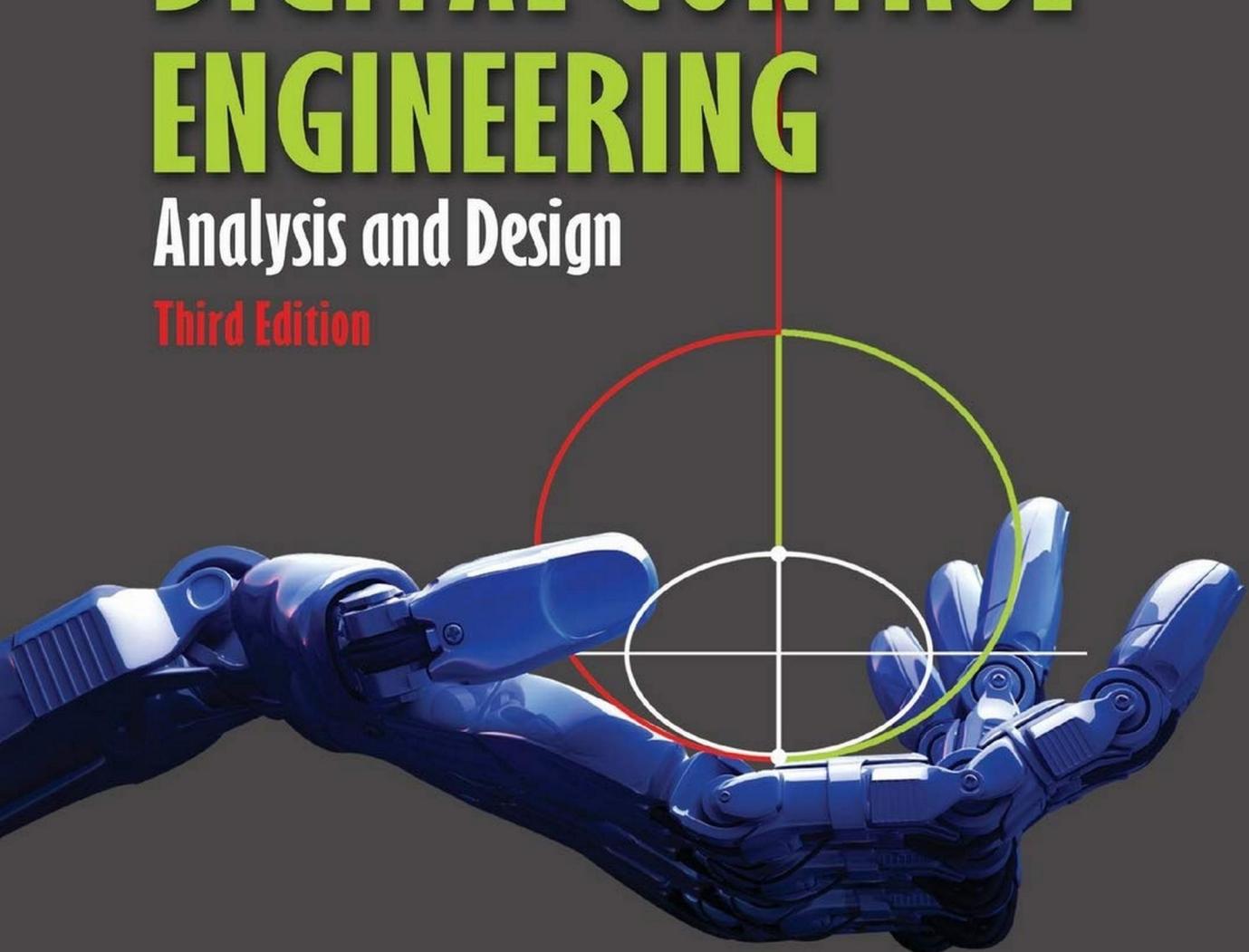


MATLAB®
and Simulink®
examples

DIGITAL CONTROL ENGINEERING

Analysis and Design

Third Edition



M. SAMI FADALI

ANTONIO VISIOLI



Digital Control Engineering

Analysis and Design

Third Edition

M. Sami Fadali
University of Nevada, Reno
Reno, NV, United States

Antonio Visioli
University of Brescia
Brescia, Italy



ACADEMIC PRESS

An imprint of Elsevier

Academic Press is an imprint of Elsevier
125 London Wall, London EC2Y 5AS, United Kingdom
525 B Street, Suite 1650, San Diego, CA 92101, United States
50 Hampshire Street, 5th Floor, Cambridge, MA 02139, United States
The Boulevard, Langford Lane, Kidlington, Oxford OX5 1GB, United Kingdom

Copyright © 2020 Elsevier Inc. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or any information storage and retrieval system, without permission in writing from the publisher. Details on how to seek permission, further information about the Publisher's permissions policies and our arrangements with organizations such as the Copyright Clearance Center and the Copyright Licensing Agency, can be found at our website: www.elsevier.com/permissions.

This book and the individual contributions contained in it are protected under copyright by the Publisher (other than as may be noted herein).

Notices

Knowledge and best practice in this field are constantly changing. As new research and experience broaden our understanding, changes in research methods, professional practices, or medical treatment may become necessary.

Practitioners and researchers must always rely on their own experience and knowledge in evaluating and using any information, methods, compounds, or experiments described herein. In using such information or methods they should be mindful of their own safety and the safety of others, including parties for whom they have a professional responsibility.

To the fullest extent of the law, neither the Publisher nor the authors, contributors, or editors, assume any liability for any injury and/or damage to persons or property as a matter of products liability, negligence or otherwise, or from any use or operation of any methods, products, instructions, or ideas contained in the material herein.

Library of Congress Cataloging-in-Publication Data

A catalog record for this book is available from the Library of Congress

British Library Cataloguing-in-Publication Data

A catalogue record for this book is available from the British Library

ISBN: 978-0-12-814433-6

For information on all Academic Press publications visit our website
at <https://www.elsevier.com/books-and-journals>

Publisher: Mara Conner

Acquisition Editor: Sonnini R. Yura

Editorial Project Manager: Naomi Robertson

Production Project Manager: Surya Narayanan Jayachandran

Cover Designer: Greg Harris

Typeset by TNQ Technologies



Working together
to grow libraries in
developing countries

www.elsevier.com • www.bookaid.org

Preface

Approach

Control systems are an integral part of everyday life in today's society. They control our appliances, our entertainment centers, our cars, and our office environments; they control our industrial processes and our transportation systems; they control our exploration of land, sea, air, and space. Almost all of these applications use digital controllers implemented with computers, microprocessors, or digital electronics. Every electrical, chemical, or mechanical engineering senior or graduate student should therefore be familiar with the basic theory of digital controllers.

This text is designed for a senior or combined senior/graduate-level course in digital controls in departments of mechanical, electrical, or chemical engineering. Although other texts are available on digital controls, most do not provide a satisfactory format for a senior/graduate-level class. Some texts have very few examples to support the theory, and some were written before the wide availability of computer-aided design (CAD) packages. Others use CAD packages in certain ways but do not fully exploit their capabilities. Most available texts are based on the assumption that students must complete several courses in systems and control theory before they can be exposed to digital control. We disagree with this assumption, and we firmly believe that students can learn digital control after a one-semester course covering the basics of analog control. As with other topics that started at the graduate level—linear algebra and Fourier analysis to name a few—the time has come for digital control to become an integral part of the undergraduate curriculum.

Features

To meet the needs of the typical senior/graduate-level course, this text includes the following features.

Numerous examples

The book includes a large number of examples. Typically, only one or two examples can be covered in the classroom because of time limitations. The student can use the remaining examples for self-study. The experience of the authors is that students need more examples

to experiment with so as to gain a better understanding of the theory. The examples are varied to bring out subtleties of the theory that students may overlook.

Extensive use of CAD packages

The book makes extensive use of CAD packages. It goes beyond the occasional reference to specific commands to the integration of these commands into the modeling, design, and analysis of digital control systems. For example, root locus design procedures given in most digital control texts are not CAD procedures and instead emphasize paper-and-pencil design. The use of CAD packages, such as MATLAB®, frees students from the drudgery of mundane calculations and allows them to ponder more subtle aspects of control system analysis and design. The availability of a simulation tool like Simulink allows the student to simulate closed-loop control systems, including aspects neglected in design such as non-linearities and disturbances.

Coverage of background material

The book itself contains review material from linear systems and classical control. Some background material is included in the appendices that could either be reviewed in class or consulted by the student as necessary. The review material, which is often neglected in digital control texts, is essential for the understanding of digital control system analysis and design. For example, the behavior of discrete-time systems in the time domain and in the frequency domain is a standard topic in linear systems texts but often receives brief coverage. Root locus design is almost identical for analog systems in the s -domain and digital systems in the z -domain. The topic is covered much more extensively in classical control texts and inadequately in digital control texts. The digital control student is expected to recall this material or rely on other sources. Often, instructors are obliged to compile their own review materials, and the continuity of the course is adversely affected.

Inclusion of advanced topics

In addition to the basic topics required for a one-semester senior/graduate class, the text includes some advanced material to make it suitable for an introductory graduate-level class or for two quarters at the senior/graduate level. We would also hope that the students in a single-semester course would acquire enough background and interest to read the additional chapters on their own. Examples of optional topics are state-space methods, which may receive brief coverage in a one-semester course, and nonlinear discrete-time systems, which may not be covered.

Standard mathematics prerequisites

The mathematics background required for understanding most of the book does not exceed what can be reasonably expected from the average electrical, chemical, or mechanical

engineering senior. This background includes three semesters of calculus, differential equations, and basic linear algebra. Some texts on digital control require more mathematical maturity and are therefore beyond the reach of the typical senior. On the other hand, the text does include optional topics for the more advanced student. The rest of the text does not require knowledge of this optional material so that it can be easily skipped if necessary.

Senior system theory prerequisites

The control and system theory background required for understanding the book does not exceed material typically covered in one semester of linear systems and one semester of control systems. Thus, students should be familiar with Laplace transforms, the frequency domain, and the root locus. They need not be familiar with the behavior of discrete-time systems in the frequency and time domain or have extensive experience with compensator design in the s -domain. For an audience with an extensive background in these topics, some topics can be skipped and the material can be covered at a faster rate.

Coverage of theory and applications

The book has two authors: the first is primarily interested in control theory and the second is primarily interested in practical applications and hardware implementation. Although some control theorists have sufficient familiarity with practical issues such as hardware implementation and industrial applications to touch on the subject in their texts, the material included is often deficient because of the rapid advances in the area and the limited knowledge that theorists have of the subject.

It became clear to the first author that to have a suitable text for his course and similar courses, he needed to find a partner to satisfactorily complete the text. He gradually collected material for the text and started looking for a qualified and interested partner. Finally, he found a co-author who shared his interest in digital control and the belief that it can be presented at a level amenable to the average undergraduate engineering student.

For many years, Dr. Antonio Visioli has been teaching an introductory and a laboratory course on automatic control, as well as a course on control systems technology. Furthermore, his research interests are in the fields of industrial regulators and robotics. Although he contributed to the material presented throughout the text, his major contribution was adding material related to the practical design and implementation of digital control systems. This material is rarely covered in control systems texts but is an essential prerequisite for applying digital control theory in practice.

The text is written to be as self-contained as possible. However, the reader is expected to have completed a semester of linear systems and classical control. Throughout the text, extensive use is made of the numerical computation and CAD package MATLAB. As with all computational tools, the enormous capabilities of MATLAB are no substitute for a sound understanding of the theory presented in the text. As an example of the inappropriate use of

supporting technology, we recall the story of the driver who followed the instructions of his GPS system and drove into the path of an oncoming train!¹ The reader must use MATLAB as a tool to support the theory without blindly accepting its computational results.

New to this edition

We made several important changes and added material to the second edition:

1. We added a new section on sensitivity.
2. We added a new section on the return difference equality of the linear quadratic regulator.
3. We added a new section on model predictive control.
4. We added a new chapter on linear matrix inequalities.
5. We added over 60 new problems including 20 for the new chapter.
6. We added explanations and revised the presentation for several sections.
7. We rewrote the section on deadbeat control.
8. We corrected many minor errors in the second edition.

Organization of text

The text begins with an *introduction to digital control* and the reasons for its popularity. It also provides a few examples of applications of digital control from the engineering literature.

Chapter 2 considers discrete-time models and their analysis using the z -transform. We review the z -transform, its properties, and its use to solve difference equations. The chapter also reviews the properties of the *frequency response* of discrete-time systems. After a brief discussion of the *sampling theorem*, we are able to provide rules of thumb for *selecting the sampling rate* for a given signal or for given system dynamics. This material is often covered in linear systems courses, and much of it can be skipped or covered quickly in a digital control course. However, the material is included because it serves as a foundation for much of the material in the text.

Chapter 3 derives simple mathematical models for *linear discrete-time systems*. We derive models for the *analog-to-digital converter* (ADC), the *digital-to-analog converter* (DAC), and an analog system with a DAC and an ADC. We include systems with time delays that are not an integer multiple of the sampling period. These transfer functions are particularly important because many applications include an analog plant with DAC and ADC.

Nevertheless, there are situations where different configurations are used. We therefore include an analysis of a variety of configurations with samplers. We also characterize the *steady-state tracking error* of discrete-time systems and define error constants for the unity

¹ The story was reported in the *Chicago Sun-Times*, on January 4, 2008. The driver, a computer consultant, escaped just in time before the train slammed into his car at 60 mph in Bedford Hills, New York.

feedback case. These error constants play an analogous role to the error constants for analog systems. Using our analysis of more complex configurations, we are able to obtain the *error due to a disturbance input*. We include an introduction to Simulink, the MATLAB toolbox that uses a graphical language to build mathematical models and simplify their simulation. We also discuss *sensitivity analysis* and its use to asses the effect of uncertainty or error in the parameters of mathematical models.

In Chapter 4, we present stability tests for input-output systems. We examine the definitions of *input-output stability* and *internal stability* and derive conditions for each. By transforming the characteristic polynomial of a discrete-time system, we are able to test it using the standard *Routh–Hurwitz criterion* for analog systems. We use the *Jury criterion*, which allows us to directly test the stability of a discrete-time system. Finally, we present the *Nyquist criterion* for the z -domain and use it to determine closed-loop stability of discrete-time systems.

Chapter 5 introduces analog s -domain design of proportional (P), proportional-integral (PI), proportional-plus-derivative (PD), and *proportional-integral-derivative (PID) control* using MATLAB. We use MATLAB as an integral part of the design process, although many steps of the design can be completed using a scientific calculator. It would seem that a chapter on analog design does not belong in a text on digital control. This is false. Analog control can be used as a first step toward obtaining a digital control. In addition, direct digital control design in the z -domain is similar in many ways to s -domain design.

Digital controller design is the topic of Chapter 6. It begins with proportional control design and then examines digital controllers based on analog design. The direct design of digital controllers is considered next. We consider *root locus design* in the z -plane for PI and PID controllers. We also consider a synthesis approach due to Ragazzini that allows us to specify the desired closed-loop transfer function. As a special case, we consider the design of *deadbeat controllers* that allow us to exactly track an input at the sampling points after a few sampling points. For completeness, we also examine *frequency response design in the w -plane*. This approach requires more experience because values of the stability margins must be significantly larger than in the more familiar analog design. As with analog design, MATLAB is an integral part of the design process for all digital control approaches.

