

We then use the discrete LQR design method to find the state feedback. The weight matrices of the associated quadratic performance index are chosen to be $\mathbf{Q} = \mathbf{I}_2$ and $r = 1$. Solving the discrete algebraic Riccati equation yields

$$\mathbf{P} = \left[\begin{array}{cc} 3.9351 & -1.5008 \\ -1.5008 & 1.7974 \end{array} \right] \quad \text{and} \quad \mathbf{k} = [-1.5008 \quad 0.7974].$$

Hence,

$$\mathbf{A}_c = \mathbf{A} - \mathbf{b}\mathbf{k} = \left[\begin{array}{cc} -0.2670 & 0.2026 \\ 0.3000 & 0 \end{array} \right].$$

Solving the discrete Lyapunov equation, $\mathbf{A}_c^T \mathbf{P}_l \mathbf{A}_c - \mathbf{P}_l = -\mathbf{I}_2$, gives

$$\mathbf{P}_l = \left[\begin{array}{cc} 1.1902 & -0.0686 \\ -0.0686 & 1.0489 \end{array} \right].$$

The next step is to calculate V_{\max} corresponding to the largest level curve such that the Lyapunov first difference $\Delta V < 0$ on this level curve for the nonlinear model. We obtain level curves by trial and error as follows. We first evaluate

$$\begin{aligned} V &= (\mathbf{x} - \mathbf{x}_e)^T \mathbf{P}_l (\mathbf{x} - \mathbf{x}_e) \\ &= [x_1 - 0.6314 \quad x_2 - 0.1894] \left[\begin{array}{cc} 1.1902 & -0.0686 \\ -0.0686 & 1.0489 \end{array} \right] \left[\begin{array}{c} x_1 - 0.6314 \\ x_2 - 0.1894 \end{array} \right] \\ &= 1.1902(x_1 - 0.6314)^2 - 0.1372(x_1 - 0.6314)(x_2 - 0.1894) \\ &\quad + 1.0489(x_2 - 0.1894)^2. \end{aligned} \tag{11.101}$$

We next choose an initial state $\mathbf{x}(0)$ such that $V(\mathbf{x}(0)) = V_0$, where V_0 is a constant selected by the designer. Let $x_2(0) = 0.1894$, which is the value of the second component of the equilibrium state of interest. Substituting this value into (11.101) allows us to compute the first component of $\mathbf{x}(0)$ such that $V(\mathbf{x}(0)) = V_0$, where

$$\mathbf{x}(0) = [\sqrt{0.8402 V_0} + 0.6314 \quad 0.1894]^T.$$

Then, a level curve is constructed by solving the differential equation

$$\nabla^T V \dot{\mathbf{x}} = 0,$$

subject to the initial condition $\mathbf{x}(0)$. To satisfy the above equation, we can select

$$\begin{aligned} \dot{x}_1 &= \frac{\partial V}{\partial x_2}, \\ \dot{x}_2 &= -\frac{\partial V}{\partial x_1} \end{aligned} \tag{11.102}$$

After constructing a level curve, we evaluate ΔV at these points and check if

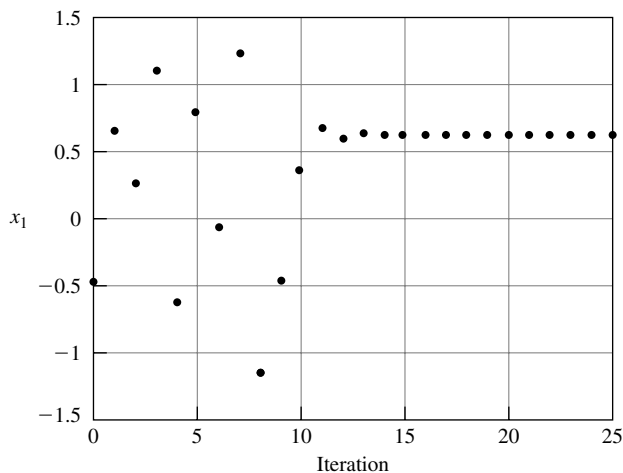


Figure 11.24 A plot of x_1 versus time for the controlled Hénon map of Example 11.5.

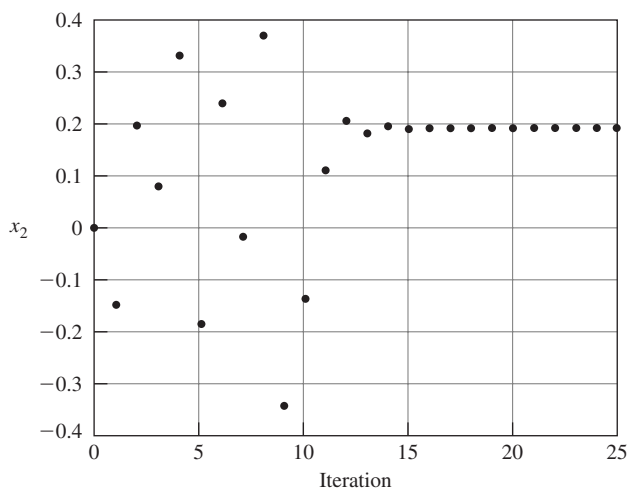


Figure 11.25 A plot of x_2 versus time for the controlled Hénon map of Example 11.5.

$\Delta V < 0$ at every point on the level curve. For the control law

$$u(k) = -0.3 \operatorname{sat}(\mathbf{k}(\mathbf{x}(t) - \mathbf{x}_e)/0.3), \quad (11.103)$$

applied to the nonlinear system (11.100), we found $V_{\max} = 0.0109$. Plots of x_1 and x_2 versus time for the closed-loop system (11.100), (11.103) are shown in Figures 11.24 and 11.25, respectively, for the initial condition $\mathbf{x}(0) = [-0.5 \ 0]^T$. The orbit of $\mathbf{x}(0) = [-0.5 \ 0]^T$ of the controlled Hénon map is shown in Figure 11.26, where the controllable target is also depicted as the set of points inside the ellipse $V(\mathbf{x} - \mathbf{x}_e) = 0.0109$.

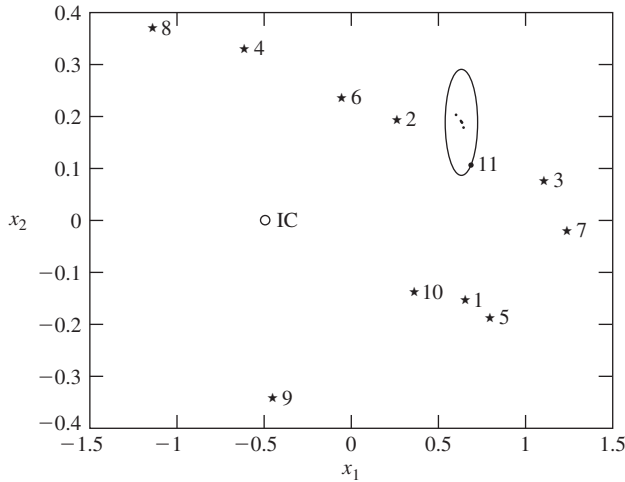


Figure 11.26 An orbit of the controlled Hénon map of Example 11.5, where the initial condition is denoted using “o” and labeled “IC”. The subsequent iterations prior to hitting the controllable target are denoted as “*”. The iterations in the controllable target are denoted using “.”.

◆ Example 11.6

Controlling the Bouncing Ball on a Vibrating Plate

The bouncing ball on a vibrating plate is another example of a chaotic system. In this example, we use the chaotic algorithm to design a control law for the vibrating plate, as shown in Figure 11.27, so that the ball, starting at rest on the plate, can be bounced up to and then maintained bouncing at a prescribed height. The amplitude of the vibrating plate is kept fixed. The control input is the frequency of the plate vibrations. To proceed, we need a mathematical model of the ball-plate system, which we derive next.

Modeling a Bouncing Ball on a Vibrating Plate

Let the displacement of the plate in the vertical direction be

$$y(t) = A \sin(\omega t),$$

where A is the amplitude of the plate displacement, ω is the angular frequency of vibrations and t is time. Then, the velocity of the plate is

$$W(t) = \frac{dy(t)}{dt} = A\omega \cos(\omega t). \quad (11.104)$$

Let U be the velocity of the ball just before impact and let V denote the velocity of the ball just after impact. In our modeling process, we use the principle of conservation of momentum, which states that momentum of the system before and after collisions is constant, that is,

$$\sum (\text{masses} \times \text{velocities})_{\text{before}} = \sum (\text{masses} \times \text{velocities})_{\text{after}}.$$

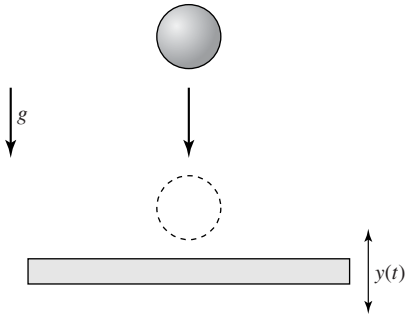


Figure 11.27 Bouncing ball on a vibrating plate.

Applying the conservation of momentum principle to our system consisting of the bouncing ball and the vibrating plate yields

$$m_{\text{ball}} U + m_{\text{plate}} W_{\text{before}} = m_{\text{ball}} V + m_{\text{plate}} W_{\text{after}}. \quad (11.105)$$

To proceed, we define the so-called *coefficient of restitution*

$$e = -\frac{V - W_{\text{after}}}{U - W_{\text{before}}}. \quad (11.106)$$

Performing simple manipulations on (11.106) yields

$$W_{\text{after}} = Ue - W_{\text{before}}e + V. \quad (11.107)$$

Substituting the above into (11.105) and rearranging gives

$$\begin{aligned} m_{\text{ball}} U + m_{\text{plate}} W_{\text{before}} \\ = (m_{\text{ball}} + m_{\text{plate}}) V + em_{\text{plate}} U - em_{\text{plate}} W_{\text{before}}. \end{aligned} \quad (11.108)$$

Dividing both sides of the above by $(m_{\text{ball}} + m_{\text{plate}})$ and rearranging, we get

$$V = \frac{m_{\text{ball}} - em_{\text{plate}}}{m_{\text{ball}} + m_{\text{plate}}} U + \frac{m_{\text{plate}} + em_{\text{plate}}}{m_{\text{ball}} + m_{\text{plate}}} W_{\text{before}}. \quad (11.109)$$

Let

$$M = \frac{m_{\text{ball}}}{m_{\text{plate}}}.$$

Substituting the above into (11.109) gives

$$V = \frac{M - e}{M + 1} U + \frac{e + 1}{M + 1} W_{\text{before}}. \quad (11.110)$$

The ball is acted upon by a constant acceleration g and therefore $V(t) = gt$. Hence, from one bounce to the next, we have

$$t(k+1) = t(k) + \frac{2V(k)}{g}. \quad (11.111)$$

We assume that the amplitude, A , of the plate vibrations is negligible compared

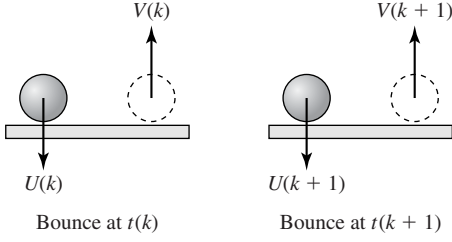


Figure 11.28 Velocities of the bouncing ball from one bounce to the next.

with the height of the bounce. Referring to Figure 11.28, we obtain

$$U(k+1) = -V(k). \quad (11.112)$$

Let

$$a_1 = \frac{e+1}{M+1} \quad \text{and} \quad a_2 = \frac{M-e}{M+1}.$$

Taking into account the above notation and (11.110), we obtain a relation between velocities of the bouncing ball and the vibrating plate at the bounce at a time $t(k+1)$:

$$V(k+1) = a_2 U(k+1) + a_1 W_{\text{before}}(k+1). \quad (11.113)$$

We assume that W_{before} is given by (11.104). Then, using this assumption and (11.112), we represent (11.113) as

$$V(k+1) = -a_2 V(k) + a_1 W(k+1). \quad (11.114)$$

Furthermore, by (11.104) and (11.111),

$$W(k+1) = A\omega \cos(\omega t(k+1)) = A\omega \cos\left(\omega \left(t(k) + \frac{2V(k)}{g}\right)\right).$$

Hence, we can represent (11.114) as

$$V(k+1) = -a_2 V(k) + a_1 A\omega \cos\left(\omega \left(t(k) + \frac{2V(k)}{g}\right)\right). \quad (11.115)$$

Let

$$\phi(k) = \omega t(k) \quad \text{and} \quad \psi(k) = 2\omega V(k)/g, \quad (11.116)$$

where $\phi(k)$ is the current phase angle of the plate and $\psi(k)$ is the phase angle change between the current bounce of the ball and the next. Then, we can write

$$\begin{aligned} \phi(k+1) &= \phi(k) + \psi(k), \\ \psi(k+1) &= -a_2 \psi(k) + \hat{a}_1 \omega^2 \cos(\phi(k) + \psi(k)), \end{aligned} \quad (11.117)$$

where $\hat{a}_1 = 2Aa_1/g$. The above dynamical system is known as *the ball map*. In our simulations, we use the same values of the ball map parameters as those used by Vincent [290], which are given in Table 11.1. Using these parameters, we compute $a_1 = 1.73333$, $\hat{a}_1 = 0.004594 \text{ sec}^2$, and $a_2 = -0.733333$.

Table 11.1 Simulation Parameters of the Ball Map

M	1/26
e	0.8
A	0.013 m

There are many possible periodic equilibrium solutions to the ball map (11.117). For instance, the ball can bounce to a fixed height at every n cycles of the plate, or it can bounce to m different heights at nm cycles of the plate. Here we are concerned with the case when $m = 1$. There are bands of frequencies of the vibrating plate that will not support periodic motion for $m = 1$. Examples of such frequency bands, as listed by Vincent [290], are

$$\begin{aligned} 32.3966 < \omega < 33.0782, \\ 36.4575 < \omega < 38.1954, \\ 40.5278 < \omega < 42.7038, \\ 44.4249 < \omega < 46.7796. \end{aligned}$$

A frequency of the plate vibration from one of the above ranges will result in either a chaotic motion of the ball or a chaotic transient to a periodic motion with $m > 1$. In our example, we use $\omega = 45$ rad/sec to generate chaotic motion.

The ball map (11.117) may produce negative values for ψ corresponding to negative bounce heights, which is physically impossible. To overcome this problem, whenever the map produces a negative value for ψ , this value is multiplied by -1 to simulate the fact that the ball cannot go through the plate. In Figure 11.29, an orbit of the modified ball map for the initial conditions $\phi(0) = 0$ and $\psi(0) = 9.1743$ rads and $\omega = 45$ rad/sec is shown. Note that the values of ψ are divided by 2π before plotting. The chaotic attractor is located between the upper and lower bounds as

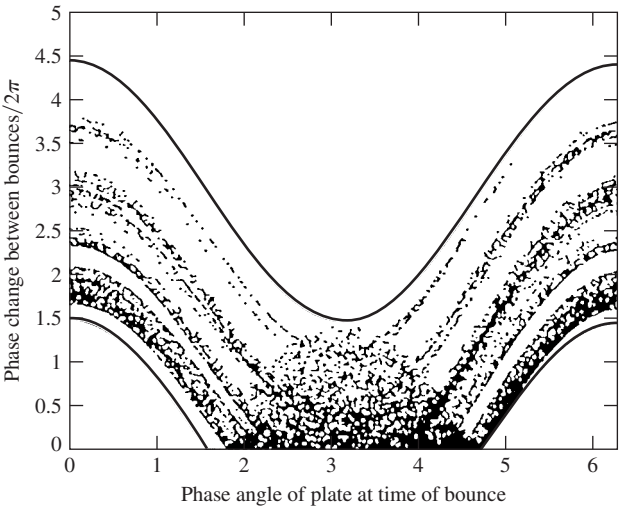


Figure 11.29 An orbit of the ball map.

depicted in the figure. The upper bound is given by

$$\psi_{\text{upper}} = 2\hat{a}_1\omega^2 + \hat{a}_1\omega^2 \cos(\phi(k) + \psi(k)),$$

while the lower bound is

$$\psi_{\text{lower}} = \hat{a}_1\omega^2 \cos(\phi(k) + \psi(k)).$$

When ψ_{lower} goes negative, it is set equal to zero.

Controller Construction for a Bouncing Ball on a Vibrating Plate

To proceed with the controller design, we first find fixed points of the ball map. They are found by solving the equations,

$$\begin{aligned}\bar{\phi} &= \bar{\phi} + \bar{\psi}, \\ \bar{\psi} &= -a_2\bar{\psi} + \hat{a}_1\omega^2 \cos(\bar{\phi} + \bar{\psi}).\end{aligned}\tag{11.118}$$

To obtain fixed points, we evaluate the right-hand side of the first of the above equations modulo $2n\pi$, where n is the number of cycles of the plate between bounces. The result of the above calculations is

$$\begin{aligned}\bar{\psi} &= 2n\pi, \\ \cos(\bar{\phi}) &= \frac{2n\pi(1 + a_2)}{\hat{a}_1\omega^2}.\end{aligned}\tag{11.119}$$

Linearizing the ball map about (11.119) gives

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix},\tag{11.120}$$

where

$$\begin{aligned}x_1(k) &= \phi(k) - \bar{\phi}, \\ x_2(k) &= \psi(k) - \bar{\psi}, \\ a_{21} &= -\hat{a}_1\omega^2 \sin(\bar{\phi}), \\ a_{22} &= -a_2 + a_{21}.\end{aligned}$$

We next construct a control law for selecting frequencies of the vibrating plate in such a way that the steady-state constant height of the bouncing ball is reached. Vincent [290] showed that for $n=1$ the range of frequencies for which there exists an asymptotically stable equilibrium is [19.0977 rad/sec, 28.9505 rad/sec]. The corresponding range of the bounce heights is [0.0578 m, 0.1327 m]. Thus, for example the height, h , of 0.1 m is in the above range and from the energy conservation law,

$$\frac{1}{2}m_{\text{ball}}V^2 = m_{\text{ball}}gh,$$

we obtain

$$V = 1.3994 \text{ m/sec.}$$

Using the above, (11.116), and (11.119), we get

$$\bar{\psi} = 2\pi = \frac{2\omega V}{g} \quad \text{and} \quad \cos(\bar{\phi}) = \frac{2\pi(1 + a_2)}{\hat{a}_1\omega^2},$$

from which we compute

$$\omega = 22 \text{ rad/sec} \quad \text{and} \quad \bar{\phi} = 41.1^\circ = 0.7173 \text{ rad.}$$

The linearized ball map about the above equilibrium state is

$$\mathbf{x}(k+1) = \begin{bmatrix} 1 & 1 \\ -1.4617 & -0.7283 \end{bmatrix} \mathbf{x}(k),$$

which has its poles in the open unit circle. Solving the discrete Lyapunov equation, $\mathbf{A}^T \mathbf{P} \mathbf{A} - \mathbf{P} = -\mathbf{I}_2$, gives

$$\mathbf{P} = \begin{bmatrix} 8.6414 & 4.6978 \\ 4.6978 & 5.9598 \end{bmatrix}.$$

Employing the same method as in the previous example, we obtain $V_{\max} = 5.006$. The controllable target is

$$\left\{ \begin{bmatrix} \phi \\ \psi \end{bmatrix} : V(\phi - \bar{\phi}, \psi - \bar{\psi}) = [\phi - 0.7173 \quad \psi - 2\pi] \mathbf{P} \begin{bmatrix} \phi - 0.7173 \\ \psi - 2\pi \end{bmatrix}^T \leq V_{\max} \right\}.$$

The controllable target is depicted in Figure 11.30. We use $\omega = 45 \text{ rad/sec}$ when the system trajectory is outside of the controllable target, and we apply $\omega = 22 \text{ rad/sec}$ when the system trajectory hits the controllable target; that is, we use the following algorithm:

IF $V(\phi - \bar{\phi}, \psi - \bar{\psi}) > V_{\max}$
 THEN $\omega = 45 \text{ rad/sec}$
 ELSE $\omega = 22 \text{ rad/sec}$
 END

An orbit of $\phi(0) = 0$, $\psi(0) = 9.1743$ of the controlled modified ball map using

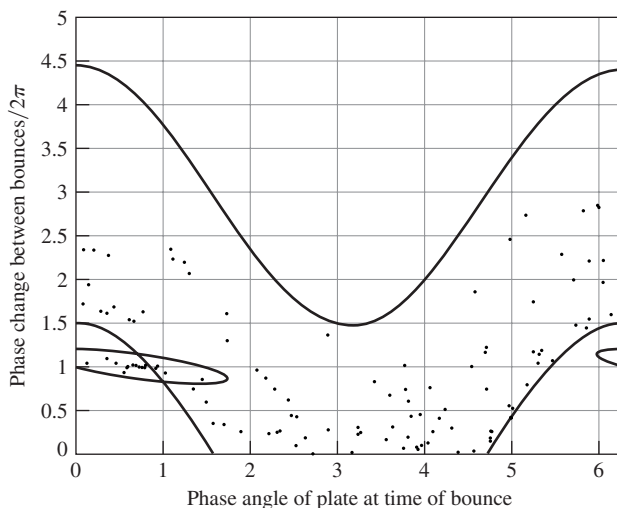


Figure 11.30 An orbit of the controlled modified ball map.

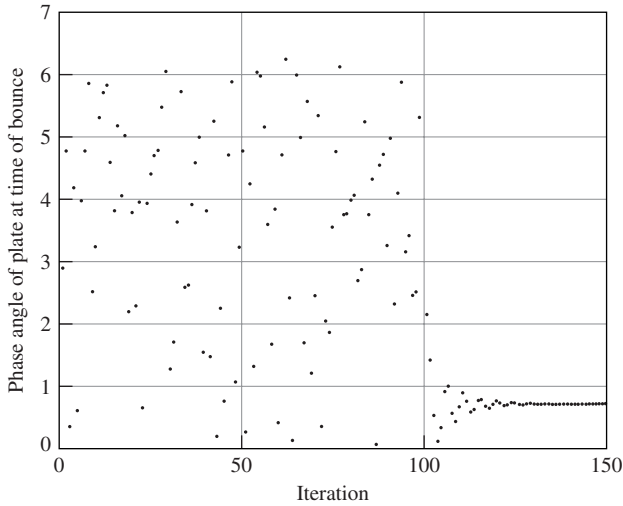


Figure 11.31 A plot of $\phi(k)$ versus time of the controlled modified ball map.

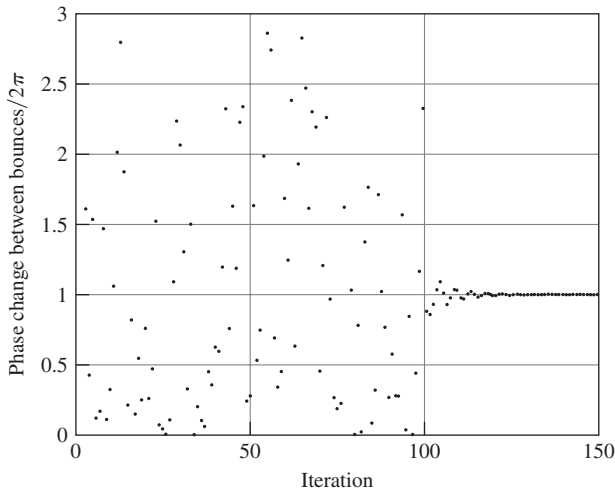


Figure 11.32 A plot of $\psi(k)$ versus time of the controlled modified ball map.

the above algorithm is shown in Figure 11.30. A plot of $\phi(k)$ versus time is shown in Figure 11.31. A plot of $\psi(k)$ versus time is shown in Figure 11.32.

Notes

Sandefur [253] is a lucid elementary introduction to chaotic discrete-time dynamical systems. Peitgen, Jürgens, and Saupe [227] and Alligood, Sauer, and Yorke [7] are superb introductions to chaos and dynamical systems. Papers by May [198] and Feigenbaum [81] are recommended for their originality and clarity. Introductory texts to nonlinear dynamics and chaos for scientists

and engineers are the books by Moon [209] and by Thompson and Stewart [281]. Numerical algorithms for chaotic systems can be found in the book by Parker and Chua [222]. Fractals are covered in Mandelbrot [195], Peitgen and Saupe [228], Devaney and Keen [63], and Bunde and Havlin [38]. Barnsley [21] and Peitgen and Saupe [228] are recommended for their coverage of iterated function systems. A reader who is interested in a mathematical approach to chaotic behavior of dynamical systems should consult the book by Arrowsmith and Place [13]. Peterson [232] gives a remarkable account of the development of models of celestial mechanics leading to the discovery of chaos in the solar system. The use of fractals in education, art, music, fashion, chess, and medicine is described in an edited volume by Pickover [233]. Other, more recent papers on controlling chaotic systems, in addition to the ones already referenced in the previous section, are references 305, 308, and 317. Applications of chaos in information processing are discussed in the papers by Dmitriev, Panas, and Starkov [65] and by Andreyev et al. [10]. Aihara, Takabe, and Toyoda [5] and Adachi and Aihara [3] discuss the presence of chaos in biological neurons and propose mathematical models of chaotic neural networks. Applications of chaotic neural networks in the complex memory search for patterns with a sought feature are discussed by Nara, Davis, and Totsuji [214] and by Chan and Žak [46].

Edward Lorenz was a meteorologist at MIT. Lorenz's discovery of chaotic behavior of his mathematical model of weather is a fascinating one. As reported by Gleick [100, Chapter 1] and Alligood et al. [7, Chapter 9], to speed up his computer output which was simulating a weather model, Lorenz used just three significant digits to print out the results of simulations although the calculations were performed using several more digits of accuracy. One day in the winter of 1961, Lorenz came across a simulation that he decided to repeat. He typed in an initial condition for the new simulation from the printout using the three digit format. He then performed the simulation expecting the same behavior as before of the repeated simulation. Instead, "he saw something unexpected, something that planted a seed for a new science" [100, p. 16]. The resulting simulation was totally different from the previous one. The difference between the two solutions was the result of the Lorenz's weather model sensitive dependence on initial conditions. In the computer memory, according to Gleick [100], the initial condition for the first run was 0.506127, while to save space, Lorenz entered 0.506 for the second run. This seemingly inconsequential difference caused the two simulations to flow apart until there was no resemblance at all between them. Later, Lorenz greatly simplified his weather model to just three nonlinear equations that are given in Exercise 11.2. The phase portrait of these famous equations resemble butterfly's wings, which became a symbol for the early activity in chaotic systems.

EXERCISES

- 11.1** Design a 1-D map f with a stable limit cycle of period 5. Simulate the dynamical behavior of the map.

Hint A method of constructing 1-D maps with cycles of any period is described in Andreyev et al. [11].

- 11.2** Simulate the system of differential equations

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} a(-x_1 + x_2) \\ bx_1 - x_2 - x_1x_3 \\ -cx_3 + x_1x_2 \end{bmatrix}, \quad \mathbf{x}(0) = \begin{bmatrix} -8 \\ -8 \\ 24 \end{bmatrix}, \quad (11.121)$$

where $a = 10$, $b = 28$, and $c = 8/3$. The set of differential equations (11.121) is commonly referred to as the *Lorenz equations*—after the American meteorologist E. N. Lorenz, who first published these equations in 1963. They model fluid flows in the atmosphere. Plot x_2 versus x_1 , x_3 versus x_1 , and x_3 versus x_2 . Use MATLAB's command `plot3` to obtain a 3-D plot of the solution of the Lorenz equations.

11.3 Simulate the system of the difference equations

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 1 - x_2(k) + |x_1(k)| \\ x_1(k) \end{bmatrix}, \quad \mathbf{x}(0) = \begin{bmatrix} -0.1 \\ 0 \end{bmatrix}. \quad (11.122)$$

The solution “fills” a region that looks like a “gingerbread man.” For this reason, Devaney [228, p. 149] refers to this example of a chaotic set in the plane as a *chaotic gingerbreadman*.

11.4 Simulate the system of differential equations

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} -x_2 - x_3 \\ x_1 + ax_2 \\ b + x_3(x_1 - c) \end{bmatrix}, \quad (11.123)$$

where $a = 0.398$, $b = 2$, and $c = 4$. The set of differential equations (11.123) is referred to as the *Rössler equations*—after O. E. Rössler, who reported on their chaotic behavior in 1976. Rössler called the chaotic attractor of the above system the “spiral chaos.” This type of chaotic attractor is also called a *folded band*. Select an initial condition, integrate the equations, then plot x_2 versus x_1 , x_3 versus x_1 , and x_3 versus x_2 . Use MATLAB's command `plot3` to obtain a 3-D plot of the folded band attractor.

11.5 For the dynamical system

$$x(k+1) = x(k) + r(0.8 - \text{atan}(x(k))), \quad (11.124)$$

plot a bifurcation diagram for the parameter r in the interval $[3, 12]$. The map (11.124) was analyzed in Hui and Žak [134].

11.6 For the dynamical system

$$x(k+1) = x(k) + r \left(0.8 - \frac{1 - e^{-x(k)}}{1 + e^{-x(k)}} \right), \quad (11.125)$$

plot a bifurcation diagram for the parameter r in the interval $[10, 55]$.

11.7 For the dynamical system

$$x(k+1) = x(k) + r \left(0.8 - \frac{1}{1 + e^{-x(k)}} \right), \quad (11.126)$$

plot a bifurcation diagram for the parameter r in the interval $[10, 45]$.

11.8 Write a MATLAB program of the following “chaos game”: Select three points that are vertices of an equilateral triangle. Choose at random one of the vertices as a starting point. Then, choose at random another vertex and move halfway to it from the starting point. Mark this halfway point with a dot. This halfway point is your new starting point. Choose at random one of the vertices and move halfway to the chosen vertex from the new starting point. Mark the new halfway point with the dot. Repeat the process. The outcome of this chaos game is a fractal called the Sierpiński gasket. Generate a plot of

the Sierpiński gasket. Other methods of constructing the Sierpiński gasket are described in Peitgen et al. [227] and in Exercise 11.9. The term “chaos game” was coined by M. F. Barnsley [21].

- 11.9** The result of the chaos game described in Exercise 11.8 is the Sierpiński gasket, which can also be constructed using the following procedure. Start with an equilateral triangle. Divide it into four equal subtriangles (quarters) and throw away the middle one—that is, the upside-down triangle. This step leaves three equilateral triangles. In the second step, divide each of the remaining three triangles into four equal quarter-triangles and throw away the middle, upside-down, triangle in each of the three triangles. This step leaves nine triangles. Repeat the process. The remaining area, as the number of iterations tends to infinity, is called the Sierpiński gasket.
- (a) Show that the area of the Sierpiński gasket is zero.
- (b) Determine the fractal dimension of the Sierpiński gasket.
- 11.10** In this exercise, we concern ourselves with another fractal, the Sierpiński carpet. Start with a square. Divide it into nine equal subsquares. Remove the central square. Repeat the procedure. The remaining area as the number of iterations goes to infinity is called the Sierpiński carpet.
- (a) Show that the area of the Sierpiński carpet is zero.
- (b) Determine the fractal dimension of the Sierpiński carpet.
- 11.11** For the double pendulum system shown in Figure 11.33, show that the Lagrangian has the form

$$L = \frac{1}{2}m_1 l_1^2 \omega_1^2 + \frac{1}{2}m_2 (l_1^2 \omega_1^2 + l_2^2 \omega_2^2 + 2l_1 l_2 \cos(\Delta\theta) \omega_1 \omega_2) \\ + m_1 g l_1 \cos(\theta_1) + m_2 g (l_1 \cos(\theta_1) + l_2 \cos(\theta_2)),$$

where $\omega_i = \dot{\theta}_i$, $i = 1, 2$, are the angular velocities of the pendulum arm, $\Delta\theta = \theta_1 - \theta_2$, and m_i and l_i are the masses and arm lengths, respectively. Assume that the input torques, T_1 and T_2 are applied at each joint. The parameter values are given in Table 11.2. Write the Lagrange equations of motion for the system and represent them in state-space format, where the state variables can be defined as $x_1 = \theta_1$, $x_2 = \dot{\theta}_1$, $x_3 = \theta_2$,

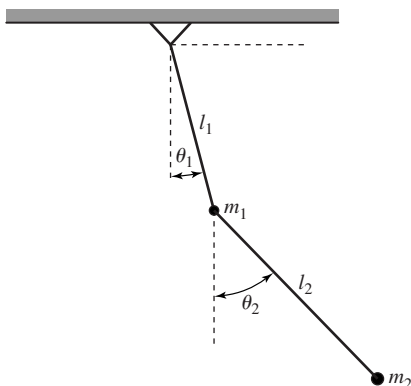


Figure 11.33 Double pendulum of Exercise 11.11.

Table 11.2 Parameter Numerical Values of the Double Pendulum System of Exercise 11.11

$m_1 = m_2$	1 kg
$l_1 = l_2$	0.5 m
g	9.8 m/sec ²

and $x_4 = \dot{\theta}_2$. Set the input torques to zero and simulate the system behavior starting from the initial state $\mathbf{x}(0) = [0 \ 0 \ \pi \ .001]^T$. Verify that one of the equilibrium states corresponding to zero input torques is $\mathbf{x}_e = [0 \ 0 \ \pi \ 0]^T$. Linearize the state-space model about this equilibrium point. Then, use the LQR approach to compute the state feedback matrix \mathbf{K} for the weight matrices

$$\mathbf{Q} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 10 & 0 \\ 0 & 0 & 0 & 10 \end{bmatrix} \quad \text{and} \quad \mathbf{R} = \begin{bmatrix} 1 & 0 \\ 0 & 10 \end{bmatrix}.$$

Solve the Lyapunov equation

$$(\mathbf{A} - \mathbf{BK})^T \mathbf{P}_l + \mathbf{P}_l (\mathbf{A} - \mathbf{BK}) = -\mathbf{Q}.$$

Suppose that the bound on the input torques are $|T_1| \leq 5$ and $|T_2| \leq 2$. Use the method of Example 11.5 to find V_{\max} and apply the chaotic control law to stabilize the closed-loop system with respect to $[0 \ 0 \ \pi \ 0]^T$. Plot x_1 , x_2 , and $V = (\mathbf{x} - \mathbf{x}_e)^T \mathbf{P}_l (\mathbf{x} - \mathbf{x}_e)$ versus time.



APPENDIX

Math Review

Do not worry about your difficulties in Mathematics. I can assure you mine are still greater.

—Albert Einstein

The last word, when all is heard: Fear God and keep his commandments, for this is man's all; because God will bring to judgment every work, with all its hidden qualities, whether good or bad.

—Ecclesiastes 12:13–14

The appendix reviews mathematical techniques that are used in the book. The appendix also serves as a glossary of symbols, notation, and concepts employed in the text. For more information on the topics covered here, the reader should consult references supplied at the end of the appendix.

In an attempt to keep the presentation as clear as possible, we employ the theorem-proof style whenever appropriate. Major results of independent interest are encapsulated as *theorems*. A *lemma* is used to denote an auxiliary result that is used as a stepping stone toward a theorem. We use the term *proposition* as a way to announce a minor theorem. A *corollary* is a direct consequence of a theorem or proposition. A *hypothesis* is a statement that is taken for further reasoning in a proof. For more information on the subject of writing for mathematical or technical sciences, we recommend a very clear and useful book by Higham [122].

A.1 Notation and Methods of Proof

We often use the term *set*. The great German mathematician G. Cantor (1845–1918), the creator of the theory of sets, defines a set as any collection of definite, distinguishable objects of our thought, to be conceived as a whole [193, p. 1]. In the definition the objects are the elements or members of the set. Usually, we use capital letters to denote sets and use lowercase letters to denote elements of the set. Thus, if X is a set, then we write $x \in X$ to mean that x is an element of X . When an object x is not an element of a set X , we write $x \notin X$. Often, we use the curly bracket notation for sets. For example, $\{x : x < 7\}$ denotes the set of all x less than 7.

The colon, :, following x is read “such that.” The following sets are used in the book:

$\mathbb{N} = \{1, 2, 3, \dots\}$, the set of positive integers,

$\mathbb{Z} = \{0, 1, -1, 2, -2, \dots\}$, the set of integers,

\mathbb{R} , the set of real numbers,

\mathbb{C} , the set of complex numbers.

Sometimes, we take a fixed set and then carry out manipulations with reference to the fixed set. We call this fixed set the *universe of discourse*, or just the *universe*, for short.

Throughout the text, we work with statements that are combinations of individual statements. It is necessary to distinguish between sentences and the statements. A *sentence* is a part of a language. Sentences are used to make different statements. In classical logic, statements are either true or false. Consider the following two statements:

$$A : x > 7;$$

$$B : x^2 > 49.$$

We can combine the above two statements into one statement, called a *conditional*, that has the form “IF A , THEN B .” Thus, a conditional is a compound statement obtained by placing the word “IF” before the first statement and inserting the word “THEN” before the second statement. We use the symbol “ \Rightarrow ” to represent the conditional “IF A , THEN B ” as $A \Rightarrow B$. The statement $A \Rightarrow B$ also reads as “ A implies B ,” or “ A only if B ,” or “ A is sufficient for B ,” or “ B is necessary for A .” Statements A and B may be either true or false. The relationship between the truth or falsity of A and B and the conditional $A \Rightarrow B$ can be illustrated by means of a diagram called a *truth table* shown in Table A.1. In the table, T stands for “true” and F stands for “false.” It is intuitively clear that if A is true, then B must also be true for the statement $A \Rightarrow B$ to be true. If A is not true, then the sentence $A \Rightarrow B$ does not have an obvious meaning in everyday language. We can interpret $A \Rightarrow B$ to mean that we cannot have A true and B not true. One can check that the truth values of the statements $A \Rightarrow B$ and “not (A and not B)” are the same. We can also check that the truth values of $A \Rightarrow B$ and “not $B \Rightarrow$ not A ” are the same. The above is the basis for three methods, inference processes, of proving statements of the form $A \Rightarrow B$. *Inference* is a process by which one statement is arrived at, using reasoning, on the basis of one or more other statements accepted as a starting point. *Reasoning* is a special kind of thinking in which inference takes place using the laws of thought. By a *proof* we mean the process of establishing the truth of a statement by arriving at it from other statements using principles of reasoning. Most of the time we are concerned with the statements of the form $A \Rightarrow B$ and use

Table A.1 Truth Table for Conditional

A	B	$A \Rightarrow B$
F	F	T
F	T	T
T	F	F
T	T	T

the following types of proofs of such statements:

1. The direct method
2. The proof by contraposition
3. The proof by contradiction

In the first case, the direct method, one starts with A and then tries to deduce B . In the proof by contraposition, one starts with “not B ” and then deduces a chain of consequences to end with “not A .” This method of proving is based on the fact that the truth values of $A \Rightarrow B$ and “not $B \Rightarrow$ not A ” are the same. The proof by contradiction is based on the fact that the truth values of $A \Rightarrow B$ and “not (A and not B)” are the same. In this method, one starts with “ A and not B ” and tries to arrive at a contradiction, which is a statement that is always false.

It is possible to combine two conditionals to form a *biconditional*, which is a statement of the form “ A if and only if B .” Occasionally, we use the symbol “ \Leftrightarrow ” to write $A \Leftrightarrow B$ to mean “ A if and only if B .” Another way of expressing $A \Leftrightarrow B$ is “ A is equivalent to B ”, or “ A is necessary and sufficient for B .”

A.2 Vectors

We define an n -dimensional column vector to be an n -tuple of numbers:

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad (\text{A.1})$$

where n is a positive integer. Thus a vector can be viewed as an ordered set of scalars. We use lowercase bold letters—for example, \mathbf{x} —to denote such n -tuples. We call the numbers x_1, x_2, \dots, x_n the *coordinates* or *components* of the vector \mathbf{x} . The number of components of a vector is its dimension. When referring to a particular component of a vector, we use its component—for example, x_i —or the subscripted name of the vector—for example $(\mathbf{x})_i$ —to denote x_i . The set of all n -tuples, that is, n -vectors, with real coefficients is denoted \mathbb{R}^n and called the real n -space. The transpose of the vector \mathbf{x} , denoted \mathbf{x}^T , is an n -dimensional row vector

$$\mathbf{x}^T = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}^T = [x_1 \quad x_2 \quad \cdots \quad x_n]. \quad (\text{A.2})$$

The sum of two n -vectors \mathbf{x} and \mathbf{y} is an n -vector, denoted $\mathbf{x} + \mathbf{y}$, whose i th component is $x_i + y_i$. The following rules are satisfied:

1. $(\mathbf{x} + \mathbf{y}) + \mathbf{z} = \mathbf{x} + (\mathbf{y} + \mathbf{z})$, that is, vector addition is associative.
2. $\mathbf{x} + \mathbf{y} = \mathbf{y} + \mathbf{x}$, which means that vector addition is commutative.
3. Let $\mathbf{0} = [0 \ 0 \ \cdots \ 0]^T$ denote the zero vector. Then,

$$\mathbf{0} + \mathbf{x} = \mathbf{x} + \mathbf{0} = \mathbf{x}.$$

4. Let $\mathbf{x} = [x_1 \ x_2 \ \cdots \ x_n]^T$ and let $-\mathbf{x} = [-x_1 \ -x_2 \ \cdots \ -x_n]^T$. Then,

$$\mathbf{x} + (-\mathbf{x}) = \mathbf{0}.$$

The product of an n -vector \mathbf{x} by a scalar c is an n th vector, denoted $c\mathbf{x}$, whose i th component is cx_i .