Chapter 7 covers *state-space models* and state–space realizations. First, we discuss analog state–space equations and their solutions. We include nonlinear analog equations and their linearization to obtain linear state–space equations. We then show that the solution of the analog state equations over a sampling period yields a discrete-time state–space model. Properties of the solution of the analog state equation can thus be used to analyze the discrete-time state equation. The discrete-time state equation is a recursion for which we obtain a solution by induction. In Chapter 8, we consider important properties of state–space models: *stability*, *controllability*, and *observability*. As in Chapter 4, we consider internal stability and input-output stability, but the treatment is based on the properties of the state–space model rather than those of the transfer function. Controllability is a property that characterizes our ability to drive the system from an arbitrary initial

state to an arbitrary final state in finite time. Observability characterizes our ability to calculate the initial state of the system using its input and output measurements. Both are structural properties of the system that are independent of its stability. Next, we consider *realizations* of discrete-time systems. These are ways of implementing discrete-time systems through their state-space equations using summers and delays.

Chapter 9 covers the design of controllers for state-space models. We show that the system dynamics can be arbitrarily chosen using state feedback if the system is controllable. If the state is not available for feedback, we can design a state estimator or *observer* to estimate it from the output measurements. These are dynamic systems that mimic the system but include corrective feedback to account for errors that are inevitable in any implementation. We give two types of observers. The first is a simpler but more computationally costly full-order observer that estimates the entire state vector. The second is a reduced-order observer with the order reduced by virtue of the fact that the measurements are available and need not be estimated. Either observer can be used to provide an estimate of the state for feedback control, or for other purposes. Control schemes based on state estimates are said to use *observer state feedback*.

Chapter 10 deals with the *optimal control* of digital control systems. We consider the problem of unconstrained optimization, followed by constrained optimization, and then generalize to dynamic optimization as constrained by the system dynamics. We are particularly interested in the linear quadratic regulator where optimization results are easy to interpret and the prerequisite mathematics background is minimal. We consider both the finite time and steady-state regulator and discuss conditions for the existence of the steady-state solution. We discuss the return difference equality of the linear quadratic regulator and how it can be used to assess its robustness. We provide an introduction to *model predictive control (MPC)*, a design strategy that optimizes the control law by predicting the system's behavior using its mathematical model.

MPC is a discrete-time optimal control strategy that exploits knowledge of a system model to calculate the control law to minimize a given cost function.

The first 10 chapters are mostly restricted to linear discrete-time systems.

Chapter 11 examines the far more complex behavior of *nonlinear discrete-time systems*. It begins with equilibrium points and their stability. It shows how equivalent discrete-time models can be easily obtained for some forms of nonlinear analog systems using *global* or *extended linearization*. For the classes of nonlinear systems for which extended linearization is straightforward, linear design methodologies can yield nonlinear controllers. It provides stability theorems and instability theorems using *Lyapunov stability theory*. The theory gives sufficient conditions for nonlinear systems, and failure of either the stability or instability tests is inconclusive. For linear systems, Lyapunov stability yields necessary and sufficient conditions. Lyapunov stability theory also allows us to design controllers by selecting a control that yields a closed-loop system that meets the Lyapunov stability conditions. The chapter also discusses *input-output stability* of nonlinear systems and the small gain theorem.

Chapter 12 deals with practical issues that must be addressed for the successful implementation of digital controllers. In particular, the hardware and software requirements for the correct implementation of a digital control system are analyzed. We discuss the choice of the sampling frequency in the presence of *antialiasing filters* and the effects of *quantization*, rounding, and truncation errors. We also discuss *bumpless switching* from automatic to manual control, avoiding discontinuities in the control input. Our discussion naturally leads to approaches for the effective implementation of a PID controller. Finally, we consider nonuniform sampling, where the sampling frequency is changed during control operation, and *multiprate sampling*, where samples of the process outputs are available at a slower rate than the controller sampling rate.

Chapter 13 deals with *linear matrix inequalities* (LMIs) and the MATLAB LMI commands available as part of the Robust Control Toolbox. The material presented is introductory in nature and does not include the numerical algorithms underlying computer solution of LMIs. It explains how some of the important problems covered in the book can be formulated as LMIs and provides some results that can be used to simplify the LMIs. It shows how the MATLAB toolbox can be used to obtain solutions for LMIs.

Supporting material

The following resources are available to instructors adopting this text for use in their courses. Please visit www.textbooks.elsevier.com to register for access to these materials:

Instructor Solutions Manual. Fully typeset solutions to the end-of-chapter problems in the text.

PowerPoint Images. Electronic images of the figures and tables from the book, useful for creating lectures.

PowerPoint Presentations.

Acknowledgments

We would like to thank the anonymous reviewers who provided excellent suggestions for improving the text. We would also like to thank Dr. Qing-Chang Zhong of the Illinois Institute of Technology, who suggested the cooperation between the two authors that led to the completion of this text. We would also like to thank Joseph P. Hayton, Michael Joyce, Lisa Lamenzo, Naomi Robertson, and the Elsevier staff for their help in producing the text. Finally, we would like to thank our wives Betsy and Silvia for their support and love throughout the months of writing this book.

Introduction to digital control

Objectives

After completing this chapter, the reader will be able to do the following:

1. Explain the reasons for the popularity of digital control systems.
2. Draw a block diagram for digital control of a given analog control system.
3. Explain the structure and components of a typical digital control system.

In most modern engineering systems, it is necessary to control the evolution with time of one or more of the system variables. Controllers are required to ensure satisfactory transient and steady-state behavior for these engineering systems. To guarantee satisfactory performance in the presence of disturbances and model uncertainty, most controllers in use today employ some form of negative feedback. A sensor is needed to measure the controlled variable and compare its behavior to a reference signal. Control action is based on an error signal defined as the difference between the reference and the actual values.

The controller that manipulates the error signal to determine the desired control action has classically been an analog system, which includes electrical, fluid, pneumatic, or mechanical components. These systems all have *analog* inputs and outputs (i.e., their input and output signals are defined over a continuous time interval and have values that are defined over a continuous range of amplitudes). In the past few decades, analog controllers have often been replaced by *digital* controllers whose inputs and outputs are defined at discrete time instances. The digital controllers are in the form of digital circuits, digital computers, or microprocessors.

Intuitively, one would think that controllers that continuously monitor the output of a system would be superior to those that base their control on sampled values of the output. It would seem that control variables (controller outputs) that change continuously would achieve better control than those that change periodically. This is in fact true! Had all other factors been identical for digital and analog control, analog control would be superior to digital control. What, then, is the reason behind the change from analog to digital that has occurred over the past few decades?

Chapter Outline

1.1 Why digital control?	2
1.2 The structure of a digital control system	3
1.3 Examples of digital control systems	3
1.3.1 Closed-loop drug delivery system	3
1.3.2 Computer control of an aircraft turbojet engine	4
1.3.3 Control of a robotic manipulator	4
Resources	6
Problems	6

1.1 Why digital control?

Digital control offers distinct advantages over analog control that explains its popularity. Here are some of its many advantages:

Accuracy: Digital signals are represented in terms of zeros and ones with typically 12 bits or more to represent a single number. This involves a very small error as compared to analog signals, where noise and power supply drift are always present.

Implementation errors: Digital processing of control signals involves addition and multiplication by stored numerical values. The errors that result from digital representation and arithmetic are negligible. By contrast, the processing of analog signals is performed using components such as resistors and capacitors with actual values that vary significantly from the nominal design values.

Flexibility: An analog controller is difficult to modify or redesign once implemented in hardware. A digital controller is implemented in firmware or software and its modification is possible without a complete replacement of the original controller. Furthermore, the structure of the digital controller need not follow one of the simple forms that are typically used in analog control. More complex controller structures involve a few extra arithmetic operations and are easily realizable.

Speed: The speed of computer hardware has increased exponentially since the 1980s. This increase in processing speed has made it possible to sample and process control signals at very high speeds. Because the interval between samples, the sampling period, can be made very small, digital controllers achieve performance that is essentially the same as that based on continuous monitoring of the controlled variable.

Cost: Although the prices of most goods and services have steadily increased, the cost of digital circuitry continues to decrease. Advances in semiconductor technology have made it possible to manufacture better, faster, and more reliable integrated circuits and to offer

them to the consumer at a lower price. This has made the use of digital controllers more economical even for small, low-cost applications.

1.2 The structure of a digital control system

To control a physical system or process using a digital controller, the controller must receive measurements from the system, process them, and then send control signals to the actuator that effects the control action. In almost all applications, both the plant and the actuator are analog systems. This is a situation where the controller and the controlled do not “speak the same language,” and some form of translation is required. The translation from controller language (digital) to physical process language (analog) is performed by a digital-to-analog converter or DAC. The translation from process language to digital controller language is performed by an analog-to-digital converter or ADC. A sensor is needed to monitor the controlled variable for feedback control. The combination of the elements discussed here in a control loop is shown in Fig. 1.1. Variations on this control configuration are possible. For example, the system could have several reference inputs and controlled variables, each with a loop similar to that of Fig. 1.1. The system could also include an inner loop with digital or analog control.

1.3 Examples of digital control systems

In this section, we briefly discuss examples of control systems where digital implementation is now the norm. There are many other examples of industrial processes that are digitally controlled, and the reader is encouraged to seek other examples from the literature.

1.3.1 Closed-loop drug delivery system

Several chronic diseases require the regulation of the patient’s blood levels of a specific drug or hormone. For example, some diseases involve the failure of the body’s natural closed-loop control of blood levels of nutrients. Most prominent among these is the

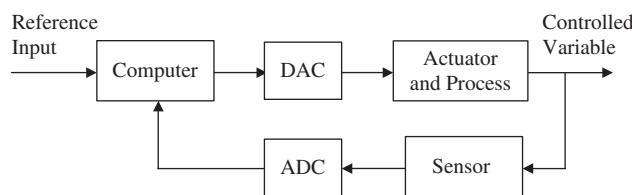


Figure 1.1

Configuration of a digital control system. *ADC*, analog-to-digital converter; *DAC*, digital-to-analog converter.

disease diabetes, where the production of the hormone insulin that controls blood glucose levels is impaired.

To design a closed-loop drug delivery system, a sensor is utilized to measure the levels of the regulated drug or nutrient in the blood. This measurement is converted to digital form and fed to the control computer, which drives a pump that injects the drug into the patient's blood. A block diagram of the drug delivery system is shown in Fig. 1.2. See Carson and Deutsch (1992) for a more detailed example of a drug delivery system.

1.3.2 Computer control of an aircraft turbojet engine

To achieve the high performance required for today's aircraft, turbojet engines employ sophisticated computer control strategies. A simplified block diagram for turbojet computer control is shown in Fig. 1.3. The control requires feedback of the engine state (speed, temperature, and pressure), measurements of the aircraft state (speed and direction), and pilot command.

1.3.3 Control of a robotic manipulator

Robotic manipulators are capable of performing repetitive tasks at speeds and accuracies that far exceed those of human operators. They are now widely used in manufacturing

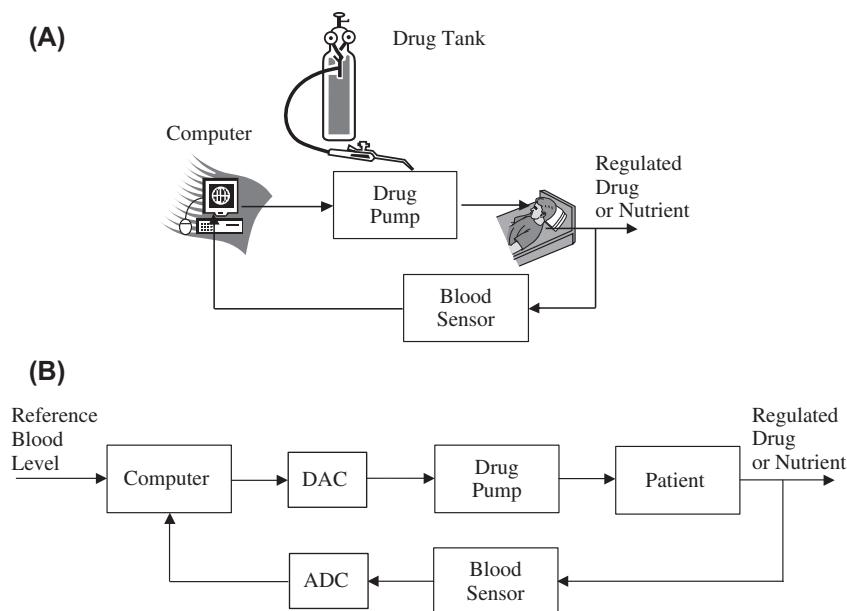


Figure 1.2

Drug delivery digital control system. (A) Schematic of a drug delivery system. (B) Block diagram of a drug delivery system. *ADC*, analog-to-digital converter; *DAC*, digital-to-analog converter.

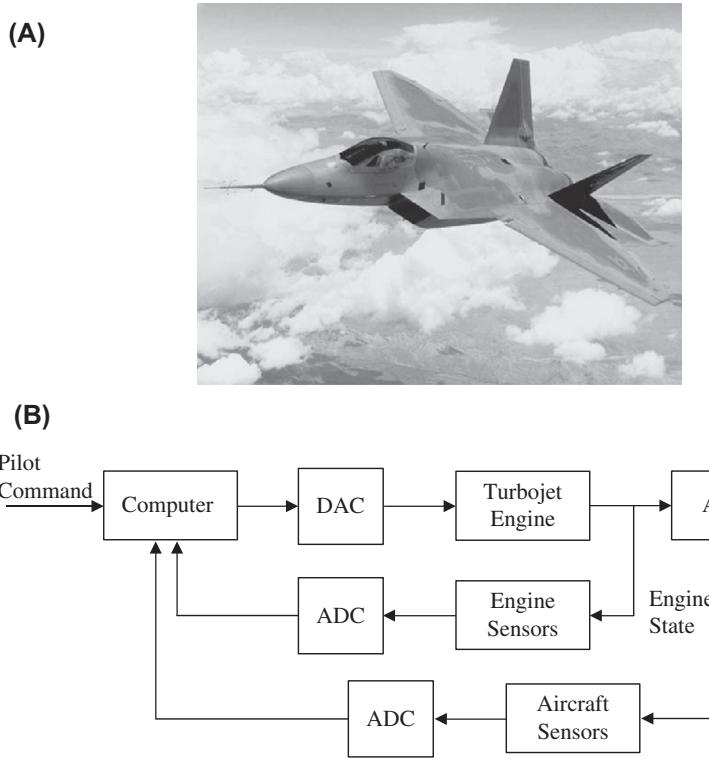
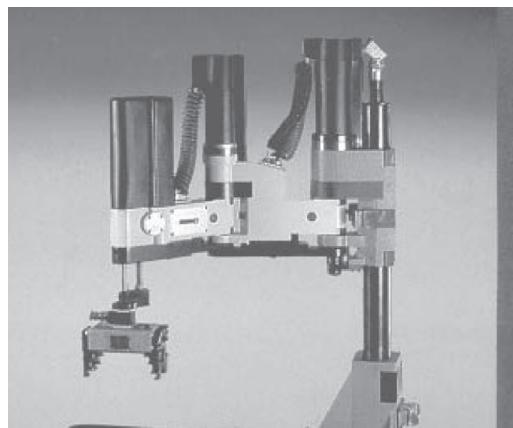


Figure 1.3

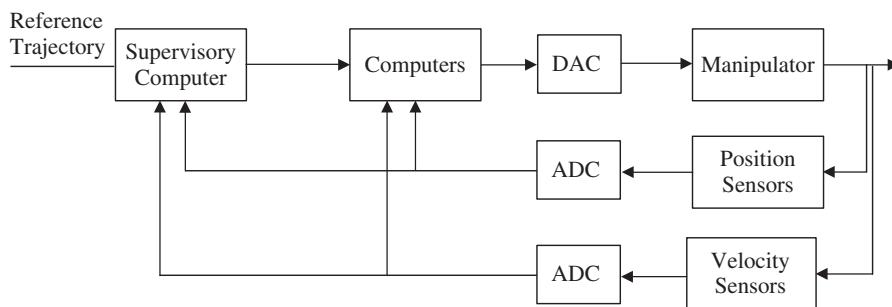
Turbojet engine control system. (A) F-22 military fighter aircraft. (B) Block diagram of an engine control system. *ADC*, analog-to-digital converter; *DAC*, digital-to-analog converter.

processes such as spot welding and painting. To perform their tasks accurately and reliably, manipulator hand (or end-effector) positions and velocities are controlled digitally. Each motion or degree of freedom (DOF) of the manipulator is positioned using a separate position control system. All the motions are coordinated by a supervisory computer to achieve the desired speed and positioning of the end effector. The computer also provides an interface between the robot and the operator that allows programming the lower-level controllers and directing their actions. The control algorithms are downloaded from the supervisory computer to the control computers, which are typically specialized microprocessors known as digital signal processing (DSP) chips. The DSP chips execute the control algorithms and provide closed-loop control for the manipulator. A simple robotic manipulator is shown in Fig. 1.4A, and a block diagram of its digital control system is shown in Fig. 1.4B. For simplicity, only one motion control loop is shown in Fig. 1.4, but there are actually n loops for an n -DOF manipulator.