For two vectors \mathbf{x} and \mathbf{y} from \mathbb{R}^n , we define their *scalar product*, also called the *inner product*, to be

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = x_1 y_1 + x_2 y_2 + \cdots + x_n y_n. \quad (\text{A.3})$$

The scalar product of vectors from \mathbb{R}^n is commutative, that is,

$$\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle,$$

and distributivity over vector addition holds:

$$\langle \mathbf{x}, \mathbf{y} + \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{z} \rangle.$$

Two vectors \mathbf{x} and \mathbf{y} are said to be *orthogonal* if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$.

A *vector norm*, denoted $\|\cdot\|$, satisfies the following properties:

VN 1. $\|\mathbf{x}\| > 0$ for any nonzero vector.

VN 2. For any scalar c , we have $\|c\mathbf{x}\| = |c|\|\mathbf{x}\|$.

VN 3. For any two vectors \mathbf{x} and \mathbf{y} , the *triangle inequality* holds, that is, $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$.

The *Euclidean norm*, or the *2-norm*, of a vector $\mathbf{x} \in \mathbb{R}^n$, denoted $\|\mathbf{x}\|_2$, is defined as

$$\|\mathbf{x}\|_2 = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{x_1^2 + x_2^2 + \cdots + x_n^2}. \quad (\text{A.4})$$

Note that $\|\mathbf{x}\|_2^2 = \langle \mathbf{x}, \mathbf{x} \rangle = \mathbf{x}^T \mathbf{x}$. The Euclidean norm is a special class of norms called the *Hölder* or *p-norms* defined by

$$\|\mathbf{x}\|_p = (|x_1|^p + \cdots + |x_n|^p)^{1/p}, \quad p \geq 1. \quad (\text{A.5})$$

Other often-used norms that are special cases of the above class of norms are the *1-norm* and the *infinity norm*, defined respectively as

$$\|\mathbf{x}\|_1 = |x_1| + \cdots + |x_n| \quad (\text{A.6})$$

and

$$\|\mathbf{x}\|_\infty = \max_i |x_i|. \quad (\text{A.7})$$

Most of the time, we use the Euclidean norm for vectors. So for the sake of brevity, we write $\|\mathbf{x}\|$ to mean $\|\mathbf{x}\|_2$.

Theorem A.1 (Pythagoras) If \mathbf{x} and \mathbf{y} are orthogonal, then

$$\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2.$$

Proof We have

$$\begin{aligned} \|\mathbf{x} + \mathbf{y}\|^2 &= (\mathbf{x} + \mathbf{y})^T (\mathbf{x} + \mathbf{y}) \\ &= \|\mathbf{x}\|^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle + \|\mathbf{y}\|^2 \\ &= \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2, \end{aligned}$$

because by assumption $\langle \mathbf{x}, \mathbf{y} \rangle = 0$. The proof of the theorem is complete.

We now introduce the notion of *projection*. Let \mathbf{x}, \mathbf{y} be two vectors, where $\mathbf{y} \neq \mathbf{0}$. Let $\mathbf{z} = c\mathbf{y}$, $c \in \mathbb{R}$, be the vector such that the vector $\mathbf{x} - \mathbf{z}$ is orthogonal to \mathbf{y} —see Figure A.1 for

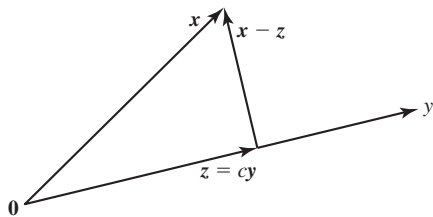


Figure A.1 Computing the projection of \mathbf{x} along \mathbf{y} .

an illustration of the above construction. Thus,

$$(\mathbf{x} - \mathbf{z})^T \mathbf{y} = 0. \quad (\text{A.8})$$

Substituting $\mathbf{z} = c\mathbf{y}$ into (A.8) yields

$$(\mathbf{x} - \mathbf{z})^T \mathbf{y} = (\mathbf{x} - c\mathbf{y})^T \mathbf{y} = 0. \quad (\text{A.9})$$

Hence,

$$c = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{y}\|^2}. \quad (\text{A.10})$$

The number c given by (A.10) is called the component of \mathbf{x} along \mathbf{y} . The projection of \mathbf{x} along \mathbf{y} is the vector

$$\mathbf{z} = c\mathbf{y} = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{y}\|^2} \mathbf{y}. \quad (\text{A.11})$$

Theorem A.2 (Cauchy-Schwarz Inequality) Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. Then,

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|. \quad (\text{A.12})$$

Proof If $\mathbf{y} = \mathbf{0}$, then both sides of (A.12) are equal to 0, and the assertion holds. Suppose that $\mathbf{y} \neq \mathbf{0}$, and let c denote the component of \mathbf{x} along \mathbf{y} . Hence, c is given by (A.10). We then write

$$\mathbf{x} = \mathbf{x} - \mathbf{z} + \mathbf{z} = \mathbf{x} - c\mathbf{y} + c\mathbf{y}. \quad (\text{A.13})$$

For an illustration of (A.13), consult Figure A.1. Applying the Pythagoras theorem to (A.13), we obtain

$$\|\mathbf{x}\|^2 = \|\mathbf{x} - c\mathbf{y}\|^2 + \|c\mathbf{y}\|^2 = \|\mathbf{x} - c\mathbf{y}\|^2 + c^2 \|\mathbf{y}\|^2. \quad (\text{A.14})$$

Hence,

$$c^2 \|\mathbf{y}\|^2 \leq \|\mathbf{x}\|^2. \quad (\text{A.15})$$

On the other hand, by (A.10),

$$c^2 \|\mathbf{y}\|^2 = \frac{\langle \mathbf{x}, \mathbf{y} \rangle^2}{\|\mathbf{y}\|^4} \|\mathbf{y}\|^2 = \frac{\langle \mathbf{x}, \mathbf{y} \rangle^2}{\|\mathbf{y}\|^2}. \quad (\text{A.16})$$

Combining (A.15) and (A.16) yields

$$\frac{\langle \mathbf{x}, \mathbf{y} \rangle^2}{\|\mathbf{y}\|^2} \leq \|\mathbf{x}\|^2. \quad (\text{A.17})$$

Multiplying both sides of (A.17) by $\|\mathbf{y}\|^2$ and taking the square root, concludes the proof.

Theorem A.3 (Triangle Inequality) Let \mathbf{x}, \mathbf{y} be two vectors in \mathbb{R}^n . Then,

$$\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|. \quad (\text{A.18})$$

Proof Both sides of (A.18) are nonnegative. Hence, to prove that (A.18) holds it suffices to show that the squares of the both sides satisfy the desired inequality,

$$(\mathbf{x} + \mathbf{y})^T(\mathbf{x} + \mathbf{y}) \leq (\|\mathbf{x}\| + \|\mathbf{y}\|)^2. \quad (\text{A.19})$$

Expanding $(\mathbf{x} + \mathbf{y})^T(\mathbf{x} + \mathbf{y})$ and applying the Schwarz inequality yields

$$\begin{aligned} (\mathbf{x} + \mathbf{y})^T(\mathbf{x} + \mathbf{y}) &= \|\mathbf{x}\|^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle + \|\mathbf{y}\|^2 \\ &\leq \|\mathbf{x}\|^2 + 2\|\mathbf{x}\|\|\mathbf{y}\| + \|\mathbf{y}\|^2 \\ &= (\|\mathbf{x}\| + \|\mathbf{y}\|)^2, \end{aligned} \quad (\text{A.20})$$

which completes the proof of the theorem.

Theorem A.4 Let \mathbf{x}, \mathbf{y} be two vectors in \mathbb{R}^n . Then,

$$|\|\mathbf{x}\| - \|\mathbf{y}\|| \leq \|\mathbf{x} - \mathbf{y}\|.$$

Proof We have

$$\|\mathbf{x}\| = \|(\mathbf{x} - \mathbf{y}) + \mathbf{y}\| \leq \|\mathbf{x} - \mathbf{y}\| + \|\mathbf{y}\|.$$

Hence,

$$\|\mathbf{x}\| - \|\mathbf{y}\| \leq \|\mathbf{x} - \mathbf{y}\|. \quad (\text{A.21})$$

The above implies that

$$\|\mathbf{x} - \mathbf{y}\| = \|\mathbf{y} - \mathbf{x}\| \geq \|\mathbf{y}\| - \|\mathbf{x}\| = -(\|\mathbf{x}\| - \|\mathbf{y}\|). \quad (\text{A.22})$$

Combining (A.21) and (A.22) yields

$$|\|\mathbf{x}\| - \|\mathbf{y}\|| \leq \|\mathbf{x} - \mathbf{y}\|,$$

thus completing the proof.

A.3 Matrices and Determinants

Let m, n be two integers such that $m, n \geq 1$. An array of numbers

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

is called a *matrix*. Thus a matrix is a set of scalars subject to two orderings. We say that it is an m by n matrix, or $m \times n$ matrix. This matrix has m rows and n columns. If $m = n$, then we say that it is a square matrix. If the elements of an $m \times n$ matrix \mathbf{A} are real numbers, then we write $\mathbf{A} \in \mathbb{R}^{m \times n}$. If the elements of the matrix \mathbf{A} are complex numbers, then we write $\mathbf{A} \in \mathbb{C}^{m \times n}$. We define the zero matrix, denoted \mathbf{O} , to be the matrix such that $a_{ij} = 0$ for all i, j . We define the

$n \times n$ identity matrix, denoted \mathbf{I}_n , as the matrix having diagonal components all equal to 1 and all other components equal to 0.

Given an $m \times n$ matrix \mathbf{A} . The $n \times m$ matrix \mathbf{B} such that $b_{ji} = a_{ij}$ is called the transpose of \mathbf{A} and is denoted \mathbf{A}^T . Thus, the transpose of a matrix is obtained by changing rows into columns and vice versa. A matrix \mathbf{A} that is equal to its transpose—that is, $\mathbf{A} = \mathbf{A}^T$ —is called symmetric.

We now define the sum of two matrices. For two matrices \mathbf{A} and \mathbf{B} of the same size, we define $\mathbf{A} + \mathbf{B}$ to be the matrix whose elements are $a_{ij} + b_{ij}$; that is, we add matrices of the same size componentwise.

Let c be a number and let \mathbf{A} be a matrix. We define $c\mathbf{A}$ to be the matrix whose components are ca_{ij} . Thus, to obtain $c\mathbf{A}$, we multiply each component of \mathbf{A} by c .

We now define the product of two matrices. Let \mathbf{A} be an $m \times n$ matrix and let \mathbf{B} be an $n \times s$ matrix. Let \mathbf{a}_i^T denote the i th row of \mathbf{A} , and let \mathbf{b}_j denote the j th column of \mathbf{B} . Then, the product of \mathbf{AB} is an $m \times s$ matrix \mathbf{C} defined as

$$\mathbf{C} = \mathbf{AB} = \begin{bmatrix} \mathbf{a}_1^T \mathbf{b}_1 & \mathbf{a}_1^T \mathbf{b}_2 & \cdots & \mathbf{a}_1^T \mathbf{b}_s \\ \mathbf{a}_2^T \mathbf{b}_1 & \mathbf{a}_2^T \mathbf{b}_2 & \cdots & \mathbf{a}_2^T \mathbf{b}_s \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{a}_m^T \mathbf{b}_1 & \mathbf{a}_m^T \mathbf{b}_2 & \cdots & \mathbf{a}_m^T \mathbf{b}_s \end{bmatrix}. \quad (\text{A.23})$$

It is not always true that $\mathbf{AB} = \mathbf{BA}$; that is, matrix multiplication is not necessarily commutative.

Let \mathbf{A} , \mathbf{B} , and \mathbf{C} be matrices such that \mathbf{A} , \mathbf{B} can be multiplied, \mathbf{A} , \mathbf{C} can be multiplied, and \mathbf{B} , \mathbf{C} can be added. Then,

$$\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}. \quad (\text{A.24})$$

The above property is called the distributive law. If x is a number, then

$$\mathbf{A}(x\mathbf{B}) = x(\mathbf{AB}). \quad (\text{A.25})$$

Furthermore, matrix multiplication is associative because

$$(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC}). \quad (\text{A.26})$$

Suppose now that \mathbf{A} , \mathbf{B} can be multiplied. Then,

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T. \quad (\text{A.27})$$

Thus the transpose of a product is the product of transposed matrices in reverse order.

Let \mathbf{A} be an $n \times n$ matrix. An inverse of \mathbf{A} , if it exists, is an $n \times n$ matrix \mathbf{B} such that

$$\mathbf{AB} = \mathbf{BA} = \mathbf{I}_n. \quad (\text{A.28})$$

If an inverse exists, then it is unique. A square matrix that is invertible is said to be *nonsingular*. We denote the inverse of \mathbf{A} by \mathbf{A}^{-1} . Hence, we have

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{AA}^{-1} = \mathbf{I}_n.$$

The transpose of an inverse is the inverse of the transpose,

$$(\mathbf{A}^{-1})^T = (\mathbf{A}^T)^{-1}, \quad (\text{A.29})$$

and we denote it as \mathbf{A}^{-T} .

We now present the Sherman–Morrison–Woodbury formula for inverting a special kind of matrices.

Lemma A.1 If A , D , and $D^{-1} + CA^{-1}B$ are invertible, then

$$(A + BDC)^{-1} = A^{-1} - A^{-1}B(D^{-1} + CA^{-1}B)^{-1}CA^{-1}.$$

Proof We prove the Sherman–Morrison–Woodbury formula by verification. We have

$$\begin{aligned} (A + BDC)(A^{-1} - A^{-1}B(D^{-1} + CA^{-1}B)^{-1}CA^{-1}) \\ &= (A + BDC)A^{-1} \\ &\quad - (A + BDC)A^{-1}B(D^{-1} + CA^{-1}B)^{-1}CA^{-1} \\ &= I + BDCA^{-1} \\ &\quad - (B + BDCA^{-1}B)(D^{-1} + CA^{-1}B)^{-1}CA^{-1} \\ &= I + BDCA^{-1} \\ &\quad - BD(D^{-1} + CA^{-1}B)(D^{-1} + CA^{-1}B)^{-1}CA^{-1} \\ &= I + BDCA^{-1} - BDCA^{-1} \\ &= I, \end{aligned}$$

as desired.

In the case when $D = I$, the Sherman–Morrison–Woodbury formula takes the form

$$(A + BC)^{-1} = A^{-1} - A^{-1}B(I + CA^{-1}B)^{-1}CA^{-1}. \quad (\text{A.30})$$

Formula (A.30) is sometimes referred to as the *matrix inversion formula*.

A set of nonzero vectors $\{q_1, q_2, \dots, q_r\}$ is said to be orthogonal if

$$\langle q_i, q_j \rangle = 0 \quad \text{for } i \neq j.$$

If each vector in an orthogonal set is of unity length—that is, $\|q_i\| = 1$ for all $i = 1, 2, \dots, r$ —then this set of vectors is called *orthonormal*. A square matrix with orthonormal columns is called a *unitary matrix*. Sometimes in the literature, a unitary matrix is called an *orthogonal matrix*. If $Q \in \mathbb{R}^{n \times n}$ is a unitary matrix, then

$$Q^T Q = Q Q^T = I_n. \quad (\text{A.31})$$

For a square, $n \times n$, matrix A , its trace, denoted $\text{trace } A$, is the sum of its diagonal elements,

$$\text{trace } A = \sum_{i=1}^n a_{ii}. \quad (\text{A.32})$$

If B is also an $n \times n$ matrix, then

$$\text{trace}(AB) = \text{trace}(BA). \quad (\text{A.33})$$

The *determinant* of a square matrix A , denoted $\det A$, can be evaluated using the Laplace expansion as follows. Let c_{ij} denote the cofactor corresponding to the element a_{ij} of an $n \times n$ matrix A . The cofactor c_{ij} is $(-1)^{i+j}$ times the determinant of the $(n-1) \times (n-1)$ submatrix of A obtained by deleting the i th row and j th column from A . Then, for any fixed i , $1 \leq i \leq n$, we have

$$\det A = \sum_{j=1}^n a_{ij} c_{ij}. \quad (\text{A.34})$$

This is the expansion of the determinant along the i th row. In a similar fashion, we can expand the determinant along a column. The determinant has the following properties:

D 1. $\det A = \det A^T$.

D 2. Let A and B be square matrices of the same size, then $\det(AB) = \det A \det B = \det(BA)$.

D 3. $\det O = 0$.

D 4. $\det I_n = 1$.

We now list additional properties of determinants. Let \mathbf{a}_i denote the i th column of A . Then, we write

$$\det A = \det(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n).$$

D 5. Suppose that the i th column \mathbf{a}_i is written as a sum

$$\mathbf{a}_i = \mathbf{b} + \mathbf{c}.$$

Then, we have

$$\begin{aligned} \det(\mathbf{a}_1, \dots, \mathbf{a}_i, \dots, \mathbf{a}_n) &= \det(\mathbf{a}_1, \dots, \mathbf{b} + \mathbf{c}, \dots, \mathbf{a}_n) \\ &= \det(\mathbf{a}_1, \dots, \mathbf{b}, \dots, \mathbf{a}_n) + \det(\mathbf{a}_1, \dots, \mathbf{c}, \dots, \mathbf{a}_n). \end{aligned}$$

D 6. If x is a number, then

$$\det(\mathbf{a}_1, \dots, x\mathbf{a}_i, \dots, \mathbf{a}_n) = x \det(\mathbf{a}_1, \dots, \mathbf{a}_i, \dots, \mathbf{a}_n).$$

D 7. If two columns of the matrix are equal, then the determinant is equal to zero.

D 8. If we add a multiple of one column to another, then the value of the determinant does not change. In other words, if x is a number, then

$$\det(\mathbf{a}_1, \dots, \mathbf{a}_i, \dots, \mathbf{a}_j, \dots, \mathbf{a}_n) = \det(\mathbf{a}_1, \dots, \mathbf{a}_i + x\mathbf{a}_j, \dots, \mathbf{a}_j, \dots, \mathbf{a}_n).$$

Properties D 5–D 8 hold also for rows instead of columns because $\det A = \det A^T$.

Using the cofactors of A and its determinant, we can give a formula for A^{-1} . First, we define the adjugate of A , denoted $\text{adj } A$, to be the matrix whose (i, j) th entry is the cofactor c_{ji} ; that is, $\text{adj } A$ is the transpose of the cofactor matrix of A . Then,

$$A^{-1} = \left(\frac{1}{\det A} \right) \text{adj } A = \frac{\text{adj } A}{\det A}. \quad (\text{A.35})$$

If A is $n \times n$ and \mathbf{v} is a nonzero $n \times 1$ vector such that for some scalar λ ,

$$A\mathbf{v} = \lambda\mathbf{v}, \quad (\text{A.36})$$

then \mathbf{v} is an *eigenvector* corresponding to the *eigenvalue* λ . In order that λ be an eigenvalue of A , it is necessary and sufficient for the matrix $\lambda I_n - A$ to be singular; that is, $\det(\lambda I_n - A) = 0$. This leads to an n th-order polynomial equation in λ of the form

$$\det(\lambda I_n - A) = \lambda^n + a_{n-1}\lambda^{n-1} + \dots + a_1\lambda + a_0 = 0. \quad (\text{A.37})$$

This equation, called the *characteristic polynomial equation*, must have n , possibly nondisjoint, complex roots that are the eigenvalues of A . The *spectral radius* of A , denoted $\rho(A)$, is

defined as

$$\rho(\mathbf{A}) = \max_{1 \leq i \leq n} |\lambda_i(\mathbf{A})|,$$

where $\lambda_i(\mathbf{A})$ denotes the i th eigenvalue of \mathbf{A} .

Any $m \times n$ matrix \mathbf{A} can be represented as

$$\mathbf{A} = \mathbf{U}[\mathbf{S}, \quad \mathbf{O}]\mathbf{V}^T \quad \text{if } m \leq n, \quad (\text{A.38})$$

or

$$\mathbf{A} = \mathbf{U} \begin{bmatrix} \mathbf{S} \\ \mathbf{O} \end{bmatrix} \mathbf{V}^T \quad \text{if } m > n, \quad (\text{A.39})$$

where \mathbf{U} is an $m \times m$ unitary matrix, \mathbf{V} is an $n \times n$ unitary matrix, and \mathbf{S} is a diagonal matrix of the form

$$\mathbf{S} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_p), \quad \text{where } p = \min(m, n) \quad \text{and} \quad \sigma_i \geq 0, \quad i = 1, 2, \dots, p.$$

The nonnegative numbers $\{\sigma_i\}$ are called the *singular values* of \mathbf{A} ; and (A.38), or (A.39), is called the *singular value decomposition* (SVD) of \mathbf{A} . Following the convention, we write

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p,$$

so that $\sigma_1 = \sigma_1(\mathbf{A})$ denotes the largest singular value of \mathbf{A} . The singular values of \mathbf{A} are the square roots of the eigenvalues of $\mathbf{A}\mathbf{A}^T$ if $m \leq n$ or of $\mathbf{A}^T\mathbf{A}$ if $m > n$. If \mathbf{A} is symmetric, its singular values are the absolute values of its eigenvalues.

A.4 Quadratic Forms

A real *quadratic form* is the product $\mathbf{x}^T \mathbf{Q} \mathbf{x}$, where \mathbf{Q} is a specified $n \times n$ matrix with real coefficients and \mathbf{x} is an $n \times 1$ real, column vector. Without loss of generality, the matrix \mathbf{Q} can be taken to be symmetric, that is, $\mathbf{Q} = \mathbf{Q}^T$. Indeed, if \mathbf{Q} is not symmetric, we can always replace it with a symmetric matrix without changing the quadratic form values. Indeed, because $z = \mathbf{x}^T \mathbf{Q} \mathbf{x}$ is a scalar, hence $z^T = z$, and we have

$$(\mathbf{x}^T \mathbf{Q} \mathbf{x})^T = \mathbf{x}^T \mathbf{Q}^T \mathbf{x} = \mathbf{x}^T \mathbf{Q} \mathbf{x}.$$

Therefore, we can write

$$\mathbf{x}^T \mathbf{Q} \mathbf{x} = \mathbf{x}^T \left(\frac{1}{2} \mathbf{Q} + \frac{1}{2} \mathbf{Q}^T \right) \mathbf{x} = \mathbf{x}^T \left(\frac{\mathbf{Q} + \mathbf{Q}^T}{2} \right) \mathbf{x},$$

where $\frac{1}{2}(\mathbf{Q} + \mathbf{Q}^T)$ is a symmetric matrix. Thus, the quadratic form values are unchanged if \mathbf{Q} is replaced with the symmetric matrix $\frac{1}{2}(\mathbf{Q} + \mathbf{Q}^T)$. From now on, in our study of quadratic forms, we assume that \mathbf{Q} is symmetric.

A quadratic form $\mathbf{x}^T \mathbf{Q} \mathbf{x}$ —or, equivalently, the matrix \mathbf{Q} —is called *positive semidefinite* if $\mathbf{x}^T \mathbf{Q} \mathbf{x} \geq 0$ for all \mathbf{x} . It is *positive definite* if $\mathbf{x}^T \mathbf{Q} \mathbf{x} > 0$ for all $\mathbf{x} \neq \mathbf{0}$. If \mathbf{Q} is positive semidefinite, then we write $\mathbf{Q} \geq 0$; if it is positive definite, then we write $\mathbf{Q} > 0$. We write $\mathbf{Q}_1 \geq \mathbf{Q}_2$ to mean that $\mathbf{Q}_1 - \mathbf{Q}_2 \geq 0$. We say that \mathbf{Q} is *negative semidefinite* if $\mathbf{x}^T \mathbf{Q} \mathbf{x} \leq 0$ for all \mathbf{x} , and it is *negative definite* if $\mathbf{x}^T \mathbf{Q} \mathbf{x} < 0$ for all $\mathbf{x} \neq \mathbf{0}$. Note that $\mathbf{Q} \leq 0$ if and only if $-\mathbf{Q} \geq 0$.

We now present tests for definiteness properties of symmetric matrices. We state these tests in terms of *principal minors* of a given real, symmetric, $n \times n$ matrix \mathbf{Q} with elements q_{ij} . We also give tests for definiteness in terms of eigenvalues of \mathbf{Q} . The principal minors of \mathbf{Q} are $\det \mathbf{Q}$ itself and the determinants of submatrices of \mathbf{Q} obtained by removing successively an i th row and i th column. Specifically, for $p = 1, 2, \dots, n$, the principal minors are

$$\Delta_p \begin{pmatrix} i_1 & i_2 & \dots & i_p \\ i_1 & i_2 & \dots & i_p \end{pmatrix} = \det \begin{bmatrix} q_{i_1 i_1} & q_{i_1 i_2} & \dots & q_{i_1 i_p} \\ q_{i_2 i_1} & q_{i_2 i_2} & \dots & q_{i_2 i_p} \\ \vdots & \vdots & \ddots & \vdots \\ q_{i_p i_1} & q_{i_p i_2} & \dots & q_{i_p i_p} \end{bmatrix}, \quad 1 \leq i_1 \leq i_2 \leq \dots \leq i_p.$$

The above are called the principal minors of order p .

In the following, we use the fact that if $\mathbf{Q} = \mathbf{Q}^T$, then its eigenvalues are all real.

Theorem A.5 The following are equivalent:

1. The symmetric matrix \mathbf{Q} is positive semidefinite.
2. Principal minors of \mathbf{Q} are nonnegative; that is, for $p = 1, 2, \dots, n$,

$$\Delta_p \begin{pmatrix} i_1 & i_2 & \dots & i_p \\ i_1 & i_2 & \dots & i_p \end{pmatrix} \geq 0, \quad 1 \leq i_1 \leq i_2 \leq \dots \leq i_p.$$

3. Eigenvalues of \mathbf{Q} are nonnegative, that is,

$$\lambda_i(\mathbf{Q}) \geq 0, \quad i = 1, 2, \dots, n.$$

We now give two tests for $\mathbf{Q} = \mathbf{Q}^T$ to be positive definite: one in terms of its leading principal minors, and the other in terms of its eigenvalues. The *leading principal minors* of \mathbf{Q} are

$$\begin{aligned} \Delta_1 \begin{pmatrix} 1 \\ 1 \end{pmatrix} &= q_{11}, & \Delta_2 \begin{pmatrix} 1 & 2 \\ 1 & 2 \end{pmatrix} &= \det \begin{bmatrix} q_{11} & q_{12} \\ q_{21} & q_{22} \end{bmatrix}, \\ \Delta_3 \begin{pmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \end{pmatrix} &= \det \begin{bmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{bmatrix}, & \dots, & \Delta_p \begin{pmatrix} 1 & 2 & \dots & p \\ 1 & 2 & \dots & p \end{pmatrix} = \det \mathbf{Q}. \end{aligned}$$

Theorem A.6 The following are equivalent:

1. The symmetric matrix \mathbf{Q} is positive definite.
2. Leading principal minors of \mathbf{Q} are positive, that is,

$$\Delta_p \begin{pmatrix} 1 & 2 & \dots & p \\ 1 & 2 & \dots & p \end{pmatrix} > 0, \quad p = 1, 2, \dots, n.$$

3. Eigenvalues of \mathbf{Q} are positive, that is,

$$\lambda_i(\mathbf{Q}) > 0, \quad i = 1, 2, \dots, n.$$

◆ Example A.1

We will find the range of the parameter γ for which the quadratic form

$$\mathbf{x}^T \begin{bmatrix} -5 - \gamma & 6 \\ 0 & -2 \end{bmatrix} \mathbf{x}$$

is negative definite. We first symmetrize the above quadratic form to get

$$f = \mathbf{x}^T \begin{bmatrix} -5 - \gamma & 6 \\ 0 & -2 \end{bmatrix} \mathbf{x} = \mathbf{x}^T \begin{bmatrix} -5 - \gamma & 3 \\ 3 & -2 \end{bmatrix} \mathbf{x}.$$

A quadratic form is negative definite if and only if its negative is positive definite. So instead of working with f , we will find the range of the parameter γ for which $g = -f$ is positive definite, where

$$g = \mathbf{x}^T \begin{bmatrix} 5 + \gamma & -3 \\ -3 & 2 \end{bmatrix} \mathbf{x}.$$

To finish to problem, we will investigate when the matrix

$$\begin{bmatrix} 5 + \gamma & -3 \\ -3 & 2 \end{bmatrix}$$

is positive definite. The first-order leading principal minor of the above matrix is $\Delta_1 = 5 + \gamma$, which is positive if and only if $\gamma > -5$. The second-order leading principal minor is

$$\Delta_2 = \det \begin{bmatrix} 5 + \gamma & -3 \\ -3 & 2 \end{bmatrix} = 2\gamma + 1,$$

which positive for $\gamma > -1/2$. Thus the quadratic form f is negative definite if and only if $\gamma > -1/2$.

If a given matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is symmetric, then its eigenvalues are real and there exists a unitary matrix \mathbf{U} such that

$$\mathbf{U}^T \mathbf{Q} \mathbf{U} = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix} = \mathbf{\Lambda}. \quad (\text{A.40})$$

Using the above, one can show that a symmetric positive semidefinite matrix \mathbf{Q} has a positive semidefinite symmetric square root, denoted $\mathbf{Q}^{1/2}$ or $\sqrt{\mathbf{Q}}$, that satisfies

$$\mathbf{Q}^{1/2} \mathbf{Q}^{1/2} = \mathbf{Q}. \quad (\text{A.41})$$

Indeed, it follows from Theorem A.5 that if $\mathbf{Q} = \mathbf{Q}^T \geq 0$ then its eigenvalues are all nonnegative. Hence,

$$\mathbf{\Lambda}^{1/2} = \begin{bmatrix} \lambda_1^{1/2} & 0 & \cdots & 0 \\ 0 & \lambda_2^{1/2} & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n^{1/2} \end{bmatrix} \quad (\text{A.42})$$

is well-defined. Using the fact that $\mathbf{U}^T \mathbf{U} = \mathbf{I}_n$, we can define

$$\mathbf{Q}^{1/2} = \mathbf{U} \mathbf{\Lambda}^{1/2} \mathbf{U}^T, \quad (\text{A.43})$$

which has the desired properties.

If $\mathbf{Q} = \mathbf{Q}^T$ is positive semidefinite and not positive definite, then some of its eigenvalues are equal to zero. Suppose that \mathbf{Q} has r nonzero eigenvalues. Then, there is a unitary matrix, say \mathbf{U} , such that

$$\mathbf{U}^T \mathbf{Q} \mathbf{U} = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & & \ddots & \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_r & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & & & & & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \end{bmatrix}. \quad (\text{A.44})$$

Let $\mathbf{V} \in \mathbb{R}^{n \times r}$ be a matrix that consists of the first r columns of the matrix \mathbf{U} , and let

$$\mathbf{C} = \begin{bmatrix} \lambda_1^{1/2} & 0 & \cdots & 0 \\ 0 & \lambda_2^{1/2} & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_r^{1/2} \end{bmatrix} \mathbf{V}^T. \quad (\text{A.45})$$

Note that the matrix \mathbf{C} is of full rank, and

$$\mathbf{Q} = \mathbf{C}^T \mathbf{C}. \quad (\text{A.46})$$

We refer to (A.46) as a full-rank factorization of \mathbf{Q} because \mathbf{C} is of full rank. We can generalize the above, using (A.38), to obtain a full-rank factorization of a nonsquare matrix. Specifically, given an $m \times n$ matrix \mathbf{A} of rank r , we can represent it as

$$\mathbf{A} = \mathbf{B} \mathbf{C},$$

where \mathbf{B} is an $m \times r$ matrix and \mathbf{C} is a $r \times n$ matrix such that

$$\text{rank}(\mathbf{B}) = \text{rank}(\mathbf{C}) = \text{rank}(\mathbf{A}) = r.$$

A.5 The Kronecker Product

Let A be an $m \times n$ and B be a $p \times q$ matrices. The *Kronecker product* of A and B , denoted $A \otimes B$, is an $mp \times nq$ matrix defined as

$$A \otimes B = \begin{bmatrix} a_{11}B & a_{12}B & \cdots & a_{1n}B \\ a_{21}B & a_{22}B & \cdots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mn}B \end{bmatrix}. \quad (\text{A.47})$$

Thus, the matrix $A \otimes B$ consists of mn blocks. Using the definition of the Kronecker product, we can verify the following two properties:

$$(A \otimes B)(C \otimes D) = AC \otimes BD, \quad (\text{A.48})$$

$$(A \otimes B)^T = A^T \otimes B^T. \quad (\text{A.49})$$

Let

$$X = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_m]$$

be a $n \times m$ matrix, where \mathbf{x}_i , $i = 1, 2, \dots, m$, are the columns of X . Each \mathbf{x}_i consists of n elements. Then, the *stacking operator* defined as

$$\begin{aligned} s(X) &= [\mathbf{x}_1^T \quad \mathbf{x}_2^T \quad \cdots \quad \mathbf{x}_m^T]^T \\ &= [x_{11} \quad x_{21} \quad \cdots \quad x_{n1} \quad x_{12} \quad x_{22} \quad \cdots \quad x_{n2} \quad \cdots \quad x_{1m} \quad x_{2m} \quad \cdots \quad x_{nm}]^T \end{aligned}$$

yields the column nm vector formed from the columns of X taken in order. Let now A be an $n \times n$, and let C and X be $n \times m$ matrices. Then, the matrix equation

$$AX = C \quad (\text{A.50})$$

can be written as

$$(I_m \otimes A)s(X) = s(C). \quad (\text{A.51})$$

If B is an $m \times m$ matrix, then the matrix equation

$$XB = C \quad (\text{A.52})$$

can be represented as

$$(B^T \otimes I_n)s(X) = s(C). \quad (\text{A.53})$$

Using the above two facts, we can verify that the matrix equation

$$AX + XB = C, \quad (\text{A.54})$$

where A is $n \times n$, B is $m \times m$, and C is $n \times m$, can be written as

$$(I_m \otimes A + B^T \otimes I_n)s(X) = s(C). \quad (\text{A.55})$$

◆ Example A.2

Consider the matrix equation (A.54), where

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} -2 & 0 \\ -3 & 1 \end{bmatrix}, \quad \text{and} \quad \mathbf{C} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}.$$

We will represent the above equation in the matrix–vector format as given by (A.55), that is,

$$(\mathbf{I}_2 \otimes \mathbf{A} - \mathbf{B}^T \otimes \mathbf{I}_2) \mathbf{s}(\mathbf{X}) = \mathbf{s}(\mathbf{C}).$$

Substituting into the above equation the given matrices yields

$$\left(\begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix} + \begin{bmatrix} -2 & 0 & -3 & 0 \\ 0 & -2 & 0 & -3 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \right) \begin{bmatrix} x_{11} \\ x_{21} \\ x_{12} \\ x_{22} \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}.$$

Performing simple manipulations yields

$$\begin{bmatrix} -2 & 1 & -3 & 0 \\ 0 & -1 & 0 & -3 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} x_{11} \\ x_{21} \\ x_{12} \\ x_{22} \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}.$$

The above system of linear equations has a unique solution because the matrix

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

is nonsingular.

Let λ_i , \mathbf{v}_i be the eigenvalues and eigenvectors, respectively of \mathbf{A} , and μ_j and \mathbf{w}_j the eigenvalues and eigenvectors of $m \times m$ matrix \mathbf{B} . Then,

$$\begin{aligned} (\mathbf{A} \otimes \mathbf{B})(\mathbf{v}_i \otimes \mathbf{w}_j) &= \mathbf{A}\mathbf{v}_i \otimes \mathbf{B}\mathbf{w}_j \\ &= \lambda_i \mathbf{v}_i \otimes \mu_j \mathbf{w}_j \\ &= \lambda_i \mu_j (\mathbf{v}_i \otimes \mathbf{w}_j). \end{aligned} \tag{A.56}$$

Thus, the eigenvalues of $\mathbf{A} \otimes \mathbf{B}$ are $\lambda_i \mu_j$, and their respective eigenvectors are $\mathbf{v}_i \otimes \mathbf{w}_j$ for $i = 1, 2, \dots, n$, $j = 1, 2, \dots, m$.

◆ Example A.3

For the given two matrices

$$\mathbf{A} = \begin{bmatrix} 0 & 0 \\ 2 & -2 \end{bmatrix} \quad \text{and} \quad \mathbf{B} = \begin{bmatrix} 4 & 6 & -8 \\ 0 & -1 & 0 \\ 0 & 2 & 7 \end{bmatrix},$$

we will find the eigenvalues of the matrix $(\mathbf{I} + \mathbf{A} \otimes \mathbf{B}^T)$, where \mathbf{I} is the identity matrix of appropriate dimensions. We note that

$$\text{eig}(\mathbf{I} + \mathbf{A} \otimes \mathbf{B}^T) = 1 + \text{eig}(\mathbf{A})\text{eig}(\mathbf{B}),$$

because for any square matrix \mathbf{M} , $\text{eig}(\mathbf{I} + \mathbf{M}) = 1 + \text{eig}(\mathbf{M})$ and $\text{eig}(\mathbf{M}^T) = \text{eig}(\mathbf{M})$. By inspection,

$$\text{eig}(\mathbf{A}) = \{0, -2\} \quad \text{and} \quad \text{eig}(\mathbf{B}) = \{4, -1, 7\}.$$

Hence,

$$\text{eig}(\mathbf{I} + \mathbf{A} \otimes \mathbf{B}^T) = \{1, 1, 1, -7, 3, -13\}.$$

We will now use (A.56) to study the matrix equation (A.54), which we represent as

$$\mathbf{M} s(\mathbf{X}) = s(\mathbf{C}), \tag{A.57}$$

where

$$\mathbf{M} = \mathbf{I}_m \otimes \mathbf{A} + \mathbf{B}^T \otimes \mathbf{I}_n. \tag{A.58}$$

The solution to (A.57) is unique if, and only if, the $mn \times mn$ matrix \mathbf{M} is nonsingular. To find the condition for this to hold, consider the matrix

$$(\mathbf{I}_m + \varepsilon \mathbf{B}^T) \otimes (\mathbf{I}_n + \varepsilon \mathbf{A}) = \mathbf{I}_m \otimes \mathbf{I}_n + \varepsilon \mathbf{M} + \varepsilon^2 \mathbf{B}^T \otimes \mathbf{A} \tag{A.59}$$

whose eigenvalues are

$$(1 + \varepsilon \mu_j)(1 + \varepsilon \lambda_i) = 1 + \varepsilon(\mu_j + \lambda_i) + \varepsilon^2 \mu_j \lambda_i \tag{A.60}$$

because for a square matrix \mathbf{Q} we have

$$\lambda_i(\mathbf{I}_n + \varepsilon \mathbf{Q}) = 1 + \varepsilon \lambda_i(\mathbf{Q}).$$

Comparing the terms in ε in (A.59) and (A.60), we conclude that the eigenvalues of \mathbf{M} are $\lambda_i + \mu_j$, $i = 1, 2, \dots, n$, $j = 1, 2, \dots, m$. Hence, \mathbf{M} is nonsingular if and only if

$$\lambda_i + \mu_j \neq 0. \tag{A.61}$$

The above is a necessary and sufficient condition for the solution \mathbf{X} of the matrix equation $\mathbf{A}\mathbf{X} + \mathbf{X}\mathbf{B} = \mathbf{C}$ to be unique.

Another useful identity involving the Kronecker product is

$$s(\mathbf{ABC}) = (\mathbf{C}^T \otimes \mathbf{A})s(\mathbf{B}). \tag{A.62}$$

For further information on the subject of the Kronecker product, we refer the reader to Brewer [33].

A.6 Upper and Lower Bounds

Consider a set \mathcal{S} of real numbers. A number M is called an *upper bound* of the set \mathcal{S} if

$$x \leq M \quad \text{for all } x \in \mathcal{S}.$$

A set of real numbers that has an upper bound is said to be *bounded above*.

The Least Upper Bound Axiom Every nonempty set \mathcal{S} of real numbers that has an upper bound has a *least upper bound*, also called *supremum*, and is denoted

$$\sup\{x : x \in \mathcal{S}\}.$$

Examples

1. $\sup\left\{\frac{1}{2}, \frac{2}{3}, \dots, \frac{n}{n+1}, \dots\right\} = 1;$
2. $\sup\left\{-\frac{1}{2}, -\frac{1}{8}, -\frac{1}{27}, \dots, -\frac{1}{n^3}, \dots\right\} = 0;$
3. $\sup\{x : x^2 < 3\} = \sup\{x : -\sqrt{3} < x < \sqrt{3}\} = \sqrt{3}.$

Theorem A.7 If $M = \sup\{x : x \in \mathcal{S}\}$ and $\varepsilon > 0$, then there is at least one number x in \mathcal{S} such that

$$M - \varepsilon < x \leq M.$$

Proof The condition $x \leq M$ is satisfied by all numbers x in \mathcal{S} by virtue of the least upper bound axiom. We have to show that for any $\varepsilon > 0$ there is a number $x \in \mathcal{S}$ such that

$$M - \varepsilon < x.$$

We prove the statement by contradiction. We suppose that there is no such number in \mathcal{S} . Then,

$$x \leq M - \varepsilon \quad \text{for all } x \in \mathcal{S}$$

and hence $M - \varepsilon$ would be an upper bound of \mathcal{S} that is less than M , which is the least upper bound. This contradicts the assumption, and thus the proof is complete.

To illustrate the above theorem, we consider the following set of real numbers:

$$\mathcal{S} = \left\{\frac{1}{2}, \frac{2}{3}, \frac{3}{4}, \dots, \frac{n}{n+1}, \dots\right\}.$$

Note that

$$\sup\{x : x \in \mathcal{S}\} = 1.$$

Let $\varepsilon = 0.1$. Then, for example, the element of the set \mathcal{S} that is equal to $99/100$ satisfies

$$1 - \varepsilon < \frac{99}{100} \leq 1.$$

A number m is called a *lower bound* of \mathcal{S} if

$$m \leq x \quad \text{for all } x \in \mathcal{S}.$$

Sets that have lower bounds are said to be *bounded below*.

Theorem A.8 Every nonempty set of real numbers that has a lower bound has a *greatest lower bound*, also called *infimum*, and is denoted as

$$\inf\{x : x \in \mathcal{S}\}.$$

Proof By assumption, \mathcal{S} is nonempty and has a lower bound that we call s . Thus

$$s \leq x \quad \text{for all } x \in \mathcal{S}.$$

Hence

$$-x \leq -s \quad \text{for all } x \in \mathcal{S};$$

that is, the set $\{-x : x \in \mathcal{S}\}$ has an upper bound $-s$. From the least upper bound axiom, we conclude that the set $\{-x : x \in \mathcal{S}\}$ has a least upper bound, supremum. We call it $-s_0$. Because

$$-x \leq -s_0 \quad \text{for all } x \in \mathcal{S},$$

we have

$$s_0 \leq x \quad \text{for all } x \in \mathcal{S},$$

and thus s_0 is a lower bound for \mathcal{S} . We now prove, by contradiction, that s_0 is the greatest lower bound, infimum, of the set \mathcal{S} . Indeed, if there was a number \tilde{x} satisfying

$$s_0 < \tilde{x} \leq x \quad \text{for all } x \in \mathcal{S},$$

then we would have

$$-x \leq -\tilde{x} < -s_0 \quad \text{for all } x \in \mathcal{S},$$

which contradicts the fact that $-s_0$ is the supremum of $\{-x : x \in \mathcal{S}\}$.

Using arguments similar to those we used in the proof of Theorem A.8, we can prove the following theorem:

Theorem A.9 If $m = \inf\{x : x \in \mathcal{S}\}$ and $\varepsilon > 0$, then there is at least one number x in \mathcal{S} such that

$$m \leq x < m + \varepsilon.$$

A.7 Sequences

A *sequence of real numbers* can be viewed as a function whose domain is the set of natural numbers $1, 2, \dots, n, \dots$ and whose range is contained in \mathbb{R} . Thus, a sequence of real numbers can be viewed as a set of ordered pairs $(1, a_1), (2, a_2), \dots, (n, a_n), \dots$. We denote a sequence

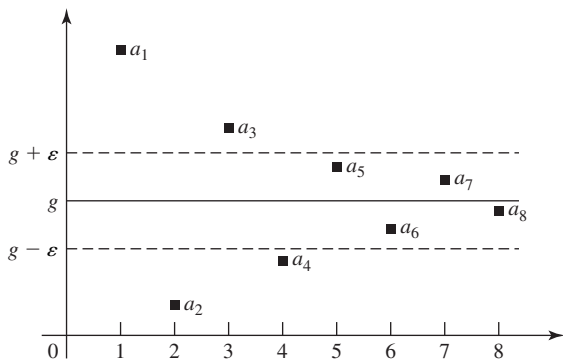


Figure A.2 An illustration of the notion of the limit of a sequence.

of real numbers

$$a_1, a_2, \dots, a_n, \dots$$

as $\{a_n\}$.

A sequence $\{a_n\}$ is increasing if $a_1 < a_2 < \dots < a_n \dots$. In general, a sequence is increasing if $a_n < a_{n+1}$. If $a_n \leq a_{n+1}$, then we say that the sequence is nondecreasing. Similarly, we can define decreasing and nonincreasing sequences. Nonincreasing or nondecreasing sequences are called *monotone* sequences.

A number g is called the *limit* of the infinite sequence $a_1, a_2, \dots, a_n, \dots$ if for any positive ε there is a natural number $k = k(\varepsilon)$ such that for all $n > k$ we have

$$|a_n - g| < \varepsilon;$$

that is, a_n lies between $g - \varepsilon$ and $g + \varepsilon$ for all $n > k$. In other words, for any $\varepsilon > 0$ we have $|a_n - g| < \varepsilon$ for sufficiently large n 's (see Figure A.2), and we write

$$g = \lim_{n \rightarrow \infty} a_n,$$

or

$$a_n \rightarrow g.$$

A sequence that has a limit is called a *convergent* sequence. A sequence that has no limit is called *divergent*.

The notion of the sequence can be extended to sequences with elements in \mathbb{R}^p . Specifically, a sequence in \mathbb{R}^p is a function whose domain is the set of natural numbers $1, 2, \dots, n, \dots$ and whose range is contained in \mathbb{R}^p .

Suppose that $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n, \dots$, denoted $\{\mathbf{a}_n\}$, is a sequence in \mathbb{R}^p . Then, $\mathbf{g} \in \mathbb{R}^p$ is called the limit of this sequence if for each $\varepsilon > 0$ there exists a natural number $k = k(\varepsilon)$ such that for all $n > k$ we have

$$\|\mathbf{a}_n - \mathbf{g}\| < \varepsilon,$$

and we write

$$\mathbf{g} = \lim_{n \rightarrow \infty} \mathbf{a}_n,$$

or

$$\mathbf{a}_n \rightarrow \mathbf{g}.$$

Theorem A.10 If a limit of a sequence in \mathbb{R}^p exists, then it is unique.

Proof We prove this theorem by contradiction. Let us assume that a sequence $\{\mathbf{a}_n\}$ has two different limits, say \mathbf{g}_1 and \mathbf{g}_2 . Then, we have

$$\|\mathbf{g}_1 - \mathbf{g}_2\| > 0.$$

Let

$$\varepsilon = \frac{1}{2}\|\mathbf{g}_1 - \mathbf{g}_2\|.$$

From the definition of a limit, there exist $k_1 = k_1(\varepsilon)$ and $k_2 = k_2(\varepsilon)$ such that for $n > k_1$ we have $\|\mathbf{a}_n - \mathbf{g}_1\| < \varepsilon$, and for $n > k_2$ we have $\|\mathbf{a}_n - \mathbf{g}_2\| < \varepsilon$. Let $m = \max(k_1, k_2)$. If $n > m$, then $\|\mathbf{a}_n - \mathbf{g}_1\| < \varepsilon$ and $\|\mathbf{a}_n - \mathbf{g}_2\| < \varepsilon$. Adding $\|\mathbf{a}_n - \mathbf{g}_1\| < \varepsilon$ and $\|\mathbf{a}_n - \mathbf{g}_2\| < \varepsilon$ yields

$$\|\mathbf{a}_n - \mathbf{g}_1\| + \|\mathbf{a}_n - \mathbf{g}_2\| < 2\varepsilon. \quad (\text{A.63})$$

It follows from the triangle inequality that

$$\|\mathbf{a} - \mathbf{b}\| \leq \|\mathbf{a}\| + \|\mathbf{b}\|.$$

Using the above inequality, we obtain

$$\begin{aligned} \|\mathbf{g}_1 - \mathbf{g}_2\| &= \|\mathbf{a}_n - \mathbf{g}_1 - \mathbf{a}_n + \mathbf{g}_2\| \\ &= \|(\mathbf{a}_n - \mathbf{g}_1) - (\mathbf{a}_n - \mathbf{g}_2)\| \\ &\leq \|\mathbf{a}_n - \mathbf{g}_1\| + \|\mathbf{a}_n - \mathbf{g}_2\|. \end{aligned} \quad (\text{A.64})$$

Applying (A.63) to (A.64) yields

$$\|\mathbf{g}_1 - \mathbf{g}_2\| = \|\mathbf{g}_1 - \mathbf{g}_2\| < 2\varepsilon.$$

However, the above contradicts the assumption that $\|\mathbf{g}_1 - \mathbf{g}_2\| = 2\varepsilon$, which completes the proof.

A sequence in \mathbb{R}^p is bounded if there exists a number $M \geq 0$ such that $\|\mathbf{a}_n\| \leq M$ for all $n = 1, 2, \dots$. The set of all bounded sequences is sometimes denoted as l_∞ .

Theorem A.11 Every convergent sequence is bounded.

Proof Let

$$\mathbf{g} = \lim_{n \rightarrow \infty} \mathbf{a}_n.$$

Choose $\varepsilon = 1$. Then, there exists a natural number $k = k(\varepsilon)$ such that

$$\|\mathbf{a}_n - \mathbf{g}\| < 1 \quad \text{for all } n > k.$$

By Theorem A.4,

$$\|\mathbf{a}\| - \|\mathbf{b}\| \leq \|\mathbf{a} - \mathbf{b}\|.$$

Using the above inequality, we get

$$\|\mathbf{a}_n\| - \|\mathbf{g}\| \leq \|\mathbf{a}_n - \mathbf{g}\| < 1 \quad \text{for all } n > k.$$

Therefore,

$$\|\mathbf{a}_n\| < \|\mathbf{g}\| + 1 \quad \text{for all } n > k.$$

Let

$$M = \max(\|\mathbf{a}_1\|, \|\mathbf{a}_2\|, \dots, \|\mathbf{a}_k\|, \|\mathbf{g}\| + 1),$$

then we have

$$M \geq \|\mathbf{a}_n\| \quad \text{for all } n,$$

which means that a convergent sequence is bounded, and the proof is complete.

Suppose we are given a sequence $\mathbf{a}_1, \mathbf{a}_2, \dots$ in \mathbb{R}^p and an increasing sequence of positive natural numbers

$$m_1, m_2, \dots, m_n, \dots$$

The sequence

$$\mathbf{a}_{m_1}, \mathbf{a}_{m_2}, \dots, \mathbf{a}_{m_n}, \dots$$

is called a *subsequence* of the sequence $\mathbf{a}_1, \mathbf{a}_2, \dots$ in \mathbb{R}^p , and we have

$$m_n \geq n.$$

Indeed, this is obvious for $n = 1$; that is, $m_1 \geq 1$ because m_1 is a positive integer. We now apply the principle of induction. We assume that $m_n \geq n$ for a given n . Then,

$$m_{n+1} > m_n \geq n.$$

Hence,

$$m_{n+1} \geq n + 1.$$

Thus, indeed $m_n \geq n$ for all n .

We can say that a subsequence is obtained from the given sequence by neglecting a number of elements in the sequence. We denote a subsequence $\mathbf{a}_{m_1}, \mathbf{a}_{m_2}, \dots, \mathbf{a}_{m_n}, \dots$ as

$$\{\mathbf{a}_{m_n}\}.$$

Theorem A.12 A subsequence of a convergent sequence in \mathbb{R}^p is convergent to the same limit; that is, if

$$\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{g},$$

and if

$$m_1 < m_2 < \dots < m_n < \dots,$$

then

$$\lim_{n \rightarrow \infty} \mathbf{a}_{m_n} = \mathbf{g}.$$

Proof Let $\varepsilon > 0$ be given. Then, there exists a natural number $k = k(\varepsilon)$ such that $\|\mathbf{a}_n - \mathbf{g}\| < \varepsilon$ for any $n > k$. From the fact that $m_n \geq n$, it follows that

$$m_n > n > k,$$

and thus

$$\|\mathbf{a}_{m_n} - \mathbf{g}\| < \varepsilon$$

for any $m_n > n > k$, which means that

$$\lim_{n \rightarrow \infty} \mathbf{a}_{m_n} = \mathbf{g}.$$

The proof is complete.

We now state, without a proof, an important result due to Bolzano and Weierstrass.

Theorem A.13 (Bolzano–Weierstrass) Every bounded sequence in \mathbb{R}^p contains a convergent subsequence.

Theorem A.14 (Cauchy) A sequence $\{\mathbf{a}_n\}$ in \mathbb{R}^p is convergent if and only if for every $\varepsilon > 0$ there exists an $r = r(\varepsilon)$ such that

$$\|\mathbf{a}_n - \mathbf{a}_r\| < \varepsilon$$

holds for all $n > r$.

Proof (\Rightarrow) Let

$$\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{g},$$

and let $\varepsilon > 0$ be given. Then, for some $l = l(\varepsilon)$,

$$\|\mathbf{a}_n - \mathbf{g}\| < \frac{1}{2}\varepsilon \quad \text{for all } n > l.$$

Let $r = l + 1$, then

$$\|\mathbf{a}_r - \mathbf{g}\| < \frac{1}{2}\varepsilon.$$

Adding the above inequalities, we obtain

$$\|\mathbf{a}_n - \mathbf{g}\| + \|\mathbf{a}_r - \mathbf{g}\| < \varepsilon \quad \text{for all } n > r.$$

On the other hand,

$$\|\mathbf{a}_n - \mathbf{a}_r\| \leq \|\mathbf{a}_n - \mathbf{g}\| + \|\mathbf{a}_r - \mathbf{g}\|.$$

Hence,

$$\|\mathbf{a}_n - \mathbf{a}_r\| < \varepsilon \quad \text{for all } n > r.$$

(\Leftarrow) We now assume that for every $\varepsilon > 0$ there exists an $r = r(\varepsilon)$ such that

$$\|\mathbf{a}_n - \mathbf{a}_r\| < \varepsilon \quad \text{for all } n > r.$$

We then show that the above condition implies

$$\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{g}.$$

First, we show that $\{\mathbf{a}_n\}$ is a bounded sequence. Let $\varepsilon = 1$. Then, an $r = r(\varepsilon)$ exists such that $\|\mathbf{a}_n - \mathbf{a}_r\| < 1$ for all $n > r$. Hence,

$$\|\mathbf{a}_n\| - \|\mathbf{a}_r\| \leq \|\mathbf{a}_n - \mathbf{a}_r\| < 1,$$

which means that

$$\|\mathbf{a}_n\| < \|\mathbf{a}_r\| + 1.$$

Let

$$M = \max\{\|\mathbf{a}_1\|, \|\mathbf{a}_2\|, \dots, \|\mathbf{a}_{r-1}\|, \|\mathbf{a}_r\| + 1\},$$

then we have

$$\|\mathbf{a}_n\| \leq M \quad \text{for all } n,$$

which implies that the sequence $\{\mathbf{a}_n\}$ is bounded. By the Bolzano–Weierstrass theorem the sequence $\{\mathbf{a}_n\}$ contains a convergent subsequence. Let

$$\lim_{n \rightarrow \infty} \mathbf{a}_{m_n} = \mathbf{g}, \quad m_1 < m_2 < \dots$$

We will show that

$$\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{g}.$$

From the hypothesis of the theorem it follows that for some $r = r(\varepsilon/3)$ we have

$$\|\mathbf{a}_n - \mathbf{a}_r\| < \frac{1}{3}\varepsilon \quad \text{for all } n > r. \quad (\text{A.65})$$

On the other hand, the assumption that the subsequence $\{\mathbf{a}_{m_n}\}$ is convergent implies that for some $k = k(\varepsilon/3)$ we have

$$\|\mathbf{a}_{m_n} - \mathbf{g}\| < \frac{1}{3}\varepsilon \quad \text{for all } n > k. \quad (\text{A.66})$$

We can select k so that $k > r$. In such a case, (A.65) and (A.66) are satisfied for all $n > k$. Because $m_n \geq n > r$, we have for all $n > k$ the following:

$$\|\mathbf{a}_{m_n} - \mathbf{a}_r\| < \frac{1}{3}\varepsilon,$$

or, equivalently,

$$\|\mathbf{a}_r - \mathbf{a}_{m_n}\| < \frac{1}{3}\varepsilon. \quad (\text{A.67})$$

Adding (A.65), (A.66), and (A.67) yields

$$\|\mathbf{a}_n - \mathbf{a}_r\| + \|\mathbf{a}_{m_n} - \mathbf{g}\| + \|\mathbf{a}_r - \mathbf{a}_{m_n}\| < \varepsilon.$$

Hence,

$$\begin{aligned} \|\mathbf{a}_n - \mathbf{g}\| &= \|(\mathbf{a}_n - \mathbf{a}_r) + (\mathbf{a}_r - \mathbf{a}_{m_n}) + (\mathbf{a}_{m_n} - \mathbf{g})\| \\ &\leq \|\mathbf{a}_n - \mathbf{a}_r\| + \|\mathbf{a}_{m_n} - \mathbf{g}\| + \|\mathbf{a}_r - \mathbf{a}_{m_n}\| \\ &< \varepsilon \quad \text{for all } n > k, \end{aligned}$$

which implies that

$$\lim_{n \rightarrow \infty} \mathbf{a}_n = \mathbf{g}.$$

The proof is complete.

Theorem A.15 Every monotone bounded sequence in \mathbb{R} is convergent.

Proof We will prove the theorem for nondecreasing sequences. The proof for non-increasing sequences is analogous.

We assume that we are given a sequence $\{a_n\}$ in \mathbb{R} that is bounded and nondecreasing. Let Z be the set of numbers belonging to the sequence, and let g be the least upper bound of Z . Hence

$$g \geq a_n, \quad n = 1, 2, \dots$$

From the properties of the least upper bound (see Theorem A.7), it follows that for any $\varepsilon > 0$ there exists a k such that

$$g - \varepsilon < a_k.$$

Because the sequence $\{a_n\}$ is nondecreasing, the inequality $n > k$ implies that

$$a_k \leq a_n,$$

and thus

$$g - \varepsilon < a_n.$$

We then conclude that for any $\varepsilon > 0$, there exists a k such that for $n > k$,

$$g - \varepsilon < a_n < g + \varepsilon.$$

Hence,

$$|a_n - g| < \varepsilon,$$

which means that

$$\lim_{n \rightarrow \infty} a_n = g.$$

Thus, a monotone bounded sequence in \mathbb{R} is convergent. The proof is complete.

Given a sequence $\{a_n\}$ of real numbers. Then the sequence $\{s_N\}$ defined by

$$s_N = \sum_{n=1}^N a_n = a_1 + a_2 + \dots + a_N$$

is called the sequence of *partial sums* of the series

$$\sum_{n=1}^{\infty} a_n. \tag{A.68}$$

If the sequence of the partial sums has a limit—that is, $s_N \rightarrow s$ as $N \rightarrow \infty$ —then we say that the series converges to the sum s , and write

$$s = \sum_{n=1}^{\infty} a_n.$$

The above formula can only make sense when the series (A.68) converges.

Theorem A.16 If the series $\sum_{n=1}^{\infty} a_n$ converges, then

$$a_n \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

Proof Let s be the sum of the series $\sum_{n=1}^{\infty} a_n$. Then,

$$s_N = \sum_{n=1}^N a_n \rightarrow s \quad \text{as } N \rightarrow \infty. \quad (\text{A.69})$$

We also have

$$s_{N-1} \rightarrow s \quad \text{as } N \rightarrow \infty, \quad (\text{A.70})$$

and

$$\begin{aligned} a_N &= (a_1 + a_2 + \cdots + a_N) - (a_1 + a_2 + \cdots + a_{N-1}) \\ &= s_N - s_{N-1}. \end{aligned}$$

Therefore, by (A.69) and (A.70),

$$a_N \rightarrow (s - s) = 0 \quad \text{as } N \rightarrow \infty,$$

and the proof is complete.

The converse of Theorem A.16 is false. The terms of the series

$$\sum_{n=1}^{\infty} \frac{1}{n}$$

tend to zero, but the series diverges.

A.8 Functions

We discuss functions of one and several variables and their basic properties. A function f from a set \mathcal{A} to a set \mathcal{B} , denoted

$$f : \mathcal{A} \rightarrow \mathcal{B},$$

is a rule that assigns to each $x \in \mathcal{A}$ a unique element $y \in \mathcal{B}$. The set \mathcal{A} is the domain of f and the range of f is a subset of \mathcal{B} , not necessarily the whole of \mathcal{B} .

Let $f : \mathcal{S} \rightarrow \mathbb{R}$, where \mathcal{S} is a subset of \mathbb{R} . Then, f is said to be continuous at $z_0 \in \mathcal{S}$ if for every $\varepsilon > 0$, there exists $\delta = \delta(\varepsilon, z_0)$ such that if $z \in \mathcal{S}$, $|z - z_0| < \delta$, then $|f(z) - f(z_0)| < \varepsilon$. This is also written as $f(z) \rightarrow f(z_0)$ as $z \rightarrow z_0$. If f is continuous at every point of \mathcal{S} , then we say that f is continuous on \mathcal{S} .

We say that the function f is bounded on a set \mathcal{S} if there is $M \geq 0$ such that for any $z \in \mathcal{S}$ we have

$$|f(z)| \leq M.$$

Lemma A.2 If f is continuous on $[a, b] \subset \mathbb{R}$, then f is bounded on $[a, b]$.

Proof Consider the set

$$\{x : x \in [a, b], \text{ and } f \text{ is bounded on } [a, x]\}.$$

This set is nonempty and bounded above by b . Let

$$c = \sup\{x : f \text{ is bounded on } [a, x]\}.$$

We now show that $c = b$. Suppose that $c < b$. From the continuity of f at c , it follows that f is bounded on $[c - \delta, c + \delta]$ for some sufficiently small $\delta > 0$. Being bounded on $[a, c - \delta]$ and on $[c - \delta, c + \delta]$, it is bounded on $[a, c + \delta]$. This contradicts our choice of c . We therefore conclude that $c = b$. This means that f is bounded on $[a, x]$ for all $x < b$. From the continuity of f , we know that f is bounded on some interval $[b - \delta, b]$. Hence, f is bounded on $[a, b - \delta]$ and on $[b - \delta, b]$, which means that f is bounded on $[a, b]$.

With the help of the above lemma, we can now prove the maximum–minimum theorem of Weierstrass.

Theorem A.17 (The Maximum–Minimum Theorem) If f is continuous on an interval $[a, b]$, then it takes on both its supremum and infimum values on $[a, b]$.

Proof By Lemma A.2, the function f is bounded on $[a, b]$. Let

$$M = \sup\{f(x) : x \in [a, b]\}.$$

We now show that there is a $c \in [a, b]$ such that $f(c) = M$. Let

$$g(x) = \frac{1}{M - f(x)}.$$

If f does not take on the value of M , then g is continuous on $[a, b]$, and by Lemma A.2 the function g is bounded on $[a, b]$. But g cannot be bounded on $[a, b]$. Thus, the assumption that f does not take on the value M has led us to a contradiction. The other case that f takes on its infimum value can be proved in a similar manner.

A *sphere*, ball, around \mathbf{x} in \mathbb{R}^p is the set of the form $\{\mathbf{y} : \|\mathbf{x} - \mathbf{y}\| < \varepsilon\}$ for some $\varepsilon > 0$. A set Ω in \mathbb{R}^p is said to be *open* if around every point \mathbf{x} in Ω there is a sphere that is contained in Ω .

If $\mathbf{x} \in \mathbb{R}^p$, then any set that contains an open set containing \mathbf{x} is called a *neighborhood* of \mathbf{x} . A sphere around \mathbf{x} is a neighborhood of \mathbf{x} .

A point $\mathbf{x} \in \mathbb{R}^p$ is called a *boundary point* of a set $\Omega \subseteq \mathbb{R}^p$ if every neighborhood of \mathbf{x} contains a point of Ω and a point that is not in Ω . The set of all boundary points of Ω is called the *boundary* of Ω . A set is said to be *closed* if it contains its boundary. Equivalently, a set Ω is closed if $\mathbf{x}_k \rightarrow \mathbf{g}$ with $\mathbf{x}_k \in \Omega$ implies $\mathbf{g} \in \Omega$. A set is *bounded* if it is contained within some sphere of finite radius. A set is *compact* if it is both closed and bounded. Any finite set $\Omega = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$ in \mathbb{R}^p is compact.

We now generalize Lemma A.2 for real-valued functions of several variables.

Lemma A.3 Let $f : \mathbb{R}^P \rightarrow \mathbb{R}$ be a continuous function on a compact set $\Omega \subset \mathbb{R}^P$. Then, f is bounded on Ω .

Proof We prove the lemma by contradiction. We assume that f is continuous on a compact set Ω , but f is not bounded on Ω . Without loss of generality, suppose that f is unbounded above on Ω . Since f is unbounded above on Ω , we can find a sequence $\{\mathbf{x}_n\} \subset \Omega$ such that $f(\mathbf{x}_n) > r(n)$, where $r(n) \rightarrow \infty$ as $n \rightarrow \infty$, that is, $f(\mathbf{x}_n) \rightarrow \infty$ as $n \rightarrow \infty$. The sequence $\{\mathbf{x}_n\}$ must be bounded because, by assumption, the set Ω is compact. Therefore, by the Bolzano-Weierstrass theorem, on page 645, there is a convergent subsequence, say $\{\mathbf{x}_{n_k}\}$, that converges to a point $\mathbf{x}^* \in \Omega$. The function is continuous on Ω . Hence,

$$\lim_{k \rightarrow \infty} f(\mathbf{x}_{n_k}) = f\left(\lim_{k \rightarrow \infty} \mathbf{x}_{n_k}\right) = f(\mathbf{x}^*). \quad (\text{A.71})$$

We arrived at a contradiction because, by assumption, the left-hand side of (A.71) is divergent, while $f(\mathbf{x}^*) < \infty$ because f is continuous on Ω . Therefore, a continuous function on a compact set in \mathbb{R}^P must be bounded.

We next prove the maximum–minimum theorem of Weierstrass for real-valued functions of several variables.

Theorem A.18 (Weierstrass) If $f : \mathbb{R}^P \rightarrow \mathbb{R}$ is a continuous function on a compact set $\Omega \subset \mathbb{R}^P$, then f achieves its supremum and infimum.

Proof The proof follows that of Theorem A.17, with the interval $[a, b]$ being replaced with the compact set Ω and by referring to Lemma A.3 rather than to Lemma A.2.

Definition A.1 (Uniform Continuity) Let $f : \mathcal{X} \rightarrow \mathcal{Y}$. Then, f is called uniformly continuous on \mathcal{X} if for every $\varepsilon > 0$ there exists $\delta = \delta(\varepsilon) > 0$, depending only on ε , such that $|x - \tilde{x}| < \delta$ implies $|f(x) - f(\tilde{x})| < \varepsilon$, where x, \tilde{x} are in \mathcal{X} .

Every uniformly continuous function is continuous, but not necessarily conversely. Indeed, $f(x) = 1/x$ is continuous on $(0, 1)$ but not uniformly continuous.

Another type of continuous functions that we use in the book are absolutely continuous functions.

Definition A.2 (Absolutely Continuous Function) Let f be a function defined on the closed interval $[a, b]$. Then, the function f is said to be absolutely continuous on $[a, b]$ if for every $\varepsilon > 0$, there exists a $\delta > 0$ such that

$$\sum_{k=1}^n (b_k - a_k) < \delta \quad (\text{A.72})$$

for all $a_1, b_1, \dots, a_n, b_n$ such that $a_1 < b_1 \leq a_2 < b_2 \leq \dots \leq a_n < b_n$, we have

$$\sum_{k=1}^n |f(b_k) - f(a_k)| < \varepsilon. \quad (\text{A.73})$$

Note that every absolutely continuous function is continuous, because the case $n = 1$ is not excluded in the above definition. A function f that satisfies the Lipschitz condition,

$$|f(x) - f(y)| \leq L|x - y|,$$

is absolutely continuous.

A.9 Linear Operators

A function

$$\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

is said to be *linear* if it satisfies the following two properties:

LO 1. For any \mathbf{x}, \mathbf{y} in \mathbb{R}^n , we have

$$\mathbf{F}(\mathbf{x} + \mathbf{y}) = \mathbf{F}(\mathbf{x}) + \mathbf{F}(\mathbf{y}).$$

LO 2. If c is a number, then

$$\mathbf{F}(c\mathbf{x}) = c\mathbf{F}(\mathbf{x}).$$

Linear functions are sometimes referred to as *linear maps* or *linear operators*.

Linear functions from \mathbb{R}^n into \mathbb{R}^m can be identified with the $m \times n$ real matrices as follows:

$$\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ is linear if and only if } \mathbf{F}(\mathbf{x}) = \mathbf{A}\mathbf{x},$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$.

If $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is linear, then there exists a constant $L \geq 0$ such that

$$\|\mathbf{F}(\mathbf{x})\| \leq L\|\mathbf{x}\| \tag{A.74}$$

for all $\mathbf{x} \in \mathbb{R}^n$. The smallest such constant L for which this is true is denoted $\|\mathbf{F}\|$ and called the norm of the linear map \mathbf{F} . If \mathbf{F} is the linear map induced by the matrix \mathbf{A} , then we define $\|\mathbf{A}\| = \|\mathbf{F}\|$. For any matrix \mathbf{A} and vector norm $\|\cdot\|$, we can compute the matrix norm using the formula

$$\|\mathbf{A}\| = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|}. \tag{A.75}$$

The matrix norm—that is, the norm of the linear map \mathbf{F} given by (A.75)—is equivalent to

$$\|\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\| \tag{A.76}$$

A matrix norm produced using (A.75) or (A.76) is called the *induced norm*. Note that

$$\|\mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{y}\| = \|\mathbf{A}(\mathbf{x} - \mathbf{y})\| \leq \|\mathbf{A}\|\|\mathbf{x} - \mathbf{y}\|;$$

that is, the linear map, \mathbf{F} , represented by the matrix \mathbf{A} is uniformly continuous.

We now comment on the notation used in (A.75) and (A.76). First, the same symbol is used for the matrix norm as for the norm of a vector. Second, the norms appearing on the right-hand side of (A.75) and (A.76) are the norms of the vectors $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{A}\mathbf{x} \in \mathbb{R}^m$. Thus, if \mathbb{R}^n is equipped with the vector norm $\|\cdot\|_s$ while \mathbb{R}^m is equipped with the vector norm $\|\cdot\|_t$, then (A.76)

should be interpreted as

$$\|\mathbf{A}\| = \max_{\|\mathbf{x}\|_s=1} \|\mathbf{Ax}\|_t.$$

Vector norms and a matrix norm are said to be *compatible*, or *consistent*, if they satisfy

$$\|\mathbf{Ax}\|_t \leq \|\mathbf{A}\| \|\mathbf{x}\|_s. \quad (\text{A.77})$$

We will show that (A.76) is well-defined. For this, we need the following lemma.

Lemma A.4 A vector norm $\|\cdot\| : \mathbb{R}^p \rightarrow \mathbb{R}$ is an uniformly continuous function.

Proof By Theorem A.4 on page 629, we have

$$|\|\mathbf{x}\| - \|\mathbf{y}\|| \leq \|\mathbf{x} - \mathbf{y}\|.$$

Hence, for any $\varepsilon > 0$, if we take $\delta = \varepsilon$, then for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$, if $\|\mathbf{x} - \mathbf{y}\| < \delta$, then $|\|\mathbf{x} - \mathbf{y}\|| < \varepsilon$, as was to be shown.

By the continuity of the vector norm and the theorem of Weierstrass on page 650, a vector \mathbf{x}_0 can be found such that $\|\mathbf{x}_0\| = 1$ and $\|\mathbf{A}\| = \|\mathbf{Ax}_0\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{Ax}\|$. Thus, the expression (A.76) is well-defined.

Next, we will show that (A.76) satisfies (A.77). Suppose that we are given a vector $\mathbf{y} \neq \mathbf{0}$. Then,

$$\mathbf{x} = \frac{1}{\|\mathbf{y}\|} \mathbf{y}$$

satisfies the condition $\|\mathbf{x}\| = 1$. Hence,

$$\|\mathbf{Ay}\| = \|\mathbf{A}(\|\mathbf{y}\|\mathbf{x})\| = \|\mathbf{y}\| \|\mathbf{Ax}\|.$$

Applying (A.76) to the right-hand side of the above yields

$$\|\mathbf{Ay}\| \leq \|\mathbf{y}\| \|\mathbf{A}\|.$$

Hence, any matrix norm produced by (A.76) is consistent.

We will now prove that a matrix norm obtained using (A.76) satisfies the following conditions:

MN 1. $\|\mathbf{A}\| > 0$ if $\mathbf{A} \neq \mathbf{O}$ and $\|\mathbf{O}\| = 0$.

MN 2. $\|c\mathbf{A}\| = |c| \|\mathbf{A}\|$, where c is a constant.

MN 3. $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$.

MN 4. $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$.

Proof of MN 1. Let $\mathbf{A} \neq \mathbf{O}$. Then a vector \mathbf{x} , $\|\mathbf{x}\| = 1$, can be found such that $\mathbf{Ax} \neq \mathbf{0}$, and thus $\|\mathbf{Ax}\| \neq 0$. Hence, $\|\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{Ax}\| \neq 0$. If, on the other hand, $\mathbf{A} = \mathbf{O}$, then $\|\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{Ox}\| = 0$.

Proof of MN 2. By definition, $\|c\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} \|c\mathbf{A}\mathbf{x}\|$. Using a property of the vector norm, we can write $\|c\mathbf{A}\mathbf{x}\| = |c|\|\mathbf{A}\mathbf{x}\|$, and hence

$$\|c\mathbf{A}\| = \max_{\|\mathbf{x}\|=1} |c|\|\mathbf{A}\mathbf{x}\| = |c| \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\| = |c|\|\mathbf{A}\|.$$

Proof of MN 3. For the sum of two matrices, $\mathbf{A} + \mathbf{B}$, there exists a vector \mathbf{x}_0 , $\|\mathbf{x}_0\| = 1$, such that

$$\|\mathbf{A} + \mathbf{B}\| = \|(\mathbf{A} + \mathbf{B})\mathbf{x}_0\|.$$

Then, using properties of the vector norm, we obtain

$$\|\mathbf{A} + \mathbf{B}\| = \|\mathbf{A}\mathbf{x}_0 + \mathbf{B}\mathbf{x}_0\| \leq \|\mathbf{A}\mathbf{x}_0\| + \|\mathbf{B}\mathbf{x}_0\|.$$

Using now (A.76) yields

$$\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\|\|\mathbf{x}_0\| + \|\mathbf{B}\|\|\mathbf{x}_0\| = \|\mathbf{A}\| + \|\mathbf{B}\|.$$

Proof of MN 4. For the product of two matrices $\mathbf{A}\mathbf{B}$, there exists a vector \mathbf{x}_0 , $\|\mathbf{x}_0\| = 1$, such that $\|\mathbf{A}\mathbf{B}\| = \|\mathbf{A}\mathbf{B}\mathbf{x}_0\|$. Then, using (A.76), we have

$$\|\mathbf{A}\mathbf{B}\| = \|\mathbf{A}(\mathbf{B}\mathbf{x}_0)\| \leq \|\mathbf{A}\|\|\mathbf{B}\mathbf{x}_0\|.$$

Employing (A.76) again yields

$$\|\mathbf{A}\mathbf{B}\| \leq \|\mathbf{A}\|\|\mathbf{B}\|\|\mathbf{x}_0\| = \|\mathbf{A}\|\|\mathbf{B}\|.$$

The most frequently used matrix norms are the p -norms,

$$\|\mathbf{A}\|_p = \max_{\|\mathbf{x}\|_p=1} \|\mathbf{A}\mathbf{x}\|_p.$$

The induced matrix norms corresponding to the one-, two-, and infinity-norm for vectors are:

- $\|\mathbf{A}\|_1 = \max_j \|\mathbf{a}_j\|_1$, that is, $\|\mathbf{A}\|_1$ is the maximum of absolute column sum;
- $\|\mathbf{A}\|_2 = \sigma_1(\mathbf{A})$, which is the largest singular value of \mathbf{A} ;
- $\|\mathbf{A}\|_\infty = \max_i \|\mathbf{a}_i^T\|_1$, which is the maximum of absolute row sum.

A matrix norm that is not induced by a vector norm is the *Frobenius norm*, defined for an $m \times n$ matrix \mathbf{A} as

$$\|\mathbf{A}\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2 \right)^{1/2}. \quad (\text{A.78})$$

The Frobenius norm satisfies the relation

$$\|\mathbf{A}\|_F^2 = \text{trace}(\mathbf{A}\mathbf{A}^T) = \text{trace}(\mathbf{A}^T\mathbf{A}). \quad (\text{A.79})$$

The Frobenius norm is compatible with the Euclidean vector norm, that is,

$$\|\mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{A}\|_F \|\mathbf{x}\|_2. \quad (\text{A.80})$$

A.10 Vector Spaces

A given set of vectors of the same dimension $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r\}$ is *linearly dependent* if the zero vector, $\mathbf{0}$, can be written as a nontrivial *linear combination* of the given vectors:

$$\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \cdots + \alpha_r \mathbf{x}_r = \mathbf{0}, \quad \text{with } \alpha_i \neq 0 \text{ for at least one } i.$$

In other words, the set of vectors $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r\}$ is *linearly independent* if

$$\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \cdots + \alpha_r \mathbf{x}_r = \mathbf{0}, \quad \text{implies } \alpha_1 = \alpha_2 = \cdots = \alpha_r = 0.$$

A set \mathcal{S} containing vectors of dimension p is a *subspace* of \mathbb{R}^p if for any scalars α and β ,

$$\mathbf{x}, \mathbf{y} \in \mathcal{S} \quad \text{implies } \alpha \mathbf{x} + \beta \mathbf{y} \in \mathcal{S}. \quad (\text{A.81})$$

It follows from the above that a subspace must contain the zero vector, $\mathbf{0}$. The set containing only the zero vector is a subspace. This set, $\{\mathbf{0}\}$, is called a *trivial* subspace.

A matrix–vector product

$$\mathbf{Ax} = \begin{bmatrix} a_1 & a_2 & \cdots & a_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \cdots + x_n \mathbf{a}_n \quad (\text{A.82})$$

is a linear combination of the columns of the matrix. Thus, linear dependence of the columns of \mathbf{A} is equivalent to the condition

$$\mathbf{Az} = \mathbf{0} \quad \text{for some vector } \mathbf{z} \neq \mathbf{0}. \quad (\text{A.83})$$

Linear independence of the columns of \mathbf{A} is equivalent to the condition

$$\mathbf{Az} = \mathbf{0} \quad \text{implies } \mathbf{z} = \mathbf{0}. \quad (\text{A.84})$$

The set of all vectors that are linear combinations of the columns of \mathbf{A} is called the *range*, or *image*, of \mathbf{A} and is denoted $\text{range}(\mathbf{A})$. Thus,

$$\text{range}(\mathbf{A}) = \{\mathbf{v} : \mathbf{v} = \mathbf{Ax} \text{ for some } \mathbf{x}\}. \quad (\text{A.85})$$

For a given $m \times n$ matrix \mathbf{A} and an m -vector \mathbf{b} , we have $\mathbf{b} \in \text{range}(\mathbf{A})$ if and only if there exists an n -vector \mathbf{x} such that $\mathbf{Ax} = \mathbf{b}$. We say that a system of linear equations represented by $\mathbf{Ax} = \mathbf{b}$ is compatible if $\mathbf{b} \in \text{range}(\mathbf{A})$.