(A)



(B)

**Figure 1.4**

Robotic manipulator control system. (A) 3 Degrees of freedom robotic manipulator. (B) Block diagram of a manipulator control system. *ADC*, analog-to-digital converter; *DAC*, digital-to-analog converter.

Resources

- Carson, E.R., Deutsch, T., 1992. A spectrum of approaches for controlling diabetes. *Control Syst. Mag.* 12 (6), 25–31.
- Chen, C.T., 1993. *Analog and Digital Control System Design*. SaundereHBJ.
- Koivo, A.J., 1989. *Fundamentals for Control of Robotic Manipulators*. Wiley.
- Shaffer, P.L., 1990. A multiprocessor implementation of a real-time control of turbojet engine. *Control Syst. Mag.* 10 (4), 38–42.

Problems

- 1.1 A fluid level control system includes a tank, a level sensor, a fluid source, and an actuator to control fluid inflow. Consult any classical control text¹ to obtain a block

¹ See, for example, Van deVegte, J., 1994. *Feedback Control Systems*. Prentice Hall.

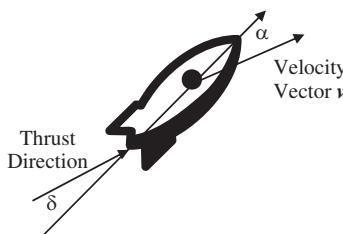


Figure P1.5
Missile angle-of-attack control.

diagram of an analog fluid control system. Modify the block diagram to show how the fluid level could be digitally controlled.

- 1.2 If the temperature of the fluid in Problem 1.1 is to be regulated together with its level, modify the analog control system to achieve the additional control. (*Hint:* An additional actuator and sensor are needed.) Obtain a block diagram for the two-input–two-output control system with digital control.
- 1.3 Position control servos are discussed extensively in classical control texts. Draw a block diagram for a direct current motor position control system after consulting your classical control text. Modify the block diagram to obtain a digital position control servo.
- 1.4 Repeat Problem 1.3 for a velocity control servo.
- 1.5 A ballistic missile (see Fig. P1.5) is required to follow a predetermined flight path by adjusting its angle of attack α (the angle between its axis and its velocity vector v). The angle of attack is controlled by adjusting the thrust angle δ (angle between the thrust direction and the axis of the missile). Draw a block diagram for a digital control system for the angle of attack, including a gyroscope to measure the angle α and a motor to adjust the thrust angle δ .
- 1.6 A system is proposed to remotely control a missile from an earth station. Because of cost and technical constraints, the missile coordinates would be measured every 20 s for a missile speed of up to 0.5 km/s. Is such a control scheme feasible? What would the designers need to do to eliminate potential problems?
- 1.7 The control of the recording head of a dual actuator hard disk drive (HDD) requires two types of actuators to achieve the required high real density. The first is a coarse voice coil motor with a large stroke but slow dynamics, and the second is a fine piezoelectric transducer (PZT) with a small stroke and fast dynamics. A sensor measures the head position, and the position error is fed to a separate controller for each actuator. Draw a block diagram for a dual actuator digital control system for the HDD.²

² Ding, J., Marcassa, F., Wu, S.C., Tomizuka, M., 2006. Multirate control for computational saving. IEEE Trans. Control Systems Tech. 14 (1), 165–169.

- 1.8 In a planar contour tracking task performed by a robot manipulator, the robot end effector is required to track the contour of an unknown object with a given reference tangential velocity and by applying a given force to the object in the normal direction. For this purpose, a force sensor can be applied on the end effector, while the end-effector velocity can be determined by means of the joint velocities. Draw a block diagram of the digital control system.³
- 1.9 A typical main irrigation canal consists of several pools separated by gates that are used for regulating the water distribution from one pool to the next. In automatically regulated canals, the controlled variables are the water levels, the manipulated variables are the gate positions, and the fundamental perturbation variables are the unknown offtake discharges.⁴ Draw a block diagram of the control scheme.

³ Jatta, F., Legnani, G., Visioli, A., Ziliani, G., 2006. On the use of velocity feedback in hybrid force/velocity control of industrial manipulators. *Control Engineering Practice* 14, 1045–1055.

⁴ Feliu-Battle, V., Rivas Perez, R., Sanchez Rodriguez, L., 2007. Fractional robust control of main irrigation canals with variable dynamic parameters. *Control Engineering Practice* 15, 673–686.

Discrete-time systems

Objectives

After completing this chapter, the reader will be able to do the following:

1. Explain why difference equations result from digital control of analog systems.
2. Obtain the z -transform of a given time sequence and the time sequence corresponding to a function of z .
3. Solve linear time-invariant (LTI) difference equations using the z -transform.
4. Obtain the z -transfer function of an LTI system.
5. Obtain the time response of an LTI system using its transfer function or impulse response sequence.
6. Obtain the modified z -transform for a sampled time function.
7. Select a suitable sampling period for a given LTI system based on its dynamics.

Digital control involves systems whose control is updated at discrete time instants.

Discrete-time models provide mathematical relations between the system variables at these time instants. In this chapter, we develop the mathematical properties of discrete-time models that are used throughout the remainder of the text. For most readers, this material provides a concise review of material covered in basic courses on control and system theory. However, the material is self-contained, and familiarity with discrete-time systems is not required. We begin with an example that illustrates how discrete-time models arise from analog systems under digital control.

Chapter Outline

- 2.1 Analog systems with piecewise constant inputs 10
- 2.2 Difference equations 12
- 2.3 The z -transform 13
 - 2.3.1 z -transforms of standard discrete-time signals 14
 - 2.3.2 Properties of the z -transform 17
 - 2.3.2.1 Linearity 17
 - 2.3.2.2 Time delay 17

2.3.2.3 Time advance	18
2.3.2.4 Multiplication by exponential	19
2.3.2.5 Complex differentiation	20
2.3.3 Inversion of the z-transform	21
2.3.3.1 Long division	21
2.3.3.2 Partial fraction expansion	22
2.3.4 The final value theorem	31
2.4 Computer-aided design	33
2.5 z-transform solution of difference equations	34
2.6 The time response of a discrete-time system	35
2.6.1 Convolution summation	36
2.6.2 The convolution theorem	38
2.7 The modified z-transform	41
2.8 Frequency response of discrete-time systems	44
2.8.1 Properties of the frequency response of discrete-time systems	47
2.8.2 MATLAB commands for the discrete-time frequency response	48
2.9 The sampling theorem	50
2.9.1 Selection of the sampling frequency	52
Resources	55
Problems	55
Computer exercises	59

2.1 Analog systems with piecewise constant inputs

In most engineering applications, it is necessary to control a physical system or **plant** so that it behaves according to given design specifications. Typically, the plant is analog, the control is piecewise constant, and the control action is updated periodically. The output is sampled every T s as shown in Fig. 2.1 where k is an integer. This arrangement results in an overall system that is conveniently described by a discrete-time model. We demonstrate this concept using a simple example.

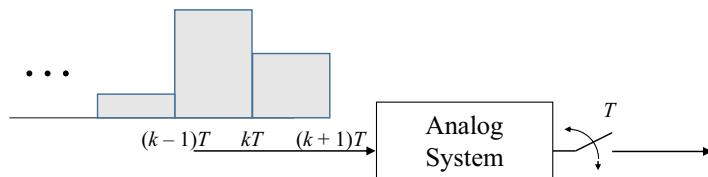


Figure 2.1

Analog system with piecewise constant input and sampled output.

Example 2.1

Consider the tank control system in Fig. 2.2. In the figure, lowercase letters denote perturbations from fixed steady-state values. The variables are defined as:

- H = steady-state fluid height in the tank
- h = height perturbation from the nominal value
- Q = steady-state flow rate through the tank
- q_i = inflow perturbation from the nominal value
- q_o = outflow perturbation from the nominal value

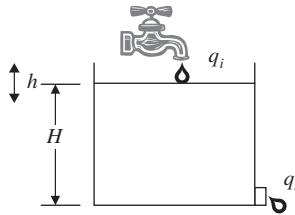


Figure 2.2
Fluid level control system.

It is necessary to maintain a constant fluid level H by adjusting the fluid flow rate into the tank around Q . Obtain an analog mathematical model of the tank, and use it to obtain a discrete-time model for the system with piecewise constant inflow perturbation q_i and output h .

Solution

Although the fluid system is nonlinear, a linear model can satisfactorily describe the system under the assumption that fluid level is regulated around a constant value. The linearized model for the outflow valve is analogous to an electrical resistor and is given by

$$h = R q_o$$

where h is the perturbation in tank level from nominal, q_o is the perturbation in the outflow from the tank from a nominal level Q , and R is the fluid resistance of the valve.

Assuming an incompressible fluid, the principle of conservation of mass reduces to the volumetric balance: rate of fluid volume increase = rate of fluid volume in – rate of fluid volume out:

$$\frac{dC(h+H)}{dt} = (q_i + Q) - (q_o + Q)$$

where C is the area of the tank or its fluid capacitance. The term H is a constant and its derivative is zero, and the term Q cancels so that the remaining terms only involve perturbations. Substituting for the outflow q_o from the linearized valve equation into the volumetric fluid balance gives the analog mathematical model

$$\frac{dh}{dt} + \frac{h}{\tau} = \frac{q_i}{C}$$

where $\tau = RC$ is the fluid time constant for the tank. The solution of this differential equation is

$$h(t) = e^{-(t-t_0)/\tau} h(t_0) + \frac{1}{C} \int_{t_0}^t e^{-(t-\lambda)/\tau} q_i(\lambda) d\lambda$$

Example 2.1—cont'd

Let q_i be constant over each sampling period T —that is, $q_i(t) = q_i(k) = \text{constant}$ for t in the interval $[kT, (k+1)T]$. Then we can solve the analog equation over any sampling period to obtain

$$h(k+1) = e^{-T/\tau} h(k) + R \left[1 - e^{-T/\tau} \right] q_i(k), \quad k = 0, 1, 2, \dots$$

where the variables at time kT are denoted by the argument k . This is the desired discrete-time model describing the system with piecewise constant control. Details of the solution are left as an exercise (Problem 2.1).

The discrete-time model obtained in Example 2.1 is known as a difference equation. Because the model involves a linear time-invariant analog plant, the equation is linear time invariant. Next, we briefly discuss difference equations, and then we introduce a transform used to solve them.

2.2 Difference equations

Difference equations arise in problems where the independent variable, usually time, is assumed to have a discrete set of possible values. The nonlinear difference equation

$$\begin{aligned} y(k+n) &= f[y(k+n-1), y(k+n-2), \dots, y(k+1), y(k), u(k+n), \\ &\quad u(k+n-1), \dots, u(k+1), u(k)] \end{aligned} \quad (2.1)$$

with forcing function $u(k)$ is said to be of order n because the difference between the highest and lowest time arguments of $y(\cdot)$ and $u(\cdot)$ is n . The equations we deal with in this text are almost exclusively linear and are of the form

$$\begin{aligned} y(k+n) + a_{n-1}y(k+n-1) + \dots + a_1y(k+1) + a_0y(k) \\ = b_nu(k+n) + b_{n-1}u(k+n-1) + \dots + b_1u(k+1) + b_0u(k) \end{aligned} \quad (2.2)$$

We further assume that the coefficients $a_i, b_i, i = 0, 1, 2, \dots$, are constant. The difference equation is then referred to as linear time invariant, or LTI. If the forcing function $u(k)$ is equal to zero, the equation is said to be *homogeneous*.

Example 2.2

For each of the following difference equations, determine the order of the equation. Is the equation (a) linear, (b) time invariant, or (c) homogeneous?

1. $y(k+2) + 0.8y(k+1) + 0.07y(k)u(k)$
2. $y(k+4) + \sin(0.4k)y(k+1) + 0.3y(k) = 0$
3. $y(k+1) = -0.1y^2(k)$

Example 2.2—cont'd**Solution**

1. The equation is second order. All terms enter the equation linearly and have constant coefficients. The equation is therefore LTI. A forcing function appears in the equation, so it is nonhomogeneous.
2. The equation is fourth order. The second coefficient is time dependent, but all the terms are linear and there is no forcing function. The equation is therefore linear time varying and homogeneous.
3. The equation is first order. The right-hand side (RHS) is a nonlinear function of $y(k)$ but does not include a forcing function or terms that depend on time explicitly. The equation is therefore nonlinear, time invariant, and homogeneous.

Difference equations can be solved using classical methods analogous to those available for differential equations. Alternatively, z -transforms provide a convenient approach for solving LTI equations, as discussed in the next section.

2.3 The z -transform

The z -transform is an important tool in the analysis and design of discrete-time systems. It simplifies the solution of discrete-time problems by converting LTI difference equations to algebraic equations and convolution to multiplication. Thus, it plays a role similar to that served by Laplace transforms in continuous-time problems. Because we are primarily interested in application to digital control systems, this brief introduction to the z -transform is restricted to **causal signals** (i.e., signals with zero values for negative time) and the one-sided z -transform.

The following are two alternative definitions of the z -transform.

Definition 2.1

Given the causal sequence $\{u_0, u_1, u_2, \dots, u_k, \dots\}$, its z -transform is defined as

$$\begin{aligned} U(z) &= u_0 + u_1 z^{-1} + u_2 z^{-2} + \dots + u_k z^{-k} \\ &= \sum_{k=0}^{\infty} u_k z^{-k} \end{aligned} \tag{2.3}$$

The variable z^{-1} in the preceding equation can be regarded as a time delay operator. The z -transform of a given sequence can be easily obtained as in the following example.

Definition 2.2

Given the impulse train representation of a discrete-time signal,

$$\begin{aligned} u^*(t) &= u_0\delta(t) + u_1\delta(t-T) + u_2\delta(t-2T) + \dots + u_k\delta(t-kT) + \dots \\ &= \sum_{k=0}^{\infty} u_k\delta(t-kT) \end{aligned} \quad (2.4)$$

the Laplace transform of (2.4) is

$$\begin{aligned} U^*(s) &= u_0 + u_1e^{-sT} + u_2e^{-2sT} + \dots + u_ke^{-ksT} + \dots \\ &= \sum_{k=0}^{\infty} u_k(e^{-sT})^k \end{aligned} \quad (2.5)$$

Let z be defined by

$$z = e^{sT} \quad (2.6)$$

Then, substituting from (2.6) in (2.5) yields the z -transform expression (2.3).

Example 2.3

Obtain the z -transform of the sequence $\{u_k\}_{k=0}^{\infty} = \{1, 3, 2, 0, 4, 0, 0, 0, \dots\}$.

Solution

Applying Definition 2.1 gives $U(z) = 1 + 3z^{-1} + 2z^{-2} + 4z^{-4}$.