The range of \mathbf{A} is a subspace of \mathbb{R}^m . Indeed, for any $\mathbf{b}_1, \mathbf{b}_2 \in \text{range}(\mathbf{A})$, we have $\mathbf{b}_1 = \mathbf{Ax}_1$ and $\mathbf{b}_2 = \mathbf{Ax}_2$ for some \mathbf{x}_1 and \mathbf{x}_2 . Using the fact that \mathbf{A} represents a linear map, for any scalars α_1 and α_2 , we obtain

$$\alpha_1 \mathbf{b}_1 + \alpha_2 \mathbf{b}_2 = \alpha_1 \mathbf{Ax}_1 + \alpha_2 \mathbf{Ax}_2 = \mathbf{A}(\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2). \quad (\text{A.86})$$

Hence, if $\mathbf{b}_1, \mathbf{b}_2 \in \text{range}(\mathbf{A})$, then $\alpha_1 \mathbf{b}_1 + \alpha_2 \mathbf{b}_2 \in \text{range}(\mathbf{A})$, which means that the range of \mathbf{A} is a subspace of \mathbb{R}^m .

Similarly, we define the range of \mathbf{A}^T to be the set

$$\text{range}(\mathbf{A}^T) = \{\mathbf{y} : \mathbf{y} = \mathbf{A}^T \mathbf{z} \text{ for some } \mathbf{z}\}. \quad (\text{A.87})$$

For a nontrivial subspace \mathcal{S} there is a unique smallest positive integer called the *dimension* of the subspace \mathcal{S} and denoted $\dim \mathcal{S}$. The dimension of the subspace is equal to the smallest number of linearly independent vectors of the subspace such that every vector of the subspace

can be expressed as a linear combination of these vectors. Any such set of vectors is called a *basis* of the subspace. The vectors forming a basis are not unique.

The number of linearly independent columns of a matrix A is called the *rank* of A and is denoted $\text{rank}(A)$. The rank of A is also equal to the number of linearly independent rows of A .

For the product of an $m \times n$ matrix A and an $n \times p$ matrix B , we have *Sylvester's inequalities*:

$$\text{rank}(A) + \text{rank}(B) - n \leq \text{rank}(AB) \leq \min\{\text{rank}(A), \text{rank}(B)\}. \quad (\text{A.88})$$

For an $m \times n$ matrix A , the set of all n -vectors orthogonal to the rows of A is called the *kernel*, or *null space*, of A and is denoted $\ker(A)$, that is,

$$\ker(A) = \{\mathbf{x} : A\mathbf{x} = \mathbf{0}\}. \quad (\text{A.89})$$

The kernel of A is a subspace of \mathbb{R}^n . Indeed, if $\mathbf{x}_1, \mathbf{x}_2 \in \ker(A)$, then $A\mathbf{x}_1 = \mathbf{0}$ and $A\mathbf{x}_2 = \mathbf{0}$. Hence, for any scalars α_1 and α_2 , we have

$$A(\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2) = \alpha_1 A\mathbf{x}_1 + \alpha_2 A\mathbf{x}_2 = \mathbf{0}. \quad (\text{A.90})$$

Thus, if $\mathbf{x}_1, \mathbf{x}_2 \in \ker(A)$, then $\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2 \in \ker(A)$, which means that $\ker(A)$ is a subspace of \mathbb{R}^n . Similarly, we define the kernel of A^T to be the set

$$\ker(A^T) = \{\mathbf{w} : A^T\mathbf{w} = \mathbf{0}\}. \quad (\text{A.91})$$

The intersection of $\text{range}(A)$ and $\ker(A^T)$ contains only the zero vector, that is,

$$\text{range}(A) \cap \ker(A^T) = \{\mathbf{0}\}. \quad (\text{A.92})$$

For a given subspace S of \mathbb{R}^n , the orthogonal complement S^\perp of S is defined as

$$S^\perp = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x}^T \mathbf{y} = 0 \text{ for all } \mathbf{y} \in S\}. \quad (\text{A.93})$$

The vectors in $\ker(A^T)$ are orthogonal to the vectors in $\text{range}(A)$, and vice versa. We say that $\text{range}(A)$ and $\ker(A^T)$ are *orthogonal complements* of one another and write

$$\text{range}(A)^\perp = \ker(A^T) \quad (\text{A.94})$$

and

$$\ker(A^T)^\perp = \text{range}(A). \quad (\text{A.95})$$

Similarly,

$$\text{range}(A^T)^\perp = \ker(A) \quad (\text{A.96})$$

and

$$\ker(A)^\perp = \text{range}(A^T). \quad (\text{A.97})$$

If A is an $m \times n$ matrix with linearly independent columns, which implies that $m \geq n$, then

$$\text{rank}(A) = n = \dim \text{range}(A)$$

and

$$\dim \ker(A^T) = m - n.$$

We now define a normed space.

Definition A.3 A vector space with a norm is called a *normed space*.

By the normed vector space \mathbb{R}^n , we usually mean the Euclidean space, which is \mathbb{R}^n with the norm $\|\mathbf{x}\|_2 = \sqrt{x_1^2 + \cdots + x_n^2}$. The vector space \mathbb{R}^n equipped with any of the p -norm forms a normed vector space. The absolute value is a norm in \mathbb{R} .

Let Ω be a compact subset of \mathbb{R}^n and let $\mathcal{C}(\Omega)$ denote the set of continuous functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ on Ω . Thus, we view functions as vectors in the function space $\mathcal{C}(\Omega)$. Then, the function $\|\cdot\| : \mathcal{C}(\Omega) \rightarrow \mathbb{R}$ defined as

$$\|f\| = \max_{x \in \Omega} |f(x)| \quad (\text{A.98})$$

satisfies the norm axioms. The function space $\mathcal{C}(\Omega)$, with scalars from \mathbb{R} and with the norm (A.98), is a linear normed space. In particular, the space of all real-valued functions, $f(x)$, continuous on some interval $[a, b]$ of the real line with the norm

$$\|f\| = \max_{x \in [a, b]} |f(x)| \quad (\text{A.99})$$

is an example of a linear normed space. This space is denoted $\mathcal{C}([a, b])$.

We are also interested in sequences of elements of $\mathcal{C}([a, b])$. Let f_1, f_2, \dots be a sequence of elements of $\mathcal{C}([a, b])$. Denote this sequence as $\{f_k\}$. In the following, $f_k(x)$, refers to the value of the function f_k at x , while f_k denotes the function itself. We say that the sequence $\{f_k\}$ converges (pointwise) on the interval $[a, b]$ if there exists a function f such that for every $x \in [a, b]$ the sequence $\{f_k(x)\}$ converges to $f(x)$. In other words, the sequence $\{f_k\}$ converges (pointwise) if there exists a function f such that for a given $x \in [a, b]$ and a given $\varepsilon > 0$, we can find $N = N(x, \varepsilon)$ such that

$$|f(x) - f_k(x)| < \varepsilon \quad \text{for } k > N. \quad (\text{A.100})$$

The sequence $\{f_k\}$ converges uniformly, on the interval $[a, b]$, if there is a function f such that for any $\varepsilon > 0$ there exists $N = N(\varepsilon)$, depending on ε but not on x , such that (A.100) holds for every $x \in [a, b]$. The distinction between convergence and uniform convergence is that in the later case the same N can be used for any value of x to show the convergence of the sequence $\{f_k(x)\}$. One can show that the limit of an uniformly convergent sequence of continuous functions is continuous. It can be verified that the norm (A.99) defines the uniform convergence because $\{f_k\}$ converges uniformly to f if and only if $\max_{x \in [a, b]} |f(x) - f_k(x)| \rightarrow 0$. This is the reason the norm (A.99) is sometimes called the *norm of uniform convergence*. It can be shown that there is no norm on $\mathcal{C}([a, b])$ that defines pointwise convergence to an element in $\mathcal{C}([a, b])$ —see Debnath and Mikusiński [58, p. 12].

Definition A.4 (Banach Space) A normed linear space is complete if every Cauchy sequence with elements from the space converges to an element in the space. A complete normed linear space is called a Banach space.

A sequence that satisfies conditions of Theorem A.14 is called a *Cauchy sequence*. The Euclidean space is an example of a Banach space, as is $\mathcal{C}([a, b])$ with the norm (A.99).

A.11 Least Squares

The theory of least squares was developed by Carl Friedrich Gauss (1777–1855) and published by him in 1809. Gauss first applied the least squares to devise a method of calculating orbits of celestial bodies. For more information on the history of the least squares development, we recommend an account by Campbell and Meyer [40, Section 2.5].

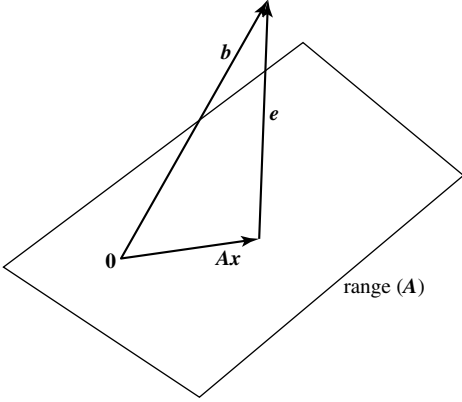


Figure A.3 Geometrical interpretation of the least squares solution to a system of linear equations.

Consider the system of equations

$$Ax = b, \quad (\text{A.101})$$

where $A \in \mathbb{R}^{m \times n}$, $x \in \mathbb{R}^n$, and $b \in \mathbb{R}^m$. A necessary and sufficient condition for the system (A.101) to have a solution is

$$\text{rank}(A) = \text{rank}[A : b]. \quad (\text{A.102})$$

In particular, if the matrix A is invertible, then there is a unique solution to (A.101) given by

$$x^* = A^{-1}b.$$

Here, we analyze the case when $m > n$ and $\text{rank}(A) < \text{rank}[A : b]$, that is, $b \notin \text{range}(A)$, which means that the system (A.101) does not have an exact solution. We modify (A.101) by incorporating an error vector e to obtain

$$Ax + e = b. \quad (\text{A.103})$$

This is illustrated in Figure A.3. Because

$$e = b - Ax, \quad (\text{A.104})$$

we have

$$\begin{aligned} \|e\|^2 &= e^T e \\ &= (b - Ax)^T (b - Ax) \\ &= x^T A^T A x - 2x^T A^T b + b^T b. \end{aligned} \quad (\text{A.105})$$

The gradient of $\|e\|^2$ with respect to x is

$$\nabla_x \|e\|^2 = 2A^T A x - 2A^T b. \quad (\text{A.106})$$

Solving

$$\nabla_x \|e\|^2 = 0$$

yields candidate minimizer points of (A.105). If the matrix A has full column rank, then the matrix $A^T A$ is invertible and the only critical point is

$x^* = (A^T A)^{-1} A^T b$

(A.107)

The above point satisfies the second-order sufficient condition to be the minimizer of (A.105) because the Hessian of (A.105) is

$$2\mathbf{A}^T \mathbf{A},$$

which is positive definite. The solution given by (A.107) is called the least squares solution of the inconsistent linear system of equations, $\mathbf{A}\mathbf{x} = \mathbf{b}$, because it minimizes the sum of squares $\|\mathbf{e}\|^2 = \mathbf{e}^T \mathbf{e} = \sum_{i=1}^m e_i^2$. For a graphical interpretation of the least squares solution, we refer to Figure A.3.

Suppose now that we interpret the i th row of the matrix \mathbf{A} and the i th component of the vector \mathbf{b} as the i th data pair. Thus, the i th data pair has the form $(\mathbf{a}_i^T; b_i)$. Suppose that the k data pairs are represented by the pair $(\mathbf{A}; \mathbf{b})$. We know that the least squares solution to the system of equations represented by the data pair $(\mathbf{A}; \mathbf{b})$ is given by (A.107). If the new data pair, $(\mathbf{a}_{k+1}^T; b_{k+1})$, becomes available, then we could find the least squares solution to the system of equations

$$\begin{bmatrix} \mathbf{A} \\ \mathbf{a}_{k+1}^T \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{b} \\ b_{k+1} \end{bmatrix} \quad (\text{A.108})$$

by applying the formula (A.107) to (A.108). Let \mathbf{x}_{k+1} denote the least squares solution to the above problem. Then,

$$\mathbf{x}_{k+1} = \left(\begin{bmatrix} \mathbf{A} \\ \mathbf{a}_{k+1}^T \end{bmatrix}^T \begin{bmatrix} \mathbf{A} \\ \mathbf{a}_{k+1}^T \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{A} \\ \mathbf{a}_{k+1}^T \end{bmatrix}^T \begin{bmatrix} \mathbf{b} \\ b_{k+1} \end{bmatrix}. \quad (\text{A.109})$$

In the case of large k the above method of recomputing the new least squares solution from scratch may be time-consuming. We present a recursive method for computing least squares solutions as the data pairs become available. To proceed, let

$$\mathbf{P}_k = (\mathbf{A}^T \mathbf{A})^{-1}. \quad (\text{A.110})$$

In view of the above, we can write

$$\begin{aligned} \mathbf{P}_{k+1}^{-1} &= \begin{bmatrix} \mathbf{A} \\ \mathbf{a}_{k+1}^T \end{bmatrix}^T \begin{bmatrix} \mathbf{A} \\ \mathbf{a}_{k+1}^T \end{bmatrix} \\ &= [\mathbf{A}^T \quad \mathbf{a}_{k+1}^T] \begin{bmatrix} \mathbf{A} \\ \mathbf{a}_{k+1}^T \end{bmatrix} \\ &= \mathbf{A}^T \mathbf{A} + \mathbf{a}_{k+1} \mathbf{a}_{k+1}^T \\ &= \mathbf{P}_k^{-1} + \mathbf{a}_{k+1} \mathbf{a}_{k+1}^T. \end{aligned} \quad (\text{A.111})$$

Hence,

$$\mathbf{P}_{k+1} = (\mathbf{P}_k^{-1} + \mathbf{a}_{k+1} \mathbf{a}_{k+1}^T)^{-1}.$$

Applying to the above the matrix inversion formula (A.30), we obtain

$$\boxed{\mathbf{P}_{k+1} = \mathbf{P}_k - \frac{\mathbf{P}_k \mathbf{a}_{k+1} \mathbf{a}_{k+1}^T \mathbf{P}_k}{1 + \mathbf{a}_{k+1}^T \mathbf{P}_k \mathbf{a}_{k+1}}}. \quad (\text{A.112})$$

Using (A.110), we write

$$\begin{aligned}\mathbf{x}_k &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b} \\ &= \mathbf{P}_k \mathbf{A}^T \mathbf{b}.\end{aligned}\tag{A.113}$$

Employing (A.111), we obtain

$$\begin{aligned}\mathbf{x}_{k+1} &= \left(\begin{bmatrix} \mathbf{A} \\ \mathbf{a}_{k+1}^T \end{bmatrix}^T \begin{bmatrix} \mathbf{A} \\ \mathbf{a}_{k+1}^T \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{A} \\ \mathbf{a}_{k+1}^T \end{bmatrix}^T \begin{bmatrix} \mathbf{b} \\ b_{k+1} \end{bmatrix} \\ &= \mathbf{P}_{k+1} (\mathbf{A}^T \mathbf{b} + \mathbf{a}_{k+1} b_{k+1}).\end{aligned}\tag{A.114}$$

Computing $\mathbf{A}^T \mathbf{b}$ from (A.113), we get

$$\mathbf{A}^T \mathbf{b} = \mathbf{P}_k^{-1} \mathbf{x}_k.\tag{A.115}$$

Substituting the above into (A.114) yields

$$\begin{aligned}\mathbf{x}_{k+1} &= \mathbf{P}_{k+1} (\mathbf{A}^T \mathbf{b} + \mathbf{a}_{k+1} b_{k+1}) \\ &= \mathbf{P}_{k+1} (\mathbf{P}_k^{-1} \mathbf{x}_k + \mathbf{a}_{k+1} b_{k+1}).\end{aligned}\tag{A.116}$$

Computing \mathbf{P}_k^{-1} from (A.111) gives

$$\mathbf{P}_k^{-1} = \mathbf{P}_{k+1}^{-1} - \mathbf{a}_{k+1} \mathbf{a}_{k+1}^T$$

Substituting the above into (A.116), we have

$$\mathbf{x}_{k+1} = \mathbf{P}_{k+1} ((\mathbf{P}_{k+1}^{-1} - \mathbf{a}_{k+1} \mathbf{a}_{k+1}^T) \mathbf{x}_k + \mathbf{a}_{k+1} b_{k+1}).\tag{A.117}$$

Performing simple manipulations yields

$$\boxed{\mathbf{x}_{k+1} = \mathbf{x}_k + (b_{k+1} - \mathbf{a}_{k+1}^T \mathbf{x}_k) \mathbf{P}_{k+1} \mathbf{a}_{k+1}.}\tag{A.118}$$

The two equations (A.112) and (A.118) constitute the *recursive least squares* (RLS) algorithm. To initialize the RLS, we can set

$$\mathbf{P}_0 = (\mathbf{A}_n^T \mathbf{A}_n)^{-1} \quad \text{and} \quad \mathbf{x}_0 = \mathbf{P}_n \mathbf{A}_n^T \mathbf{b}_n,$$

where the pair $(\mathbf{A}_n; \mathbf{b}_n)$ represents the first n data pairs.

A.12 Contraction Maps

We now discuss the contraction maps that found interesting applications in control and stability analyses of nonlinear systems. In particular, in 1964, Leibovic [180] presented an interesting idea of how to apply contraction maps in nonlinear and adaptive control. The proposed application required a significant computational power not available at that time. Leibovic's idea looks like a good challenge for today's computers. More recently, Kelly [152] used the contraction map theorem as well as Lyapunov functions to investigate stability properties of a class of neural networks.

An operator \mathbf{F} that maps Euclidean space \mathbb{R}^p into itself is called a *contraction map* if for any two elements \mathbf{x} and \mathbf{y} of \mathbb{R}^p ,

$$\|\mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|, \quad 0 < L < 1.\tag{A.119}$$

A point $\mathbf{x}^* \in \mathbb{R}^p$ is called a *fixed point* of \mathbf{F} if

$$\mathbf{F}(\mathbf{x}^*) = \mathbf{x}^*. \quad (\text{A.120})$$

The importance of the contraction operators lies in the fact that there is a unique fixed point. Furthermore, this point may be found by successive iterations as stated in the following theorem.

Theorem A.19 Let $\mathbf{F} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ be a contraction map with the contraction constant $0 < L < 1$. Then \mathbf{F} has a unique fixed point, and the sequence

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$$

generated by $\mathbf{x}^{(n+1)} = \mathbf{F}(\mathbf{x}^{(n)})$, $n = 0, 1, 2, \dots$, converges to \mathbf{x}^* for any initial choice of $\mathbf{x}^{(0)} \in \mathbb{R}^p$.

Proof To show the uniqueness of the fixed point \mathbf{x}^* , suppose that there exist two distinct fixed points \mathbf{x}^* and \mathbf{y}^* . Then,

$$\begin{aligned} \|\mathbf{x}^* - \mathbf{y}^*\| &= \|\mathbf{F}(\mathbf{x}^*) - \mathbf{F}(\mathbf{y}^*)\| \\ &\leq L \|\mathbf{x}^* - \mathbf{y}^*\|. \end{aligned}$$

Thus,

$$\|\mathbf{x}^* - \mathbf{y}^*\| \leq L \|\mathbf{x}^* - \mathbf{y}^*\|, \quad 0 < L < 1,$$

which is impossible unless $\|\mathbf{x}^* - \mathbf{y}^*\| = 0$. Hence, $\mathbf{x}^* = \mathbf{y}^*$ which means that \mathbf{F} has at most one fixed point.

We will now show that a fixed point exists and that the sequence

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$$

generated by $\mathbf{x}^{(n+1)} = \mathbf{F}(\mathbf{x}^{(n)})$, $n = 0, 1, 2, \dots$, converges to \mathbf{x}^* for any initial choice of $\mathbf{x}^{(0)} \in \mathbb{R}^p$. If $\mathbf{F}(\mathbf{x}^{(0)}) = \mathbf{x}^{(0)}$, then the theorem is proved. In general, $\mathbf{F}(\mathbf{x}^{(0)}) \neq \mathbf{x}^{(0)}$. For any n, m with $m > n$, we have

$$\begin{aligned} \|\mathbf{x}^{(m)} - \mathbf{x}^{(n)}\| &= \|\mathbf{F}^m(\mathbf{x}^{(0)}) - \mathbf{F}^n(\mathbf{x}^{(0)})\| \\ &= \|\mathbf{F}(\mathbf{F}^{m-1}(\mathbf{x}^{(0)})) - \mathbf{F}(\mathbf{F}^{n-1}(\mathbf{x}^{(0)}))\| \\ &\leq L \|\mathbf{F}^{m-1}(\mathbf{x}^{(0)}) - \mathbf{F}^{n-1}(\mathbf{x}^{(0)})\| \\ &\vdots \\ &\leq L^n \|\mathbf{x}^{(m-n)} - \mathbf{x}^{(0)}\| \\ &= L^n \|\mathbf{x}^{(0)} - \mathbf{x}^{(m-n)}\|. \end{aligned} \quad (\text{A.121})$$

We write

$$\begin{aligned} \|\mathbf{x}^{(0)} - \mathbf{x}^{(m-n)}\| &= \|\mathbf{x}^{(0)} - \mathbf{x}^{(1)} + \mathbf{x}^{(1)} - \mathbf{x}^{(2)} + \mathbf{x}^{(2)} \\ &\quad - \dots - \mathbf{x}^{(m-n-1)} + \mathbf{x}^{(m-n-1)} - \mathbf{x}^{(m-n)}\|. \end{aligned} \quad (\text{A.122})$$

Substituting (A.122) into (A.121) and then applying the triangle inequality yields

$$\|\mathbf{x}^{(m)} - \mathbf{x}^{(n)}\| \leq L^n (\|\mathbf{x}^{(0)} - \mathbf{x}^{(1)}\| + \|\mathbf{x}^{(1)} - \mathbf{x}^{(2)}\| + \dots + \|\mathbf{x}^{(m-n-1)} - \mathbf{x}^{(m-n)}\|). \quad (\text{A.123})$$

Observe that

$$\begin{aligned}
 \| \mathbf{x}^{(j)} - \mathbf{x}^{(j+1)} \| &= \| \mathbf{F}(\mathbf{x}^{(j-1)}) - \mathbf{F}(\mathbf{x}^{(j)}) \| \\
 &\leq L \| \mathbf{x}^{(j-1)} - \mathbf{x}^{(j)} \| \\
 &\vdots \\
 &\leq L^j \| \mathbf{x}^{(0)} - \mathbf{x}^{(1)} \|,
 \end{aligned} \tag{A.124}$$

We use the above to further evaluate the right-hand side of (A.123) to obtain

$$\begin{aligned}
 \| \mathbf{x}^{(m)} - \mathbf{x}^{(n)} \| &\leq L^n \| \mathbf{x}^{(0)} - \mathbf{x}^{(1)} \| (1 + L + \dots + L^{m-n-1}) \\
 &\leq L^n \| \mathbf{x}^{(0)} - \mathbf{x}^{(1)} \| (1 + L + \dots + L^{m-n-1} + L^{m-n} + \dots) \\
 &= \frac{L^n}{1-L} \| \mathbf{x}^{(0)} - \mathbf{x}^{(1)} \|
 \end{aligned}$$

because $0 < L < 1$. Recall that a sequence $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots$ is a Cauchy sequence if for each $\varepsilon > 0$ there exists a positive integer $n = n(\varepsilon)$ such that $\| \mathbf{x}^{(m)} - \mathbf{x}^{(n)} \| < \varepsilon$ for all $m > n$. Given an $\varepsilon > 0$, let n be such that

$$L^n < \frac{\varepsilon(1-L)}{\| \mathbf{x}^{(0)} - \mathbf{x}^{(1)} \|}.$$

Then,

$$\| \mathbf{x}^{(m)} - \mathbf{x}^{(n)} \| \leq \frac{L^n}{1-L} \| \mathbf{x}^{(0)} - \mathbf{x}^{(1)} \| < \varepsilon \quad \text{for all } m > n. \tag{A.125}$$

Therefore, the sequence $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots$ is a Cauchy sequence, and we will now show that the limit \mathbf{g} of this sequence is the fixed point \mathbf{x}^* of the map \mathbf{F} . Note that \mathbf{F} is uniformly continuous. Indeed, for each $\varepsilon > 0$ there exists a $\delta = \delta(\varepsilon)$ such that if $\| \mathbf{x} - \mathbf{y} \| < \delta$, then $\| \mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{y}) \| < \varepsilon$. Since $\| \mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{y}) \| \leq L \| \mathbf{x} - \mathbf{y} \|$, given an $\varepsilon > 0$, we can choose $\delta = \varepsilon/L$. Because \mathbf{F} is continuous, we have

$$\lim_{n \rightarrow \infty} \mathbf{F}(\mathbf{x}^{(n)}) = \mathbf{F}\left(\lim_{n \rightarrow \infty} \mathbf{x}^{(n)}\right) = \mathbf{F}(\mathbf{g}).$$

However, $\mathbf{x}^{(n+1)} = \mathbf{F}(\mathbf{x}^{(n)})$. Therefore,

$$\begin{aligned}
 \mathbf{g} &= \lim_{n \rightarrow \infty} \mathbf{x}^{(n+1)} \\
 &= \lim_{n \rightarrow \infty} \mathbf{F}(\mathbf{x}^{(n)}) \\
 &= \mathbf{F}\left(\lim_{n \rightarrow \infty} \mathbf{x}^{(n)}\right) \\
 &= \mathbf{F}(\mathbf{g}).
 \end{aligned}$$

Hence, $\mathbf{g} = \mathbf{x}^*$ is a fixed point of \mathbf{F} , and the proof is complete.

Using the above proof we can give an estimate of the norm of $\| \mathbf{x}^{(n)} - \mathbf{x}^* \|$. Indeed, if we fix n in (A.125) and allow $m \rightarrow \infty$, then

$$\| \mathbf{x}^{(n)} - \mathbf{x}^* \| \leq \frac{L^n}{1-L} \| \mathbf{x}^{(0)} - \mathbf{F}(\mathbf{x}^{(0)}) \|. \tag{A.126}$$

A.13 First-Order Differential Equation

Let \mathbb{R} denote the set of real numbers and let $D \subseteq \mathbb{R}^2$ be a domain, that is, an open connected nonempty subset of \mathbb{R}^2 . Let f be a real-valued function defined and continuous on D . Let $\dot{x} = dx/dt$ denote the derivative of x with respect to t . We call

$$\dot{x} = f(t, x) \quad (\text{A.127})$$

an *ordinary differential equation of the first order*. By a *solution* of the differential equation (A.127) on an open interval $I = \{t \in \mathbb{R} : t_1 < t < t_2\}$, we mean a real-valued continuously differentiable function ϕ defined on I such that $(t, \phi(t)) \in D$ for all $t \in I$ and

$$\dot{\phi}(t) = f(t, \phi(t)) \quad \text{for all } t \in I.$$

Definition A.5 Given $(t_0, x_0) \in D$, the *initial value problem* for (A.127) is

$$\dot{x} = f(t, x), \quad x(t_0) = x_0. \quad (\text{A.128})$$

A function ϕ is a solution to (A.128) if ϕ is a solution of the differential equation (A.127) on some interval I containing t_0 and $x(t_0) = x_0$.

We can represent the initial value problem (A.128) equivalently by an integral equation,

$$\phi(t) = x_0 + \int_{t_0}^t f(s, \phi(s)) ds. \quad (\text{A.129})$$

To prove this equivalence, let ϕ be a solution of the initial value problem (A.128). Then $x(t_0) = x_0$ and

$$\dot{\phi}(t) = f(t, \phi(t)) \quad \text{for all } t \in I.$$

Integrating the above from t_0 to t , we obtain

$$\int_{t_0}^t \dot{\phi}(s) ds = \int_{t_0}^t f(s, \phi(s)) ds.$$

Hence,

$$\phi(t) - x_0 = \int_{t_0}^t f(s, \phi(s)) ds.$$

Therefore, ϕ is a solution of the integral equation (A.129).

Conversely, let ϕ be a solution of the integral equation (A.129). Then $\phi(t_0) = x_0$, and differentiating both sides of (A.129) with respect to t gives $\dot{\phi}(t) = f(t, \phi(t))$. Therefore, ϕ is also a solution of the initial value problem (A.128).

The general *ordinary differential linear equation of the first order* has the form

$$\dot{x}(t) = a(t)x(t) + b(t), \quad t \in (t_1, t_2), \quad (\text{A.130})$$

where $a(t)$ and $b(t)$ are assumed to be continuous on the interval (t_1, t_2) . We associate with equation (A.130) the initial condition

$$x(t_0) = x_0, \quad x_0 \in \mathbb{R}.$$

To find the solution to (A.130), we write it in the form

$$\dot{x}(t) - a(t)x(t) = b(t). \quad (\text{A.131})$$

We then multiply (A.131) by the “integrating factor”

$$e^{-\int_{t_0}^t a(s) ds} \quad (\text{A.132})$$

to get

$$\dot{x}(t)e^{-\int_{t_0}^t a(s) ds} - a(t)x(t)e^{-\int_{t_0}^t a(s) ds} = e^{-\int_{t_0}^t a(s) ds} b(t). \quad (\text{A.133})$$

The above equation can be represented in an equivalent form as

$$\frac{d}{dt} \left\{ x(t)e^{-\int_{t_0}^t a(s) ds} \right\} = e^{-\int_{t_0}^t a(s) ds} b(t). \quad (\text{A.134})$$

Integrating both sides of (A.134) yields

$$x(t)e^{-\int_{t_0}^t a(s) ds} - x_0 = \int_{t_0}^t e^{-\int_{t_0}^{\tau} a(s) ds} b(\tau) d\tau. \quad (\text{A.135})$$

Hence,

$$x(t) = e^{\int_{t_0}^t a(s) ds} \left[x_0 + \int_{t_0}^t e^{-\int_{t_0}^{\tau} a(s) ds} b(\tau) d\tau \right], \quad (\text{A.136})$$

or, in an equivalent form,

$$x(t) = e^{\int_{t_0}^t a(s) ds} x_0 + \int_{t_0}^t e^{\int_{\tau}^t a(s) ds} b(\tau) d\tau \quad (\text{A.137})$$

A.14 Integral and Differential Inequalities

A.14.1 The Bellman-Gronwall Lemma

Lemma A.5 Let $v(t)$ and $w(t)$ be continuous real functions of t . Let $w(t) \geq 0$, and c be a real constant. If

$$v(t) \leq c + \int_0^t w(s) v(s) ds, \quad (\text{A.138})$$

then

$$v(t) \leq ce^{\int_0^t w(s) ds}. \quad (\text{A.139})$$

Proof Denote the right-hand side of (A.138) by $x(t)$; that is, let

$$x(t) = c + \int_0^t w(s) v(s) ds \quad (\text{A.140})$$

and let

$$z(t) = x(t) - v(t). \quad (\text{A.141})$$

Note that $z(t) \geq 0$, and $x(t)$ is differentiable. Differentiating both sides of (A.140) and using (A.141) yields

$$\begin{aligned}\dot{x}(t) &= w(t)v(t) \\ &= w(t)x(t) - w(t)z(t).\end{aligned}\tag{A.142}$$

We use (A.137) to represent (A.142) in the equivalent integral form as

$$x(t) = e^{\int_0^t w(s) ds} x_0 - \int_0^t e^{\int_\tau^t w(s) ds} w(\tau) z(\tau) d\tau.\tag{A.143}$$

Because $z(t) \geq 0$ and $w(t) \geq 0$, we have

$$x(t) \leq e^{\int_0^t w(s) ds} x_0.\tag{A.144}$$

From the definition of $x(t)$ given by (A.140), we have $x(0) = c$. On the other hand,

$$v(t) = x(t) - z(t) \leq x(t),\tag{A.145}$$

because $z(t) \geq 0$. Hence,

$$v(t) \leq ce^{\int_0^t w(s) ds},$$

and the proof is completed.

The Bellman–Gronwall lemma allows us to convert the investigation of integral equations into integral inequalities.

A.14.2 A Comparison Theorem

We consider the differential inequality

$$\dot{x}(t) \leq \omega(t, x(t)), \quad x(t_0) = x_0,$$

where $\omega : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a continuous function in a domain $\Delta \subset \mathbb{R}^2$. We assume that in Δ , the function ω satisfies the Lipschitz condition,

$$|\omega(t, x) - \omega(t, y)| \leq L|x - y|,$$

where $L \geq 0$ is a Lipschitz constant.

Consider now an associated differential equation

$$\dot{\tilde{x}}(t) = \omega(t, \tilde{x}(t)), \quad \tilde{x}(t_0) = \tilde{x}_0.$$

We refer to the above equation as the *equation of comparison*. Denote the solution to the equation of comparison by

$$\tilde{x}(t) = \tilde{x}(t; t_0, \tilde{x}_0).$$

The following theorem, which can be found, for example, in Corduneanu [55, p. 49], relates the solutions of the differential inequality and the associated differential equation.

Theorem A.20 If $x(t)$ is a differentiable solution to the differential inequality and $\tilde{x}(t)$ is the solution to the associated differential equation with the corresponding initial conditions such that

$$x_0 \leq \tilde{x}_0,$$

then

$$x(t) \leq \tilde{x}(t)$$

for all $t \geq t_0$.

Proof We prove the theorem by contradiction. We suppose that $x_0 \leq \tilde{x}_0$ but that $x(t) \leq \tilde{x}(t)$ does not hold. Thus, there exists $t_1 > t_0$ such that

$$x(t_1) > \tilde{x}(t_1).$$

Let

$$y(t) = x(t) - \tilde{x}(t).$$

Note that $y(t_0) = x(t_0) - \tilde{x}(t_0) \leq 0$ and that $y(t_1) = x(t_1) - \tilde{x}(t_1) > 0$. Denote by τ the largest $t \in [t_0, t_1]$ such that $y(\tau) = 0$. Such a value of t exists by virtue of the continuity of y . Thus,

$$y(t) > 0 \quad \text{for all } t \in (\tau, t_1].$$

On the other hand, for $t \in [\tau, t_1]$,

$$\begin{aligned} \dot{y}(t) &= \dot{x}(t) - \dot{\tilde{x}}(t) \\ &\leq \omega(t, x(t)) - \omega(t, \tilde{x}(t)) \\ &\leq |\omega(t, x(t)) - \omega(t, \tilde{x}(t))| \\ &\leq L|x(t) - \tilde{x}(t)|. \end{aligned}$$

In the interval $[\tau, t_1]$,

$$y(t) = x(t) - \tilde{x}(t) \geq 0.$$

Hence,

$$\begin{aligned} \dot{y}(t) &\leq L|x(t) - \tilde{x}(t)| \\ &= L(x(t) - \tilde{x}(t)) \\ &= Ly(t). \end{aligned}$$

Therefore, in the interval $[\tau, t_1]$, the function y satisfies

$$y(t) \leq y(\tau)e^{L(t-\tau)} = 0$$

because $y(\tau) = 0$. But this is a contradiction to the previous conclusion that $y(t) > 0$ for all $t \in (\tau, t_1]$, which completes the proof.

A.15 Solving the State Equations

A.15.1 Solution of Uncontrolled System

Given a dynamical system's time-invariant model of the form

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) \tag{A.146}$$

subject to an initial condition

$$\mathbf{x}(0) = \mathbf{x}_0, \quad (\text{A.147})$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$. The solution to (A.146) subject to (A.147) is

$$\mathbf{x}(t) = e^{\mathbf{A}t} \mathbf{x}_0$$

The above is a generalization of the situation when \mathbf{A} is a scalar. The exponential matrix $e^{\mathbf{A}t}$ is defined as

$$e^{\mathbf{A}t} = \mathbf{I}_n + t\mathbf{A} + \frac{t^2}{2!}\mathbf{A}^2 + \frac{t^3}{3!}\mathbf{A}^3 + \cdots$$

Observe that

$$\frac{d}{dt}e^{\mathbf{A}t} = \mathbf{A}e^{\mathbf{A}t} = e^{\mathbf{A}t}\mathbf{A}.$$

Suppose now that the initial condition $\mathbf{x}(0) = \mathbf{x}_0$ is replaced with a more general condition

$$\mathbf{x}(t_0) = \mathbf{x}_0,$$

where t_0 is a given initial time instant. Then, the solution to $\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t)$ subject to $\mathbf{x}(t_0) = \mathbf{x}_0$ is

$$\mathbf{x}(t) = e^{\mathbf{A}(t-t_0)} \mathbf{x}_0$$

The matrix $e^{\mathbf{A}(t-t_0)}$ is often written as

$$e^{\mathbf{A}(t-t_0)} = \Phi(t, t_0)$$

and is called the *state transition matrix* because it relates the state at any instant of time t_0 to the state at any other time t . There are many methods for calculating the state transition matrix. For an overview of some of them, we refer the reader to the classic paper on the subject by Moler and Van Loan [208]. In the case when the matrix \mathbf{A} is of low dimensions, we can use the formula that results from applying the Laplace transform to the time-invariant system $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$, $\mathbf{x}(0) = \mathbf{x}_0$. The Laplace transform of the above matrix differential equation is

$$s\mathbf{X}(s) - \mathbf{x}_0 = \mathbf{A}\mathbf{X}(s).$$

Hence

$$\mathbf{X}(s) = (s\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{x}_0,$$

and the inverse transform yields

$$e^{\mathbf{A}t} = \mathcal{L}^{-1}\{(s\mathbf{I}_n - \mathbf{A})^{-1}\} \quad (\text{A.148})$$

where \mathcal{L}^{-1} denotes the inverse Laplace transform operator.

A.15.2 Solution of Controlled System

Consider a model of a controlled dynamical system,

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) \quad (\text{A.149})$$

subject to a given initial condition $\mathbf{x}(0) = \mathbf{x}_0$, where \mathbf{B} is an $n \times m$ matrix, and $m < n$ as is usually in practice. Premultiplying both sides of (A.149) by the integrating factor, $e^{-\mathbf{A}t}$, yields

$$e^{-\mathbf{A}t} \dot{\mathbf{x}}(t) = e^{-\mathbf{A}t} \mathbf{A} \mathbf{x}(t) + e^{-\mathbf{A}t} \mathbf{B} \mathbf{u}(t).$$

We then rearrange the above equation as

$$e^{-\mathbf{A}t} \dot{\mathbf{x}}(t) - e^{-\mathbf{A}t} \mathbf{A} \mathbf{x}(t) = e^{-\mathbf{A}t} \mathbf{B} \mathbf{u}(t). \quad (\text{A.150})$$

Noting that

$$\frac{d}{dt}(e^{-\mathbf{A}t} \mathbf{x}(t)) = -\mathbf{A} e^{-\mathbf{A}t} \mathbf{x}(t) + e^{-\mathbf{A}t} \dot{\mathbf{x}}(t),$$

we can represent (A.150) in the form

$$\frac{d}{dt}(e^{-\mathbf{A}t} \mathbf{x}(t)) = e^{-\mathbf{A}t} \mathbf{B} \mathbf{u}(t).$$

Integrating both sides of the resulting equation gives

$$e^{-\mathbf{A}t} \mathbf{x}(t) - \mathbf{x}(0) = \int_0^t e^{-\mathbf{A}\tau} \mathbf{B} \mathbf{u}(\tau) d\tau.$$

Hence,

$$\mathbf{x}(t) = e^{\mathbf{A}t} \mathbf{x}(0) + \int_0^t e^{\mathbf{A}(t-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau.$$

If the initial condition is $\mathbf{x}(t_0) = \mathbf{x}_0$, then integration from t_0 to t yields

$$\boxed{\mathbf{x}(t) = e^{\mathbf{A}(t-t_0)} \mathbf{x}(t_0) + \int_{t_0}^t e^{\mathbf{A}(t-\tau)} \mathbf{B} \mathbf{u}(\tau) d\tau.} \quad (\text{A.151})$$

As in the case of the uncontrolled system, we can also use the Laplace transform technique to obtain the solution to (A.149) subject to the initial condition $\mathbf{x}(0) = \mathbf{x}_0$. Indeed, taking the Laplace transform of (A.149) and performing some manipulations yields

$$\mathbf{X}(s) = (s\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{x}(0) + (s\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{B} \mathbf{U}(s).$$

Hence,

$$\mathbf{x}(t) = \mathcal{L}^{-1}\{(s\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{x}(0)\} + \mathcal{L}^{-1}\{(s\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{B} \mathbf{U}(s)\}.$$

If $\mathbf{x}_0 = \mathbf{x}(t_0)$, then

$$\mathbf{x}(t) = \mathcal{L}^{-1}\{e^{-t_0 s} (s\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{x}(t_0)\} + \mathcal{L}^{-1}\{e^{-t_0 s} (s\mathbf{I}_n - \mathbf{A})^{-1} \mathbf{B} \mathbf{U}(s)\}.$$

A.16 Curves and Surfaces

Consider a point moving along a curve in n -dimensional space. The position of the point at time t can be described as

$$\mathbf{x}(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \\ \vdots \\ x_n(t) \end{bmatrix}. \quad (\text{A.152})$$

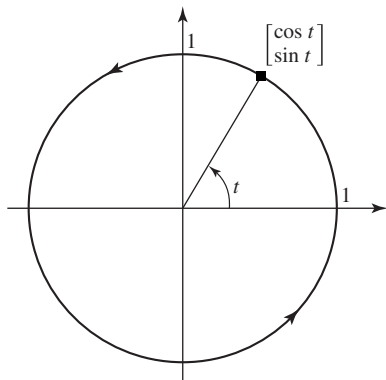


Figure A.4 An example of a curve in two-dimensional space.

For example, if the point is moving along a straight line passing through fixed points \mathbf{y} and \mathbf{z} in \mathbb{R}^n , then

$$\mathbf{x}(t) = \mathbf{y} + t\mathbf{z}. \quad (\text{A.153})$$

As another example, consider

$$\mathbf{x}(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} = \begin{bmatrix} \cos t \\ \sin t \end{bmatrix} \in \mathbb{R}^2. \quad (\text{A.154})$$

Then, the point moves around a circle of radius 1 in counterclockwise direction as shown in Figure A.4. We used t as the variable corresponding to the angle corresponding to the position of the point on the circle. We are now ready for a formal definition of a curve.

Definition A.6 Let $I \in \mathbb{R}$ be an interval. A parameterized curve defined on the interval I is an association which to each point of I associates a vector, where $\mathbf{x}(t)$ denotes the vector associated to $t \in I$ by \mathbf{x} . We write the association $t \mapsto \mathbf{x}(t)$ as

$$\mathbf{x} : I \rightarrow \mathbb{R}^n.$$

We also call this association the parameterization of a curve. We call $\mathbf{x}(t)$ the position vector at time t .

The position vector can be represented in terms of its coordinates as in (A.152), where each component $x_i(t)$ is a function of time t . We say that this curve is differentiable if each function $x_i(t)$ is a differentiable function of time t . We define the derivative $d\mathbf{x}(t)/dt$ to be

$$\dot{\mathbf{x}}(t) = \frac{d\mathbf{x}(t)}{dt} = \begin{bmatrix} \frac{dx_1(t)}{dt} \\ \frac{dx_2(t)}{dt} \\ \vdots \\ \frac{dx_n(t)}{dt} \end{bmatrix}$$

and call it the velocity vector of the curve at time t . The velocity vector is located at the origin. However, when we translate it to the point $\mathbf{x}(t)$, as in Figure A.5, then we visualize it as tangent to the curve $\mathbf{x}(t)$.

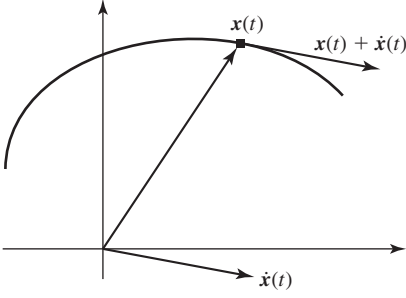


Figure A.5 Velocity vector of a curve at time t is parallel to a tangent vector at time t .

We define the tangent line to a curve \mathbf{x} at time t to be the line passing through $\mathbf{x}(t)$ in the direction of $\dot{\mathbf{x}}(t)$, provided that $\dot{\mathbf{x}}(t) \neq \mathbf{0}$. If $\dot{\mathbf{x}}(t) = \mathbf{0}$, we do not define a tangent line. The speed of the curve $\mathbf{x}(t)$, denoted $v(t)$, is $v(t) = \|\dot{\mathbf{x}}(t)\|$. To illustrate the above, we consider the point moving on the circle. The point position at time t is

$$\mathbf{x}(t) = [\cos t \quad \sin t]^T.$$

The velocity vector of the curve at time t is

$$\dot{\mathbf{x}}(t) = [-\sin t \quad \cos t]^T,$$

and the speed of the curve $\mathbf{x}(t)$ at time t is

$$v(t) = \|\dot{\mathbf{x}}(t)\| = \sqrt{(-\sin t)^2 + \cos^2 t} = 1.$$

The set of points in \mathbb{R}^3 such that for some $F : \mathbb{R}^3 \rightarrow \mathbb{R}$ we have

$$F(\mathbf{x}) = F(x_1, x_2, x_3) = 0$$

and

$$DF(x_1, x_2, x_3) = \nabla F(x_1, x_2, x_3)^T = \begin{bmatrix} \frac{\partial F}{\partial x_1} & \frac{\partial F}{\partial x_2} & \frac{\partial F}{\partial x_3} \end{bmatrix} \neq \mathbf{0}^T$$

is called a surface in \mathbb{R}^3 , where $\nabla F = \text{grad } F$. For example, let

$$F(x_1, x_2, x_3) = x_1^2 + x_2^2 + x_3^2 - 1;$$

then, the surface is the sphere of radius 1, centered at the origin of \mathbb{R}^3 .

Now let

$$\mathbf{C}(t) = [x_1(t) \quad x_2(t) \quad \cdots \quad x_n(t)]^T$$

be a differentiable curve. We say that the curve lies on the surface described by

$$\begin{aligned} F_1(x_1, x_2, \dots, x_n) &= 0, \\ &\vdots \\ F_m(x_1, x_2, \dots, x_n) &= 0 \end{aligned}$$

if, for all t , we have

$$\begin{aligned} F_1(x_1(t), x_2(t), \dots, x_n(t)) &= 0, \\ &\vdots \\ F_m(x_1(t), x_2(t), \dots, x_n(t)) &= 0. \end{aligned}$$

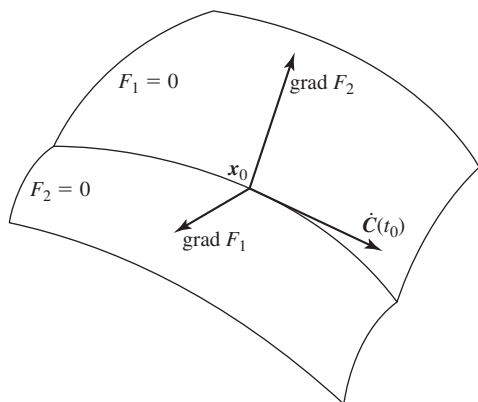


Figure A.6 Tangent space to a surface is the set of vectors orthogonal to the gradients of the functions defining the surface.

Differentiating the above equations, using the chain rule, we obtain

$$DF_i(C(t))\dot{C}(t) = \nabla F_i(C(t))^T \dot{C}(t) = 0, \quad i = 1, 2, \dots, m.$$

Let \mathbf{x}_0 be a point of the surface, and let $C(t)$ be a curve on the surface passing through \mathbf{x}_0 ; that is, for some t_0 ,

$$C(t_0) = \mathbf{x}_0.$$

For this t_0 we have

$$\nabla F_i(\mathbf{x}_0)^T \dot{C}(t_0) = 0, \quad i = 1, 2, \dots, m,$$

where we assumed that $\dot{C}(t_0) \neq \mathbf{0}$. Thus, the gradients ∇F_i at \mathbf{x}_0 are orthogonal to the tangent vector of the curve at the point \mathbf{x}_0 . This is true for every differentiable curve on the surface passing through \mathbf{x}_0 . For an illustration of the above properties in three-dimensional space, we refer to Figure A.6.

Definition A.7 The tangent space to the surface $\{\mathbf{x} : \mathbf{F}(\mathbf{x}) = \mathbf{0}\}$ at the point \mathbf{x}_0 is the set of vectors at \mathbf{x}_0 that are tangent to differentiable curves on the surface passing through \mathbf{x}_0 .

A.17 Vector Fields and Curve Integrals

Definition A.8 Let \mathcal{U} be an open subset of \mathbb{R}^n . By a vector field on \mathcal{U} we mean an association

$$\mathbf{F} : \mathcal{U} \rightarrow \mathbb{R}^n,$$

which to every point of \mathcal{U} associates a vector of the same dimension.

We can interpret a vector field as a field of forces. We can visualize a vector field as a field of arrows, which to each point associates an arrow as shown in Figure A.7. If $f : \mathcal{U} \rightarrow \mathbb{R}$ is a function, then

$$\mathbf{F} = \nabla f$$

is a vector field, and we have

$$\mathbf{F}(\mathbf{x}) = \nabla f(\mathbf{x})$$

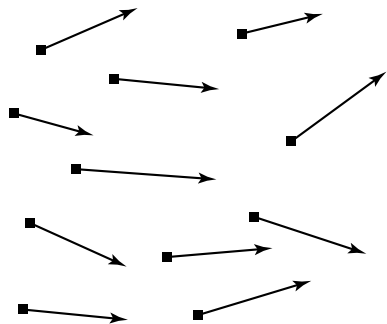


Figure A.7 A vector field.

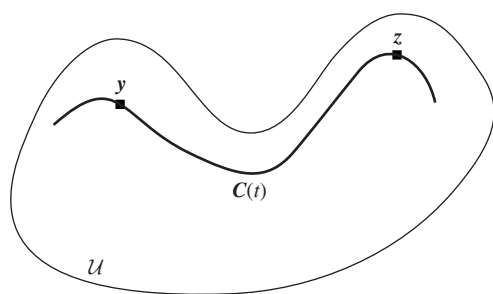


Figure A.8 A connected open set in which there are points that cannot be joined by a straight line.

for all $\mathbf{x} \in \mathcal{U}$. We now consider an inverse problem: When is a given vector field in \mathbb{R}^n a gradient?

Definition A.9 Let \mathbf{F} be a vector field on an open set \mathcal{U} . If f is a differentiable function on \mathcal{U} such that

$$\mathbf{F} = \nabla f,$$

then we say that f is a potential function for \mathbf{F} .

We then can ask the question, Do the potential functions exist? To answer the question, we recast the problem into the one that involves solving n first-order partial differential equations:

$$\frac{\partial f}{\partial x_i} = F_i(\mathbf{x}), \quad i = 1, 2, \dots, n,$$

with the initial condition

$$\mathbf{F}(\mathbf{x}_0) = \mathbf{F}_0.$$

To answer the above question we need one more definition.

Definition A.10 An open set \mathcal{U} is said to be connected if, given two points in \mathcal{U} , there is a differentiable curve in \mathcal{U} that joins the two points.

If $\mathcal{U} = \mathbb{R}^n$, then any two points can be joined by a straight line. Of course, it is not always the case that two points of an open set can be joined by a straight line, as illustrated in Figure A.8. We say that the function $\mathbf{F} = [F_1 \ F_2 \ \cdots \ F_n]^T$ is a C^1 function if all partial derivatives of all functions F_i exist and are continuous.

Theorem A.21 Let $\mathbf{F}(\mathbf{x})$ be \mathcal{C}^1 function defined on an open connected region \mathcal{U} of \mathbb{R}^n . Then, there is a solution to the differential equation

$$\nabla f(\mathbf{x}) = \mathbf{F}(\mathbf{x}), \quad \mathbf{F}(\mathbf{x}_0) = \mathbf{F}_0,$$

if and only if

$$\frac{\partial F_i}{\partial x_j} = \frac{\partial F_j}{\partial x_i};$$

that is, the Jacobian matrix of \mathbf{F} is symmetric.

Proof We prove the theorem for $n = 2$, where the connected region \mathcal{U} is a rectangle. We select a point $\mathbf{x}_0 = [x_{10} \ x_{20}]^T$ in the rectangle. We assume that $\mathbf{F} = [F_1 \ F_2]^T$ and that

$$D_2 F_1 = D_1 F_2; \quad \text{that is,} \quad \frac{\partial F_1}{\partial x_2} = \frac{\partial F_2}{\partial x_1}.$$

Our goal is to find a potential function $f(x_1, x_2)$. We assume that the potential function has the form

$$f(x_1, x_2) = \int_{x_{10}}^{x_1} F_1(s, x_2) ds + u(x_2),$$

where $u(x_2)$ is a function of x_2 only and will be determined later in the proof. We illustrate the integration path in Figure A.9. Applying the fundamental theorem of calculus, we obtain

$$\begin{aligned} \frac{\partial f}{\partial x_1} &= \frac{\partial}{\partial x_1} \int_{x_{10}}^{x_1} F_1(s, x_2) ds + \frac{\partial u(x_2)}{\partial x_1} \\ &= F_1(x_1, x_2) \end{aligned}$$

because

$$\frac{\partial u(x_2)}{\partial x_1} = 0.$$

Hence, $D_1 f(x_1, x_2) = F_1(x_1, x_2)$ as required.

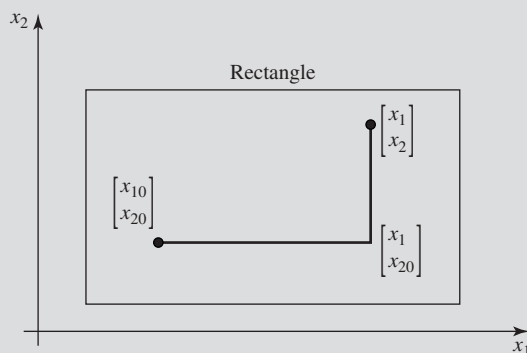


Figure A.9 Integration path used in the proof of Theorem A.21.

We now differentiate $f(x_1, x_2)$ with respect to x_2 to get

$$\begin{aligned}\frac{\partial f}{\partial x_2} &= \int_{x_{10}}^{x_1} \frac{\partial F_1(s, x_2)}{\partial x_2} ds + \frac{\partial u(x_2)}{\partial x_2} \\ &= \int_{x_{10}}^{x_1} \frac{\partial F_2(s, x_2)}{\partial x_1} ds + \frac{\partial u(x_2)}{\partial x_2}\end{aligned}$$

because, by assumption, $D_2 F = D_1 F$. Therefore,

$$\begin{aligned}\frac{\partial f}{\partial x_2} &= F_2(s, x_2)|_{x_{10}}^{x_1} + \frac{\partial u(x_2)}{\partial x_2} \\ &= F_2(x_1, x_2) - F_2(x_{10}, x_2) + \frac{\partial u(x_2)}{\partial x_2}.\end{aligned}$$

Let

$$-F_2(x_{10}, x_2) + \frac{\partial u(x_2)}{\partial x_2} = 0,$$

that is, $\frac{\partial u(x_2)}{\partial x_2} = F_2(x_{10}, x_2)$. Then, $D_2 f = F_2$ as required. To conclude the proof, we let

$$u(x_2) = \int_{x_{20}}^{x_2} F_2(x_{10}, x_2) dx_2,$$

which yields

$$f(x_1, x_2) = \int_{x_{10}}^{x_1} F_1(s, x_2) ds + \int_{x_{20}}^{x_2} F_2(x_{10}, t) dt \quad (\text{A.155})$$

The proof is complete for $n = 2$.

We note that if the connected region \mathcal{U} is an arbitrary open set, then the above proof may not work. This is because we needed to integrate over intervals that were contained in the region \mathcal{U} .

Suppose we are given a vector field \mathbf{F} on an open set U in the plane. We can interpret \mathbf{F} as a field of forces. We wish to calculate the work done when moving a particle from a point $\mathbf{C}(t_1)$ to a point $\mathbf{C}(t_2)$ along the curve $\mathbf{C}(t)$ in U , as illustrated in Figure A.10. Let

$$\mathbf{F}(\mathbf{x}) = [F_1(x_1, x_2) \quad F_2(x_1, x_2)]^T.$$

We assume that each component of \mathbf{F} is differentiable; that is, \mathbf{F} is a differentiable vector field. We next assume that the curve $\mathbf{C}(t)$ is defined on a closed interval $[t_1, t_2]$ with $t_1 < t_2$. For each $t \in [t_1, t_2]$ the value $\mathbf{C}(t)$ is a point in \mathbb{R}^2 . The curve \mathbf{C} lies in U if $\mathbf{C}(t)$ is a point in U for all $t \in [t_1, t_2]$. We further assume that the curve \mathbf{C} is continuously differentiable and that its derivative $\dot{\mathbf{C}}(t) = d\mathbf{C}(t)/dt$ exists and is continuous; that is, we assume that $\mathbf{C} \in \mathcal{C}^1$.

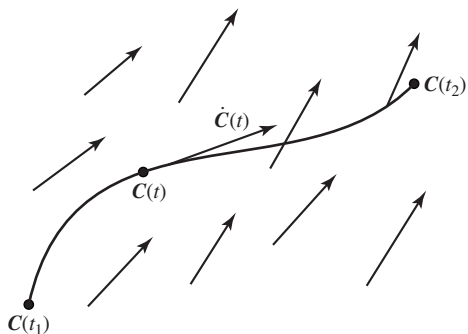


Figure A.10 A particle moving along a curve in vector field.

We define the integral of \mathbf{F} along \mathbf{C} to be

$$\begin{aligned}\int_{\mathbf{C}} \mathbf{F} &= \int_{t_1}^{t_2} \mathbf{F}(\mathbf{C}(t))^T \frac{d\mathbf{C}}{dt} dt \\ &= \int_{t_1}^{t_2} \mathbf{F}(\mathbf{C}(t))^T d\mathbf{C},\end{aligned}$$

which is the work done against the force field \mathbf{F} along the curve \mathbf{C} . Let

$$\mathbf{C}(t) = [x_1(t) \quad x_2(t)]^T.$$

Then, we can write

$$\int_{\mathbf{C}} \mathbf{F} = \int_{\mathbf{C}} F_1 dx_1 + F_2 dx_2. \quad (\text{A.156})$$

Theorem A.22 (Green) Let the vector field \mathbf{F} be given on a region A that is the interior of a closed path \mathbf{C} parameterized counterclockwise. Then,

$$\int_{\mathbf{C}} F_1 dx_1 + F_2 dx_2 = \iint_A \left(\frac{\partial F_2}{\partial x_1} - \frac{\partial F_1}{\partial x_2} \right) dx_1 dx_2.$$

A.18 Matrix Calculus Formulas

Below, we list useful formulas from matrix calculus that involve computing the gradient (with respect to \mathbf{x}). Recall that in our convention the gradient is a column vector. We have

$$\nabla(\mathbf{x}^T \mathbf{y}) = \nabla(\mathbf{y}^T \mathbf{x}) = \mathbf{y}, \quad (\text{A.157})$$

$$\nabla(\mathbf{x}^T \mathbf{x}) = 2\mathbf{x}, \quad (\text{A.158})$$

$$\nabla(\mathbf{x}^T \mathbf{A} \mathbf{y}) = \mathbf{A} \mathbf{y}, \quad (\text{A.159})$$

$$\nabla(\mathbf{y}^T \mathbf{A} \mathbf{x}) = \mathbf{A}^T \mathbf{y}, \quad (\text{A.160})$$

$$\nabla(\mathbf{x}^T \mathbf{A} \mathbf{x}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}. \quad (\text{A.161})$$

Let $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be two maps, and let $\partial \mathbf{f} / \partial \mathbf{x}$ and $\partial \mathbf{g} / \partial \mathbf{x}$ be the Jacobian matrices of \mathbf{f} and \mathbf{g} , respectively. Then,

$$\nabla(\mathbf{f}(\mathbf{x})^T \mathbf{g}(\mathbf{x})) = \nabla(\mathbf{g}(\mathbf{x})^T \mathbf{f}(\mathbf{x})) = \left(\frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right)^T \mathbf{g} + \left(\frac{\partial \mathbf{g}}{\partial \mathbf{x}} \right)^T \mathbf{f}. \quad (\text{A.162})$$

For an $n \times n$ symmetric matrix \mathbf{Q} , we have

$$\nabla(\mathbf{f}(\mathbf{x})^T \mathbf{Q} \mathbf{f}(\mathbf{x})) = 2 \left(\frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right)^T \mathbf{Q} \mathbf{f}. \quad (\text{A.163})$$

Notes

For more on the subject of proof techniques, we refer to Velleman [288]. For further reading on linear algebra, and in particular matrix theory, we recommend classical texts of Gel'fand [95] and Gantmacher [93]. Analyses of numerical methods for matrix computations can be found in Faddeeva [78], Householder [127], Golub and Van Loan [104], and Gill, Murray, and Wright [97]. Lang [175] and Seeley [255] are recommended for a review of calculus of several variables. A straightforward approach to mathematical analysis is presented in Binmore [28]. Principles of real analysis can be found in Natanson [215], Bartle [22], and Rudin [244]. Some facts from basic functional analysis that we touched upon in this appendix are discussed, for example, by Maddox [193] and by Debnath and Mikusiński [58]. Lucid expositions of differential equations are in Driver [69], Miller and Michel [205], and Boyce and DiPrima [32].

G. H. Hardy (1877–1947) argues in reference 114 (p. 113) that in truly great theorems and their proofs, “there is a very high degree of *unexpectedness*, combined with *inevitability* and *economy*. The arguments take so odd and surprising a form; the weapons used seem so childishly simple when compared with the far-reaching results; but there is no escape from the conclusions.” He further states that “A mathematical proof should resemble a simple and clear-cut constellation, not a scattered cluster in the Milky Way.”

EXERCISES

A.1 The open unit disc with respect to the norm $\|\cdot\|_p$ is the set

$$\Delta_p = \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x}\|_p < 1\}.$$

The unit discs are quite different for different choices of the norm $\|\cdot\|_p$. Sketch Δ_1 , Δ_2 , Δ_3 , and Δ_∞ .

A.2 Given a linear combination of matrices,

$$\mathbf{M} = \alpha_1 \mathbf{M}_1 + \cdots + \alpha_r \mathbf{M}_r, \quad (\text{A.164})$$

where for $i = 1, \dots, r$, $\alpha_i \in \mathbb{R}$, and $\mathbf{M}_i \in \mathbb{R}^{n \times n}$. Let \mathbf{m}_i^j be the j th row of the matrix

M_i . Thus,

$$M_i = \begin{bmatrix} m_i^1 \\ m_i^2 \\ \vdots \\ m_i^n \end{bmatrix}.$$

Show that

$$\det \mathbf{M} = \sum_{s_1=1}^r \cdots \sum_{s_n=1}^r \alpha_{s_1} \cdots \alpha_{s_n} \det \begin{bmatrix} m_{s_1}^1 \\ \vdots \\ m_{s_n}^n \end{bmatrix}. \quad (\text{A.165})$$

A.3 Let $\mathbf{X} \in \mathbb{R}^{n \times m}$ and $\mathbf{Y} \in \mathbb{R}^{m \times n}$. Show that

$$\det(\mathbf{I}_n + \mathbf{XY}) = \det(\mathbf{I}_m + \mathbf{YX}).$$

A.4 Let $\lambda_i(\mathbf{A})$ denote the i th eigenvalue of a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ and let α be a scalar. Show that

$$(a) \quad \lambda_i(\alpha \mathbf{A}) = \alpha \lambda_i(\mathbf{A});$$

$$(b) \quad \lambda_i(\mathbf{I}_n \pm \alpha \mathbf{A}) = 1 \pm \alpha \lambda_i(\mathbf{A}).$$

A.5 Show that the formulas (A.75) and (A.76) for matrix norms are equivalent.

A.6 Show that for any square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, we have

$$\|\mathbf{A}\| \geq \max_{1 \leq i \leq n} |\lambda_i(\mathbf{A})|. \quad (\text{A.166})$$

A.7 Show that if $\mathbf{P} \in \mathbb{R}^{n \times n}$ is symmetric, then

$$\lambda_{\min}(\mathbf{P}) \|\mathbf{x}\|_2^2 \leq \mathbf{x}^T \mathbf{P} \mathbf{x} \leq \lambda_{\max}(\mathbf{P}) \|\mathbf{x}\|_2^2.$$

A.8 Show that a sufficient condition for

$$\lim_{k \rightarrow \infty} \mathbf{A}^k \rightarrow \mathbf{O}$$

is

$$\|\mathbf{A}\| < 1.$$

A.9 Let $\mathbf{A} \in \mathbb{R}^{n \times n}$. Show that a necessary and sufficient condition for

$$\lim_{k \rightarrow \infty} \mathbf{A}^k \rightarrow \mathbf{O}.$$

is

$$|\lambda_i(\mathbf{A})| < 1, \quad i = 1, 2, \dots, n.$$

A.10 Let $\text{eig}(\mathbf{M})$ denote the set of eigenvalues of the square matrix \mathbf{M} and let Δ be the open unit disk in the complex plane. Prove that if $\text{eig}(\mathbf{M}) \subset \Delta \cup \{1\}$, then $\lim_{k \rightarrow \infty} \mathbf{M}^k$ exists if and only if the dimension of the eigenspace corresponding to the eigenvalue 1 is equal to the multiplicity of the eigenvalue 1; that is, the geometric multiplicity of the eigenvalue 1 is equal to its algebraic multiplicity.

A.11

- (a) Show that the series

$$\sum_{n=1}^{\infty} \frac{1}{n}$$

diverges to $+\infty$.

- (b) Show that if
- α
- is a rational number such that
- $\alpha > 1$
- , then the series

$$\sum_{n=1}^{\infty} \frac{1}{n^{\alpha}}$$

converges.

- (c) Use MATLAB's Symbolic Math Toolbox to compute

$$\sum_{n=1}^{\infty} \frac{1}{n^2}.$$

A.12

- (a) Let
- $0 < x < 1$
- . Show that

$$\ln(1-x) \leq -x.$$

- (b) Let
- $0 < x \leq 1/2$
- . Show that

$$\ln(1-x) \geq -2x.$$

- A.13**
- Let
- $0 < a_k < 1$
- for
- $k \geq 0$
- . Use a proof by contraposition, Theorem A.16, and Exercise A.12 to show that

$$\prod_{k=0}^{\infty} (1 - a_k) = 0 \quad \text{if and only if} \quad \sum_{k=0}^{\infty} a_k = \infty.$$

- A.14**
- Prove Sylvester's inequalities (A.88).

- A.15**
- Compute
- e^{At}
- for

- (a)

$$A = \begin{bmatrix} 0 & -3 & 0 \\ 3 & 0 & 0 \\ 0 & 0 & -1 \end{bmatrix},$$

- (b)

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -2 & -3 \end{bmatrix}.$$

- A.16**
- Find
- $\mathbf{x}(t)$
- for the dynamical system model,

$$\dot{\mathbf{x}}(t) = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{x}(t) + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u(t), \quad \mathbf{x}(0) = \begin{bmatrix} 0 \\ 1 \end{bmatrix},$$

where $u(t)$ is the unit step function; that is, $u(t) = 1$ for $t \geq 0$ and $u(t) = 0$ for $t < 0$.

A.17 For

$$\mathbf{A} = \begin{bmatrix} -2 & -6 \\ 2 & -10 \end{bmatrix},$$

find constants $M > 0$ and $\alpha > 0$ such that

$$\|e^{\mathbf{A}t}\|_2 \leq M e^{-\alpha t}, \quad t \geq 0.$$

A.18 Show that the function space $\mathcal{C}([a, b])$ of real-valued functions, $f : \mathbb{R} \rightarrow \mathbb{R}$, with the norm defined by (A.99) is a normed linear space. Show that this space is a Banach space.

A.19 Derive a sufficient condition for the convergence of the iteration process

$$\mathbf{x}^{(k+1)} = \mathbf{A}\mathbf{x}^{(k)} + \mathbf{b}, \quad k = 0, 1, \dots,$$

for any initial $\mathbf{x}(0)$, where $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{b} \in \mathbb{R}^n$.

BIBLIOGRAPHY

1. P. P. Acarnley. *Stepping Motors: A Guide to Modern Theory and Practice*. IEE Control Engineering Series 19. Peter Peregrinus, London, 1982.
2. J. Ackermann. Robust control prevents car skidding. *IEEE Control Systems*, 17(3):23–31, June 1997.
3. M. Adachi and K. Aihara. Associative dynamics in a chaotic neural network. *Neural Networks*, 10(1):83–98, January 1997.
4. J. A. Adam and N. Bellomo, editors. *A Survey of Models for Tumor-Immune System Dynamics*. Birkhäuser, Boston, 1997.
5. K. Aihara, T. Takabe, and M. Toyoda. Chaotic neural networks. *Physics Letters A*, 144(6, 7):333–340, March 12, 1990.
6. A. W. Al-Khafaji and J. R. Tooley. *Numerical Methods in Engineering Practice*. Holt, Rinehart and Winston, New York, 1986.
7. K. T. Alligood, T. D. Sauer, and J. A. Yorke. *CHAOS An Introduction to Dynamical Systems*. Springer-Verlag, New York, 1997.
8. J. A. Anderson. *An Introduction to Neural Networks*. A Bradford Book, MIT Press, Cambridge, MA, 1995.
9. J. A. Anderson, J. W. Silverstein, S. A. Ritz, and R. S. Jones. Distinctive features, categorical perception, probability learning: Some applications of a neural model. In J. A. Anderson and E. Rosenfeld, editors, *Neurocomputing: Foundations of Research*, Chapter 22, pages 283–325. MIT Press, Cambridge, MA, 1989. Reprint from *Psychological Review*, 84:413–451, 1977.
10. Y. V. Andreyev, Y. L. Belsky, A. S. Dmitriev, and D. A. Kuminov. Information processing using dynamical chaos: Neural networks implementation. *IEEE Transactions on Neural Networks*, 7(2):290–299, March 1996.
11. Yu. V. Andreyev, A. S. Dmitriev, L. O. Chua, and C. W. Wu. Associative and random access memory using one-dimensional maps. *International Journal of Bifurcation and Chaos*, 2(3):483–504, September 1992.
12. P. J. Antsaklis and A. N. Michel. *Linear Systems*. McGraw-Hill, New York, 1997.
13. D. K. Arrowsmith and C. M. Place. *An Introduction to Dynamical Systems*. Cambridge University Press, Cambridge, Great Britain, 1990.
14. J.-P. Aubin and A. Cellina. *Differential Inclusions: Set-Valued Maps and Viability Theory*. Springer-Verlag, Berlin, 1984.
15. T. Bäck, U. Hammel, and H.-P. Schwefel. Evolutionary computation: Comments on the history and current state. *IEEE Transactions on Evolutionary Computation*, 1(1):3–17, April 1997.
16. F. N. Bailey. The application of Lyapunov's second method to interconnected systems. *Journal of SIAM on Control*, 3(3):443–462, 1966.
17. G. J. Bakker and H. Nijmeijer. Comments on 'Control design for multivariable nonlinear time-varying systems.' *International Journal of Control*, 41(6):1627–1629, 1985.
18. P. Ball, editor. *A Biography of Water: Life's Matrix*. Farra, Straus and Giroux, New York, 1999.

19. S. P. Banks. *Control Systems Engineering: Modelling and Simulation, Control Theory and Microprocessor Implementation*. Prentice-Hall International, Englewood Cliffs, NJ, 1986.
20. S. Barnett. *Introduction to Mathematical Control Theory*. Clarendon Press, Oxford, 1975.
21. M. Barnsley. *Fractals Everywhere*. Academic Press, San Diego, 1988.
22. R. G. Bartle. *The Elements of Real Analysis*, second edition. John Wiley & Sons, New York, 1976.
23. E. F. Beckenbach and R. Bellman. *Inequalities*. Springer-Verlag, New York, second revised printing, 1965.
24. R. Bellman. *Introduction to Matrix Analysis*. McGraw-Hill, New York, 1960.
25. R. Bellman. Vector Lyapunov functions. *Journal of SIAM on Control*, 1(1):32–34, 1962.
26. D. P. Bertsekas. Necessary and sufficient conditions for a penalty method to be exact. *Mathematical Programming*, 9(1):87–99, August 1975.
27. D. Bestle and M. Zeitz. Canonical form observer design for non-linear time-variable systems. *International Journal of Control*, 38(2):419–431, 1983.
28. K. G. Binmore. *Mathematical Analysis: A Straightforward Approach*, second edition. Cambridge University Press, Cambridge, Great Britain, 1982.
29. D. L. Boley. On Kalman's procedure for the computation of the controllable/observable canonical form. *SIAM Journal of Control and Optimization*, 18(6):624–626, November 1980.
30. V. G. Boltyanskii. *Mathematical Methods of Optimal Control*. Holt, Rinehart and Winston, New York, 1971.
31. T. L. Boullion and P. L. Odell. *Generalized Inverse Matrices*. Wiley-Interscience, New York, 1971.
32. W. E. Boyce and R. C. DiPrima. *Elementary Differential Equations and Boundary Value Problems*, sixth edition. John Wiley & Sons, New York, 1997.
33. J. W. Brewer. Kronecker products and matrix calculus in system theory. *IEEE Transactions on Circuits and Systems*, CAS-25(9):772–781, September 1978.
34. R. W. Brockett. *Finite Dimensional Linear Systems*. John Wiley & Sons, New York, 1970.
35. W. L. Brogan. *Modern Control Theory*, third edition. Prentice-Hall, Englewood Cliffs, NJ, 1991.
36. A. E. Bryson and Y. C. Ho. *Applied Optimal Control*. Blaisdell, Waltham, MA, 1975.
37. A. E. Bryson, Jr. *Dynamic Optimization*. Addison Wesley Longman, Menlo Park, CA, 1999.
38. A. Bunde and S. Havlin, editors. *Fractals in Science*. Springer-Verlag, Berlin, 1994.
39. J. B. Burl. *Linear Optimal Control; \mathcal{H}_2 and \mathcal{H}_∞ Methods*. Addison Wesley Longman, Menlo Park, CA, 1999.
40. S. L. Campbell and C. D. Meyer, Jr. *Generalized Inverses of Linear Transformations*. Dover Publications, New York, 1991.
41. R. H. Cannon, Jr. *Dynamics of Physical Systems*. McGraw-Hill, New York, 1967.
42. S. G. Cao, N. W. Rees, and G. Feng. Quadratic stability analysis and design of continuous-time fuzzy control systems. *International Journal of System Science*, 27(2):193–203, 1996.
43. H. Y. Chan, S. Hui, and S. H. Žak. Synchronous and asynchronous Brain-State-in-a-Box information system neural models. *International Journal of Neural Systems*, 9(1):61–74, February 1999.
44. H. Y. Chan and S. H. Žak. On neural networks that design neural associative memories. *IEEE Transactions on Neural Networks*, 8(2):360–372, March 1997.
45. H. Y. Chan and S. H. Žak. Real-time synthesis of sparsely interconnected neural associative memories. *Neural Networks*, 11(4):749–759, June 1998.
46. H. Y. Chan and S. H. Žak. Chaotic neural fuzzy associative memory. *International Journal of Bifurcation and Chaos*, 9(8):1597–1617, August 1999.
47. K. C. C. Chan, V. Lee, and H. Leung. Generating fuzzy rules for target tracking using a steady-state genetic algorithm. *IEEE Transactions on Evolutionary Computation*, 1(3):189–200, September 1997.
48. C.-T. Chen. *Linear Systems Theory and Design*, third edition. Oxford University Press, New York, 1999.

49. E. K. P. Chong, S. Hui, and S. H. Žak. An analysis of a class of neural networks for solving linear programming problems. *IEEE Transactions on Automatic Control*, 44(11):1995–2006, November 1999.
50. E. K. P. Chong and S. H. Žak. *An Introduction to Optimization*, second edition. Wiley-Interscience, New York, 2001.
51. L. O. Chua and G.-N. Lin. Nonlinear programming without computation. *IEEE Transactions on Circuits and Systems*, CAS-31(2):182–188, February 1984.
52. A. Cichocki and R. Unbehauen. *Neural Networks for Optimization and Signal Processing*. John Wiley, Chichester, England, 1993.
53. S. D. Conte and C. de Boor. *Elementary Numerical Analysis: An Algorithmic Approach*, third edition. McGraw-Hill, New York, 1980.
54. I. M. Copi. *Introduction to Logic*, fifth edition. Macmillan, New York, 1978.
55. C. Corduneanu. *Principles of Differential and Integral Equations*. Allyn and Bacon, Boston, MA, 1971.
56. M. J. Corless and G. Leitmann. Continuous state feedback guaranteeing uniform ultimate boundedness for uncertain dynamic systems. *IEEE Transactions on Automatic Control*, AC-26(5):1139–1144, October 1981.
57. J. J. D'Azzo and C. H. Houpis. *Linear Control System Analysis and Design: Conventional and Modern*, fourth edition. McGraw-Hill, New York, 1995.
58. L. Debnath and P. Mikusiński. *Introduction to Hilbert Spaces with Applications*. Academic Press, San Diego, 1990.
59. R. A. DeCarlo. *Linear Systems: A State Variable Approach with Numerical Implementation*. Prentice-Hall, Englewood Cliffs, NJ, 1989.
60. R. A. DeCarlo, S. H. Žak, and G. P. Matthews. Variable structure control of nonlinear multivariable systems: A tutorial. *Proceedings of the IEEE*, 76(3):212–232, March 1988.
61. V. Del Toro. *Electric Machines and Power Systems*. Prentice-Hall, Englewood Cliffs, NJ, 1985.
62. M. J. Denham. Canonical forms for the identification of multivariable linear systems. *IEEE Transactions on Automatic Control*, AC-19(6):646–656, December 1974.
63. R. L. Devaney and L. Keen, editors. *Chaos and Fractals; The Mathematics Behind the Computer Graphics*. Proceedings of Symposia in Applied Mathematics, Volume 39. American Mathematical Society, Providence, RI, 1989.
64. B. M. Diong and J. V. Medanic. Dynamic output feedback variable structure control for system stabilization. *International Journal of Control*, 56(3):607–630, September 1992.
65. A. S. Dmitriev, A. I. Panas, and S. O. Starkov. Storing and recognizing information based on stable cycles of one-dimensional maps. *Physics Letters A*, 155(8, 9):494–499, May 27, 1991.
66. R. C. Dorf and R. H. Bishop. *Modern Control Systems*, eighth edition. Addison-Wesley, Menlo Park, CA, 1998.
67. D. Driankov, H. Hellendoorn, and M. Reinfrank. *An Introduction to Fuzzy Control*. Springer-Verlag, Berlin, 1993.
68. M. Driels. *Linear Control Systems Engineering*. McGraw-Hill, New York, 1996.
69. R. D. Driver. *Ordinary and Delay Differential Equations*. Springer-Verlag, New York, 1977.
70. W. Dunham. *Journey Through Genius: The Great Theorems of Mathematics*. Wiley Science Editions, John Wiley & Sons, New York, 1990.
71. C. Edwards and S. K. Spurgeon. On the development of discontinuous observers. *International Journal of Control*, 59(5):1211–1229, May 1994.
72. C. Edwards and S. K. Spurgeon. Compensator based output feedback sliding mode controller design. *International Journal of Control*, 71(4):601–614, November 10, 1998.
73. C. Edwards and S. K. Spurgeon. *Sliding Mode Control: Theory and Applications*. Systems and Control Book Series. Taylor & Francis, London, England, 1998.
74. G. D. Elseth and K. D. Baumgardner. *Genetics*. Addison-Wesley, Reading, MA, 1984.