Although the preceding two definitions yield the same transform, each has its advantages and disadvantages. The first definition allows us to avoid the use of impulses and the Laplace transform. The second allows us to treat z as a complex variable and to use some of the familiar properties of the Laplace transform, such as linearity.

Clearly, it is possible to use Laplace transformation to study discrete time, continuous time, and mixed systems. However, the z -transform offers significant simplification in notation for discrete-time systems and greatly simplifies their analysis and design.

2.3.1 ***z**-transforms of standard discrete-time signals*

Having defined the z -transform, we now obtain the z -transforms of commonly used discrete-time signals such as the sampled step, exponential, and the discrete-time impulse. The following identities are used repeatedly to derive several important results:

$$\sum_{k=0}^n a^k = \frac{1-a^{n+1}}{1-a}, \quad a \neq 1 \quad (2.7)$$

$$\sum_{k=0}^{\infty} a^k = \frac{1}{1-a}, \quad |a| < 1$$

Example 2.4**Unit impulse**

Consider the discrete-time impulse (Fig. 2.3)

$$u(k) = \delta(k) = \begin{cases} 1, & k = 0 \\ 0, & k \neq 0 \end{cases}$$

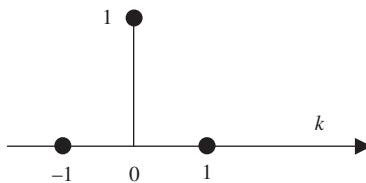


Figure 2.3
Discrete-time impulse.

Applying Definition 2.1 gives the z-transform

$$U(z) = 1$$

Alternatively, one may consider the impulse-sampled version of the delta function $u^*(t) = \delta(t)$. This has the Laplace transform

$$U^*(s) = 1$$

Substitution from (2.6) has no effect. Thus, the z-transform obtained using Definition 2.2 is identical to that obtained using Definition 2.1.

Example 2.5**Sampled step**

Consider the sequence $\{u_k\}_{k=0}^{\infty} = \{1, 1, 1, 1, 1, 1, \dots\}$ of Fig. 2.4. Definition 2.1 gives the z-transform

$$U(z) = 1 + z^{-1} + z^{-2} + z^{-3} + \dots + z^{-k} + \dots$$

$$= \sum_{k=0}^{\infty} z^{-k}$$

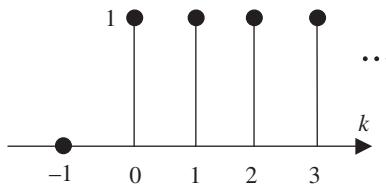
Example 2.5—cont'd

Figure 2.4
Sampled unit step.

Using the identity (2.7) gives the following closed-form expression for the z-transform:

$$\begin{aligned} U(z) &= \frac{1}{1 - z^{-1}} \\ &= \frac{z}{z - 1} \end{aligned}$$

Note that (2.7) is only valid for $|z| < 1$. This implies that the z-transform expression we obtain has a region of convergence outside which it is not valid. The region of convergence must be clearly given when using the more general two-sided transform with functions that are nonzero for negative time. However, for the one-sided z-transform and time functions that are zero for negative time, we can essentially extend regions of convergence and use the z-transform in the entire z-plane.¹

¹ The idea of extending the definition of a complex function to the entire complex plane is known as **analytic continuation**. For a discussion of this topic, consult any text on complex analysis.

Example 2.6**Exponential**

Let

$$u(k) = \begin{cases} a^k, & k \geq 0 \\ 0, & k < 0 \end{cases}$$

Fig. 2.5 shows the case where $0 < a < 1$. Definition 2.1 gives the z-transform

$$U(z) = 1 + az^{-1} + a^2z^{-2} + \dots + a^kz^{-k} + \dots$$

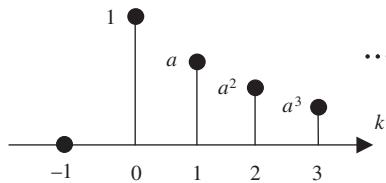


Figure 2.5
Sampled exponential.

Example 2.6—cont'd

Using (2.7), we obtain

$$\begin{aligned} U(z) &= \frac{1}{1 - (\alpha/z)} \\ &= \frac{z}{z - \alpha} \end{aligned}$$

As in Example 2.5, we can use the transform in the entire z -plane in spite of the validity condition for (2.7) because our time function is zero for negative time.

2.3.2 Properties of the z -transform

The z -transform can be derived from the Laplace transform as shown in Definition 2.2. Hence, it shares several useful properties with the Laplace transform, which can be stated without proof. These properties can also be easily proved directly, and the proofs are left as an exercise for the reader. Proofs are provided for properties that do not obviously follow from the Laplace transform.

2.3.2.1 Linearity

This equation follows directly from the linearity of the Laplace transform.

$$\mathcal{Z}\{\alpha f_1(k) + \beta f_2(k)\} = \alpha F_1(z) + \beta F_2(z) \quad (2.8)$$

Example 2.7

Find the z -transform of the causal sequence

$$f(k) = 2 \times 1(k) + 4\delta(k), \quad k = 0, 1, 2, \dots$$

Solution

Using linearity, the transform of the sequence is

$$F(z) = \mathcal{Z}\{2 \times 1(k) + 4\delta(k)\} = 2\mathcal{Z}\{1(k)\} + 4\mathcal{Z}\{\delta(k)\} = \frac{2z}{z-1} + 4 = \frac{6z-4}{z-1}$$

2.3.2.2 Time delay

This equation follows from the time delay property of the Laplace transform and Eq. (2.6).

$$\mathcal{Z}\{f(k-n)\} = z^{-n}F(z) \quad (2.9)$$

Example 2.8

Find the z-transform of the causal sequence

$$f(k) = \begin{cases} 4, & k = 2, 3, \dots \\ 0, & \text{otherwise} \end{cases}$$

Solution

The given sequence is a sampled step starting at $k = 2$ rather than $k = 0$ (i.e., it is delayed by two sampling periods). Using the delay property, we have

$$F(z) = \mathcal{Z}\{4 \times 1(k-2)\} = 4z^{-2} \mathcal{Z}\{1(k)\} = z^{-2} \frac{4z}{z-1} = \frac{4}{z(z-1)}$$

2.3.2.3 Time advance

$$\begin{aligned} \mathcal{Z}\{f(k+1)\} &= zF(z) - zf(0) \\ \mathcal{Z}\{f(k+n)\} &= z^n F(z) - z^n f(0) - z^{n-1} f(1) - \cdots - zf(n-1) \end{aligned} \quad (2.10)$$

Proof

Only the first part of the theorem is proved here. The second part can be easily proved by induction. We begin by applying the z-transform Definition 2.1 to a discrete-time function advanced by one sampling interval. This gives

$$\begin{aligned} \mathcal{Z}\{f(k+1)\} &= \sum_{k=0}^{\infty} f(k+1)z^{-k} \\ &= z \sum_{k=0}^{\infty} f(k+1)z^{-(k+1)} \end{aligned}$$

Now add and subtract the initial condition $f(0)$ to obtain

$$\mathcal{Z}\{f(k+1)\} = z \left\{ \left[f(0) + \sum_{k=0}^{\infty} f(k+1)z^{-(k+1)} \right] - f(0) \right\}$$

Next, change the index of summation to $m = k + 1$ and rewrite the z-transform as

$$\begin{aligned} \mathcal{Z}\{f(k+1)\} &= z \left\{ \left[\sum_{m=0}^{\infty} f(m)z^{-m} \right] - f(0) \right\} \\ &= zF(z) - zf(0) \end{aligned}$$

Example 2.9

Using the time advance property, find the z-transform of the causal sequence

$$\{f(k)\} = \{4, 8, 16, \dots\}$$

Solution

The sequence can be written as

$$f(k) = 2^{k+2} = g(k+2), \quad k = 0, 1, 2, \dots$$

where $g(k)$ is the exponential time function

$$g(k) = 2^k, \quad k = 0, 1, 2, \dots$$

Using the time advance property, we write the transform

$$F(z) = z^2 G(z) - z^2 g(0) - zg(1) = z^2 \frac{z}{z-2} - z^2 - 2z = \frac{4z}{z-2}$$

Clearly, the solution can be obtained directly by rewriting the sequence as

$$\{f(k)\} = 4\{1, 2, 4, \dots\}$$

and using the linearity of the z-transform.

2.3.2.4 Multiplication by exponential

$$\mathcal{Z}\{a^{-k}f(k)\} = F(az) \quad (2.11)$$

Proof

$$\text{LHS} = \sum_{k=0}^{\infty} a^{-k} f(k) z^{-k} = \sum_{k=0}^{\infty} f(k) (az)^{-k} = F(az)$$

Example 2.10

Find the z-transform of the exponential sequence

$$f(k) = e^{-\alpha k T}, \quad k = 0, 1, 2, \dots$$

Solution

Recall that the z-transform of a sampled step is

$$F(z) = \frac{z}{z-1} = (1 - z^{-1})^{-1}$$

Example 2.10—cont'd

and observe that $f(k)$ can be rewritten as

$$f(k) = (e^{\alpha T})^{-k} \times 1, \quad k = 0, 1, 2, \dots$$

Then apply the multiplication by exponential property to obtain

$$\mathcal{Z}\left\{(e^{\alpha T})^{-k} f(k)\right\} = \left[1 - (e^{\alpha T} z)^{-1}\right]^{-1} = \frac{z}{z - e^{-\alpha T}}$$

Example 2.6 gives the same answer with $a = e^{-\alpha T}$.

2.3.2.5 Complex differentiation

$$\mathcal{Z}\{k^m f(k)\} = \left(-z \frac{d}{dz}\right)^m F(z) \quad (2.12)$$

Proof

To prove the property by induction, we first establish its validity for $m = 1$. Then we assume its validity for any m and prove it for $m + 1$. This establishes its validity for $1 + 1 = 2$, then $2 + 1 = 3$, and so on.

For $m = 1$, we have

$$\begin{aligned} \mathcal{Z}\{kf(k)\} &= \sum_{k=0}^{\infty} kf(k)z^{-k} = \sum_{k=0}^{\infty} f(k) \left(-z \frac{d}{dz}\right) z^{-k} \\ &= \left(-z \frac{d}{dz}\right) \sum_{k=0}^{\infty} f(k)z^{-k} = \left(-z \frac{d}{dz}\right) F(z) \end{aligned}$$

Next, let the statement be true for any m and define the sequence

$$f_m(k) = k^m f(k), \quad k = 0, 1, 2, \dots$$

and obtain the transform

$$\begin{aligned} \mathcal{Z}\{kf_m(k)\} &= \sum_{k=0}^{\infty} kf_m(k)z^{-k} \\ &= \sum_{k=0}^{\infty} f_m(k) \left(-z \frac{d}{dz}\right) z^{-k} \\ &= \left(-z \frac{d}{dz}\right) \sum_{k=0}^{\infty} f_m(k)z^{-k} = \left(-z \frac{d}{dz}\right) F_m(z) \end{aligned}$$

Proof—cont'd

Substituting for $F_m(z)$, we obtain the result

$$\mathcal{Z}\{k^{m+1}f(k)\} = \mathcal{Z}\{kf_m(k)\} = \left(-z \frac{d}{dz}\right)^{m+1} F(z)$$

Example 2.11

Find the z-transform of the sampled ramp sequence

$$f(k) = k, \quad k = 0, 1, 2, \dots$$

Solution

Recall that the z-transform of a sampled step is

$$F(z) = \frac{z}{z - 1}$$

and observe that $f(k)$ can be rewritten as

$$f(k) = k \times 1, \quad k = 0, 1, 2, \dots$$

Then apply the complex differentiation property to obtain

$$\mathcal{Z}\{k \times 1\} = \left(-z \frac{d}{dz}\right) \left(\frac{z}{z - 1}\right) = (-z) \frac{(z - 1) - z}{(z - 1)^2} = \frac{z}{(z - 1)^2}$$

2.3.3 Inversion of the z-transform

Because the purpose of z-transformation is often to simplify the solution of time domain problems, it is essential to inverse-transform z-domain functions. As in the case of Laplace transforms, a complex integral can be used for inverse transformation. This integral is difficult to use and is rarely needed in engineering applications. Two simpler approaches for inverse z-transformation are discussed in this section.

2.3.3.1 Long division

This approach is based on Definition 2.1, which relates a time sequence to its z-transform directly. We first use long division to obtain as many terms as desired of the z-transform

expansion; then we use the coefficients of the expansion to write the time sequence. The following two steps give the inverse z -transform of a function $F(z)$:

1. Using long division, expand $F(z)$ as a series to obtain

$$F_t(z) = f_0 + f_1 z^{-1} + \dots + f_i z^{-i} = \sum_{k=0}^i f_k z^{-k}$$

2. Write the inverse transform as the sequence

$$\{f_0, f_1, \dots, f_i, \dots\}$$

The number of terms i obtained by long division is selected to yield a sufficient number of points in the time sequence.

Example 2.12

Obtain the inverse z -transform of the function

$$F(z) = \frac{z+1}{z^2 + 0.2z + 0.1}$$

Solution

1. Long Division

$$\begin{array}{r} z^{-1} + 0.8z^{-2} - 0.26z^{-3} + \dots \\ z^2 + 0.2z + 0.1 \overline{)z + 1} \\ \hline z + 0.2 + 0.1z^{-1} \\ \hline 0.8 - 0.10z^{-1} \\ \hline 0.8 + 0.16z^{-1} + 0.08z^{-2} \\ \hline -0.26z^{-1} - \dots \end{array}$$

Thus, $F_t(z) = 0 + z^{-1} + 0.8z^{-2} - 0.26z^{-3}$

2. Inverse Transformation

$$\{f_k\} = \{0, 1, 0.8, -0.26, \dots\}$$

2.3.3.2 Partial fraction expansion

This method is almost identical to that used in inverting Laplace transforms. However, because most z -functions have the term z in their numerator, it is often convenient to expand $F(z)/z$ rather than $F(z)$. As with Laplace transforms, partial fraction expansion

allows us to write the function as the sum of simpler functions that are the z -transforms of known discrete-time functions. The time functions are available in z -transform tables such as the table provided in Appendix I.

The procedure for inverse z -transformation is

1. Find the partial fraction expansion of $F(z)/z$ or $F(z)$.
2. Obtain the inverse transform $f(k)$ using the z -transform tables.

We consider three types of z -domain functions $F(z)$: functions with simple (nonrepeated) real poles, functions with complex conjugate and real poles, and functions with repeated poles. We discuss examples that demonstrate partial fraction expansion and inverse z -transformation in each case.

Case 1 Simple real roots

The most convenient method to obtain the partial fraction expansion of a function with simple real roots is the method of residues. The residue of a complex function $F(z)$ at a simple pole z_i is given by

$$A_i = (z - z_i)F(z)]_{z \rightarrow z_i} \quad (2.13)$$

This is the partial fraction coefficient of the i th term of the expansion

$$F(z) = \sum_{i=1}^n \frac{A_i}{z - z_i} \quad (2.14)$$

Because most terms in the z -transform tables include a z in the numerator (see Appendix I), it is often convenient to expand $F(z)/z$ and then to multiply both sides by z to obtain an expansion whose terms have a z in the numerator. Except for functions that already have a z in the numerator, this approach is slightly longer but has the advantage of simplifying inverse transformation. Both methods are examined through Example 2.13.

Example 2.13

Obtain the inverse z -transform of the function

$$F(z) = \frac{z + 1}{z^2 + 0.3z + 0.02}$$

Solution

It is instructive to solve this problem using two different methods. First, we divide by z ; then we obtain the partial fraction expansion.