75. L. E. Elsgolc. *Calculus of Variations*. Pergamon Press, London, 1961.
76. K. D. Elworthy, W. N. Everitt, and E. B. Lee, editors. *Differential Equations, Dynamical Systems, and Control Science: A Festschrift in Honor of Lawrence Markus*. Lecture Notes in Pure and Applied Mathematics, Volume 152. Marcel Dekker, New York, 1994.
77. S. V. Emelyanov, editor. *Titles in Theory of Variable Structure Control Systems (1957–1989)*. International Research Institute for Management Sciences, Moscow, 1989.
78. V. N. Faddeeva. *Computational Methods of Linear Algebra*. Dover Publications, New York, 1959.
79. L. Fausett. *Fundamentals of Neural Networks; Architectures, Algorithms, and Applications*. Prentice-Hall, Englewood Cliffs, NJ, 1994.
80. L. Faybusovich. Dynamical systems that solve linear programming problems. In *Proceedings of the 31st Conference on Decision and Control, Tucson, Arizona*, pages 1626–1631, December 1992.
81. M. J. Feigenbaum. Universal behavior in nonlinear systems. *Physica D*, 7D(1–3):16–39, May 1983.
82. G. Feng, S. G. Cao, N. W. Rees, and C. K. Chak. Design of fuzzy control systems with guaranteed stability. *Fuzzy Sets and Systems*, 85(1):1–10, January 1, 1997.
83. A. F. Filippov. Differential equations with discontinuous right-hand side. In *American Mathematical Society Translations*, pages 199–231. American Mathematical Society, Providence, RI, 1964.
84. A. F. Filippov. *Differential Equations with Discontinuous Righthand Sides*. Mathematics and Its Applications (Soviet Series). Kluwer Academic Publishers, Dordrecht, The Netherlands, 1988.
85. K. Fischle and D. Schröder. An improved stable adaptive fuzzy control method. *IEEE Transactions on Fuzzy Systems*, 7(1):27–40, February 1999.
86. S. Flåm and J. Zowe. A primal-dual differential method for convex programming. In A. Ioffe, M. Marcus, and S. Reich, editors, *Optimization and Nonlinear Analysis*, pages 119–129. Longman Scientific & Technical, Essex, England, 1992.
87. S. D. Flåm. Solving convex programs by means of ordinary differential equations. *Mathematics of Operations Research*, 17(2):290–302, May 1992.
88. D. B. Fogel. What is evolutionary computation? *IEEE Spectrum*, 37(2):26–32, February 2000.
89. G. F. Franklin and J. D. Powell. *Digital Control of Dynamic Systems*. Addison-Wesley, Reading, MA, 1980.
90. G. F. Franklin, J. D. Powell, and A. Emami-Naeini. *Feedback Control of Dynamic Systems*, fourth edition. Prentice-Hall, Upper Saddle River, NJ, 2002.
91. B. Friedland. *Advanced Control System Design*. Prentice-Hall, Englewood Cliffs, NJ, 1996.
92. A. T. Fuller. Maxwell's treatment of governors fitted with differential gears. *International Journal of Control*, 65(3):385–408, October 1996.
93. F. R. Gantmacher. *Applications of the Theory of Matrices*. Interscience Publishers, New York, 1959.
94. S. S. Ge, T. H. Lee, and G. Zhu. Genetic algorithm tuning of Lyapunov-based controllers: An application to a single-link flexible robot system. *IEEE Transactions on Industrial Electronics*, 43(5):567–574, October 1996.
95. I. M. Gel'fand. *Lectures on Linear Algebra*. Interscience Publishers, New York, 1961.
96. M. Gen and R. Cheng. *Genetic Algorithms and Engineering Design*. Wiley-Interscience, New York, 1997.
97. P. E. Gill, W. Murray, and M. H. Wright. *Numerical Linear Algebra and Optimization*, Volume 1. Addison-Wesley, Redwood City, CA, 1991.
98. M. P. Glazos, S. Hui, and S. H. Žak. Sliding modes in solving convex programming problems. *SIAM Journal of Control and Optimization*, 36(2):680–697, March 1998.
99. M. P. Glazos and S. H. Žak. Practical stabilization of nonlinear/uncertain dynamic systems with bounded controllers. *International Journal of Control*, 62(1):153–171, July 1995.
100. J. Gleick. *Chaos: Making a New Science*. Penguin Books, New York, 1988.
101. I. Glynn, editor. *An Anatomy of Thought: The Origin and Machinery of Mind*. Oxford University Press, New York, 1999.

102. R. M. Golden. The "Brain-State-in-a-Box" neural model is a gradient descent algorithm. *Journal of Mathematical Psychology*, 30(1):73–80, 1986.
103. R. M. Golden. Stability and optimization analyses of the generalized Brain-State-in-a-Box neural network model. *Journal of Mathematical Psychology*, 37(2):282–298, 1993.
104. G. H. Golub and C. F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, Baltimore, MD, 1983.
105. D. Graham and D. McRuer. *Analysis of Nonlinear Control Systems*. John Wiley & Sons, New York, 1961.
106. W. J. Grantham and A. M. Athalye. Discretization chaos: Feedback control and transistion to chaos. In C. T. Leondes, editor, *Control and Dynamic Systems; Advances in Theory and Applications, Volume 34: Advances in Control Mechanics, Part 1 of 2*, pages 205–277. Academic Press, San Diego, CA, 1990.
107. H. J. Greenberg. Equilibria of the Brain-State-in-a-Box (BSB) neural model. *Neural Networks*, 1(4):323–324, 1988.
108. S. Grossberg. Nonlinear neural networks: Principles, mechanisms, and architectures. *Neural Networks*, 1(1):17–61, 1988.
109. V. B. Haas. EE 675 Class Notes. Purdue University, School of Electrical Engineering, West Lafayette, IN, Fall 1980.
110. M. T. Hagan, H. B. Demuth, and M. Beale. *Neural Network Design*. PWS Publishing Co., Boston, MA, 1996.
111. W. Hahn. Über die Anwendung der Methode von Ljapunov auf Differenzengleichungen. *Mathematische Annalen*, 136(1):430–441, 1958.
112. W. Hahn. *Theory and Application of Liapunov's Direct Method*. Prentice-Hall International Series in Applied Mathematics. Prentice-Hall, Englewood Cliffs, NJ, 1963.
113. F. Hakl. Recall strategy of B-S-B neural network. *Neural Network World*, (1):59–82, 1992.
114. G. H. Hardy. *A Mathematician's Apology*. Cambridge University Press, Cambridge, Great Britain, 1992. Canto edition with a foreword by C. P. Snow.
115. R. L. Harvey. *Neural Network Principles*. Prentice-Hall, Englewood Cliffs, NJ, 1994.
116. M. H. Hassoun. *Fundamentals of Artificial Neural Networks*. A Bradford Book, MIT Press, Cambridge, MA, 1995.
117. S. W. Hawking. *A Brief History of Time: From the Big Bang to Black Holes*. Bantam Books, New York, 1988.
118. S. Haykin. *Neural Networks: A Comprehensive Foundation*, second edition. Prentice-Hall, Upper Saddle River, NJ, 1999.
119. U. Helmke and J. B. Moore. *Optimization and Dynamical Systems*. Communications and Control Engineering Series. Springer-Verlag, London, 1994.
120. J. Hertz, A. Krogh, and R. G. Palmer. *Introduction to the Theory of Neural Computation*. Lecture Notes Volume I in the Santa Fe Institute Studies in the Sciences of Complexity. Addison-Wesley, Redwood City, CA, 1991.
121. M. R. Hestenes. *Calculus of Variations and Optimal Control Theory*. John Wiley & Sons, New York, 1966.
122. N. J. Higham. *Handbook of Writing for the Mathematical Sciences*. SIAM, Philadelphia, 1993.
123. H. Hochstadt. *Differential Equations: A Modern Approach*. Dover Publications, New York, 1975.
124. J. H. Holland. *Adaptation in Natural and Artificial Systems; An Introductory Analysis with Applications to Biology, Control, and Artificial Intelligence*. A Bradford Book, MIT Press, Cambridge, MA, 1992.
125. J. J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences of the USA, Biophysics*, 79:2554–2558, April 1982.

126. J. J. Hopfield. Neurons with graded response have collective computational properties like those of two-state neurons. *Proceedings of the National Academy of Sciences of the USA, Biophysics*, 81:3088–3092, May 1984.
127. A. S. Householder. *The Theory of Matrices in Numerical Analysis*. Dover Publications, New York, 1975.
128. J. C. Hsu and A. U. Meyer. *Modern Control Principles and Applications*. McGraw-Hill, New York, 1968.
129. S. Hui, W. E. Lillo, and S. H. Žak. Dynamics and stability analysis of the Brain-State-in-a-Box (BSB) neural models. In M. H. Hassoun, editor, *Associative Neural Memories; Theory and Implementation*, Chapter 11, pages 212–224. Oxford University Press, New York, 1993.
130. S. Hui and S. H. Žak. Dynamical analysis of the Brain-State-in-a-Box (BSB) neural models. *IEEE Transactions on Neural Networks*, 3(1):86–94, January 1992.
131. S. Hui and S. H. Žak. Robust control synthesis for uncertain/nonlinear dynamical systems. *Automatica*, 28(2):289–298, 1992.
132. S. Hui and S. H. Žak. Low-order state estimators and compensators for dynamical systems with unknown inputs. *Systems & Control Letters*, 21(6):493–502, December 1993.
133. S. Hui and S. H. Žak. Robust output feedback stabilization of uncertain dynamic systems with bounded controllers. *International Journal of Robust and Nonlinear Control*, 3(2):115–132, August 1993.
134. S. Hui and S. H. Žak. The Widrow–Hoff algorithm for McCulloch–Pitts type neurons. *IEEE Transactions on Neural Networks*, 5(6):924–929, November 1994.
135. J. Y. Hung, W. Gao, and J. C. Hung. Variable structure control: A survey. *IEEE Transactions on Industrial Electronics*, 40(1):2–22, February 1993.
136. L. R. Hunt, R. Su, and G. Meyer. Global transformations of nonlinear systems. *IEEE Transactions on Automatic Control*, AC-28(1):24–31, January 1983.
137. D. R. Hush and B. G. Horne. Progress in supervised neural networks: What’s new since Lippmann? *IEEE Signal Processing Magazine*, 10(1):8–39, January 1993.
138. H. S. Hwang. Automatic design of fuzzy rule base for modelling and control using evolutionary programming. *IEE Proceedings—Control Theory Appl.*, 146(1):9–16, January 1999.
139. A. Isidori. *Nonlinear Control Systems*, third edition. Springer-Verlag, London, 1995.
140. U. Itkis. *Control Systems of Variable Structure*. A Halsted Press Book, a division of John Wiley & Sons, New York, 1976.
141. M. Jamshidi, M. Tarokh, and B. Shafai. *Computer-Aided Analysis and Design of Linear Control Systems*. Prentice-Hall, Englewood Cliffs, NJ, 1992.
142. J.-S. R. Jang, C.-T. Sun, and E. Mizutani. *Neuro-Fuzzy and Soft Engineering: A Computational Approach to Learning and Machine Intelligence*. MATLAB Curriculum Series. Prentice-Hall, Upper Saddle River, NJ, 1997.
143. T. Kaczorek. *Teoria Sterowania i Systemów*. Wydawnictwo Naukowe PWN—Polish Scientific Publishers, Warszawa, 1993.
144. T. Kaczorek. *Wektory i Macierze w Automatyce i Elektrotechnice*, second edition. Wydawnictwa Naukowe-Techniczne—Scientific-Technical Publishers, Warszawa, 1998.
145. P. B. Kahn. *Mathematical Methods for Scientists and Engineers: Linear and Nonlinear Systems*. Wiley-Interscience, New York, 1990.
146. T. Kailath. *Linear Systems*. Prentice-Hall, Englewood Cliffs, NJ, 1980.
147. R. E. Kalman. Mathematical description of linear dynamical systems. *SIAM Journal of Control, Series A*, 1(2):152–192, 1963.
148. R. E. Kalman. On the computation of the reachable/observable canonical form. *SIAM Journal of Control and Optimization*, 20(2):258–260, March 1982.
149. S.-J. Kang, C.-H. Woo, H.-S. Hwang, and K. B. Woo. Evolutionary design of fuzzy rule base for nonlinear system modeling and control. *IEEE Transactions on Fuzzy Systems*, 8(1):37–45, February 2000.

150. N. N. Karpinskaya. Method of "penalty" functions and the foundations of Pyne's method. *Automation and Remote Control*, 28(1):124–129, January 1967.
151. A. Kaufmann and M. M. Gupta. *Introduction to Fuzzy Arithmetic: Theory and Applications*. Van Nostrand Reinhold Electrical/Computer Science and Engineering Series. Van Nostrand Reinhold, New York, 1985.
152. D. G. Kelly. Stability in contractive nonlinear neural networks. *IEEE Transactions on Biomedical Engineering*, 37(3):231–242, March 1990.
153. T. Kenjo. *Stepping Motors and Their Microprocessor Controls*. Clarendon Press, Oxford, 1984.
154. T. Kenjo and S. Nagamori. *Permanent-Magnet and Brushless DC Motors*. Clarendon Press, Oxford, 1985.
155. M. P. Kennedy and L. O. Chua. Neural networks for nonlinear programming. *IEEE Transactions on Circuits and Systems*, 35(5):554–562, May 1988.
156. H. K. Khalil. Adaptive output feedback control of nonlinear systems represented by input–output models. *IEEE Transactions on Automatic Control*, 41(2):177–188, February 1996.
157. H. K. Khalil. *Nonlinear Systems*, second edition. Prentice-Hall, Upper Saddle River, NJ, 1996.
158. T. Khanna. *Foundations of Neural Networks*. Addison-Wesley, Reading, MA, 1990.
159. J.-H. Kim and H. Myung. Evolutionary programming techniques for constrained optimization problems. *IEEE Transactions on Evolutionary Computation*, 1(2):129–140, July 1997.
160. D. E. Kirk. *Optimal Control Theory: An Introduction*. Prentice-Hall, Englewood Cliffs, NJ, 1970.
161. G. J. Klir, U. St. Clair, and B. Yuan. *Fuzzy Set Theory: Foundations and Applications*. Prentice-Hall PTR, Upper Saddle River, NJ, 1997.
162. D. E. Knuth. *The Art of Computer Programming*, Volume 1. Fundamental Algorithms, second edition. Addison-Wesley, Reading, MA, 1973.
163. C. Koch. Computation and the single neuron. *Nature*, 385(6613):207–210, January 16, 1997.
164. S. K. Korovin and V. I. Utkin. Using sliding modes in static optimization and nonlinear programming. *Automatica*, 10(5):525–532, September 1974.
165. B. Kosko. *Fuzzy Engineering*. Prentice-Hall, Upper Saddle River, NJ, 1997.
166. J. R. Koza. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. A Bradford Book, MIT Press, Cambridge, MA, 1992.
167. J. R. Koza, F. H. Bennett, III, D. Andre, M. A. Keane, and F. Dunlap. Automated synthesis of analog electrical circuits by means of genetic programming. *IEEE Transactions on Evolutionary Computation*, 1(2):109–128, July 1997.
168. T. Kozek, T. Roska, and L. O. Chua. Genetic algorithm for CNN template learning. *IEEE Transactions on Circuits and Systems—I: Fundamental Theory and Applications*, 40(6):392–402, June 1993.
169. P. C. Krause and O. Wasynczuk. *Electromechanical Motion Devices*. McGraw-Hill, New York, 1989.
170. A. J. Krener and W. Respondek. Nonlinear observers with linearizable error dynamics. *SIAM Journal of Control and Optimization*, 23(2):197–216, March 1985.
171. M. Krstić, I. Kanellakopoulos, and P. V. Kokotović. *Nonlinear and Adaptive Control Design*. John Wiley & Sons, New York, 1995.
172. V. Kučera. A contribution to matrix quadratic equations. *IEEE Transactions on Automatic Control*, AC-17(3):344–347, June 1972.
173. H. Kwakernaak and R. Sivan. *Linear Optimal Control Systems*. Wiley-Interscience, New York, 1972.
174. C. M. Kwan. On variable structure output feedback controllers. *IEEE Transactions on Automatic Control*, 41(11):1691–1693, November 1996.
175. S. Lang. *Calculus of Several Variables*, third edition. Springer-Verlag, New York, 1987.
176. J. P. LaSalle. *The Stability and Control of Discrete Processes*. Springer-Verlag, New York, 1986.
177. E. B. Lee and L. Markus. *Foundations of Optimal Control Theory*. Robert E. Krieger Publishing Company, Malabar, FL, 1986.

178. Y. Lee, J. Q. Gong, B. Yao, and S. H. Žak. Fuzzy adaptive robust tracking control of a class of nonlinear systems. In *Proceedings of 2001 American Control Conference*, pages 4040–4045. Arlington, VA, June 25–27, 2001.
179. Y. Lee and S. H. Žak. Genetic fuzzy tracking controllers for autonomous ground vehicles. In *Proceedings of 2002 American Control Conference*, pages 2144–2149. Anchorage, AK, May 8–10, 2002.
180. K. N. Leibovic. The principle of contraction mapping in nonlinear and adaptive control systems. *IEEE Transactions on Automatic Control*, AC-9(4):393–398, October 1964.
181. J. R. Leigh. *Modelling and Simulation*. IEE Topics in Control Series 1. Peter Peregrinus, London, 1983.
182. F. L. Lewis and V. L. Syrmos. *Optimal Control*, second edition. Wiley-Interscience, New York, 1995.
183. J.-H. Li, A. N. Michel, and W. Porod. Analysis and synthesis of a class of neural networks: Linear systems operating on a closed hypercube. *IEEE Transactions on Circuits and Systems*, 36(11):1405–1422, 1989.
184. W. E. Lillo, D. C. Miller, S. Hui, and S. H. Žak. Synthesis of Brain-State-in-a-Box (BSB) based associative memories. *IEEE Transactions on Neural Networks*, 5(5):730–737, 1994.
185. C.-T. Lin and C. S. G. Lee. *Neural Fuzzy Systems: A Neuro-Fuzzy Synergism to Intelligent Systems*. Prentice-Hall PTR, Upper Saddle River, NJ, 1996.
186. R. P. Lippmann. An introduction to computing with neural nets. *IEEE ASSP Magazine*, 4(2):4–22, April 1987.
187. G. P. Liu, V. Kadirkamanathan, and S. A. Billings. Stable sequential identification of continuous nonlinear dynamical systems by growing radial basis function networks. *International Journal of Control*, 65(1):53–69, 10 September 1996.
188. D. G. Luenberger. Observers for multivariable systems. *IEEE Transactions on Automatic Control*, AC-11(2):190–197, April 1966.
189. D. G. Luenberger. An introduction to observers. *IEEE Transactions on Automatic Control*, AC-16(6):596–602, December 1971.
190. D. G. Luenberger. *Linear and Nonlinear Programming*, second edition. Addison-Wesley, Reading, MA, 1984.
191. X.-J. Ma, Z.-Q. Sun, and Y.-Y. He. Analysis and design of fuzzy controller and fuzzy observer. *IEEE Transactions on Fuzzy Systems*, 6(1):41–51, February 1998.
192. A. G. J. MacFarlane. The development of frequency-response methods in automatic control. *IEEE Transactions on Automatic Control*, AC-24(2):250–265, April 1979.
193. I. J. Maddox. *Elements of Functional Analysis*, second edition. Cambridge University Press, Cambridge, Great Britain, 1988.
194. E. H. Mamdani and S. Assilian. An experiment in linguistic synthesis with a fuzzy logic controller. *International Journal of Man–Machine Studies*, 7(1):1–13, January 1975.
195. B. B. Mandelbrot. *The Fractal Geometry of Nature*. W. H. Freeman, New York, 1983.
196. R. Marino and P. Tomei. *Nonlinear Control Design: Geometric, Adaptive and Robust*. Prentice-Hall Information and System Sciences Series. Prentice-Hall Europe, London, 1995.
197. J. C. Maxwell. On governors. *Proceedings of the Royal Society of London*, XVI:270–283, March 1868. Printed by Taylor and Francis, London, Red Lion Court, Fleet Street.
198. R. M. May. Simple mathematical models with very complicated dynamics. *Nature*, 261(5560):459–467, June 10, 1976.
199. O. Mayr. *The Origins of Feedback Control*. MIT Press, Cambridge, MA, 1970.
200. W. S. McCulloch and W. Pitts. A logical calculus of the ideas immanent in nervous activity. *Bulletin of Mathematical Biophysics*, 5:115–133, 1943.
201. J. M. Mendel. *Uncertain Rule-Based Fuzzy Logic Systems: Introduction and New Directions*. Prentice-Hall PTR, Upper Saddle River, NJ, 2001.
202. A. N. Michel and J. A. Farrell. Associative memories via artificial neural networks. *IEEE Control Systems Magazine*, 10(3):6–17, 1990.

203. A. N. Michel, J. A. Farrell, and W. Porod. Qualitative analysis of neural networks. *IEEE Transactions on Circuits and Systems*, 36(2):229–243, February 1989.
204. A. N. Michel, J. A. Farrell, and H.-F. Sun. Analysis and synthesis technique for Hopfield type synchronous discrete time neural networks with application to associative memory. *IEEE Transactions on Circuits and Systems*, 37(11):1356–1366, 1990.
205. R. K. Miller and A. N. Michel. *Ordinary Differential Equations*. Academic Press, New York, 1982.
206. M. Mitchell. *An Introduction to Genetic Algorithms*. A Bradford Book, MIT Press, Cambridge, MA, 1996.
207. R. R. Mohler. *Nonlinear Systems, Volume I: Dynamics and Control*. Prentice-Hall, Englewood Cliffs, NJ, 1991.
208. C. Moler and C. Van Loan. Nineteen dubious ways to compute the exponential of a matrix. *SIAM Review*, 20(4):801–836, October 1978.
209. F. C. Moon. *Chaotic Vibrations: An Introduction for Applied Scientists and Engineers*. Wiley-Interscience, John Wiley & Sons, New York, 1987.
210. Fr. M. Müller, C.S.S.R. *The Blessed Eucharist—Our Greatest Treasure*. TAN Books and Publishers, Rockford, IL, 1994. First published by Kelly & Piet of Baltimore, Maryland in 1868.
211. A. G. Nalecz. Sensitivity analysis of vehicle design attributes in frequency domain. *Vehicle System Dynamics*, 17:141–163, 1988.
212. A. G. Nalecz and A. C. Bindemann. Analysis of the dynamic response of four wheel steering vehicles at high speed. *International Journal of Vehicle Design*, 9(2):179–202, 1988.
213. A. G. Nalecz and A. C. Bindemann. Investigation into the stability of four wheel steering vehicles. *International Journal of Vehicle Design*, 9(2):159–178, 1988.
214. S. Nara, P. Davis, and H. Totsuji. Memory search using complex dynamics in a recurrent neural network model. *Neural Networks*, 6(7):963–973, 1993.
215. I. P. Natanson. *Theory of Functions of a Real Variable*. Frederick Ungar, New York, 1955.
216. K. Ogata. *Modern Control Engineering*, third edition. Prentice-Hall, Upper Saddle River, NJ, 1997.
217. E. Ott, C. Grebogi, and J. A. Yorke. Controlling chaos. *Physical Review Letters*, 64(11):1196–1199, March 12, 1990.
218. D. H. Owens. *Multivariable and Optimal Systems*. Academic Press, London, 1981.
219. S. K. Pal and D. Bhandari. Genetic algorithms with fuzzy fitness function for object extraction using cellular networks. *Fuzzy Sets and Systems*, 65(2/3):129–139, August 10, 1994.
220. J. Park, H. Cho, and D. Park. Design of GBSB neural associative memories using semidefinite programming. *IEEE Transactions on Neural Networks*, 10(4):946–950, July 1999.
221. T. S. Parker and L. O. Chua. Chaos: A tutorial for engineers. *Proceedings of the IEEE*, 75(8):982–1008, August 1987.
222. T. S. Parker and L. O. Chua. *Practical Numerical Algorithms for Chaotic Systems*. Springer-Verlag, New York, 1989.
223. R. V. Patel and M. Toda. Quantitative measures of robustness for multivariable systems. In *Proceedings of Joint Automatic Control Conference (JACC)*, San Francisco, CA, page TP8-A, 1980.
224. R. J. Patton and G. P. Liu. Robust control design via eigenstructure assignment, genetic algorithms and gradient-based optimization. *IEE Proceedings—Control Theory Applications*, 141(3):202–208, May 1994.
225. W. Pedrycz. *Fuzzy Control and Fuzzy Systems*, second, extended edition. Research Studies Press, Taunton, Somerset, England, 1993.
226. W. Pedrycz, editor. *Fuzzy Modeling: Paradigms and Practice*. Kluwer Academic Publishers, Boston, 1996.
227. H.-O. Peitgen, H. Jürgens, and D. Saupe, editors. *Chaos and Fractals: New Frontiers of Science*. Springer-Verlag, New York, 1992.
228. H.-O. Peitgen and D. Saupe, editors. *The Science of Fractal Images*. Springer-Verlag, New York, 1988.

229. R. Perfetti. A synthesis procedure for Brain-State-in-a-Box neural networks. *IEEE Transactions on Neural Networks*, 6(5):1071–1080, September 1995.
230. L. Personnaz, I. Guyon, and G. Dreyfus. Information storage and retrieval in spin-glass like neural networks. *Le Journal de Physique—Letters*, 46(8):L-359–L-365, 1985.
231. L. Personnaz, I. Guyon, and G. Dreyfus. Collective computational properties of neural networks: New learning mechanisms. *Physical Review A*, 34(5):4217–4228, 1986.
232. I. Peterson. *Newton's Clock: Chaos in the Solar System*. W. H. Freeman, New York, 1993.
233. C. A. Pickover, editor. *Fractal Horizons: The Future Use of Fractals*. St. Martin's Press, New York, 1996.
234. D. A. Pierre. *Optimization Theory with Applications*. Dover Publications, New York, 1986.
235. L. S. Pontryagin. *Ordinary Differential Equations*. Addison-Wesley, Reading, MA, 1962.
236. V. M. Popov. Invariant description of linear, time-invariant controllable systems. *SIAM Journal of Control*, 10(2):252–264, May 1972.
237. H. M. Power and R. J. Simpson. *Introduction to Dynamics and Control*. McGraw-Hill, London, 1978.
238. R. Pratap. *Getting Started with MATLAB 5: A Quick Introduction for Scientists and Engineers*. Oxford University Press, New York, 1999.
239. I. B. Pyne. Linear programming on an electronic analogue computer. *Transactions of the American Institute of Electrical Engineers*, 75:139–143, May 1956.
240. F. H. Raven. *Automatic Control Engineering*, fourth edition. McGraw-Hill, NY, 1987.
241. *Reader's Digest Practical Problem Solver. Substitutes, Shortcuts, and Ingenious Solutions for Making Life Easier*. The Reader's Digest Association, Pleasantville, New York, 1991.
242. A. Rodríguez-Vázquez, R. Domínguez-Castro, A. Rueda, J. L. Huertas, and E. Sánchez-Sinencio. Nonlinear switched-capacitor “neural” networks for optimization problems. *IEEE Transactions on Circuits and Systems*, 37(3):384–398, March 1990.
243. C. E. Rohrs, J. L. Melsa, and D. G. Schultz. *Linear Control Systems*. McGraw-Hill, New York, 1993.
244. W. Rudin. *Principles of Mathematical Analysis*, third edition. McGraw-Hill, New York, 1976.
245. W. J. Rugh. *Linear System Theory*, second edition. Prentice-Hall, Upper Saddle River, NJ, 1996.
246. W. J. Rugh and G. J. Murphy. A new approach to the solution of linear optimal control problems. *Transactions of the ASME: Journal of Basic Engineering*, 91(2):149–154, June 1969.
247. D. E. Rumelhart, J. L. McClelland, and the PDP Research Group. *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, Volume 1. A Bradford Book, MIT Press, Cambridge, MA, 1986.
248. M. V. Rybashov. Analog solution of algebraic and transcendental equations by the gradient method. *Automation and Remote Control*, 22(1):66–76, August 1961.
249. M. V. Rybashov. Some methods of solving systems of linear algebraic equations by electronic simulation. *Automation and Remote Control*, 23(2):223–231, July 1962.
250. M. V. Rybashov. The gradient method of solving convex programming problems on electronic analog computers. *Automation and Remote Control*, 26(11):1886–1898, November 1965.
251. M. V. Rybashov. Gradient method of solving linear and quadratic programming problems on electronic analog computers. *Automation and Remote Control*, 26(12):2079–2089, December 1965.
252. M. G. Salvadori and M. L. Baron. *Numerical Methods in Engineering*, second edition. Prentice-Hall, Englewood Cliffs, NJ, 1961.
253. J. T. Sandefur. *Discrete Dynamical Systems; Theory and Applications*. Clarendon Press, Oxford, 1990.
254. A. Schultz. Collective recall via the Brain-State-in-a-Box network. *IEEE Transactions on Neural Networks*, 4(4):580–587, July 1993.
255. R. T. Seeley. *Calculus of Several Variables: An Introduction*. Scott, Foresman and Company, Glenview, IL, 1970.
256. S. Seshagiri and H. K. Khalil. Output feedback control of nonlinear systems using RBF neural networks. *IEEE Transactions on Neural Networks*, 11(1):69–79, January 2000.

257. G. M. Shepherd, editor. *The Synaptic Organization of the Brain*, third edition. Oxford University Press, New York, 1990.
258. H. Sira-Ramírez and E. Colina-Morles. A sliding mode strategy for adaptive learning in adalines. *IEEE Transactions on Circuits and Systems—I: Fundamental Theory and Applications*, 42(12):1001–1012, December 1995.
259. H. Sira-Ramírez, E. Colina-Morles, and F. Rivas-Echeverría. Sliding mode-based adaptive learning in dynamical-filter-weights neurons. *International Journal of Control*, 73(8):678–685, May 20, 2000.
260. K.-Y. Siu, V. Roychowdhury, and T. Kailath. *Discrete Neural Computation; A Theoretical Foundation*. Prentice Hall PTR, Englewood Cliffs, NJ, 1995.
261. R. E. Skelton. *Dynamic Systems Control: Linear Systems Analysis and Synthesis*. John Wiley & Sons, New York, 1988.
262. J.-J. E. Slotine and W. Li. *Applied Nonlinear Control*. Prentice-Hall, Englewood Cliffs, NJ, 1991.
263. D. E. Smith and J. M. Starkey. Effects of model complexity on the performance of automated vehicle steering controllers: Controller development and evaluation. *Vehicle System Dynamics*, 23(8):627–645, November 1994.
264. D. E. Smith and J. M. Starkey. Effects of model complexity on the performance of automated vehicle steering controllers: Model development, validation and comparison. *Vehicle System Dynamics*, 24(2):163–181, March 1995.
265. M. C. Smith. Special issue on control education—The United Kingdom. *IEEE Control Systems*, 16(2):51–56, April 1996.
266. W. E. Snyder. *Industrial Robots: Computer Interfacing and Control*. Prentice-Hall, Englewood Cliffs, NJ, 1985.
267. O. A. Solheim. Design of optimal control systems with prescribed eigenvalues. *International Journal of Control*, 15(1):143–160, 1972.
268. R. Sommer. Control design for multivariable non-linear time-varying systems. *International Journal of Control*, 31(5):883–891, 1980.
269. E. D. Sontag. *Mathematical Control Theory: Deterministic Finite Dimensional Systems*, second edition. Springer-Verlag, New York, 1998.
270. R. Su. On the linear equivalents of nonlinear systems. *Systems & Control Letters*, 2(1):48–52, July 1982.
271. M. Sugeno and G. T. Kang. Structure identification of fuzzy model. *Fuzzy Sets and Systems*, 28(1):15–33, October 1988.
272. M. Sunwoo and K. C. Cheok. An application of explicit self-tuning controller to vehicle active suspension systems. In *Proceedings of the 29th Conference on Decision and Control, Honolulu, Hawaii*, pages 2251–2257, December 1990.
273. H. J. Sussmann and J. C. Willems. 300 years of optimal control: From the brachistochrone to the maximum principle. *IEEE Control Systems*, 17(3):32–44, June 1997.
274. K. Takaba. Robust servomechanism with preview action for polytopic uncertain systems. *International Journal of Robust and Nonlinear Control*, 10(2):101–111, February 2000.
275. T. Takagi and M. Sugeno. Fuzzy identification of systems and its applications to modeling and control. *IEEE Transactions on Systems, Man, and Cybernetics*, 15:116–132, January/February 1985.
276. K. Tanaka. Stability and stabilizability of fuzzy-neural-linear control systems. *IEEE Transactions on Fuzzy Systems*, 3(4):438–447, November 1995.
277. K. Tanaka and M. Sugeno. Stability analysis and design of fuzzy control systems. *Fuzzy Sets and Systems*, 45(2):135–156, January 24, 1992.
278. D. W. Tank and J. J. Hopfield. Simple ‘neural’ optimization networks: An A/D converter, signal decision circuit, and a linear programming circuit. *IEEE Transactions on Circuits and Systems*, CAS-36(5):533–541, 1986.

279. M. C. M. Teixeira and S. H. Žak. Stabilizing controller design for uncertain nonlinear systems using fuzzy models. *IEEE Transactions on Fuzzy Systems*, 7(2):133–142, April 1999.
280. T. Terano, K. Asai, and M. Sugeno. *Fuzzy Systems Theory and Its Applications*. Academic Press, San Diego, CA, 1992.
281. J. M. T. Thompson and H. B. Stewart. *Nonlinear Dynamics and Chaos*. John Wiley & Sons, Chichester, England, 1986.
282. Y. I. Toptsheev and A. P. Tsypliyakov. *Problem Solver in the Theory of Automatic Control*. Mashinostroyeniye, Moscow, 1977. In Russian.
283. J. G. Truxal. *Introductory System Engineering*. McGraw-Hill, New York, 1972.
284. V. I. Utkin. Variable structure systems with sliding modes. *IEEE Transactions on Automatic Control*, AC-22(2):212–222, April 1977.
285. V. I. Utkin. *Sliding Modes and Their Application in Variable Structure Systems*. Mir Publishers, Moscow, 1978.
286. V. I. Utkin. *Sliding Modes in Control and Optimization*. Springer-Verlag, Berlin, 1992.
287. I. Varga, G. Elek, and S. H. Žak. On the Brain-State-in-a-Convex-Domain neural models. *Neural Networks*, 9(7):1173–1184, October 1996.
288. D. J. Velleman. *How to Prove It: A Structured Approach*. Cambridge University Press, Cambridge, Great Britain, 1994.
289. M. Vidyasagar. *Nonlinear Systems Analysis*, second edition. Prentice-Hall, Englewood Cliffs, NJ, 1993.
290. T. L. Vincent. Control using chaos. *IEEE Control Systems*, 17(6):65–76, December 1997.
291. B. L. Walcott, M. J. Corless, and S. H. Žak. Comparative study of non-linear state-observation techniques. *International Journal of Control*, 45(6):2109–2132, 1987.
292. B. L. Walcott and S. H. Žak. State observation of nonlinear uncertain dynamical systems. *IEEE Transactions on Automatic Control*, AC-32(2):166–170, February 1987.
293. H. O. Wang, K. Tanaka, and M. F. Griffin. An approach to fuzzy control of nonlinear systems: Stability and design issues. *IEEE Transactions on Fuzzy Systems*, 4(1):14–23, February 1996.
294. L.-X. Wang. Stable adaptive fuzzy control of nonlinear systems. *IEEE Transactions on Fuzzy Systems*, 1(2):146–155, May 1993.
295. L.-X. Wang. *Adaptive Fuzzy Systems and Control; Design and Stability Analysis*. PTR Prentice-Hall, Englewood Cliffs, NJ, 1994.
296. L.-X. Wang. *A Course in Fuzzy Systems and Control*. Prentice-Hall PTR, Upper Saddle River, NJ, 1997.
297. R. Weinstock. *Calculus of Variations: With Applications to Physics and Engineering*. Dover Publications, New York, 1974.
298. J. R. Wertz, editor. *Spacecraft Attitude Determination and Control*. Astrophysics and Space Science Library, Volume 73. Kluwer Academic Publishers, Dordrecht, The Netherlands, 1997.
299. A. B. Will and S. H. Žak. Modeling and control of an automated vehicle. *Vehicle System Dynamics*, 27(3):131–155, March 1997.
300. J. L. Willems. *Stability Theory of Dynamical Systems*. John Wiley & Sons, New York, 1970.
301. W. A. Wolovich. *Linear Multivariable Systems*. Springer-Verlag, New York, 1974.
302. W. A. Wolovich. *Automatic Control Systems: Basic Analysis and Design*. Saunders College Publishing, Harcourt Brace College Publishers, Fort Worth, TX, 1994.
303. J. Y. Wong. *Theory of Ground Vehicles*, second edition. John Wiley & Sons, New York, 1993.
304. W. M. Wonham. On pole assignment in multi-input controllable linear systems. *IEEE Transactions on Automatic Control*, AC-12(6):660–665, December 1967.
305. T. Yang, C.-M. Yang, and L.-B. Yang. A detailed study of adaptive control of chaotic systems with unknown parameters. *Dynamics and Control*, 8(3):255–267, July 1998.
306. Z.-J. Yang, T. Hachino, and T. Tsuji. On-line identification of continuous time-delay systems combining least-squares techniques with a genetic algorithm. *International Journal of Control*, 66(1):23–42, January 10, 1997.

- 307. G. Yen and A. N. Michel. A learning and forgetting algorithm in associative memories: Results involving pseudo-inverses. *IEEE Transactions on Circuits and Systems*, 38(10):1193–1205, 1991.
- 308. X. Yu. Controlling chaos using input-output linearization approach. *International Journal of Bifurcation and Chaos*, 7(7):1659–1664, July 1997.
- 309. L. A. Zadeh. Fuzzy sets. *Information and Control*, 8(3):338–353, June 1965.
- 310. S. H. Žak and S. Hui. On variable structure output feedback controllers for uncertain dynamic systems. *IEEE Transactions on Automatic Control*, 38(10):1509–1512, October 1993.
- 311. S. H. Žak and S. Hui. Output feedback variable structure controllers and state estimators for uncertain/nonlinear dynamic systems. *IEE Proceedings—Control Theory Applications*, 140(1):41–50, January 1993.
- 312. S. H. Žak, W. E. Lillo, and S. Hui. Learning and forgetting in generalized Brain-State-in-a-Box (BSB) neural associative memories. *Neural Networks*, 9(5):845–854, July 1996.
- 313. S. H. Žak and C. A. MacCarley. State-feedback control of non-linear systems. *International Journal of Control*, 43(5):1497–1514, 1986.
- 314. S. H. Žak, V. Upatising, W. E. Lillo, and S. Hui. A dynamical systems approach to solving linear programming problems. In K. D. Elworthy, W. N. Everitt, and E. B. Lee, editors, *Differential Equations, Dynamical Systems, and Control Science: A Festschrift in Honor of Lawrence Markus*, Chapter 54, pages 913–946. Lecture Notes in Pure and Applied Mathematics, Volume 152. Marcel Dekker, New York, 1994.
- 315. M. Zamparelli. Genetically trained cellular neural networks. *Neural Networks*, 10(6):1143–1151, August 1997.
- 316. M. Zeitz. Comments on ‘Comparative study of non-linear state-observation techniques.’ *International Journal of Control*, 46(5):1861–1863, 1987.
- 317. Y. Zeng and S. N. Singh. Adaptive control of chaos in Lorenz system. *Dynamics and Control*, 7(2):143–154, April 1997.
- 318. A. S. I. Zinober, editor. *Deterministic Control of Uncertain Systems*. Peter Peregrinus, London, 1990.
- 319. M. Zribi and J. Chiasson. Position control of a PM stepper motor by exact linearization. *IEEE Transactions on Automatic Control*, 36(5):620–625, May 1991.
- 320. J. M. Zurada. *Introduction to Artificial Neural Systems*. West Publishing, St. Paul, MN, 1992.

Index Term	Links			
1-D map,	620			
1-norm,	627			
<i>See also</i> p -norms 2-bit A/D converter,	515			
<i>See also</i> Analog-to-digital (A/D) converter				
2-cycle,	590			
Attracting,	590			
2-D Look-Up Table block,	418			
2-norm. <i>See</i> Euclidean norm				
3-bit A/D converter,	517			
4-bit A/D converter,	517			
4-cycle,	590			
Attracting,	590			
8-cycle,	591			
A/D converter,	514			
<i>See also</i> Analog-to-digital (A/D) converter				
2-bit,	515			
3-bit,	517			
4-bit,	517			
Activation function,	516			
Cost function,	514	516		
Dynamical behavior of,	516			
Hypercube,	514			
Least significant bit (LSB),	514			
Most significant bit (MSB),	514			
N -bit,	515			
AAM. <i>See</i> Address addressable memory (AAM)				
Absolutely continuous function,	323	346	349	650
Accelerating potential,	91			
<i>See also</i> cyclotron				
Ackermann's formula,	121	122		
<i>See also</i> Pole placement				
Activation function,	509	516	518	

Index Term	Links	
Continuously differentiable,	518	519
Model of neuron,	509	
Sigmoid,	509	
Sign,	470	
Strictly monotonically increasing,	518	519
Activation nonlinearity,	526	
<i>See also</i> Activation function		
Active suspension,	223	
Adaline,	478	
<i>See also</i> Adaptive linear element		
Adaptation algorithm for weights,	479	
Adaptation strategy,	481	
Forward dynamics identification,	482	
Inverse dynamics identification,	483	
Lyapunov function candidate,	480	
Prefilter,	481	
Sira-Ramírez and Colina-Morles		
theorem,	479	
Adams–Bashforth algorithm,	71	
Adams–Moulton algorithm,	71	
Adaptation,	554	
Algorithm for weights,	479	
Law,	212	449
Strategy,	481	
Strategy of Sira-Ramírez and		
Colina-Morles,	481	
Adaptation parameter error,	449	
Adaptation law,	449	
Adaptive control,	659	
Adaptive fuzzy control,	447	
Adaptation parameter error,	449	
Bounded disturbance,	449	
Control objective,	447	
Fuzzy rules,	447	

Index Term	Links				
Fuzzy set,	447				
Plant model,	447				
Projection operator,	448				
Regressor,	448				
Sliding modes,	448				
Sliding surface,	448				
Adaptive linear element,	478				
<i>See also</i> Adaline					
Adaptive tracking controllers,	447				
Addition of fuzzy numbers,	399				
Addition of intervals of confidence,	397				
Address addressable memory					
(AAM),	532				
Adjoint equation. <i>See</i> Costate equation					
Admissible control,	143				
Admissible controllers,	280				
Affine model,	423				
Taylor series linearization,	424				
Affine transformation,	596	597			
Contractive,	597				
Expansive,	597				
Symmetry,	597				
Algebraic equations,	192				
Algebraic multiplicity,	676				
Algebraic product,	404				
Algebraic product defuzzifier,	464				
<i>See also</i> Fuzzy logic					
Algebraic Riccati equation (ARE),	244	246	267	268	611
Eigenvalues of Hamiltonian matrix,	311				
Eigenvector method,	248				
Represented as a vector equation using					
Kronecker product,	308				
Algebraic sum,	405				
Algorithm,	553				
Knuth's definition of,	553				

Index Term	Links				
Algorithms,	125				
Numerically stable,	125				
Alleles,	553				
<i>See also</i> Genetics Allelic genes,	553				
Analog-to-digital (A/D) converter,	514				
<i>See also</i> A/D converter					
Angular velocity,	22				
Aperiodic solution,	585	588			
Arc length,	227				
As a functional,	227				
Armature,	29	288			
<i>See also</i> DC motor					
Armature controlled DC motor,	287	100	135	256	296
	368				
Costate equation,	296				
Hamiltonian function,	288	296			
Optimal armature voltage,	297				
Optimal control,	296				
Optimal state-feedback,	287				
Observability matrix,	385	386			
Observer form,	386				
State variable transformation,	386				
Arrows,	670				
Field of,	670				
Artificial networks,	549				
<i>See also</i> Neural networks					
Artificial neural networks,	469				
<i>See also</i> Neural networks					
Associative memory,	509	529	530		
Correct pattern,	509				
Corrupted pattern,	509				
Crosstalk,	530				
Distorted pattern,	529				
Eigenstructure method,	550				
Forgetting capability,	533				

Index Term	Links		
High storage,	533		
Learning capability,	533		
Nonsymmetric interconnection structure,	532		
Outer product method,	550		
Recall phase,	529		
Retrieval efficiency,	533		
Spurious states,	532		
Storage phase,	529		
Stored pattern,	529		
Asymmetric trapezoid membership function,	415		
Asymptotes,	54	89	
<i>See also</i> Isocline			
Asymptotic stability in the large,	150		
Asymptotic state estimator,	387		
Asymptotically periodic trajectory,	601		
Asymptotically stable,	151		
Closed-loop system,	245		
Equilibrium states,	538		
In the large,	435		
Vertex,	537		
Variable structure system,	316		
Attracting,	581	589	590
2-cycle,	590		
4-cycle,	589	590	
Fixed points,	581		
Period-2 cycle,	586		
Attractive,	151	320	
<i>See also</i> Convergent;			
Switching surface			
Autoassociation,	529		
Automatic programming,	575		
Automotive suspension,	42	222	
Autonomous,	150		

Index Term	Links			
Average contractivity condition,	597			
<i>See also</i> Iterated function system (IFS)				
Average crossover operator,	562			
Average fitness,	559			
Axon,	469			
End branches,	469			
Synaptic knobs,	469			
Back emf, 34				
<i>See also</i> DC motor Backpropagation,	484	487	491	
Chain rule,	491			
Descent gradient method,	484			
Differentiable activation function,	484			
Forward pass,	484			
Gradient algorithm,	486			
Gradient descent algorithm,	486			
Identification algorithm,	492			
Identifier,	493			
Perceptron,	485			
Reverse pass,	484			
Three-layer neural network,	485			
Training algorithm,	487	491		
Training pairs,	485			
Backward Euler algorithm,	63			
Ball map,	615			
Bands of frequencies,	615			
Equilibrium solutions to,	616			
Orbit of,	618			
Banach space,	231	656		
Continuous functions,	231	656		
Bang-bang control,	295			
Switching curve,	295			
Barbalat's lemma,	210	213	451	454
Barnsley fern,	598	599		
Basin of attraction,	530	533	603	
Behavior on hypercube boundary,	533			

Index Term	Links			
Reduced-order gBSB-type models,	533			
Bellman's lemma,	443			
Bellman–Gronwall lemma,	332	663		
Bendixon's theorem,	58	59	89	
Bessel inequality,	9			
Bias,	470			
<i>See also</i> Threshold Bias function,	500			
Center of,	500			
Biconditional,	626			
Bicycle model,	24	563		
Ground vehicle,	563			
Bifurcation,	583			
Diagram,	587	621		
Definition of,	583			
Exchange of properties,	583			
Binary complement,	572			
Binary encoding,	556			
Binary fuzzy relation,	404			
Biological neuron	469			
Model of,	470			
Birth rate,	580			
Bit maps,	543			
Stacking operator,	543	544		
Bolzano–Weierstrass theorem,	214	645	646	650
Bouncing ball on vibrating plate,	613			
Coefficient of restitution,	614			
Collisions,	613			
Controllable target,	618			
Controller construction,	617			
Discrete Lyapunov equation,	618			
Plate displacement,	613			
Principle of conservation of				
momentum,	613			
Velocity of the plate,	613			
Bouncing ball system,	607			

Index Term	Links		
<i>See also</i> Bouncing			
ball on vibrating plate			
Boundary,	649		
Boundary conditions,	307		
Boundary layer,	201	202	208
Controller,	203	209	
Estimator,	338	339	
Tracking controller,	208		
Boundary point,	649		
Bounded above,	640		
Bounded below,	641		
Bounded disturbance,	449		
Bounded equilibrium state,	186		
<i>See also</i> Lagrange stable			
Bounded function,	648		
Bounded product,	405		
Bounded sequence,	645	647	
Bounded set,	649		
Bounded sum,	405		
Brachistochrone (shortest time) problem,	229	297	
Angle of departure,	300		
Boundary conditions,	300		
Condition for control optimality,	298		
Costate equations,	298		
Cycloid,	302		
Generalization,	313		
Hamiltonian function,	298		
Implicit feedback control law,	300		
Model,	297		
Performance index,	297		
State-space model,	301		
Time of quickest descent,	301		
Upside-down cycloid,	302		
Brachystochrone problem.			
<i>See</i> Brachistochrone problem			

Index Term	Links	
Brain,	469	549
Modeling,	549	
Brainlike computation,	469	
Brain-State-in-a-Box (BSB) neural		
network,	529	530
BSB network	533	
<i>See also</i> Brain-State-in-a-Box (BSB)		
neural network		
Interconnection matrix,	533	
BSB neural network,	529	530
<i>See also</i> Brain-State-in-a-Box (BSB)		
neural network		
Basin of attraction,	530	
Discrete dynamical system,	530	
Hypercube,	530	
Noisy patterns,	532	
Calculus of variations,	227	
CAM. <i>See</i> Content addressable memory		
Canonical GA,	555	559
<i>see also</i> Canonical		
genetic algorithm (GA)		
Analysis,	559	
Schema,	559	
Template,	559	
Canonical genetic algorithm (GA),	555	
Cantor set,	594	595
Disconnected fractal,	594	
Fractal dust,	594	
Middle third,	594	
Cantor's ternary set,	594	
<i>See also</i> Cantor set		
Carrying capacity,	38	580
<i>See also</i> Logistic equation		
Cart with an inverted pendulum.		
<i>See</i> Stick balancer		

Index Term	Links			
Cartesian coordinates,	11			
Cartesian product space,	407			
Cauchy sequence,	656	661		
Cauchy theorem,	645			
Cauchy–Schwarz inequality,	351	477	628	
Causality property,	6			
Definition of,	6			
Cayley–Hamilton theorem,	99	113	121	
Cell nucleus,	553			
Center average defuzzifier,	409	566	567	570
<i>See also</i> Fuzzy logic				
Center-of-gravity defuzzifier,	464			
<i>See also</i> Fuzzy logic				
Centrifugal governor,	21			
<i>See also</i> Flyball governor; Watt’s governor				
Centrifugal speed regulator,	4			
<i>See also</i> Flyball governor; Watt’s governor				
Centroid of support,	408			
Centroidal defuzzifier,	408			
Chain rule,	347	353	383	491
Change of state variables,	111			
Change-in-error,	412			
Chaos,	578			
Double pendulum,	622			
Game,	621	622		
In information processing,	620			
Logistic equation,	579			
Loss of information,	578			
Lyapunov exponents,	578			
Uncertainty growth,	578			
Chaotic,	578			
Algorithm,	609			
Attractor,	607			
Behavior,	592			

Index Term	Links			
Control algorithm,	607			
Discrete-time system,	619			
Gingerbreadman,	621			
Motions,	588			
Neural networks,	620			
Solutions,	588			
System,	578			
Transient to periodic motion,	616			
Chaotic motion,	588	601		
Positive largest Lyapunov exponent,	604			
Chaotic system,	578			
Random-like behavior,	578			
Chaotic trajectory,	601			
Definition of,	601	602		
Characteristic equation,	99	120	217	
Of Hamiltonian matrix,	251	253		
Characteristic polynomial,	113	129		
Characteristic polynomial equation,	632			
Chattering,	201	209	319	
Explanation of,	319			
Chromosomes,	554	567		
Fitness,	559			
Clockwise rotation,	256			
Closed ball,	201			
Closed path,	58			
Closed set,	649			
Closed-loop characteristic equation,	84			
Closed-loop control system,	3			
Closed-loop estimator,	128	129		
Closed-loop feedback system,	41			
Toilet-flushing system,	41			
Closed-loop system,	120	196	245	419
Characteristic equation of,	120			
Poles of,	120			
Stability analysis of,	444			

Index Term	Links				
Closed-loop systems,	207				
Clustering training patterns,	506				
Cobweb,	583				
Aperiodic solution,	585				
Period-four solution,	585				
Two-cycle solution,	585				
Coefficient of restitution,	614				
Coercive function,	185				
<i>See also</i> Radially unbounded function					
Cofactor,	631				
Cofactor matrix,	632				
Collisions,	613				
Colon,	625				
Notation for sets,	625				
Combined controller–estimator compensator,	120 468	131	148	388	392
Full-order estimator,	131				
Separation principle,	133				
Combined controller–full-order–estimator compensator,	131				
DC motor system,	138				
Combined fuzzy controller–estimator compensator,	442				
Combined objective function,	562				
Common positive definite matrix,	439				
Compact set,	216	649			
Compact subset,	656				
Companion forms,	111				
Controller form,	111				
Observer form,	117	382			
Comparison method,	442				
For linear systems,	442				
Comparison theorem,	185	450	454	664	

Index Term	Links		
Comparison,	664		
Equation of,	664		
Compatible norms,	652		
Compensator,	4	131	
<i>See also</i> Combined controller–estimator compensator			
Competing species,	605		
Complement,	395		
Completing the squares,	452	456	
Complex conjugate pairs,	120		
Complex conjugate transpose,	166		
Complex numbers,	629		
Complex plane,	84		
Components of vector,	626		
Composition of fuzzy relations,	405	406	
Fuzzy matrix,	406		
Max-* composition,	405	406	
<i>See also</i> Max-star composition			
Max-star composition,	405		
Compromise	226		
<i>See also</i> Engineering practice			
Computational energy function,	510	511	
Conditional,	625		
Conjugate transpose,	337		
Connected set,	671		
Conservative vector field,	11	40	375
Constrained optimization,	342		
Penalty function method,	342		
Constraints,	2		
Constraints on the control inputs,	292		
Pontryagin's minimum principle,	292		
Constructability,	107	111	
Constructable,	111		
Content addressable memory (CAM),	532		

Index Term	Links				
<i>See also</i> Address addressable memory (AAM)					
Continuous function,	79	648			
Definition of,	648				
Maximum–minimum theorem,	649				
Weierstrass theorems,	649	650			
Continuous minimum time regulator problem,	280				
Continuously differentiable,	154				
Function,	189	216			
Continuous-time linear system,	97				
Controllable,	97				
Reachable,	97				
Contraction maps,	659	660			
Adaptive control,	659				
Fixed point,	660				
Contractive affine transformation,	597				
Contradiction,	626				
Proof by,	626				
Contraposition,	99	106	626		
Proof by,	626				
Control error,	410				
Control objective,	2	447	563		
Control problem,	2				
Control strategy,	102	456			
Control surface,	414				
Control system,	3				
Closed-loop,	3				
Open-loop,	3				
Controllability,	94	96	97		
Controllability index,	114				
Controllability indices,	114	116	118	125	
Controllability matrix,	98	105	112	114	122
	124	370			
Controllability property,	145				

Index Term	Links				
Controllable pair,	114				
Controllable target,	618				
Controlled system,	666				
Solution of,	666				
Controller,	2	3			
Design algorithm,	466				
Model reference,	197				
Controller form,	111	112	113	115	369
Canonical,	126				
Companion form,	111	116	121	122	123
	124	138			
State variable transformation to,	115				
State-space transformation,	391				
Controller inputs,	410				
Controller outputs,	410				
Controlling chaotic systems,	607				
Bouncing ball on vibrating plate,	613				
Bouncing ball system,	607				
Chaotic attractor,	607	608			
Controllable target,	608				
Discrete Lyapunov equation,	611				
Ergodicity,	607				
Feedback saturating control,	609				
Hénon map,	607				
LQR technique,	607	608			
<i>See also</i> Linear					
quadratic regulator (LQR)					
Lyapunov function,	609				
Ott, Grebogi, and Yorke (OGY) method,	607				
Saturating linear optimal control,	267				
Convergence,	128				
Rate of,	128				
Convergent,	151				
<i>See also</i> Attractive					
Convergent matrix,	167				
Convergent sequence,	214	642	643	645	

Index Term	Links	
Convex combination,	226	
Convex constrained optimization,	425	
Local linear model,	425	
Convex function,	341	
Convex fuzzy set,	398	
Convex optimization,	340	341
Feasible set,	341	
Convex set,	322	341
Coordinates of vector,	626	
Cornering stiffness coefficient,	25	
Corollary,	624	
Definition of,	624	
Correct pattern,	509	
Corrector,	65	
<i>See also</i> Predictor-corrector method		
Equations,	66	
Corrupted pattern,	509	
Cost function,	514	516
As a landscape,	514	
Cost functional,	3	
Costate equation,	293	304
Solution to,	304	
Coulomb friction model,	35	
Counterclockwise rotation,	256	
Crisp set,	394	
Critical point,	657	
Crossing site,	554	
Crossover,	554	557
<i>See also</i> Genetic algorithm		
Crossing site,	554	
Offspring chromosomes,	557	561
One-point,	554	
Parent chromosomes,	554	561
Two-point,	558	

Index Term	Links	
Crossover probability,	557	
<i>See also</i> Crossover rate		
Crossover rate,	557	567
<i>See also</i> Crossover probability		
Crosstalk,	530	
Curve,	667	668
Differentiable,	668	
Parameterization of,	668	
Cycloid,	302	
Cyclotron,	90	91
Accelerating potential,	91	
Dee,	90	
Electric oscillator,	90	
Magnetic field,	90	
D'Alembert's equation,	12	
Damping term,	580	
Data pairs,	659	
DC motor,	100	368
Armature-controlled,	100	
Field-controlled,	47	
Dead zone,	92	
Describing function of,	92	
Dead zone nonlinearity,	93	
Decision boundary,	471	
Separating hyperplane,	471	
Decision-making logic,	410	
<i>See also</i> Inference engine		
Decoding,	561	
Decreasing sequence,	642	
Decrescent function,	180	
Dee,	90	
Defining length of schema,	559	
Defuzzification,	408	
Center average defuzzifier,	409	

Index Term	Links				
Centroidal defuzzifier,	408				
Dendrites,	469	471			
Derivative of a function with respect to a vector field,	370				
Derivative of a vector field with respect to a vector field,	369				
Derivative of Dh with respect to the vector field,	370				
Descent gradient algorithm,	491				
Descent gradient method,	484	506			
Describing function,	84				
Ideal relay,	84				
Normalized plot of,	86				
Describing function method,	59	77	82		
Closed-loop characteristic equation,	84				
Complex plane,	84				
Filtering hypothesis,	83				
Fundamental component,	83				
Fundamental harmonic,	87				
Higher harmonics,	87				
Low-pass filter,	83				
Mean level,	83				
Polar plot,	84				
Sinusoidal describing function,	88				
Transfer function,	83				
Design algorithm for fuzzy systems,	466				
Design model,	2	8	420	436	467
Takagi-Sugeno-Kang (TSK) fuzzy model,	422				
Design parameter,	450	451			
Desired closed-loop characteristic polynomial,	120	122	125	137	
Desired gains,	121				
Detectable,	111	250			

Index Term	Links			
Determinant,	631			
Cofactor matrix,	632			
Properties,	631			
Laplace expansion,	631			
Diagonalizability,	263			
Sufficient condition for,	263			
Diffeomorphism,	371	373	374	375
Invertibility,	374			
Inverse transformation,	373			
Difference equation,	3	213		
Differentiable activation function,	484			
Differentiable curve,	668			
Differentiable vector field,	673			
Differential equation,	661			
First-order,	661			
Integrating factor,	663			
Laplace transform of,	161			
Solution of,	662			
Differential equations,	2	320		
Filippov's approach,	321			
Generalized solution,	322			
With discontinuous right-hand sides,	320			
Differential inclusion,	323			
Differential inequalities,	663			
Digital positioning,	28			
Digital pulses,	29			
Dimension of subspace,	654			
Dimension of vector,	626			
Disconnected fractal,	594			
Discontinuous control strategy,	201			
Discontinuous controller chattering,	201			
Discontinuous estimator,	337			
Discontinuous robust controllers,	194			
Discontinuous surface,	320			
<i>See also</i> Switching surface				

Index Term	Links	
Discontinuous tracking controller,	208	
Discrete algebraic Riccati equation,	273	280
Discrete fuzzy models,	437	
Stability of,	437	
Discrete linear quadratic regulator		
problem,	277	
First-order necessary condition,	278	
Optimal feedback gain,	279	
Sequence of gain matrices,	279	
Discrete Lyapunov equation,	611	618
Discrete-time dynamical system,	220	604
Discrete-time Lyapunov equation,	165	166
Discretization chaos,	605	
Competing species,	605	
Distance from a point,	213	
Distance measure,	344	
Distorted pattern,	529	
Disturbance,	461	
Divergent sequence,	642	
DNA,	553	
Domain of function,	648	
Don't care element,	104	375
<i>See also</i> Don't cares		
Don't cares,	559	
Douady's rabbit,	596	
<i>See also</i> Filled Julia set		
Double integrator,	317	
Double pendulum,	622	
Lagrangian,	622	
LQR approach,	622	
Drastic product,	405	
Drastic sum,	405	
Dual composition,	406	
Dual system,	117	
Dynamic programming,	273	

Index Term	Links	
Principle of optimality (PO),	273	
Dynamic programming,	308	
Maximum cost path,	308	
Dynamical behavior of A/D converter,	516	
Dynamical gradient circuit analysis,	343	
Convergence phase,	348	
Distance measure,	344	
Reaching phase,	343	
Surface to which trajectory is confined,	345	
Dynamical gradient system,	342	
Filippov's method,	343	
Dynamical system,	1	5
Axiomatic definition of,	5	
Schematic of,	5	
EA tunes fuzzy rules and membership		
functions,	568	
<i>See also</i> Evolutionary algorithm		
EA tunes fuzzy rules,	565	
EA tunes membership functions,	567	
EA,	561	
<i>See also</i> Evolutionary algorithm		
EDVAC. <i>See</i> Electronic Discrete Variable		
Automatic Computer		
Eigenstructure method,	550	
Eigenvalue of a matrix,	676	
Algebraic multiplicity of,	676	
Geometric multiplicity of,	676	
Multiplicity of,	676	
Eigenvalues of a matrix,	333	
<i>See also</i> Spectrum		
Eigenvalues on the imaginary axis,	192	
Eigenvector method,	248	
Eigenvector test.		
<i>See</i> Popov-Belevitch-Hautus		
(PBH) test		

Index Term	Links			
Eigenvectors,	253			
Electric drive system,	456			
Disturbance,	461			
Fuzzy logic system,	456	458		
m-file,	458			
Reference signal,	458			
Reference system,	456			
SIMULINK block diagram,	457	460	461	
Sliding surface,	457			
Tracking error,	458			
Electric oscillator,	90			
Electromechanical motion devices,	28			
Electronic Discrete Variable Automatic				
Computer (EDVAC),	550			
Electronic Numerical Integrator and				
Computer (ENIAC),	550			
Elitism,	575			
Elitist strategy,	563	567	572	
Super chromosome,	563			
Ellipsoid semiaxes,	602			
Encoding,	555	563		
Binary,	556			
End branches,	469			
End-point conditions,	238	293		
Pontryagin's minimum principle,	293			
Energy,	9			
Energy of configuration. <i>See</i> Potential energy				
Engineering practice,	133			
Compromise,	133			
Engine-governor equations,	76			
ENIAC. <i>See</i> Electronic Discrete Variable				
Automatic Computer				
Equation of comparison,	444	664		
Equilibrium point,	150	423	518	519
Equilibrium position,	149			

Index Term	Links			
Equilibrium state,	74	150	151	213
Asymptotically stable,	151	535		
Attractive,	151			
Constant solution,	533			
Convergent,	151			
Instability of,	215			
Lyapunov exponents of,	603			
Stable,	151	534		
Translation of,	527			
Uniformly asymptotically stable,	152			
Uniformly exponentially stable,	526			
Unstable,	153			
Equivalent control,	327	330		
Equivalent system,	327			
Ergodicity,	607			
Error,	60	412		
Estimate of,	60			
Error signal,	4			
Estimate of the error,	60			
Estimate of the RAS,	268			
<i>See also</i> Region of asymptotic stability (RAS)				
Optimization problem,	268			
Estimation error,	128	498		
Estimator,	120	126		
Closed-loop,	128			
Full-order,	127			
More descriptive than “observer,”	127			
Open-loop,	127			
Pole selection,	133			
Euclidean norm,	479	627		
<i>See also</i> 2-norm				
Euler algorithm,	63			
Backward,	63			

Index Term	Links				
Forward,	63				
Euler formula,	89				
Euler method,	501				
See also Euler algorithm Euler's method,	62				
Step length,	62				
Euler–Lagrange equation,	237	238	239	240	241
Evolutionary algorithm (EA),	561	562	563		
Mating pool,	567				
Evolutionary algorithm-aided fuzzy logic controller,	565				
Evolutionary fuzzy logic controllers,	563				
Ground vehicle,	563				
Step-lane change maneuver,	564				
Evolutionary processes,	553				
Exact penalty function,	343				
Exclusion,	393				
Exclusive-OR (XOR) problem,	487	488	489		
Expansive affine transformation,	597				
Expert knowledge,	393				
Exponential matrix,	666				
Exponential stability,	522				
External input signal,	132				
Extremal,	238	237	243	307	
See also Calculus of variations;					
Euler–Lagrange equation					
Falsity,	625				
Family of curves,	235				
Feasible direction,	361				
Feedback controller,	120				
Feedback gains,	326				
Sliding mode controller,	326				
Feedback linearization,	367	380			
Feedback loop,	4				
Feedback matrix,	120				
Feedback saturating control,	609				

Index Term	Links	
Algebraic Riccati equation,	611	
Feedback system,	83	
Nonlinear,	83	
Feedforward network,	490	
Adaptation algorithm,	490	
Backpropagation,	490	
Feigenbaum number α ,	593	
Feigenbaum number δ ,	592	593
Field inductance,	47	
Field of arrows,	670	
Field of forces,	670	
Field-controlled DC motor,	47	
Filippov's approach,	321	357
Filled Julia set,	595	
Douady's rabbit,	596	
Filtering hypothesis,	83	
First difference,	165	
First-order differential equation,	662	
First-order discrete-time equation,	583	
Cobweb,	583	
First-order necessary condition,	236	
For relative minimum,	272	
For unconstrained optimization,	269	
Fitness function,	554	568
Fitness of chromosome,	559	
Fitness of schema,	559	
Fitness,	559	
Average,	559	
Fixed coordinates,	27	
Fixed point,	213	660
Shifting,	586	
Fixed set,	625	
Universe of discourse,	625	
FLC. <i>See</i> Fuzzy logic controller		
Flip bifurcations,	588	

Index Term	Links		
Flow,	602		
Flyball governor,	4	21	175
<i>See also</i> Centrifugal speed regulator; Watt's governor			
Stability of,	175		
Force,	9		
Force field,	670	674	
Forgetting	547		
Forgetting capability,	533		
Forward dynamics identification,	482		
Forward Euler algorithm,	63		
Forward loop,	4		
Forward pass,	484		
Fourier series,	77	81	83
Fourth-order Runge–Kutta method,	70		
Fractal	594		
Cantor set,	594		
Chaos game,	621	622	
Dimension,	594		
Scaling	594		
Self-similarity,	594		
Self-similarity dimension,	594		
Sierpiński carpet,	622		
Sierpiński gasket,	621	622	
Von Koch snowflake curve,	593		
Fractals,	578		
Non-Euclidean objects,	593		
Irregular,	593		
Friction,	15		
Friction coefficient,	434		
Frobenius norm,	503	653	
<i>See also</i> Matrix norm			
Front axle,	24		
Front tire slip angle,	26		
Front-wheel steering,	199	573	

Index Term	Links				
Vehicle,	45				
Vehicle bicycle model,	45				
Full rank factorization,	250	636			
Full-order asymptotic state estimator,	146				
Full-order estimator,	127	131			
DC motor system,	138				
Function,	229	323	395	648	
Absolutely continuous,	323	346	349	650	
Bounded,	648				
Continuous,	648	649			
Domain of,	648				
Odd,	85				
Ordered pairs,	395				
Range of,	648				
Support of,	407				
Uniformly continuous,	650				
Function space,	656				
Real-valued functions,	656				
Functional,	227	230	233	236	237
	307				
Variation of,	230	231	233		
Fundamental component,	83				
Fundamental harmonic,	87				
Fundamental theorem of calculus,	672				
Fuzziness,	463				
Fuzzy,	463				
Fuzzy adaptive control law,	450	453			
Barbalat's lemma,	451	454			
Comparison theorem,	450	454			
Completing the squares,	452	456			
Design parameter,	450	451			
Lyapunov function candidate,	451	454			
Time derivative of V ,	450	454			
Fuzzy arithmetic,	396				
Interval arithmetic,	396				
Interval of confidence,	396				

Index Term	Links			
Fuzzy closed-loop system,	444			
Stability analysis of,	444			
Fuzzy conjunction,	404			
<i>See also t</i> -norm,	404			
Fuzzy controller,	567			
Fuzzy controller-estimator compensator,	442			
Fuzzy disjunction,	404			
<i>See also t</i> -conorm,	404			
Fuzzy engineering,	463			
Fuzzy estimator,	440			
IF-THEN rules,	440			
Fuzzy implication,	404	405		
Fuzzy conjunction,	404			
Fuzzy disjunction,	404			
Fuzzy input,	407			
Fuzzy intersection,	404			
Fuzzy logic,	393			
Crisp set,	394			
Exclusion,	393			
Fuzzy number,	396			
Fuzzy set,	393			
Inclusion,	393			
Linguistic variable,	394			
Membership,	393	394		
Membership function,	395			
Partial membership,	394			
Standard additive model (SAM),	407			
Universe of discourse,	395			
Fuzzy logic controller (FLC),	393	410	563	568
Change-in-error,	412			
Closed-loop system,	419			
Control error,	410			
Control surface,	414			
Controller inputs,	410			

Index Term	Links	
Controller outputs,	410	
Error,	412	
Inference engine,	410	
Knowledge base,	410	
Linguistic variable,	413	414
Rules matrix,	412	
SIMULINK diagram,	419	
Vehicle model,	412	
Fuzzy logic system,	447	490
Feedforward network,	490	
Fuzzy logic target tracking,	577	
Fuzzy matrix,	404	406
Fuzzy maximum,	402	
Of fuzzy numbers,	402	
Fuzzy minimum,	402	
Of fuzzy numbers,	402	
Fuzzy minimum operator,	405	
Fuzzy model,	420	421
Asymptotic stabilizability in the		
large of,	428	
Design model,	420	421
Stabilizability of,	427	
Takagi-Sugeno-Kang (TSK) fuzzy		
model,	421	
Fuzzy number,	396	397
Convex and normal fuzzy set,	398	
Interval of confidence,	397	
Fuzzy numbers,	398	
Addition,	399	
Level of presumption,	398	
Triangular,	400	402
Fuzzy output,	407	
Fuzzy product,	404	
Fuzzy reasoning,	408	409
<i>See also</i> Defuzzification		

Index Term	Links		
Fuzzy relation,	407	396	403
Binary,	404		
Composition,	405		
Dual composition,	406		
Fuzzy implication,	404		
Fuzzy matrix,	404		
Product space,	403		
Fuzzy rule,	573		
Premise of,	571		
Fuzzy rule base,	570		
Fuzzy rules,	442	447	577
Fuzzy set,	393	447	
Algebra of,	463		
Membership function,	395		
Support of,	407		
Universe of discourse,	395		
Fuzzy sets,	395		
Complement,	395		
Intersection of,	395		
Union of,	395		
Fuzzy singletons,	567	568	
Fuzzy sliding mode controller,	464		
Time delay,	464		
Fuzzy state estimator,	441		
Fuzzy system,	407	466	
Controller design algorithm,	466		
Standard additive model (SAM),	407	409	
State feedback linear control law,	466		
Fuzzy union,	405		
Fuzzy zero,	399		
GA. <i>See</i> Genetic algorithm			
Gain matrix,	124		
Gaussian radial basis function,	494	496	500
Center of basis function,	500		
gBSB network,	533		

Index Term	Links			
<i>See also</i> generalized				
BSB network				
Asymptotically stable equilibrium				
states,	538			
Spurious states,	539			
Weight matrix construction algorithm,	540			
Weight matrix,	541	544		
gBSB neural net,	552			
Asymptotically stable equilibrium,	552			
Vertex,	552			
gBSB-based associative memory,	543			
Bit maps,	543	544	545	546
Forgetting	547			
Iterative learning	546			
Learning	546			
Gear,	22	257		
<i>See also</i> Centrifugal governor;				
One-link robot manipulator				
Gear multistep algorithm,	71			
Gear train,	256			
Gene,	553			
Alleles,	553			
General ordinary differential equation,	662			
Generalized BSB (gBSB) network,	533			
Generalized force,	15			
Generalized intersection operator,	404			
<i>See also</i> t -norm				
Generalized Lyapunov function,	325			
Sliding mode controller,	325			
Generalized solution in the sense of				
Filippov,	322			
Convex set,	322			
Differential inclusion,	323			
Generalized union operator,	404			
<i>See also</i> t -conorm				

Index Term	Links		
Generation,	558		
Genetic algorithm,	553	555	
Canonical,	555		
Combined objective function,	562		
Crossover operator,	554		
Decoding,	561		
Derivative-free optimization,	554		
Encoding,	555		
Fitness function,	554	555	
Generation,	558		
Genetic operators,	554		
Implementing,	555		
Mutation operator,	555		
On-line identification,	575		
Population of chromosomes,	555	556	
Run,	558		
Selection operator,	554		
Set of generations,	558		
Genetic operator,	554	571	572
Binary complement,	572		
Crossover,	554		
Mutation,	554		
Random variable,	572		
Selection,	554		
Genetic processes,	553		
Genetic programming,	575		
Genetics,	574		
Garden pea,	574		
Principles of,	574		
Genome,	554		
Genotype,	554		
Genotype space,	554		
Geodesics,	229		
Geodesics problem,	229		
Geometric multiplicity,	676		

Index Term	Links				
Global asymptotic stability,	150				
Global state transformation,	380				
Global transformation,	380				
Controller form,	380				
Globally asymptotically stable,	438				
Globally stable,	149				
Globally stable equilibrium state,	219				
Globally uniformly asymptotically stable,	446				
Gradient,	363	425	432	657	674
Algorithm,	486	490			
Descent algorithm,	486				
Operation,	363				
Greatest lower bound,	641				
<i>See also</i> Infimum,	641				
Green's theorem,	58	674			
Grid over the input-state space,	500				
Ground vehicle,	22	563	573		
Bicycle model of,	563				
Evolutionary fuzzy logic controllers,	563				
Front wheel steering,	563				
Fuzzy logic controller,	563				
Lateral position,	564				
Lateral velocity,	564				
Modeling of,	22				
Steering angle,	564				
Tire slip angel,	564				
Turning maneuver of,	22				
Growing RBF network,	500				
<i>See also</i> Radial-basis					
function (RBF) networks					
Grid over the input-state space,	500				
Simulating using Euler's method,	501				
State estimation error,	501				
H.O.T. <i>See</i> Higher-order terms					
Hamiltonian,	294				

Index Term	Links			
Hamiltonian function,	281	283	288	
Hamiltonian matrix,	248	250	252	255
Eigenvalues of,	250	311		
Linearly independent eigenvectors of,	253			
Hamilton-Jacobi-Bellman (HJB) equation,	283			
Boundary condition,	283			
Hamiltonian function,	283			
Hamming distance,	552			
Harmonic linearization,	77			
<i>See also</i> Describing function method				
Hénon map,	607	609	610	612
Controllable target,	612			
Orbit of,	610			
Hessian,	658			
Heun's method. <i>See</i> Predictor-corrector method				
Hidden layer,	486			
High storage,	533			
Higher harmonics,	87			
Higher-order terms (H.O.T.),	189	284		
High-frequency dynamics,	436			
High-frequency modes,	201			
HJB equation,	283			
<i>See also</i> Hamilton-Jacobi-Bellman (HJB) equation				
Hölder inequality,	432			
Hölder norms,	627			
<i>See also</i> p -norms Hooke's law,	142	144		
Hopfield network,	508	517		
Activation function,	518			
Activation nonlinearity,	526			
Analog-to-digital (A/D) converter,	514			
Associative memory,	509			
Asymptotic stability,	517			
Circuit implementation of,	509			

Index Term	Links				
Circuit realization of,	510				
Computational energy function,	510	511			
Equilibrium points,	518				
Equivalent representation of,	510				
Exponential stability,	522				
Interconnected information,	517	522			
Interconnected system,	512				
Interconnections,	524				
Lyapunov function candidate,	523				
Lyapunov's second method,	517	522			
Model of biological neural networks,	508				
Model of,	509	513			
Neuron,	508				
Nonlinear amplifier,	509				
Nonlinear interconnected systems,	508				
Phase plane,	528				
Sector condition,	524				
Single neuron of,	520				
Stability analysis of,	517				
Uniformly exponentially stable,	523	525	526		
Hurwitz and Routh stability criteria,	171				
Hurwitz matrix,	171	173	174	175	503
Hurwitz theorem,	173				
Hypercube,	465	514	530	533	535
	537				
Vertex of,	534				
Hypothesis,	624				
Definition of,	624				
Ideal relay,	84				
Describing function of,	84				
Ideal sliding mode,	331				
Identification,	492				
Algorithm,	492				
Architecture,	497				
Error,	497	499	504		

Index Term	Links			
Network,	497			
Identification of dynamical systems,	497			
Neural identifier,	497	498		
Identification of multivariable systems,	503			
Identifier,	493			
Identity matrix,	630			
Identity observer,	130			
IF,	625			
<i>See also</i> IF–THEN rules				
IFS,	598			
<i>See also</i> Iterated function system (IFS)				
Barnsley fern,	598	599		
Definition of,	597			
Roulette-wheel method,	598			
IF–THEN rules,	407	411	440	571
IF–THEN statement,	393			
Image,	654			
<i>See also</i> Range				
Imaginary part of complex number,	333			
Inclusion,	393			
Inconsistent linear equations,	658			
Least squares solution,	658			
Increasing function,	180			
Increasing sequence,	642			
Independent rows,	655			
Index set,	407			
Induced matrix norms,	653			
Induced norm,	651			
Induction,	644			
Inference,	625			
Inference engine,	410			
Infimum,	641	649		
<i>See also</i> Greatest lower bound				

Index Term	Links	
Infinity norm,	627	
<i>See also</i> p -norms		
Initial conditions,	599	
Sensitive dependence upon,	599	
Initial value problem,	213	
Inner product,	627	
<i>See also</i> Scalar product		
Input scaling,	197	
Input vectors,	472	
Linearly separable,	472	
Input voltage switchings,	31	
Inputs,	1	3
Input-state space,	500	
Grid over,	500	
Instability,	187	215
Of equilibrium states,	215	
Of zero solution,	219	
Instance of schema,	560	
Integral inequalities,	663	
Integral of the square of the error		
(ISE),	162	566
Integral of time multiplied by the square		
error (ITSE),	164	
Integrating by parts,	236	
Integrating factor,	663	667
Interacting populations,	38	
Interconnected information,	517	522
Interconnection matrix,	533	
Interconnections,	524	
Sector condition,	524	
Symmetry requirement,	524	
Intermediate point,	273	
Principle of optimality (PO),	273	
Internal energy,	15	
Interpolation matrix,	495	
Interpolation using RBF networks,	494	496

Index Term	Links				
Interrelations,	1				
Intersection,	395				
Interval arithmetic,	396				
Interval of confidence,	396	400			
Addition,	397				
Addition of fuzzy numbers,	399				
Fuzzy number,	397				
Level of presumption,	398				
Maximum,	398				
Minimum,	398				
Multiplication,	397				
Presumption level,	399				
Inverse dynamics identification,	483				
Inverse function theorem,	374	390			
Inverse transform,	666				
Inverse transformation,	373	380			
Inverted pendulum,	35	47	147	359	467
Design model of,	467				
Membership functions,	467				
Truth model,	467				
Iron end-caps,	30				
ISE criterion,	162				
<i>See also</i> Integral of the square of the error (ISE)					
Isocline,	53				
Asymptotes,	54				
Iterated function system (IFS),	596				
Iterated parameter,	590				
Iteration process,	678				
Iterative algorithm,	220				
Iterative learning,	546				
ITSE criterion,	164				
<i>See also</i> Integral of time multiplied by the square error (ITSE)					
Jacobian matrices,	75	369			

Interrelations,

Index Term	Links				
Jacobian matrix,	189	192	374	375	382
	384	386	390	672	
Julia set,	595				
Douady's rabbit,	596				
Filled,	595				
Kalman's decomposition,	141				
Karush-Kuhn-Tucker condition,	341	352			
k -cycle,	599				
Kernel. <i>See also</i> Null space,	655				
Kinetic energy,	9	15			
Kirchhoff's current law,	509				
Kirchhoff's voltage law,	34	100	135	258	288
	310	368			
Knowledge base,	410				
Kohonen learning rule,	506				
Kronecker product,	131	308	637	639	
Kučera's theorem,	250				
L'Hôpital's rule,	158	513			
Lagrange equations of motion,	11–17	46	239		
Lagrange multiplier,	425				
Lagrange multipliers method,	269				
Solving optimization problems,	269				
Lagrange stable,	185	186			
Definition of	186				
Lagrange's multiplier theorem,	309	310			
Optimal control sequence,	309	310			
Lagrangian. <i>See</i> Lagrangian function					
Lagrangian function,	11	16	46	622	
Lane-change maneuver,	27				
Laplace expansion,	631				
Laplace transform,	109	161	666	667	
Large control signals,	226				
Largest Lyapunov exponent,	604				
Largest singular value,	633				
LaSalle's invariance principle,	213	215			

Index Term	Links				
Continuous-time systems,	216				
Difference equation,	213				
Lateral force,	24				
Lateral position,	564				
Lateral velocity,	564				
Leading principal minors,	169	171	634		
First-order,	169				
Second-order,	169				
Learning	546				
Learning capability,	533				
Learning phase,	493				
Identifier,	493				
Least significant bit (LSB),	514				
Least squares,	575	656			
Least upper bound,	640				
<i>See also</i> Supremum					
Least upper bound axiom,	640				
Left inverse,	333				
Leibniz formula,	370	372	376	384	385
	389				
Lemma,	604				
Definition of,	624				
Level curve,	611				
Level of presumption,	398	400			
Liapunoff. <i>See</i> Lyapunov					
Liapunov. <i>See</i> Lyapunov					
Lie bracket,	369				
Derivative of a vector field with					
respect to a vector field,	369				
Lie derivative notation,	383				
Lie derivatives,	369				
Derivative of a function with respect to					
a vector field,	370				
Derivative of a vector field with respect to					
a vector field,	369				