1. Partial fraction expansion

Dividing the function by z , we expand as

Example 2.13—cont'd

$$\begin{aligned}\frac{F(z)}{z} &= \frac{z+1}{z(z^2 + 0.3z + 0.02)} \\ &= \frac{A}{z} + \frac{B}{z+0.1} + \frac{C}{z+0.2}\end{aligned}$$

where the partial fraction coefficients are given by

$$\begin{aligned}A &= z \frac{F(z)}{z} \Big|_{z=0} = F(0) = \frac{1}{0.02} = 50 \\ B &= (z+0.1) \frac{F(z)}{z} \Big|_{z=-0.1} = \frac{1-0.1}{(-0.1)(0.1)} = -90 \\ C &= (z+0.2) \frac{F(z)}{z} \Big|_{z=-0.2} = \frac{1-0.2}{(-0.2)(-0.1)} = 40\end{aligned}$$

Thus, the partial fraction expansion is

$$F(z) = \frac{50z}{z} - \frac{90z}{z+0.1} + \frac{40z}{z+0.2}$$

2. Table lookup

$$f(k) = \begin{cases} 50\delta(k) - 90(-0.1)^k + 40(-0.2)^k, & k \geq 0 \\ 0, & k < 0 \end{cases}$$

Note that $f(0) = 0$, so the time sequence can be rewritten as

$$f(k) = \begin{cases} -90(-0.1)^k + 40(-0.2)^k, & k \geq 1 \\ 0, & k < 1 \end{cases}$$

Now, we solve the same problem without dividing by z .

1. Partial fraction expansion

We obtain the partial fraction expansion directly

$$\begin{aligned}F(z) &= \frac{z+1}{z^2 + 0.3z + 0.02} \\ &= \frac{A}{z+0.1} + \frac{B}{z+0.2}\end{aligned}$$

where the partial fraction coefficients are given by

$$\begin{aligned}A &= (z+0.1)F(z) \Big|_{z=-0.1} = \frac{1-0.1}{0.1} = 9 \\ B &= (z+0.2)F(z) \Big|_{z=-0.2} = \frac{1-0.2}{-0.1} = -8\end{aligned}$$

Example 2.13—cont'd

Thus, the partial fraction expansion is

$$F(z) = \frac{9}{z + 0.1} - \frac{8}{z + 0.2}$$

2. Table lookup

Standard z-transform tables do not include the terms in the expansion of $F(z)$. However, $F(z)$ can be written as

$$F(z) = \frac{9z}{z + 0.1}z^{-1} - \frac{8z}{z + 0.2}z^{-1}$$

Then we use the delay theorem to obtain the inverse transform

$$f(k) = \begin{cases} 9(-0.1)^{k-1} - 8(-0.2)^{k-1}, & k \geq 1 \\ 0, & k < 1 \end{cases}$$

Verify that this is the answer obtained earlier when dividing by z written in a different form (observe the exponent in the preceding expression).

Although it is clearly easier to obtain the partial fraction expansion without dividing by z , inverse transforming requires some experience. There are situations where division by z may actually simplify the calculations, as seen in Example 2.14.

Example 2.14

Find the inverse z-transform of the function

$$F(z) = \frac{z}{(z + 0.1)(z + 0.2)(z + 0.3)}$$

Solution**1. Partial fraction expansion**

Dividing by z simplifies the numerator and gives the expansion

$$\begin{aligned} \frac{F(z)}{z} &= \frac{1}{(z + 0.1)(z + 0.2) + (z + 0.3)} \\ &= \frac{A}{z + 0.1} + \frac{B}{z + 0.2} + \frac{C}{z + 0.3} \end{aligned}$$

Example 2.14—cont'd

where the partial fraction coefficients are

$$A = (z + 0.1) \frac{F(z)}{z} \Big|_{z=-0.1} = \frac{1}{(0.1)(0.2)} = 50$$

$$B = (z + 0.2) \frac{F(z)}{z} \Big|_{z=-0.2} = \frac{1}{(-0.1)(0.1)} = -100$$

$$C = (z + 0.3) \frac{F(z)}{z} \Big|_{z=-0.3} = \frac{1}{(-0.2)(-0.1)} = 50$$

Thus, the partial fraction expansion after multiplying by z is

$$F(z) = \frac{50z}{z + 0.1} - \frac{100z}{z + 0.2} + \frac{50z}{z + 0.3}$$

2. Table lookup

$$f(k) = \begin{cases} 50(-0.1)^k - 100(-0.2)^k + 50(-0.3)^k, & k \geq 0 \\ 0, & k < 0 \end{cases}$$

Case 2 Complex conjugate and simple real roots

For a function $F(z)$ with real and complex poles, the partial fraction expansion includes terms with real roots and others with complex roots. Assuming that $F(z)$ has real coefficients, then its complex roots occur in complex conjugate pairs and can be combined to yield a function with real coefficients and a quadratic denominator. To inverse-transform such a function, use the following z-transforms (see Appendix I):

$$\mathcal{Z}\{e^{-\alpha k} \sin(k\omega_d)\} = \frac{e^{-\alpha} \sin(\omega_d)z}{z^2 - 2e^{-\alpha} \cos(\omega_d)z + e^{-2\alpha}} \quad (2.15)$$

$$\mathcal{Z}\{e^{-\alpha k} \cos(k\omega_d)\} = \frac{z[z - e^{-\alpha} \cos(\omega_d)]}{z^2 - 2e^{-\alpha} \cos(\omega_d)z + e^{-2\alpha}} \quad (2.16)$$

Note that in some tables the sampling period T is omitted, but then ω_d is given in radians and α is dimensionless. The denominators of the two transforms are identical and have complex conjugate roots. The numerators can be scaled and combined to give the desired inverse transform.

To obtain the partial fraction expansion, we use the residues method shown in Case 1. With complex conjugate poles, we obtain the partial fraction expansion

$$F(z) = \frac{Az}{z - p} + \frac{A^*z}{z - p^*} \quad (2.17)$$

Case 2 Complex conjugate and simple real roots—cont'd

We then inverse z-transform to obtain

$$\begin{aligned}f(k) &= Ap^k + A^*p^{*k} \\&= |A||p|^k \left[e^{j(\theta_p k + \theta_A)} + e^{-j(\theta_p k + \theta_A)} \right]\end{aligned}$$

where θ_p and θ_A are the angle of the pole p and the angle of the partial fraction coefficient A , respectively. We use the exponential expression for the cosine function to obtain

$$f(k) = 2|A||p|^k \cos(\theta_p k + \theta_A) \quad (2.18)$$

Most modern calculators can perform complex arithmetic, and the residues method is preferable in most cases. Alternatively, by equating coefficients, we can avoid the use of complex arithmetic entirely, but the calculations can be quite tedious. Example 2.15 demonstrates the two methods.

Example 2.15

Find the inverse z-transform of the function

$$F(z) = \frac{z^3 + 2z + 1}{(z - 0.1)(z^2 + z + 0.5)}$$

Solution: equating coefficients*1. Partial fraction expansion*

Dividing the function by z gives

$$\begin{aligned}\frac{F(z)}{z} &= \frac{z^3 + 2z + 1}{z(z - 0.1)(z^2 + z + 0.5)} \\&= \frac{A_1}{z} + \frac{A_2}{z - 0.1} + \frac{Az + B}{z^2 + z + 0.5}\end{aligned}$$

The first two coefficients can be easily evaluated as before. Thus,

$$A_1 = F(0) = -20$$

$$A_2 = (z - 0.1) \frac{F(z)}{z} \cong 19.689$$

To evaluate the remaining coefficients, we multiply the equation by the denominator and equate coefficients to obtain

$$z^3: A_1 + A_2 + A = 1$$

$$z^1: 0.4A_1 + 0.5A_2 - 0.1B = 2$$

Example 2.15—cont'd

where the coefficients of the third- and first-order terms yield separate equations in A and B . Because A_1 and A_2 have already been evaluated, we can solve each of the two equations for one of the remaining unknowns to obtain

$$A \cong 1.311 \quad B \cong -1.557$$

If we chosen to equate coefficients without first evaluating A_1 and A_2 , we would have faced that considerably harder task of solving four equations in four unknowns. The remaining coefficients can be used to check our calculations:

$$z^0: -0.05A_1 = 0.05(20) = 1$$

$$z^2: 0.9A_1 + A_2 - 0.1A + B = 0.9(-20) + 19.689 - 0.1(1.311) - 1.557 \cong 0$$

The results of these checks are approximate, because approximations were made in the calculations of the coefficients. The partial fraction expansion is

$$F(z) = -20 + \frac{19.689z}{z - 0.1} + \frac{1.311z^2 - 1.557z}{z^2 + z + 0.5}$$

2. Table lookup

The first two terms of the partial fraction expansion can be easily found in the z -transform tables. The third term resembles the transforms of a sinusoid multiplied by an exponential if rewritten as

$$\frac{1.311z^2 - 1.557z}{z^2 - 2(-0.5)z + 0.5} = \frac{1.311z[z - e^{-\alpha} \cos(\omega_d)] - Ce^{-\alpha} \sin(\omega_d)}{z^2 - 2e^{-\alpha} \cos(\omega_d)z + e^{-2\alpha}}$$

Starting with the constant term in the denominator, we equate coefficients to obtain

$$e^{-\alpha} = \sqrt{0.5} = 0.707$$

Next, the denominator z^1 term gives

$$\cos(\omega_d) = -0.5/e^{-\alpha} = -\sqrt{0.5} = -0.707$$

Thus, $\omega_d = 3\pi/4$, an angle in the second quadrant, with $\sin(\omega_d) = 0.707$.

Finally, we equate the coefficients of z^1 in the numerator to obtain

$$-1.311e^{-\alpha} \cos(\omega_d) - Ce^{-\alpha} \sin(\omega_d) = -0.5(C - 1.311) = -1.557$$

and solve for $C = 4.426$. Referring to the z -transform tables, we obtain the inverse transform

$$f(k) = -20\delta(k) + 19.689(0.1)^k + (0.707)^k [1.311 \cos(3\pi k / 4) - 4.426 \sin(3\pi k / 4)]$$

for positive time k . The sinusoidal terms can be combined using the trigonometric identities

$$\sin(A - B) = \sin(A)\cos(B) - \sin(B)\cos(A)$$

$$\sin^{-1}(1.311/4.616) = 0.288$$

and the constant $4.616 = \sqrt{(1.311)^2 + (4.426)^2}$. This gives

$$f(k) = -20\delta(k) + 19.689(0.1)^k - 4.616(0.707)^k \sin(3\pi k / 4 - 0.288)$$

Example 2.15—cont'd**Residues****1. Partial fraction expansion**

Dividing by z gives

$$\begin{aligned} \frac{F(z)}{z} &= \frac{z^3 + 2z + 1}{z(z - 0.1)[(z + 0.5)^2 + 0.5^2]} \\ &= \frac{A_1}{z} + \frac{A_2}{z - 0.1} + \frac{A_3}{z + 0.5 - j0.5} + \frac{A_3^*}{z + 0.5 + j0.5} \end{aligned}$$

The partial fraction expansion can be obtained as in the first approach

$$\begin{aligned} A_3 &= \left. \frac{z^3 + 2z + 1}{z(z - 0.1)(z + 0.5 + j0.5)} \right|_{z=-0.5+j0.5} \cong 0.656 + j2.213 \\ F(z) &= -20 + \frac{19.689z}{z - 0.1} + \frac{(0.656 + j2.213)z}{z + 0.5 - j0.5} + \frac{(0.656 - j2.213)z}{z + 0.5 + j0.5} \end{aligned}$$

We convert the coefficient A_3 from Cartesian to polar form:

$$A_3 = 0.656 + j2.213 = 2.308e^{j1.283}$$

We inverse z -transform to obtain

$$f(k) = -20\delta(k) + 19.689(0.1)^k + 4.616(0.707)^k \cos(3\pi k / 4 + 1.283)$$

This is equivalent to the answer obtained earlier because $1.283 - \pi/2 = -0.288$.

Case 3 Repeated roots

For a function $F(z)$ with a repeated root of multiplicity r , r partial fraction coefficients are associated with the repeated root. The partial fraction expansion is of the form

$$F(z) = \frac{N(z)}{(z - z_1)^r \prod_{j=r+1}^n z - z_j} = \sum_{i=1}^r \frac{A_{1i}}{(z - z_1)^{r+1-i}} + \sum_{j=r+1}^n \frac{A_j}{z - z_j} \quad (2.19)$$

The coefficients for repeated roots are governed by

$$A_{1,i} = \left. \frac{1}{(i-1)!} \frac{d^{i-1}}{dz^{i-1}} (z - z_1)^r F(z) \right|_{z \rightarrow z_1}, \quad i = 1, 2, \dots, r \quad (2.20)$$

The coefficients of the simple or complex conjugate roots can be obtained as before using (2.13).

Example 2.16

Obtain the inverse z-transform of the function

$$F(z) = \frac{1}{z^2(z - 0.5)}$$

Solution**1. Partial fraction expansion**

Dividing by z gives

$$\frac{F(z)}{z} = \frac{1}{z^3(z - 0.5)} = \frac{A_{11}}{z^3} + \frac{A_{12}}{z^2} + \frac{A_{13}}{z} + \frac{A_4}{z - 0.5}$$

where

$$\begin{aligned} A_{11} &= z^3 \frac{F(z)}{z} \Big|_{z=0} = \frac{1}{z - 0.5} \Big|_{z=0} = -2 \\ A_{12} &= \frac{1}{1!} \frac{d}{dz} z^3 \frac{F(z)}{z} \Big|_{z=0} = \frac{d}{dz} \frac{1}{z - 0.5} \Big|_{z=0} = \frac{-1}{(z - 0.5)^2} \Big|_{z=0} = -4 \\ A_{13} &= \frac{1}{2!} \frac{d^2}{dz^2} z^3 \frac{F(z)}{z} \Big|_{z=0} \\ &= \left(\frac{1}{2}\right) \frac{d}{dz} \frac{-1}{(z - 0.5)^2} \Big|_{z=0} = \left(\frac{1}{2}\right) \frac{(-1)(-2)}{(z - 0.5)^3} \Big|_{z=0} = -8 \\ A_4 &= (z - 0.5) \frac{F(z)}{z} \Big|_{z=0.5} = \frac{1}{z^3} \Big|_{z=0.5} = 8 \end{aligned}$$

Thus, we have the partial fraction expansion

$$F(z) = \frac{1}{z^2(z - 0.5)} = \frac{8z}{z - 0.5} - 2z^{-2} - 4z^{-1} - 8$$

2. Table lookup

The z-transform tables and Definition 2.1 yield

$$f(k) = \begin{cases} 8(0.5)^k - 2\delta(k - 2) - 4\delta(k - 1) - 8\delta(k), & k \geq 0 \\ 0, & k < 0 \end{cases}$$

Evaluating $f(k)$ at $k = 0, 1, 2$ yields

$$f(0) = 8 - 8 = 0$$

$$f(1) = 8(0.5) - 4 = 0$$

$$f(2) = 8(0.5)^2 - 2 = 0$$

Example 2.16—cont'd

We can therefore rewrite the inverse transform as

$$f(k) = \begin{cases} (0.5)^{k-3}, & k \geq 3 \\ 0, & k < 3 \end{cases}$$

Note that the solution can be obtained directly using the delay theorem without the need for partial fraction expansion because $F(z)$ can be written as

$$F(z) = \frac{z}{z - 0.5} z^{-3}$$

The delay theorem and the inverse transform of an exponential yield the solution obtained earlier.

2.3.4 The final value theorem

The final value theorem allows us to calculate the limit of a sequence as k tends to infinity, if one exists, from the z -transform of the sequence. If one is only interested in the final value of the sequence, this constitutes a significant shortcut. The main pitfall of the theorem is that there are important cases where the limit does not exist. The two main cases are as follows:

1. An unbounded sequence
2. An oscillatory sequence

The reader is cautioned against blindly using the final value theorem, because this can yield misleading results.

Theorem 2.1: the final value theorem

If a sequence approaches a constant limit as k tends to infinity, then the limit is given by

$$f(\infty) = \lim_{k \rightarrow \infty} f(k) = \lim_{z \rightarrow 1} \left(\frac{z-1}{z} \right) F(z) = \lim_{z \rightarrow 1} (z-1) F(z) \quad (2.21)$$

Proof

Let $f(k)$ have a constant limit as k tends to infinity; then the sequence can be expressed as the sum

$$f(k) = f(\infty) + g(k), \quad k = 0, 1, 2, \dots$$

Proof—cont'd

with $g(k)$ a sequence that decays to zero as k tends to infinity; that is,

$$\lim_{k \rightarrow \infty} f(k) = f(\infty)$$

The z-transform of the preceding expression is

$$F(z) = \frac{f(\infty)z}{z - 1} + G(z)$$

The final value $f(\infty)$ is the partial fraction coefficient obtained by expanding $F(z)/z$ as follows:

$$f(\infty) = \lim_{z \rightarrow 1} (z - 1) \frac{F(z)}{z} = \lim_{z \rightarrow 1} (z - 1) F(z)$$

Example 2.17

Verify the final value theorem using the z-transform of a decaying exponential sequence and its limit as k tends to infinity.

Solution

The z-transform pair of an exponential sequence is

$$\{e^{-akT}\} \leftrightarrow \frac{z}{z - e^{-aT}}$$

with $a > 0$. The limit as k tends to infinity in the time domain is

$$f(\infty) = \lim_{k \rightarrow \infty} e^{-akT} = 0$$

The final value theorem gives

$$f(\infty) = \lim_{z \rightarrow 1} \left(\frac{z - 1}{z} \right) \left(\frac{z}{z - e^{-aT}} \right) = 0$$

Example 2.18

Obtain the final value for the sequence whose z-transform is

$$F(z) = \frac{z^2(z - a)}{(z - 1)(z - b)(z - c)}$$

What can you conclude concerning the constants b and c if it is known that the limit exists?

Solution

Applying the final value theorem, we have

$$f(\infty) = \lim_{z \rightarrow 1} \frac{z(z - a)}{(z - b)(z - c)} = \frac{1 - a}{(1 - b)(1 - c)}$$

Example 2.18—cont'd

To inverse z-transform the given function, one would have to obtain its partial fraction expansion, which would include three terms: the transform of the sampled step, the transform of the exponential b^k , and the transform of the exponential c^k . If either b or c is unity, then the inverse transform is a ramp. Therefore, the conditions for the sequence to converge to a constant limit and for the validity of the final value theorem are $|b| < 1$ and $|c| < 1$.

2.4 Computer-aided design

In this text, we make extensive use of computer-aided design (CAD) and analysis of control systems. We use MATLAB,² a powerful package with numerous useful commands. For the reader's convenience, we list some MATLAB commands after covering the relevant theory. The reader is assumed to be familiar with the CAD package but not with the digital system commands. We adopt the notation of bolding all user commands throughout the text. Readers using other CAD packages will find similar commands for digital control system analysis and design.

MATLAB typically handles coefficients as vectors with the coefficients listed in descending order. The function $G(z)$ with numerator $5(z + 3)$ and denominator $z^3 + 0.1z^2 + 0.4z$ is represented as the numerator polynomial

$$\gg \mathbf{num} = \mathbf{5}^*[1, 3]$$

and the denominator polynomial

$$\gg \mathbf{den} = [1, 0.1, 0.4, 0]$$

Multiplication of polynomials is equivalent to the convolution of their vectors of coefficients and is performed using the command

$$\gg \mathbf{denp} = \mathbf{conv}(\mathbf{den1}, \mathbf{den2})$$

where **denp** is the product of **den1** and **den2**.

The partial fraction coefficients are obtained using the command

$$\gg [\mathbf{r}, \mathbf{p}, \mathbf{k}] = \mathbf{residue}(\mathbf{num}, \mathbf{den})$$

where **p** represents the poles, **r** their residues, and **k** the coefficients of the polynomial resulting from dividing the numerator by the denominator. If the highest power in the

² MATLAB® is a copyright of MathWorks Inc., of Natick, Massachusetts.

numerator is smaller than the highest power in the denominator, \mathbf{k} is zero. This is the usual case encountered in digital control problems.

MATLAB allows the user to sample a function and z -transform it with the commands

```
>>g = tf(num, den)
>>gd = c2d(g, 0.1, 'imp')
```

Other useful MATLAB commands are available with the symbolic manipulation toolbox.

ztrans z-transform
iztrans inverse z-transform

To use these commands, we must first define symbolic variables such as \mathbf{z} , \mathbf{g} , and \mathbf{k} with the command

```
>>syms z g k
```

Powerful commands for symbolic manipulations are also available through packages such as MAPLE, MATHEMATICA, and MACSYMA.

2.5 z-transform solution of difference equations

By a process analogous to Laplace transform solution of differential equations, one can easily solve linear difference equations. The equations are first transformed to the z -domain (i.e., both the right- and left-hand sides of the equation are z -transformed). Then the variable of interest is solved for and inverse z -transformed. To transform the difference equation, we typically use the time delay or the time advance property. Inverse z -transformation is performed using the methods of [Section 2.3](#).

Example 2.19

Solve the linear difference equation

$$x(k+2) - (3/2)x(k+1) + (1/2)x(k) = 1(k)$$

with the initial conditions $x(0) = 1$, $x(1) = 5/2$.

Solution

1. z-transform

We begin by z -transforming the difference equation using [\(2.10\)](#) to obtain

$$[z^2X(z) - z^2x(0) - zx(1)] - (3/2)[zX(z) - zx(0)] + (1/2)X(z) = z/(z-1)$$

2. Solve for $X(z)$

Then we substitute the initial conditions and rearrange terms to obtain

$$[z^2 - (3/2)z + (1/2)]X(z) = z/(z-1) + z^2 + (5/2 - 3/2)z$$

Example 2.19—cont'd

which we solve for

$$X(z) = \frac{z[1 + (z + 1)(z - 1)]}{(z - 1)(z - 1)(z - 0.5)} = \frac{z^3}{(z - 1)^2(z - 0.5)}$$

3. Partial fraction expansion

The partial fraction of $X(z)/z$ is

$$\frac{X(z)}{z} = \frac{z^2}{(z - 1)^2(z - 0.5)} = \frac{A_{11}}{(z - 1)^2} + \frac{A_{12}}{z - 1} + \frac{A_3}{z - 0.5}$$

where

$$A_{11} = (z - 1)^2 \frac{X(z)}{z} \Big|_{z=1} = \frac{z^2}{z - 0.5} \Big|_{z=1} = \frac{1}{1 - 0.5} = 2$$

$$A_3 = (z - 0.5) \frac{X(z)}{z} \Big|_{z=0.5} = \frac{z^2}{(z - 1)^2} \Big|_{z=0.5} = \frac{(0.5)^2}{(0.5 - 1)^2} = 1$$

To obtain the remaining coefficient, we multiply by the denominator and get the equation

$$z^2 = A_{11}(z - 0.5) + A_{12}(z - 0.5)(z - 1) + A_3(z - 1)^2$$

Equating the coefficient of z^2 gives

$$z^2: 1 = A_{12} + A_3 = A_{12} + 1 \quad \text{i.e., } A_{12} = 0$$

Thus, the partial fraction expansion in this special case includes two terms only. We now have

$$X(z) = \frac{2z}{(z - 1)^2} + \frac{z}{z - 0.5}$$

4. Inverse z-transformation

From the z-transform tables, the inverse z-transform of $X(z)$ is

$$x(k) = 2k + (0.5)^k$$

2.6 The time response of a discrete-time system

The time response of a discrete-time linear system is the solution of the difference equation governing the system. For the *linear time-invariant* (LTI) case, the response due to the initial conditions and the response due to the input can be obtained separately and then added to obtain the overall response of the system. The response due to the input, or

the forced response, is the convolution summation of its input and its response to a unit impulse. In this section, we derive this result and examine its implications.

2.6.1 Convolution summation

The response of a discrete-time system to a unit impulse is known as the *impulse response sequence*. The impulse response sequence can be used to represent the response of a linear discrete-time system to an arbitrary input sequence

$$\{u(k)\} = \{u(0), u(1), \dots, u(i), \dots\} \quad (2.22)$$

To derive this relationship, we first represent the input sequence in terms of discrete impulses as follows:

$$\begin{aligned} u(k) &= u(0)\delta(k) + u(1)\delta(k-1) + u(2)\delta(k-2) + \dots + u(i)\delta(k-i) + \dots \\ &= \sum_{i=0}^{\infty} u(i)\delta(k-i) \end{aligned} \quad (2.23)$$

For a linear system, the principle of superposition applies, and the system output due to the input is the following sum of impulse response sequences:

$$\{y(l)\} = \{h(l)\}u(0) + \{h(l-1)\}u(1) + \{h(l-2)\}u(2) + \dots + \{h(l-i)\}u(i) + \dots$$

Hence, the output at time k is given by

$$y(k) = h(k)^*u(k) = \sum_{i=0}^{\infty} h(k-i)u(i) \quad (2.24)$$

where $(*)$ denotes the convolution operation.

For a causal system, the response due to an impulse at time i is an impulse response starting at time i and the delayed response $h(k-i)$ satisfies (Fig. 2.6):

$$h(k-i) = 0, \quad i > k \quad (2.25)$$

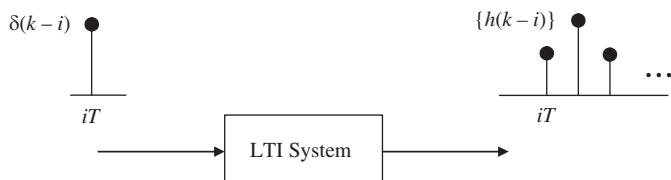


Figure 2.6

Response of a causal LTI discrete-time system to an impulse at iT .

In other words, a causal system is one whose impulse response is a causal time sequence. Thus, (2.24) reduces to

$$\begin{aligned} y(k) &= u(0)h(k) + u(1)h(k-1) + u(2)h(k-2) + \dots + u(k)h(0) \\ &= \sum_{i=0}^k u(i)h(k-i) \end{aligned} \quad (2.26)$$

A simple change of summation variable ($j = k - i$) transforms (2.26) to

$$\begin{aligned} y(k) &= u(k)h(0) + u(k-1)h(1) + u(k-2)h(2) + \dots + u(0)h(k) \\ &= \sum_{j=0}^k u(k-j)h(j) \end{aligned} \quad (2.27)$$

Eq. (2.24) is the convolution summation for a noncausal system, whose impulse response is nonzero for negative time, and it reduces to (2.26) for a causal system. The summations for time-varying systems are similar, but the impulse response at time i is $h(k, i)$. Here, we restrict our analysis to LTI systems. We can now summarize the result obtained in Theorem 2.2.

Theorem 2.2: Response of an LTI system

The response of an LTI discrete-time system to an arbitrary input sequence is given by the convolution summation of the input sequence and the impulse response sequence of the system.

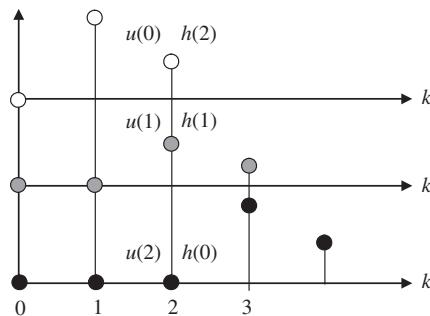
To better understand the operations involved in convolution summation, we evaluate one point in the output sequence using (2.24). For example,

$$\begin{aligned} y(2) &= \sum_{i=0}^2 u(i)h(2-i) \\ &= u(0)h(2) + u(1)h(1) + u(2)h(0) \end{aligned}$$

From Table 2.1 and Fig. 2.7, one can see the output corresponding to various components of the input of (2.23) and how they contribute to $y(2)$. Note that future input values do not contribute because the system is causal. The following example shows how the response of a linear discrete-time system to an arbitrary input sequence (2.22) can also be obtained using the coefficients of the unit step response $\{g(1), g(2), \dots, g(i), \dots\}$. Recall that for continuous-time systems the impulse response of a linear time-invariant system is the

Table 2.1: Input components and corresponding output components.

Input	Response	Fig. 2.6 Color
$u(0) \cdot \delta(k)$	$u(0) \cdot \{h(k)\}$	<i>White</i>
$u(1) \cdot \delta(k - 1)$	$u(1) \cdot \{h(k - 1)\}$	<i>Gray</i>
$u(2) \cdot \delta(k - 2)$	$u(2) \cdot \{h(k - 2)\}$	<i>Black</i>

**Figure 2.7**
Output at $k = 2$.

derivative of its step response. In the discrete time case, we have an analogous result based on the fact that a unit impulse can be expressed as the difference between two step inputs

$$\delta(k) = 1(k) - 1(k - 1)$$

By linearity, the impulse response is the difference between the two step responses

$$h(i) = g(i) - g(i - 1)$$

Substituting in (2.27) gives the convolution expression

$$\begin{aligned}
 y(k) &= \sum_{i=1}^k [g(i) - g(i - 1)]u(k - i) \\
 &= \sum_{i=0}^k g(i)[u(k - i) - u(k - i - 1)]
 \end{aligned} \tag{2.28}$$

2.6.2 The convolution theorem

The convolution summation is considerably simpler than the convolution integral that characterizes the response of linear continuous-time systems. Nevertheless, it is a fairly complex operation, especially if the output sequence is required over a long time period. Theorem 2.3 shows how the convolution summation can be avoided by z -transformation.

Theorem 2.3: The convolution theorem

The z -transform of the convolution of two time sequences is equal to the product of their z -transforms.

Proof

z -transforming (2.24) gives

$$\begin{aligned} Y(z) &= \sum_{k=0}^{\infty} y(k)z^{-k} \\ &= \sum_{k=0}^{\infty} \left[\sum_{i=0}^{\infty} u(i)h(k-i) \right] z^{-k} \end{aligned} \quad (2.29)$$

Interchange the order of summation and substitute $j = k - i$ to obtain

$$Y(z) = \sum_{i=0}^{\infty} \sum_{j=-i}^{\infty} u(i)h(j)z^{-(i+j)} \quad (2.30)$$

Using the causality property, (2.24) reduces (2.30) to

$$Y(z) = \left[\sum_{i=0}^{\infty} u(i)z^{-i} \right] \left[\sum_{j=0}^{\infty} h(j)z^{-j} \right] \quad (2.31)$$

Therefore,

$$Y(z) = H(z)U(z) \quad (2.32)$$

The function $H(z)$ is known as the **z -transfer function** or simply the **transfer function**. It plays an important role in obtaining the response of an LTI system to any input, as explained later. Note that the transfer function and impulse response sequence are z -transform pairs. The following example demonstrates this relationship.

Applying the convolution theorem to the response of an LTI system allows us to use the z -transform to find the output of a system without convolution as follows:

1. z -transform the input.
2. Multiply the z -transform of the input and the z -transfer function.
3. Inverse z -transform to obtain the output temporal sequence.

An added advantage of this approach is that the output can often be obtained in closed form. The preceding procedure is demonstrated in Example 2.21.