Index Term	Links			
Derivative of Dh with respect to the vector				
field,	370			
Leibniz formula,	370			
Limit cycle,	58	77	84	87
Amplitude,	87			
Stability of nonlinear system,	84			
Limit of sequence,	642	643		
Limiter,	92			
Describing function of,	92			
Linear approximation,	72			
Linear combination,	654			
Linear controller design algorithm,	466			
Linear equations,	657			
Least squares solution,	657			
Linear function,	651			
Linear map,	651			
Norm of,	651			
Linear operators,	651			
Linear quadratic regulator (LQR),	227	244	270	290
Discrete systems,	270			
First-order necessary condition				
for a relative minimum,	272			
HJB equation,	290			
Infinite time interval,	270			
Problem,	227			
Second-order sufficient condition for				
a relative minimizer,	272			
Linear quarter-car model,	223			
Linear saturation nonlinearity,	208			
Linear spring,	94			
Linear state-feedback controller,	120	244		
Linear system,	88	94		
Mathematical abstraction,	88			
Linear systems,	442			
Comparison method for,	442			

Index Term	Links		
Stability of,	154		
Linearization,	71		
Linearization method. <i>See</i> Lyapunov's indirect method			
Linearized model,	264	423	602
Linearizing,	74		
Differential equations,	74		
Linearizing state-feedback controller,	380	381	391
Canceling nonlinearities,	380		
Gains,	381		
One-link manipulator,	380		
Linearly dependent stored patterns,	541		
Linearly dependent vectors,	654		
Linearly independent columns,	114		
Linearly independent eigenvectors of Hamiltonian matrix,	253		
Linearly independent patterns,	541		
Maximal set of,	541		
Linearly nonseparable,	472		
Linearly separable,	472		
Linguistic values,	394		
Linguistic variable,	394	413	
Lipschitz condition,	651	664	
Lipschitz constant,	216	330	664
Liu, Kadirkamanathan, and Billings' identification scheme,	500		
Local linear model,	423		
Local models,	421		
Locally stable,	149		
Logistic equation,	38	579	
<i>See also</i> Verhulst equation			
Attracting fixed points,	581		
Bifurcation diagram of,	587		
Carrying capacity,	580		

Index Term	Links		
Damping term,	580		
Lyapunov exponent,	604		
Population,	579		
Stability analysis,	581		
Unrestricted growth rate,	580		
Logistic map. <i>See</i> Logistic equation			
Longitudinal force,	24		
Lorenz equations,	621		
Lorenz's weather model,	620		
Loss of information,	578		
Lotka–Volterra equations,	39		
<i>See also</i> Predator-prey relationship			
Lower bound,	641		
Low-order estimator,	129		
Low-pass filter,	83		
LQR,	313		
<i>See also</i> Linear quadratic regulator			
Ground vehicle model,	313		
LQR method,	608		
LQR problem,	227		
<i>See also</i> Linear quadratic regulator			
Luenberger observer,	128		
<i>See also</i> State estimator			
Lukasiewicz's logic,	463		
<i>See also</i> Multivalued logic			
Lumped uncertainties,	194	201	336
Lyapunov,	151		
Sometimes transliterated as Liapunov or			
Liapunoff,	151		
Lyapunov derivative,	430	505	
<i>See also</i> Time derivative of V			
Lyapunov difference,	215		
Lyapunov equation,	159		
Lyapunov exponent,	601		
Discrete-time dynamical system,	604		

Index Term	Links				
Trajectories' divergence,	600				
Trajectory separation,	582				
Lyapunov exponents,	578	582	599		
Basin of attraction,	603				
Ellipsoid semiaxes,	602				
Reference trajectory,	602				
Transition matrix,	603				
Lyapunov first difference,	611				
Level curve,	611				
Lyapunov function,	157	194	195	201	215
	216	245	270	439	
Tool to estimate RAS,	266				
Lyapunov function candidate,	196	198	211	212	218
	219	430	451	454	480
	504	520	523	525	
Instability,	219				
Sliding mode controller,	325				
Lyapunov matrix equation,	155	158	160	168	171
	189	198	222	503	
Lyapunov number,	599				
Measure of contraction or expansion,	602				
Rate of separation,	599				
Smooth map,	602				
Lyapunov theorem,	337	439			
Geometrical interpretation of,	157				
Lyapunov's first method,	188				
<i>See also</i> Lyapunov's indirect method					
Lyapunov's indirect method,	188	200			
Eigenvalues on the imaginary axis,	192				
Lyapunov's second method,	517	522			
Lyapunov's theorem,	155	157			
Lyapunov-based control,	431				
Lyapunov-based stabilizer,	431	463			
Lyapunov-like analysis,	210				
Lyapunov-like lemma,	211	212	499	505	

Index Term	Links				
MacFarlane and Potter method.					
<i>See</i> Eigenvector method					
Machine tool,	44				
Magnetic field,	90				
Mandelbrot set,	595				
One-page dictionary of all Julia sets,	596				
Quadratic complex dynamical system,	595				
Manifold,	364				
<i>See also</i> Sliding mode domain					
Map,	579				
Orbit of,	579				
Mass-spring-damper system,	90				
Matched uncertainty,	168	361	464		
Matching conditions,	466	468			
Mathematical abstraction,	88				
Linear system,	88				
Mathematical model,	393	436			
Mathematical modeling,	5				
Mating pool,	567				
MATLAB,	20				
MATLAB command,	105	139	415	493	457
	544	621			
MATLAB function,	58	468	586	590	
MATLAB peaks function,	575				
MATLAB script,	57	147	148	313	
MATLAB Symbolic Math Toolbox,	143	300	677		
Matrix,	629				
Definition of,	629				
Eigenvalue of,	676				
Row diagonal dominant matrix,	535				
Strongly row diagonal dominant matrix,	535				
Matrix calculus formulas,	674				
Matrix computations,	675				
Matrix differential inequality,	446				
Matrix exponential,	442				
Matrix inequality lemma,	437				

Index Term	Links				
Matrix inversion formula,	548	631	658		
Matrix Lyapunov equation,	170	207	212	224	
<i>See also</i> Lyapunov matrix equation					
Matrix multiplication,	630				
Matrix norm,	651				
Conditions of,	652				
Frobenius norm,	653				
Matrix square root,	362				
Matrix-vector product,	654				
Max-* composition,	405	406			
<i>See also</i> Max-star composition					
Maximizer,	555				
Of functional,	234				
Maximizing point,	555				
Maximum of a functional,	233				
Maximum of intervals of confidence,	398				
Maximum–minimum					
theorem of Weierstrass,	649	650			
Max-star,	405				
<i>See also</i> Max-* composition					
Composition,	405				
Fuzzy relation,	405	406			
Mean level,	83				
Mean value theorem,	582				
Mechanical energy,	15				
Membership,	393				
Exclusion,	393				
Inclusion,	393				
Partial,	393				
Membership function,	395	400	565	569	571
MATLAB function of,	415				
Pi membership function,	415				
Triangular membership function,	411	417			
Method of dynamic programming,	273				
Method of isoclines,	53				

Index Term	Links	
m-file,	458	
Adaptation rates,	458	
Minimization,	226	
Minimizer of a functional,	234	
Minimum of a functional,	233	
Minimum of intervals of confidence,	398	
Minimum operator,	405	
Minimum principle,	310	
Optimal charging voltage,	310	
Minimum principle of Pontryagin,	292	
Minimum time regulator problem,	280	
Admissible controllers,	282	
Admissible optimal controllers,	282	
Bang-bang,	303	
Costate equation,	303	
Hamiltonian function,	281	303
Multi-input systems,	303	
Optimal transfer,	280	
Piecewise continuous control,	280	
Pontryagin's minimum principle,	303	
System controllable by the i th input,	304	
Model,	8	
Biological neural networks,	508	
Biological neuron,	470	
Hopfield neural network,	509	513
Satellite,	50	
Modeling,	5	
Bouncing ball on a vibrating plate,	613	
Mathematical,	5	8
Model-reference,	197	
Signal generator,	197	
Tracking controller,	197	
Monotone sequence,	642	
Moore–Penrose inverse,	535	536
<i>See also</i> Pseudo-inverse		

Index Term	Links	
Most significant bit (MSB),	514	
Motor shaft,	257	
Clockwise rotation,	257	
Counterclockwise rotation,	257	
Motor-torque constant,	257	
Multi-input systems,	114	
Controller form for,	114	
Multiplication of fuzzy numbers,	401	
Multiplication of intervals of confidence,	397	401
Multiplicative rate,	582	
Trajectory separation,	582	
Multiplicity of eigenvalue,	676	
Multistep algorithm,	71	
Adams–Bashforth,	71	
Adams–Moulton,	71	
Gear,	71	
Multivalued logic,	463	
<i>See also</i> Lukasiewicz's logic		
Multivariable systems,	503	
Identification,	503	
Mutation operator,	555	568
Mutation probability,	557	
<i>See also</i> Mutation rate		
Mutation rate,	557	
<i>See also</i> Mutation probability		
Mutually orthogonal,	77	78
n -ary fuzzy relation,	404	
N -bit A/D converter,	515	
Necessary and sufficient condition for		
observability,	107	
Negative definite,	633	
Negative definite function,	154	
Negative semidefinite,	633	
Negative semidefinite function,	154	
Negative unity feedback,	217	

Index Term	Links			
Neighborhood,	185	191	649	
Of the zero state,	202			
Nervous pulse,	470			
Neural,	549			
Fuzzy identifier,	489			
Fuzzy on-line identification,	551			
Neural identifier,	497	498		
Identification architecture,	497			
Identification network,	497			
Neural network model in the new coordinates,	527			
Neural-network-based fuzzy logic control,	489			
Neural networks,	469			
Neurocomputing,	469			
Neuron,	508			
Axon,	469			
Biological,	469			
Dendrites,	469	470		
Model,	470			
Nervous pulse,	470			
Neurotransmitters,	469			
Receptor site,	469			
Synapse,	469			
Neurons,	469			
Neurotransmitters,	469			
Newton's second law,	9	34	257	368
Noisy patterns,	532			
Nominal system,	194			
Nonautonomous,	150			
Nonconservative forces,	15			
Noncontrollable part,	105			
Nondecreasing function,	179			
Nondecreasing sequence,	642			
Non-Euclidean objects,	593			

Index Term	Links	
Fractals,	593	
Nonideal sliding mode,	329	
Lipschitz constant,	330	
Utkin's theorem,	330	
Nonincreasing function,	181	
Nonincreasing sequence,	642	
Nonlinear amplifier,	509	
Nonlinear dynamical system,	367	
Nonlinear feedback systems,	83	
Nonlinear interconnected systems,	508	
Nonlinear projector,	468	
MATLAB's function of,	468	
Nonlinear system,	94	
Asymptotic state estimator,	387	
Combined controller-estimator		
compensator,	388	392
Controllability matrix,	370	
Observability matrix,	385	
Observer form,	382	
Stability of,	84	176
Nonlinear system controllability matrix,	370	373
Nonlinear system observability matrix,	385	
Nonlinear system observer form,	382	392
State variable transformation,	382	
Lie derivative notation,	383	
Nonlinear system state transformation,	371	
Diffeomorphism,	371	
Nonnegative real numbers,	397	
Nonobservable part,	111	
Nonreachable,	106	
Nonsymmetric interconnection structure,	532	
Norm,	78	
Equivalence,	432	
Of linear map,	651	

Index Term	Links			
Of uniform convergence,	656			
Normal fuzzy set,	399			
Normal to hyperplane,	475			
Normalized plot of describing function,	86			
Normed linear space,	678			
Normed space,	655	656		
Norms,	652			
Compatible,	652			
Notation,	624			
For sets,	624			
n-space,	626			
Null solution,	520	521		
Lyapunov function candidate,	523	525		
Uniformly exponentially stable,	523			
Null space,	655			
<i>See also</i> Kernel Numerical integration,	63			
Rectangular rule of,	63			
Simpson's formula,	67–69			
Trapezoidal rule of,	64			
Numerical methods,	53			
Numerical properties,	125			
Numerically stable algorithms,	125			
Nyquist stability criterion,	82			
Objective function,	221			
Observability,	107			
Index,	118			
Indices,	118			
Matrix,	110			
Matrix of nonlinear system,	385			
Observable,	107			
Pair,	117	119		
Part,	111			
Observer,	120	127		
Observer companion form,	129	146		
Observer form,	117	118	129	382

Index Term	Links			
State variable transformation,	382			
Odd function,	85			
Rectangular wave,	85			
Offspring,	555			
Offspring chromosomes,	557	561		
OGY method.				
<i>See</i> Ott, Grebogi, and Yorke				
(OGY) method				
Ohm's law,	288	296		
Armature circuit,	288	296		
One-link manipulator,	147	256	387	
Asymptotic state estimator,	387			
Combined controller-estimator				
compensator,	388			
Controller form,	374			
Design model,	389			
Simulation model,	389			
Sliding surface,	389			
State estimator dynamics,	387			
One-link robot,	258			
State-space model of,	258			
One-link robot manipulator,	256	368	374	464
Fuzzy sliding mode controller,	464			
Motor-torque constant,	257			
Nonlinear model control,	261			
Orthogonal projector,	464			
One-point crossover,	560			
On-line identification,	492	575		
Open left-half plane,	190	197		
Open right-half plane,	190			
Open set,	390	649	671	
Open subset,	670			
Open unit disc,	675			
Open unit disk,	166			
Open-loop,	3			

Index Term	Links					
Control system,	3					
Estimator,	127					
Transfer function,	217					
Operators over Banach spaces,	231					
Optimal closed-loop system,	247					
Optimal control with constraints on inputs,	292					
Optimal controller,	3	225	244	247	260	
Optimal gain,	256					
Optimal linear state feedback controller,	246					
Optimal policy,	273					
Optimal saturating controllers,	266					
Optimal systems with prescribed poles,	251					
Optimal trajectory,	235	273				
Orbit,	579					
Calculating using least squares,	656					
Of celestial bodies,	656					
Order of schema,	559					
Order reduction,	324					
Ordered pairs,	395					
Ordinary differential equation,	3	48	662			
Orthogonal complement,	350	655				
Orthogonal matrix,	631					
Orthogonal projection,	346	350				
Orthogonal projector,	355	464	548			
Orthonormal vectors,	631					
Ott, Grebogi, and Yorke (OGY) method,	607					
Outer product method,	550					
Output feedback,	145	328				
Output layer,	486					
Output of system,	6					
Output space,	6					
Outputs,	1	4				
p -norms,	363	627				
<i>See also</i> Hölder norms						
Parameterization of curve,	668					

Index Term	Links				
Parent chromosomes,	554	561			
Partial derivatives,	58				
Partial membership,	394				
Partial sums,	647				
Path of flux,	32				
Pattern associator,	529				
PBH rank test,	333				
<i>See also</i> Popov–Belevitch–Hautus test					
PD controller,	217				
<i>See also</i> Proportional-plus-derivative (PD)					
Penalty function,	342	343			
Penalty function method,	342	343			
Pendulum,	17				
Double,	622				
Inverted,	35				
Simple,	17				
Perceptron learning algorithm. <i>See</i> TLU learning algorithm					
Perceptron,	470	471	485		
Training pairs,	471				
Performance index (PI),	3	225	226	227	251
	252	566			
Period,	81				
Period-2 cycle,	586				
Period-four solution,	585				
Periodic cycles,	587				
Periodic function,	81				
Per-iteration changes,	602				
Permanent magnet stepper motor,	28	29	91		
Controllability matrix,	377				
Controller form,	376	380			
Electrical organization,	29				
Input voltage switchings,	31				

Index Term	Links				
Mathematical model,	31				
Step length,	31				
Tooth pitch,	31				
Persistently exciting inputs,	463				
Persistently exciting signals,	463				
Phase plane,	48	528			
Phase portrait,	48	57	89	259	260
	295	356			
MATLAB's script,	57				
Phenotype,	554				
Phenotype space,	554				
Transformation into genotype space,	556				
Pi membership function,	415				
Piecewise constant optimal control law,	295				
Piecewise continuous control,	280				
Piecewise continuous functions,	79	80			
Plant,	3				
Input,	3				
Model,	447				
Output,	3				
PM stepper motor					
See Permanent magnet					
stepper motor					
p -norm,	431	653			
See also Höldernorm,	627				
Polar plot,	84				
Pole placement,	11				
Ackermann's formula for,	121				
Algorithm for multi-input systems,	122				
Problem,	120				
Problem solvability,	126				
Poles,	120				
Complex conjugate pairs,	120				
Polyhedron,	341				
Polynomial,	171				

Index Term	Links			
Polytope,	341			
Polytopic system,	422	427		
Pontryagin's minimum principle,	292			
Costate equation,	293			
End-point condition,	293			
Generalization of the Euler–Lagrange equations,	292			
Hamiltonian function,	292			
HJB equation,	292			
Popov–Belevitch–Hautus (PBH) eigenvector test,	105			
Population,	38	579		
Dynamics,	38			
Of chromosomes,	555	556		
Positive definite,	179	633	634	
Decrescent function,	185			
Defintion of,	154	180		
Function,	154	157	180	
Quadratic function,	429			
Symmetric matrix,	189			
Positive limit point,	213	216		
Positive limit set,	213	216		
Positive linear combination (PLC),	477			
Positive point of accumulation				
<i>See</i> Positive limit point				
Positive real,	337			
Transfer function matrix,	337			
Positive semidefinite,	154	180	633	634
Definition of,	154	180		
Symmetric square root,	635			
Positively invariant,	214			
Positively invariant set,	216			
Positively invariant subset,	213			
Potential energy,	10			
Potential function,	373	374	671	

Index Term	Links				
Predator,	39				
Predator–prey relationship,	39				
Predictor,	65				
Equations of,	66				
Predictor–corrector method,	64				
Corrector,	65				
Predictor,	65				
Prefilter,	481				
Premise part of fuzzy rule,	571				
Preprocessing training pairs,	475				
Prescribed poles,	251				
Presumption level,	399				
Prey,	39				
Principal minors,	634	446			
Principle of conservation of energy,	16				
Principle of conservation of momentum,	613				
Principle of induction,	644				
Principle of optimality (PO),	273				
Application to routing network,	274	275			
Backward pass,	276	277			
Discrete linear quadratic regulator problem,	277				
Forward pass,	276	277			
Johann Bernoulli’s statement of,	274				
Richard E. Bellman’s statement of,	273				
Product of matrices,	630				
Product space,	403				
Cartesian,	407				
Projection,	627				
Proof,	625				
By contradiction	58	214	216	245	270
	477	650	665		
Direct,	626				
Proof by contraposition,	106	235	677		
Proportional-plus-derivative (PD)					

Index Term	Links				
controller,	217				
Proposition,	624				
Definition of,	624				
Pseudo-inverse,	535	536	537	538	547
	548				
<i>See also</i> Moore–Penrose inverse					
Pythagoras theorem,	627	628			
Quadratic form,	154	633			
Negative definite,	633				
Negative semidefinite,	633				
Positive definite,	633	634			
Positive semidefinite,	633	634			
Quadratic indices,	159				
Quadratic performance indices,	159				
Quadratic programming problem,	356				
Quarter-car model,	43	222	223		
Radial-basis function (RBF) networks,	494	495			
Radially unbounded function,	185	186	446		
Random variable,	572				
Random-like behavior,	578				
Range of function,	648				
Range,	654				
<i>See also</i> Image RAS,	266				
<i>See also</i> Region of asymptotic stability (RAS)					
RAS estimation,	266				
Optimization problem,	266				
Single-input linear system,	267				
Rate of convergence,	128				
<i>See also</i> State estimator					
Rational function,	161				
RBF identifier,	498				
<i>See also</i> Radial-basis function (RBF) networks					
Estimation errors,	498				

Index Term	Links		
Identification error,	499		
Lyapunov function candidate,	498		
Lyapunov-like lemma,	499		
State estimation error,	499		
RBF network,	494	500	
<i>See also</i> Radial-basis function (RBF) networks			
Growing,	500		
Identification,	497	503	
Interpolation,	494	496	
Interpolation matrix,	495		
State estimation error,	501		
RBF neural identifier,	504		
<i>See also</i> Radial-basis function (RBF) networks			
Identification error,	504		
Lyapunov derivative,	505		
Lyapunov function candidate,	504		
Lyapunov-like lemma,	505		
Multivariable systems,	503	504	
State estimation error,	505		
RBF neural identifier (<i>continued</i>)			
Weight estimation error,	504		
Weight matrix update law,	505		
Reachability,	94	96	97
Reachable,	106		
Real numbers,	397		
Nonnegative,	397		
Real-valued functions,	656	678	
Rear axle,	24		
Rear tire slip angle,	27		
Reasoning,	625		
Recall phase,	529		
Receptor sites,	469		
Rectangular wave,	85		

Index Term	Links		
Recursive least squares (RLS),	659		
Data pairs,	659		
Reduced-order estimator,	130	133	146
Reduced-order gBSB-type models,	533		
Reduced-order state estimator,	138		
DC motor system,	138		
Reduced-order system,	325	362	
Poles of,	328		
Reference,	3	199	
Input,	3	199	
Signal,	458	564	
Signal generator,	564		
Trajectory,	602		
Reference system,	197	456	
Transfer function,	456		
Region of asymptotic stability (RAS),	266		
Region of attraction,	366		
<i>See also</i> Region of asymptotic stability (RAS)			
Regressor,	448		
Regulator,	4		
Reject the uncertainty,	431		
Representative point (RP),	49	53	
Reproduction,	557		
Retrieval efficiency,	533		
Reverse pass,	484		
Riccati differential equation,	306		
Two point boundary-value problem,	306		
Riccati equation,	244	255	260
Algebraic,	244		
Differential,	306		
Road disturbances,	43	224	
Robot arm,	256		
Clockwise rotation,	256		
Counterclockwise rotation,	256		

Index Term	Links			
Robot manipulator,	18			
<i>See also</i> One-link manipulator				
Robust controllers,	194			
Discontinuous,	194			
Boundary layer controller,	203			
Robust linear controllers,	168			
Rössler equations,	621			
Rotational motion,	23			
Rotor,	29			
Roulette-wheel method,	557	567		
Routh necessary and sufficient condition,	173			
<i>See also</i> Routh test				
Routh stability criterion,	175			
Routh test,	217			
Routh–Hurwitz criterion,	174	218		
Row diagonal dominant matrix,	535			
Rule,	393			
Rules,	412			
IF-THEN rules,	407	411	440	571
Matrix,	412			
Run,	558			
Set of generations,	558			
Runge method,	67			
Runge–Kutta method,	70	89		
SAE,	23	563		
<i>See also</i> Society of Automotive Engineers (SAE)				
SAE convention,	23			
x coordinate,	23			
y coordinate,	23			
z coordinate,	23			
Sampling interval,	2			
Satellite model,	50			
Saturating linear optimal control,	267	268		
Algebraic Riccati equation,	268			

Index Term	Links
Saturation level,	38
Scalar product,	627
<i>See also</i> Inner product	
Scalar product of functions,	77
Scaling,	594
Schema,	559
Average fitness,	559
Defining length,	559
Don't cares,	559
Instance of,	560
One-point crossover effects on,	560
Order of,	559
Survives crossover,	561
Survives mutation,	561
Template,	559
Schema theorem,	561
Second-order sufficiency condition,	246
Second-order sufficient condition for	
a relative minimizer,	272
Sector condition,	524
Selection method,	568
Selection operator,	554
Reproduction,	557
Roulette-wheel method,	557
Self-organization process,	506
Self-organizing network,	506
Clustering training patterns,	506
Similarity matching,	506
Teacher,	506
Unsupervised mode,	506
Winning neuron,	506 507
Self-similarity,	594
Semigroup property,	6
Sensitive dependence upon initial	
conditions,	599

Index Term	Links		
Sentence,	625		
Defintion of,	625		
Separating hyperplane,	471	472	
Decision boundary,	471		
Separation,	599		
Rate of,	599		
Separation principle,	133	134	
<i>See also</i> Combined controller-estimator compensator			
Sequence,	213	214	641
Convergent,	642		
Decreasing,	642	647	
Divergent,	642		
Increasing,	642		
Limit of,	642		
Monotone,	642		
Nondecreasing,	642		
Nonincreasing,	642		
Series,	647		
Converges,	647	648	
Diverges,	648		
Partial sums of,	647		
Servomechanism,	5		
Set,	393		
Closed,	649		
Compact,	649		
Definition of,	624		
Of functions,	78		
Of generations,	558		
<i>See also</i>			
Run Of ordered pairs,	641		
<i>See also</i> Sequence Open,	649		
Sherman–Morrison–Woodbury formula,	630		
Sierpiński carpet,	622		
Sierpiński gasket,	621	622	
Sign activation function,	470		

Index Term	Links				
Sign function approximation,	483				
Signal generator,	414				
Transfer function of,	414				
Similarity matching,	506				
Similarity transformation,	103	104			
Simple EA,	561				
<i>See also</i> Simple evolutionary algorithm					
Simple evolutionary algorithm,	561	562			
Average crossover operator,	562				
Elitist strategy,	563				
Feasible set,	562				
Mutation operator,	562				
Simple pendulum,	17				
Simplified TLU learning algorithm,	475				
Simplified TLU learning problem,	475	476			
Simpson's formula,	67–69				
Simulation model,	2				
SIMULINK 2-D Look-Up Table block,	412	418			
SIMULINK block diagram,	419	457	460	461	565
Single neuron,	521				
Asymptotic stability condition,	521				
Exponential stability condition,	522				
Lyapunov function candidate,	520				
Null solution,	520	521			
Stability analysis of,	520				
Unstable,	521				
Single-input plants,	120				
Singular value decomposition (SVD),	633				
Sira–Ramfrez and Colina–Morles					
adaptation strategy,	481				
Sira–Ramírez and Colina–Morles theorem,	479				
Sliding,	320	324			
Controller objective,	325				
Equivalent control,	327				

Index Term	Links			
Equivalent system,	327			
Motion,	320			
Multi-input systems,	327			
Order reduction,	324			
Surface,	389	448		
Reduced set of differential equations,	327			
Reduced-order system,	325			
Sliding mode controller,	324			
Feedback gains,	326			
Selecting feedback gains,	324			
Sliding surface design,	324			
Sliding mode domain,	364			
Definition of,	364			
Sliding modes,	315	320	340	448
Definition of,	320			
Ideal,	331			
Nonideal,	329			
Solving optimization problems,	340			
System zeros,	328			
Sliding surface design,	333			
Linear,	327			
Nonlinear,	327			
Sliding surface design algorithm,	448			
<i>See also</i> Switching surface design algorithm				
Slip angle	25			
Front tire,	26			
Rear tire,	27			
Small control signals,	226			
s-map,	590	591		
Smooth map,	602			
Smooth state-space transformation,	391			
Society of Automotive Engineers (SAE),	23			
Coordinate system of,	23			
Solution in the sense of Filippov,	322	323		

Index Term	Links	
Construction of,	322	323
Solving state equations,	665	
Space,	670	
Tangent,	670	
Spectral radius,	632	
Spectrum,	333	
<i>See also</i> Eigenvalues of a matrix		
Speed of curve,	669	
Sphere,	649	
Spring coefficient,	223	
Spring-mass system,	49	54
Sprung mass,	42	223
Spurious states,	532	539
Square matrix,	676	
Square root of matrix,	635	
Stability,	149	520
Analysis of single neuron,	520	
In the sense of Lyapunov,	151	180
Of discrete fuzzy models,	437	
Of linear systems,	154	
Of logistic map,	581	
Of nonlinear system,	84	176
Stability of Hopfield network,	517	
Lyapunov's second method,	517	
Stabilizability of fuzzy models,	427	
Lyapunov function candidate,	430	
Stabilizable,	105	250
Stabilization algorithm,	435	
Stabilizing control,	420	
Stabilizing controller,	434	
Stable,	149	
Globally,	149	
Locally,	149	
Null solution,	521	
Stacking operation,	577	

Index Term	Links				
Stacking operator,	166	543	544	637	
Standard additive model (SAM),	407	409			
Defuzzification,	408				
Fuzzy system,	407				
State estimation,	336				
Discontinuous estimators,	336				
Error,	499	501	504		
State estimation of uncertain systems,	336				
Boundary layer estimator	338				
Discontinuous estimator,	336				
Matched uncertainty,	338				
State estimator,	126	387			
State estimator dynamics,	387				
State feedback,	145				
State feedback linear control law,	466				
State feedback to the inputs,	328				
State plane,	48				
State space,	5				
State transition matrix,	666				
Laplace transform,	666				
State variable transformation,	106	112	114	328	382
State variables,	1	48			
Change of,	111				
State-feedback control,	120	197			
State-feedback controller,	244				
State-feedback equivalent,	197				
State-feedback stabilizing controller,	369				
Statement,	525				
States,	1				
State-space format,	22				
State-space model,	317				
Double integrator,	317				
State-space realization,	162				
Stator circumference of,	32				
Stator,	29				

Index Term	Links	
Steady-state error,	200	
Steam engines,	4	
<i>See also</i> Centrifugal speed regulator		
Steer angle,	25	
<i>See also</i> Steering angle		
Steering angle,	564	
Step length,	31	63
Step-lane change maneuver,	563	564
Center average defuzzifier,	566	
Integral of the square of the error (ISE),	566	
Membership functions,	565	
Performance index (PI),	566	
Reference signal,	564	
Tracking error,	565	567
Tracking performance,	566	
Stepper motor. <i>See</i> Permanent magnet stepper motor		
Stick balancer,	35	
Coulomb friction model,	35	
Storage phase,	529	
Stored patterns,	541	
Linearly dependent,	541	
Straight line,	668	
Interpolation,	496	
Strictly increasing,	187	
Strongly row diagonal dominant matrix,	535	
Submatrix,	631	
Subsequence,	644	646
Subset,	656	
Compact,	656	
Subspace,	654	
Dimension of,	654	
Subsystem,	1	
Sufficient condition for diagonalizability,	263	
Sum of intervals of confidence,	400	
Sum of matrices,	630	

Index Term	Links	
Summing junction,	4	
Super chromosome,	563	
<i>See also</i> Elitist strategy		
Support,	407	
Of function,	407	
Of fuzzy sets,	407	
Supremum,	640	649
<i>See also</i> Least upper bound		
Surface,	669	
Surface of revolution,	228	
Example of a functional,	228	
Surface to which trajectory is confined,	345	
Surfaces,	667	
Suspension,	42	
Automotive,	42	
SVD. <i>See</i> Singular value decomposition		
Switched-capacitor technology,	340	
Switching,	201	
Switching curve,	295	
Minimum time regulator problem,	304	
Switching logic,	315	318
Double integrator,	318	
Switching surface,	202	358
Attractivity of,	320	
State velocity vectors directed toward,	320	
Switching surface design algorithm,	335	
Sylvester's inequalities,	655	677
Symbolic logic,	463	
Symbols,	624	
Glossary of,	624	
Symmetric matrix,	184	675
Symmetric positive definite,	429	
Matrix,	439	
Symmetric positive semidefinite,	226	
Symmetry affine transformation,	597	

Index Term	Links
Symmetry requirement,	524
Synapse,	469
Synaptic cleft,	469
Synaptic knobs,	469
Synaptic vesicles,	469
System,	1
Boundaries of,	1
Closed-loop,	3
Components of,	1
Controllable by the i th input,	304
Dynamical,	1
Inputs,	1
Matrix,	328
Objective,	2
Of linear equations,	657
Open-loop,	3
Outputs,	1
State,	110
States,	1
System zeros,	328
Invariance properties,	328
Output feedback,	328
State feedback to the inputs,	328
State variable transformation,	328
Transformation of the inputs,	328
Transformation of the outputs,	328
Systems controllable by the i th input,	304
Minimum time regulator problem,	304
System matrix,	328
Takagi–Sugeno (T-S) fuzzy model.	
<i>See</i> Takagi–Sugeno–Kang (TSK)	
fuzzy model	
Takagi–Sugeno–Kang (TSK) fuzzy model,	421
Design model,	422
Local models,	421

Index Term	Links			
Polytopic system,	422			
Weighted average of local models,	421			
Tangent line,	669			
Tangent plane,	74	323	346	
To the surfaces of discontinuity,	323			
Tangent space,	670			
Tangent to curve,	668			
Tangent vector,	670			
Taylor series,	59	89	422	
Linearization,	424			
Method,	59			
Taylor's theorem,	189			
<i>t</i> -conorm,	404	406		
Algebraic sum,	405			
Bounded sum,	405			
Drastic sum,	405			
Fuzzy union,	405			
Template,	559			
<i>See also</i> Schema Testing phase,	494			
Identifier,	493			
THEN,	625			
<i>See also</i> IF–THEN rules Theorem,	624			
Definition of,	624			
Theorem of Lyapunov,	503			
Theorem of Weierstrass,	649	652		
<i>See also</i> Maximum–minimum theorem				
Theorem-proof style,	624			
Three-layer neural network,	485			
Hidden layer,	486			
Output layer,	486			
Threshold,	470	471		
Threshold logic unit (TLU),	470			
Time delay,	464			
Time derivative of V ,	445	448	450	454

Index Term	Links			
Time-delay systems,	575			
On-line identification of,	575			
Time-invariant model,	665			
Tire slip angle,	25	564		
TLU,	470	471		
<i>See also</i> Perceptron,	470	471		
Threshold,	472	473	474	
Training pairs,	471	474		
Weight vector,	472	473	474	
TLU convergence theorem,	476			
TLU learning algorithm,	473	477	551	
Simplified,	475			
TLU learning problem,	475			
Normal to hyperplane,	475			
Preprocessing training pairs,	475			
Simplified,	475	476		
TLU training problem,	472			
t-map,	591			
<i>t</i> -norm,	404			
Algebraic product,	404			
Bounded product,	405			
Drastic product,	405			
Fuzzy intersection,	404			
Toilet-flushing system,	41			
Tooth pitch,	31			
Torsional spring,	45			
TPBVP,	304			
<i>See also</i> Two point boundary-value problem				
Trace operator,	505			
Tracking controller,	5	201		
Tracking error,	458	565	567	
Tracking performance,	566	568		
Training pairs,	471	474	485	507
Training patterns,	507			

Index Term	Links				
Trajectories' divergence,	600				
Trajectory separation,	582				
Transducer,	217				
Transfer function matrix,	337				
Positive real,	337				
Transfer function,	83	143	146	359	414
	456	564			
Reference signal generator,	564				
Transformation of the inputs,	328				
Transformation of the outputs,	328				
Transient response,	48				
Transition mapping,	5–7				
Modeling of,	21				
Transition matrix,	603				
Eigenvalues of,	603				
Translation,	527				
Transpose of matrix,	630				
Transpose of vector,	626				
Transversality conditions,	241	242	244		
Trapezoidal membership function,	569	571			
Traveling salesman problem (TSP),	576				
Triangle inequality,	331	627	629	660	
Triangular fuzzy numbers,	400	402			
Triangular membership function,	411	417			
Trivial subspace,	654				
Truth model,	2	8	467		
Truth table,	625				
TSK model,	440				
<i>See also</i> Takagi–Sugeno–Kang (TSK)					
fuzzy model					
Turbine supply pump resistance,	42				
Turning maneuver,	22				
Two point boundary-value problem					
(TPBVP),	304				
Control optimality condition,	305				

Index Term		Links
Hamiltonian function,	304	
Optimal controller,	305	
Optimal state-feedback controller,	306	
Riccati differential equation,	306	
Unconstrained control vector,	304	
Two-cycle solution,	585	
Two-dimensional affine transformation,	596	
Two-dimensional IFS,	598	
<i>See also</i> Iterated function systems		
Barnsley fern,	598	
Definition of,	597	
Two-part controller,	437	
Two-point crossover,	558	
Two-tank pumped storage,	42	
Types of proofs,	626	
Contradiction,	626	
Contraposition,	626	
Direct,	626	
Uncertain dynamical system,	360	431
Uncertain element,	201	
Uncertain nonlinear dynamical systems,	420	
Uncertain systems,	336	
State estimation,	336	
Uncertainties,	194	
Lumped,	194	
Uncertainty,	168	
Growth,	578	
Matched	168	
Unmatched	168	
Unconstrained optimization,	342	
Dynamical gradient system,	342	
First-order necessary condition for,	269	
Uncontrolled nominal system,	194	
Uncontrolled system,	665	
Solution of,	665	

Index Term	Links			
Uniform asymptotic stability in the large,	444			
Uniform boundedness,	205			
Uniform continuity,	650			
Uniform convergence,	656			
Norm of,	656			
Uniform exponential stability,	152			
Uniform stability,	151	182		
Uniformly ultimately bounded (UUB),	338			
Uniform ultimate boundedness,	201			
Definition of,	201			
Uniformly bounded,	201			
Uniformly bounded equilibrium,	338			
UUB,	338			
<i>See also</i> Uniformly ultimately bounded				
Uniformly convergent,	152			
Definition of,	152			
Uniformly exponentially stable null				
solution,	523			
Uniformly stable,	201			
Uniformly stable ball,	201	206	338	
Union,	395			
Unit disc,	675			
Unit step,	162	199		
Unitary matrix,	633	635		
Universe of discourse,	395	412	567	
Fixed set,	625			
Universe,	395			
<i>See also</i> Universe of discourse				
Unmatched uncertainty,	168	361	362	465
Unrestricted growth rate,	580			
Unspecified scalar				
<i>See</i> Don't care element				
Unsprung mass,	42	223		
Unstable null solution,	521			
Unsupervised mode,	506			
Upper bound,	640			

Index Term	Links			
Upper row triangular,	105			
Utkin's theorem,	330			
UUB closed-loop system,	366			
<i>See also</i> Uniformly ultimately bounded (UUB)				
Vague,	463			
<i>See also</i> Fuzzy,	463			
Vagueness,	463			
<i>See also</i> Fuzziness,	463			
Van der Pol equation,	57			
Phase portrait of,	57			
Variable force element,	223			
Variable structure sliding mode control systems,	315			
Variable structure system,	315	357		
Asymptotically stable,	315			
Independent structures,	315			
Switching logic,	315			
Variation,	229			
Variation of functional,	230	231	233	240
Vector,	626			
Components of,	626			
Coordinates of,	626			
Dimension of,	626			
Tangent,	670			
Transpose of,	626			
Vector field,	367	370	670	671
Conservative,	375			
Differentiable,	673			
Methods,	367			
Potential function,	373			
Vector norm,	226	627		
Vector spaces,	654			
Vehicle center of gravity (CG),	23			
Vehicle CG,	24			
<i>See also</i> Vehicle center				

Index Term	Links		
of gravity (CG)			
Vehicle model,	145	576	
Fuzzy logic controller (FLC) for,	412		
Vehicle suspension system,	144		
Vehicle system,	565		
SIMULINK block diagram,	565		
Velocity vector,	668		
Verhulst equation,	38		
<i>See also</i> Logistic equation			
Vertex,	534		
Vesicles,	469		
Vibrating plate,	615		
Bands of frequencies,	615		
v-map,	590		
Von Koch curve,	594		
Self-similarity dimension of,	594		
Von Koch snowflake curve,	593		
Walcott and Zak state estimator,	337		
Estimation error equation,	337		
Watt's governor,	4	76	
<i>See also</i> Centrifugal speed regulator;			
Flyball governor			
Engine-governor equations,	76		
Weierstrass theorem,	177	650	
<i>See also</i> Maximum–minimum			
theorem			
Weight estimation error,	504		
Weight matrix,	255	256	262
Update law,	505		
Weight matrix construction algorithm,	540		
Weight vector,	472		
Weighted average of local models,	441		
Windings,	29		
Winner-take-all (WTA) network,	506		

Index Term	Links		
Winner-take-all learning algorithm,	507		
Training patterns,	507		
Winning neuron,	507		
Winning neuron,	506	507	
Work,	9		
Work–energy theorem,	10		
WTA algorithm,	551		
<i>See also</i> Winner-take-all			
(WTA) network			
WTA learning algorithm,	507		
WTA network,	506		
<i>See also</i> Winner-take-all			
(WTA) network			
Descent gradient method,	506		
x coordinate,	23		
XOR problem,	487	488	489
<i>See also</i> Exclusive-OR			
(XOR) problem			
Threshold,	488		
Weight,	488		
y coordinate,	23		
Yaw velocity,	23		
z coordinate,	23		
Zadeh’s fuzzy set algebra,	463		
Zero matrix,	629		
Zero vector,	626	654	
z -map,	590	591	