Example 2.20

Given the discrete-time system

$$y(k+1) - 0.5y(k) = u(k), y(0) = 0$$

find the impulse response of the system $h(k)$:

1. From the difference equation
2. Using z-transformation

Solution

1. By definition, the impulse response is due to the input. From the difference equation, we obtain the output sequence

$$\begin{aligned} y(1) &= 1 \\ y(2) &= 0.5y(1) = 0.5 \\ y(3) &= 0.5y(2) = (0.5)^2 \\ \text{i.e., } h(i) &= \begin{cases} (0.5)^{i-1}, & i = 1, 2, 3, \dots \\ 0, & i < 1 \end{cases} \end{aligned}$$

2. Alternatively, z-transforming the difference equation yields the transfer function

$$H(z) = \frac{Y(z)}{U(z)} = \frac{1}{z - 0.5}$$

Inverse-transforming with the delay theorem gives the impulse response

$$h(i) = \begin{cases} (0.5)^{i-1}, & i = 1, 2, 3, \dots \\ 0, & i < 1 \end{cases}$$

Even in this simple example, it is easier to obtain the impulse response from the transfer function than from the difference equation. For higher order difference equations, the latter becomes difficult if not impossible while the former is quite simple. Observe that the response decays exponentially because the pole has magnitude less than unity. In Chapter 4, we discuss this property and relate to the stability of discrete-time systems.

Example 2.21

Given the discrete-time system

$$y(k+1) - y(k) = u(k+1)$$

find the system transfer function and its response to a sampled unit step.

Example 2.21—cont'd**Solution**

The transfer function corresponding to the difference equation is

$$H(z) = \frac{z}{z - 1}$$

We multiply the transfer function by the sampled unit step's z-transform to obtain

$$Y(z) = \left(\frac{z}{z - 1}\right) \times \left(\frac{z}{z - 1}\right) = \left(\frac{z}{z - 1}\right)^2 = z \frac{z}{(z - 1)^2}$$

The z-transform of a unit ramp is

$$F(z) = \frac{z}{(z - 1)^2}$$

Then, using the time advance property of the z-transform, we have the inverse transform

$$y(i) = \begin{cases} k + 1, & k = 0, 1, 2, 3, \dots \\ 0, & k < 0 \end{cases}$$

It is obvious from Example 2.21 that z-transforming yields the response of a system in closed form more easily than direct evaluation. For higher-order difference equations, obtaining the response in closed form directly may be impossible, whereas z-transforming to obtain the response remains a relatively simple task.

2.7 The modified z-transform

Sampling and z-transformation capture the values of a continuous-time function at the sampling points only. To evaluate the time function between sampling points, we need to delay the sampled waveform by a fraction of a sampling interval before sampling. We can then vary the sampling points by changing the delay period. The z-transform associated with the delayed waveform is known as the modified z-transform.

We consider a causal continuous-time function $y(t)$ sampled every T seconds. Next, we insert a delay $T_d < T$ before the sampler as shown in Fig. 2.8. The output of the delay element is the waveform

$$y_d(t) = \begin{cases} y(t - T_d), & t \geq 0 \\ 0, & t < 0 \end{cases} \quad (2.33)$$

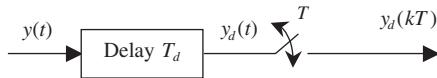


Figure 2.8
Sampling of a delayed signal.

Note that delaying a causal sequence always results in an initial zero value. To avoid inappropriate initial values, we rewrite the delay as

$$\begin{aligned} T_d &= T - mT, \quad 0 \leq m < 1 \\ m &= 1 - T_d/T \end{aligned} \tag{2.34}$$

If $y_{-1}(t + mT)$ is defined as $y(t + mT)$ delayed by one complete sampling period, then, based on (2.33), $y_d(t)$ is given by

$$y_d(t) = y(t - T + mT) = y_{-1}(t + mT) \tag{2.35}$$

We now sample the delayed waveform with sampling period T to obtain

$$y_d(kT) = y_{-1}(kT + mT), \quad k = 0, 1, 2, \dots \tag{2.36}$$

For example, a time delay of 0.03 s with a sampling period T of 0.1 s corresponds to $m = 0.7$ —that is, a time advance of 0.7 of a sampling period and a time delay of one sampling period. We also have

$$y_d(t) = y(t - 1 + 0.07) = y_{-1}(t + 0.07)$$

Sampling the delayed waveform with period 0.1 s gives

$$y_d(0.1k) = y_{-1}(0.1k + 0.07), \quad k = 0, 1, 2, \dots$$

From the delay theorem, we know the z -transform of $y_{-1}(t)$

$$Y_{-1}(z) = z^{-1} Y(z) \tag{2.37}$$

We need to determine the effect of the time advance by mT to obtain the z -transform of $y_d(t)$. We determine this effect by considering specific examples. At this point, it suffices to write

$$Y(z, m) = \mathcal{Z}_m\{y(kT)\} = z^{-1} \mathcal{Z}\{y(kT + mT)\} \tag{2.38}$$

where $\mathcal{Z}_m\{\bullet\}$ denotes the modified z -transform.

Example 2.22**Step**

The step function has fixed amplitude for all time arguments. Thus, shifting it or delaying it does not change the sampled values. We conclude that the modified z-transform of a sampled step is the same as its z-transform, times z^{-1} for all values of the time advance mT —that is, $1/(1 - z^{-1})$.

Example 2.23**Exponential**

We consider the exponential waveform

$$y(t) = e^{-pt} \quad (2.39)$$

The effect of a time advance mT on the sampled values for an exponential decay is shown in Fig. 2.9. The sampled values are given by

$$y(kT + mT) = e^{-p(k+m)T} = e^{-pmT} e^{-pkT}, \quad k = 0, 1, 2, \dots \quad (2.40)$$

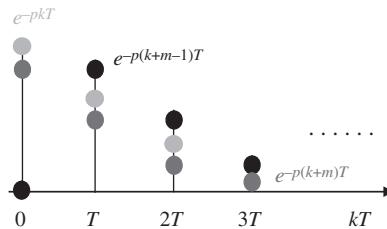


Figure 2.9

Effect of time advance on sampling an exponential decay.

We observe that the time advance results in a scaling of the waveform by the factor e^{-pmT} . By the linearity of the z-transform, we have the following:

$$\mathcal{Z}\{y(kT + mT)\} = e^{-pmT} \frac{z}{z - e^{-pT}} \quad (2.41)$$

Using (2.38), we have the modified z-transform

$$Y(z, m) = \frac{e^{-pmT}}{z - e^{-pT}} \quad (2.42)$$

For example, for $p = 4$ and $T = 0.2$ s, to delay by $0.7 T$, we let $m = 0.3$ and calculate $e^{-pmT} = e^{-0.24} = 0.787$ and $e^{-pT} = e^{-0.8} = 0.449$. We have the modified z-transform

$$Y(z, m) = \frac{0.787}{z - 0.449}$$

The modified z-transforms of other important functions, such as the ramp and the sinusoid, can be obtained following the procedure presented earlier. The derivations of these modified z-transforms are left as exercises.

2.8 Frequency response of discrete-time systems

In this section, we discuss the steady-state response of a discrete-time system to a sampled sinusoidal input.³ It is shown that, as in the continuous-time case, the response is a sinusoid of the same frequency as the input with frequency-dependent phase shift and magnitude scaling. The scale factor and phase shift define a complex function of frequency known as the frequency response.

We first obtain the frequency response using impulse sampling and the Laplace transform to exploit the well-known relationship between the transfer function $H_a(s)$ and the frequency response $H_a(j\omega)$.

$$H_a(j\omega) = H_a(s)|_{s=j\omega} \quad (2.43)$$

The impulse-sampled representation of a discrete-time waveform is

$$u^*(t) = \sum_{k=0}^{\infty} u(kT) \delta(t - kT) \quad (2.44)$$

where $u(kT)$ is the value at time kT , and $\delta(t - kT)$ denotes a Dirac delta at time kT . The representation (2.44) allows Laplace transformation to obtain

$$U^*(s) = \sum_{k=0}^{\infty} u(kT) e^{-kTs} \quad (2.45)$$

It is now possible to define a transfer function for sampled inputs as

$$H^*(s) = \left. \frac{Y^*(s)}{U^*(s)} \right|_{\text{zero initial conditions}} \quad (2.46)$$

Then, using (2.43), we obtain

$$H^*(j\omega) = H^*(s)|_{s=j\omega} \quad (2.47)$$

The expression for $H^*(j\omega)$ is the ratio of polynomials in the exponential $e^{j\omega T}$. To rewrite (2.47) in terms of the complex variable $z = e^{sT}$, we use the equation

$$H(z) = H^*(s)|_{s=\frac{1}{T}\ln(z)} \quad (2.48)$$

³ Unlike continuous sinusoids, sampled sinusoids are only periodic if the ratio of the period of the waveform and the sampling period is a rational number (equal to a ratio of integers). However, the continuous envelope of the sampled form is clearly always periodic. See the text by Oppenheim et al. (1997), p. 26, for more details.

Thus, the frequency response is given by

$$\begin{aligned} H^*(j\omega) &= H(z)|_{z=e^{j\omega T}} \\ &= H(e^{j\omega T}) \end{aligned} \quad (2.49)$$

For example, the transfer function

$$H^*(s) = \frac{1}{1 - e^{-sT}}$$

corresponds to the z -transfer function

$$H(z) = \frac{1}{1 - z^{-1}}$$

and the frequency response

$$H^*(j\omega) = \frac{1}{1 - e^{-j\omega T}}$$

We obtain the same frequency response by substituting $z = e^{j\omega T}$ in $H(z)$

$$H(e^{j\omega T}) = \frac{1}{1 - e^{-j\omega T}}$$

[Eq. \(2.49\)](#) can also be verified without the use of impulse sampling by considering the sampled complex exponential

$$\begin{aligned} u(kT) &= u_0 e^{jk\omega_0 T} \\ &= u_0 [\cos(k\omega_0 T) + j \sin(k\omega_0 T)], \quad k = 0, 1, 2, \dots \end{aligned} \quad (2.50)$$

This eventually yields the sinusoidal response while avoiding its second-order z -transform. The z -transform of the chosen input sequence is the first-order function

$$U(z) = u_0 \frac{z}{z - e^{j\omega_0 T}} \quad (2.51)$$

Assume the system z -transfer function to be

$$H(z) = \frac{N(z)}{\prod_{i=1}^n (z - p_i)} \quad (2.52)$$

where $N(z)$ is a numerator polynomial of order n or less, and the system poles p_i are assumed to lie inside the unit circle.

The system output due to the input of [\(2.50\)](#) has the z -transform

$$Y(z) = \left[\frac{N(z)}{\prod_{i=1}^n (z - p_i)} \right] u_0 \frac{z}{z - e^{j\omega_0 T}} \quad (2.53)$$

This can be expanded into the partial fractions

$$Y(z) = \frac{Az}{z - e^{j\omega_0 T}} + \sum_{i=1}^n \frac{B_i z}{z - p_i} \quad (2.54)$$

Then inverse z -transforming gives the output

$$Y(kT) = A e^{jk\omega_0 T} + \sum_{i=1}^n B_i p_i^k, \quad k = 0, 1, 2, \dots \quad (2.55)$$

The assumption of poles inside the unit circle implies that, for sufficiently large k , the output reduces to

$$y_{ss}(kT) = A e^{jk\omega_0 T}, \quad K \text{ large} \quad (2.56)$$

where $y_{ss}(kT)$ denotes the steady-state output.

The term A is the partial fraction coefficient

$$\begin{aligned} A &= \frac{Y(z)}{z} (z - e^{j\omega_0 T}) \Big|_{z=e^{j\omega_0 T}} \\ &= H(e^{j\omega_0 T}) u_0 \end{aligned} \quad (2.57)$$

Thus, we write the steady-state output in the form

$$y_{ss}(kT) = |H(e^{j\omega_0 T})| u_0 e^{j[k\omega_0 T + \angle H(e^{j\omega_0 T})]}, \quad k \text{ large} \quad (2.58)$$

The real part of this response is the response due to a sampled cosine input, and the imaginary part is the response of a sampled sine. The sampled cosine response is

$$y_{ss}(kT) = |H(e^{j\omega_0 T})| u_0 \cos[k\omega_0 T + \angle H(e^{j\omega_0 T})], \quad k \text{ large} \quad (2.59)$$

where we used the frequency response expression

$$H(e^{j\omega_0 T}) = |H(e^{j\omega_0 T})| \angle H(e^{j\omega_0 T}) \quad (2.60)$$

This is the frequency response function obtained earlier using impulse sampling. The sampled sine response is similar to (2.59) with the cosine replaced by sine.

Eqs. (2.58) and (2.59) show that the response to a sampled sinusoid is a sinusoid of the same frequency scaled by $|H(e^{j\omega_0 T})|$ and phase-shifted by $\angle H(e^{j\omega_0 T})$. Thus, one can use complex arithmetic to determine the steady-state response due to a sampled sinusoid without the need for z -transformation.

Example 2.24

Find the steady-state response of the system

$$H(z) = \frac{1}{(z - 0.1)(z - 0.5)}$$

due to the sampled sinusoid $u(kT) = 3 \cos(0.2k)$.

Solution

Using (2.59) gives the response

$$\begin{aligned} y_{ss}(kT) &= |H(e^{j0.2})| 3 \cos(0.2k + \angle H(e^{j0.2})), \quad k \text{ large} \\ &= \left| \frac{1}{(e^{j0.2} - 0.1)(e^{j0.2} - 0.5)} \right| 3 \cos\left(0.2k + \angle \frac{1}{(e^{j0.2} - 0.1)(e^{j0.2} - 0.5)}\right) \\ &= 6.4 \cos(0.2k - 0.614) \end{aligned}$$

2.8.1 Properties of the frequency response of discrete-time systems

Using (2.49), the following frequency response properties can be derived:

1. *DC gain*: The DC gain is equal to $H(1)$.

Proof

From (2.49),

$$\begin{aligned} H(e^{j\omega T})|_{\omega \rightarrow 0} &= H(z)|_{z \rightarrow 1} \\ &= H(1) \end{aligned}$$

2. *Periodic nature*: The frequency response is a periodic function of frequency with period $\omega_s = 2\pi/T$ rad/s.

Proof

The complex exponential

$$e^{j\omega T} = \cos(\omega T) + j \sin(\omega T)$$

is periodic with period $\omega_s = 2\pi/T$ rad/s. Because $H(e^{j\omega T})$ is a single-valued function of its argument, it follows that it also is periodic and that it has the same repetition frequency.

3. *Symmetry:* For transfer functions with real coefficients, the magnitude of the transfer function is an even function of frequency and its phase is an odd function of frequency.

Proof

For negative frequencies, the transfer function is

$$H(e^{-j\omega T}) = H\left(\overline{e^{j\omega T}}\right)$$

For real coefficients, we have

$$\overline{H(e^{j\omega T})} = H\left(\overline{e^{j\omega T}}\right)$$

Combining the last two equations gives

$$H(e^{-j\omega T}) = \overline{H(e^{j\omega T})}$$

Equivalently, we have

$$\begin{aligned}|H(e^{-j\omega T})| &= |H(e^{j\omega T})| \\ \angle H(e^{-j\omega T}) &= -\angle H(e^{j\omega T})\end{aligned}$$

Hence, it is only necessary to obtain $H(e^{j\omega T})$ for frequencies ω in the range from DC to $\omega_s/2$. The frequency response for negative frequencies can be obtained by symmetry, and for frequencies above $\omega_s/2$ the frequency response is periodically repeated. If the frequency response has negligible amplitudes at frequencies above $\omega_s/2$, the repeated frequency response cycles do not overlap. The overall effect of sampling for such systems is to produce a periodic repetition of the frequency response of a continuous-time system.

Because the frequency response functions of physical systems are not band-limited, overlapping of the repeated frequency response cycles, known as *folding*, occurs. The frequency $\omega_s/2$ is known as the *folding frequency*. Folding distorts the frequency response and should be minimized. This can be accomplished by proper choice of the sampling frequency $\omega_s/2$ or filtering. Fig. 2.10 shows the magnitude of the frequency response of a second-order underdamped digital system.

2.8.2 MATLAB commands for the discrete-time frequency response

The MATLAB commands **bode**, **nyquist**, and **nichols** calculate and plot the frequency response of a discrete-time system. For a sampling period of 0.2 s and a transfer function with numerator **num** and denominator **den**, the three commands have the form

```
>> g = tf(num, den, 0.2)
>> bode(g)
>> nyquist(g)
>> nichols(g)
```

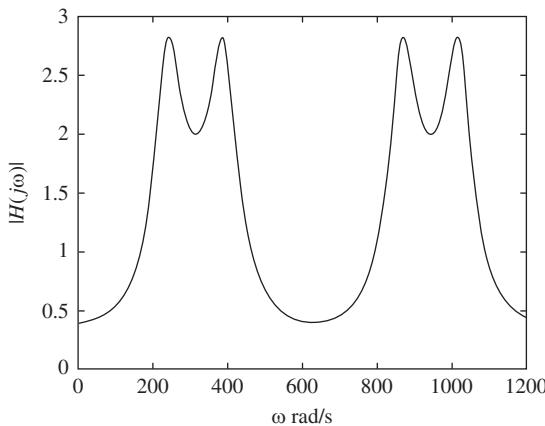


Figure 2.10
Magnitude of the frequency response of a digital system.

MATLAB limits the user's options for automatically generated plots. However, all the commands have alternative forms that allow the user to obtain the frequency response data for later plotting.

The commands **bode** and **nichols** have the alternative form

$$\begin{aligned}\gg [\mathbf{M}, \mathbf{P}, \mathbf{w}] &= \mathbf{bode}(\mathbf{g}, \mathbf{w}) \\ \gg [\mathbf{M}, \mathbf{P}, \mathbf{w}] &= \mathbf{nichols}(\mathbf{g}, \mathbf{w})\end{aligned}$$

where **w** is a predefined frequency grid, **M** is the magnitude, and **P** is the phase of the frequency response. MATLAB selects the frequency grid if none is given and returns the same list of outputs. The frequency vector can also be eliminated from the output or replaced by a scalar for single-frequency computations. The command **nyquist** can take similar forms to those just described but yields the real and imaginary parts of the frequency response as follows:

$$\gg [\mathbf{Real}, \mathbf{Imag}, \mathbf{w}] = \mathbf{nyquist}(\mathbf{g}, \mathbf{w})$$

Another useful command that calculates the frequency response is **evalfr** as in the following example

```
>> g=zpk([], [0.1, 0.5], 1, 0.02) % Transfer function
>> z1= exp(j*0.04) % w=2 rad/s, wT=0.04
>> f_resp = evalfr(g, z1) % Evaluate at z=z1
```

Unlike frequency response commands, the second argument is the complex exponential corresponding to the frequency where the frequency response is evaluated and not the frequency.

As with all MATLAB commands, printing the output is suppressed if any of the frequency response commands are followed by a semicolon. The output can then be used with the command **plot** to obtain user-selected plot specifications. For example, a plot of the actual frequency response points without connections is obtained with the command

```
>> plot(Real(:), Imag(:), '*')
```

where the locations of the data points are indicated with the '*' . The command

```
>> subplot(2,3,4)
```

creates a 2-row, 3-column grid and draw axes at the first position of the second row (the first three plots are in the first row, and **4** is the plot number). The next plot command superimposes the plot on these axes. For other plots, the subplot and plot commands are repeated with the appropriate arguments. For example, a plot in the first row and second column of the grid is obtained with the command

```
>> subplot(2,3,2)
```

2.9 The sampling theorem

Sampling is necessary for the processing of analog data using digital elements. Successful digital data processing requires that the samples reflect the nature of the analog signal and that analog signals be recoverable, at least in theory, from a sequence of samples.

[Fig. 2.11](#) shows two distinct waveforms with identical samples. Obviously, faster sampling of the two waveforms would produce distinguishable sequences. Thus, it is obvious that sufficiently fast sampling is a prerequisite for successful digital data processing. The sampling theorem gives a lower bound on the sampling rate necessary for a given **band-limited** signal (i.e., a signal with a known finite bandwidth).

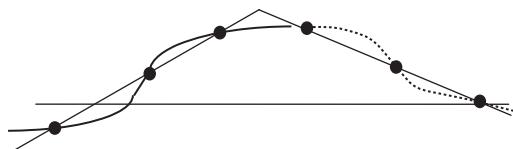


Figure 2.11
Two different waveforms with identical samples.

Theorem 2.4: The sampling theorem

The band-limited signal with bandwidth ω_m

$$\begin{aligned} f(t) &\xleftrightarrow{\mathcal{F}} F(j\omega), F(j\omega) \neq 0, & -\omega_m \leq \omega \leq \omega_m \\ F(j\omega) &= 0, & \text{elsewhere} \end{aligned} \quad (2.61)$$

Theorem 2.4: The sampling theorem—cont'd

and with \mathcal{F} denoting the Fourier transform, can be reconstructed from the discrete-time waveform

$$f^*(t) = \sum_{k=-\infty}^{\infty} f(t)\delta(t - kT) \quad (2.62)$$

if and only if the sampling angular frequency $\omega_s = 2\pi/T$ satisfies the condition

$$\omega_s > 2\omega_m \quad (2.63)$$

The spectrum of the continuous-time waveform can be recovered using an ideal low-pass filter of bandwidth ω_b in the range

$$\omega_m < \omega_b < \omega_s/2 \quad (2.64)$$

Proof

Consider the unit impulse train

$$\delta_T(t) = \sum_{k=-\infty}^{\infty} \delta(t - kT) \quad (2.65)$$

and its Fourier transform

$$\delta_T(\omega) = \frac{2\pi}{T} \sum_{n=-\infty}^{\infty} \delta(\omega - n\omega_s) \quad (2.66)$$

Impulse sampling is achieved by multiplying the waveforms $f(t)$ and $\delta_T(t)$. By the frequency convolution theorem, the spectrum of the product of the two waveforms is given by the convolution of their two spectra; that is,

$$\begin{aligned} \mathcal{F}\{\delta_T(t) \times f(t)\} &= \frac{1}{2\pi} \delta_T(j\omega)^* F(j\omega) \\ &= \left[\frac{1}{T} \sum_{n=-\infty}^{\infty} \delta(\omega - n\omega_s) \right] * F(j\omega) \\ &= \frac{1}{T} \sum_{n=-\infty}^{\infty} F(\omega - n\omega_s) \end{aligned}$$

Therefore, the spectrum of the sampled waveform is a periodic function of the sampling frequency ω_s . Assuming that $f(t)$ is a real valued function, then it is well known that the magnitude $|F(j\omega)|$ is an even function of frequency, whereas the phase $\angle F(j\omega)$ is an odd function. For a band-limited function, the amplitude and phase in the frequency range 0 to $\omega_s/2$ can be recovered by an ideal low-pass filter as shown in Fig. 2.12.

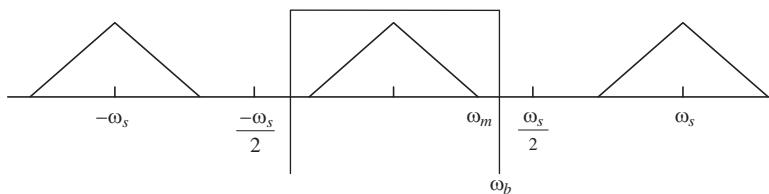


Figure 2.12
Sampling theorem.

2.9.1 Selection of the sampling frequency

In practice, finite bandwidth is an idealization associated with infinite-duration signals, whereas finite duration implies infinite bandwidth. To show this, assume that a given signal is to be band-limited. Band limiting is equivalent to multiplication by a pulse in the frequency domain. By the convolution theorem, multiplication in the frequency domain is equivalent to convolution of the inverse Fourier transforms. Hence, the inverse transform of the band-limited function is the convolution of the original time function with the sinc function, a function of infinite duration. We conclude that a band-limited function is of infinite duration.

A time-limited function is the product of a function of infinite duration and a pulse. The frequency convolution theorem states that multiplication in the time domain is equivalent to convolution of the Fourier transforms in the frequency domain. Thus, the spectrum of a time-limited function is the convolution of the spectrum of the function of infinite duration with a sinc function, a function of infinite bandwidth. Hence, the Fourier transform of a time-limited function has infinite bandwidth. Because all measurements are made over a finite time period, infinite bandwidths are unavoidable. Nevertheless, a given signal often has a finite “effective bandwidth” beyond which its spectral components are negligible. This allows us to treat physical signals as band limited and choose a suitable sampling rate for them based on the sampling theorem.

In practice, the sampling rate chosen is often larger than the lower bound specified in the sampling theorem. A rule of thumb is to choose ω_s as

$$\omega_s = k\omega_m, \quad 5 \leq k \leq 10 \quad (2.67)$$

The choice of the constant k depends on the application. In many applications, the upper bound on the sampling frequency is either prohibitively expensive or well below the capabilities of state-of-the-art hardware. A closed-loop control system cannot have a sampling period below the minimum time required for the output measurement; that is, the

sampling frequency is upper-bounded by the **sensor delay**.⁴ For example, oxygen sensors used in automotive air/fuel ratio control have a sensor delay of about 20 ms, which corresponds to a sampling frequency upper bound of 50 Hz. Another limitation is the computational time needed to update the control. This is becoming less restrictive with the availability of faster microprocessors but must be considered in sampling rate selection.

In digital control, the sampling frequency must be chosen so that samples provide a good representation of the analog physical variables. A more detailed discussion of the practical issues that must be considered when choosing the sampling frequency is given in Chapter 12. Here, we only discuss choosing the sampling period based on the sampling theorem.

For a linear system, the output of the system has a spectrum given by the product of the frequency response and input spectrum. Because the input is not known *a priori*, we must base our choice of sampling frequency on the frequency response using appropriate rules of thumb. These rules guide our choice of sampling period but, in some cases, lower sampling rates may be acceptable in practice.

The amplitude of the frequency response becomes negligibly small at about 7 times the bandwidth of the system. Thus, higher frequencies will be attenuated by the system and can be neglected. We can then design the system assuming that its output signal has a bandwidth of about 7 times the bandwidth of the system, then use the rule of thumb for a signal.

The frequency response of a first-order system is

$$H(j\omega) = \frac{K}{j\omega/\omega_b + 1} \quad (2.68)$$

where K is the DC gain and ω_b is the system bandwidth. The frequency response amplitude drops below the DC level by a factor of about 10 at the frequency $7\omega_b$. If we consider $\omega_m = 7\omega_b$, the sampling frequency is chosen as

$$\omega_s = k\omega_b, \quad 35 \leq k \leq 70 \quad (2.69)$$

For a second-order system with frequency response

$$H(j\omega) = \frac{K}{j2\zeta\omega/\omega_n + 1 - (\omega/\omega_n)^2} \quad (2.70)$$

the bandwidth of the system is approximated by the damped natural frequency

$$\omega_d = \omega_n \sqrt{1 - \zeta^2} \quad (2.71)$$

⁴ It is possible to have the sensor delay as an integer multiple of the sampling period if a state estimator is used, as discussed in Franklin et al. (1998).

Using a frequency of $7\omega_d$ as the maximum significant frequency, we choose the sampling frequency as

$$\omega_s = k\omega_d, \quad 35 \leq k \leq 70 \quad (2.72)$$

In addition, the impulse response of a second-order system is of the form

$$y(t) = Ae^{-\zeta\omega_n t} \sin(\omega_d t + \phi) \quad (2.73)$$

where A is a constant amplitude and ϕ is a phase angle. Thus, the choice of sampling frequency of (2.72) is sufficiently fast for oscillations of frequency ω_d and time to first peak π/ω_d .

Example 2.25

Given a first-order system of bandwidth 10 rad/s, select a suitable sampling frequency and find the corresponding sampling period.

Solution

A suitable choice of sampling frequency is $\omega_s = 60$, $\omega_d = 600$ rad/s. The corresponding sampling period is approximately $T = 2\pi/\omega_s \approx 0.01$ s.

Example 2.26

A closed-loop control system must be designed for a steady-state error not to exceed 5%, a damping ratio of about 0.7, and an undamped natural frequency of 10 rad/s. Select a suitable sampling period for the system if the system has a sensor delay of

1. 0.02 s
2. 0.03 s

Solution

Let the sampling frequency be

$$\begin{aligned}\omega_s &\geq 35\omega_d \\ &= 35\omega_n \sqrt{1 - \zeta^2} \\ &= 350 \sqrt{1 - 0.49} \\ &= 249.95 \text{ rad/s}\end{aligned}$$

The corresponding sampling period is $T = 2\pi/\omega_s \leq 0.025$ s.

1. A suitable choice is $T = 20$ ms because this is equal to the sensor delay.
2. We are forced to choose $T = 30$ ms, which is equal to the sensor delay.

The situation of the above example occurs in applications where two sensors are available but the faster sensor is more expensive and the less expensive sensor is preferred to reduce the cost of the product.

Resources

- Chen, C.-T., 1989. System and Signal Analysis. Saunders, New York.
- Feuer, A., Goodwin, G.C., 1996. Sampling in Digital Signal Processing and Control. Birkhauser, Boston.
- Franklin, G.F., Powell, J.D., Workman, M.L., 1998. Digital Control of Dynamic Systems. Addison-Wesley, Menlo Park, CA.
- Goldberg, S., 1986. Introduction to Difference Equations. Dover, Mineola, NY.
- Jacquot, R.G., 1981. Modern Digital Control Systems. Marcel Dekker, New York.
- Kuo, B.C., 1992. Digital Control Systems. Saunders, Ft. Worth, TX.
- Mickens, R.E., 1987. Difference Equations. Van Nostrand Reinhold, New York.
- Oppenheim, A.V., Willsky, A.S., Nawab, S.H., 1997. Signals and Systems. Prentice Hall, Englewood Cliffs, NJ.

Problems

- 2.1 Derive the discrete-time model of Example 2.1 from the solution of the system differential equation with initial time kT and final time $(k + 1)T$.
- 2.2 For each of the following equations, determine the order of the equation and then test it for (i) linearity, (ii) time invariance, and (iii) homogeneity.
 - a. $y(k + 2) = y(k + 1)y(k) + u(k)$
 - b. $y(k + 3) + 2y(k) = 0$
 - c. $y(k + 4) + y(k - 1) = u(k)$
 - d. $y(k + 5) = y(k + 4) + u(k + 1) - u(k)$
 - e. $y(k + 2) = y(k)u(k)$
- 2.3 Find the transforms of the following sequences using Definition 2.1.
 - a. $\{0, 1, 2, 4, 0, 0, \dots\}$
 - b. $\{0, 0, 0, 1, 1, 1, 0, 0, 0, \dots\}$
 - c. $\{0, 2^{-0.5}, 1, 2^{-0.5}, 0, 0, 0, \dots\}$
- 2.4 Obtain closed forms of the transforms of Problem 2.3 using the table of z -transforms and the time delay property.
- 2.5 Prove the linearity and time delay properties of the z -transform from basic principles.
- 2.6 Use the linearity of the z -transform and the transform of the exponential function to obtain the transforms of the discrete-time functions.
 - a. $\sin(k\omega T)$
 - b. $\cos(k\omega T)$

- 2.7 Use the multiplication by exponential property to obtain the transforms of the discrete-time functions.
- $e^{-\alpha k} T \sin(k\omega T)$
 - $e^{-\alpha k} T \cos(k\omega T)$
- 2.8 Find the inverse transforms of the following functions using Definition 2.1 and, if necessary, long division.
- $F(z) = 1 + 3z^{-1} + 4z^{-2}$
 - $F(z) = 5z^{-1} + 4z^{-5}$
 - $F(z) = \frac{z}{z^2 + 0.3z + 0.02}$
 - $F(z) = \frac{z - 0.1}{z^2 + 0.04z + 0.25}$
- 2.9 For Problems 2.8(c) and (d), find the inverse transforms of the functions using partial fraction expansion and table lookup.
- 2.10 Use the complex differentiation property and the transform of the unit ramp to obtain the z -transform of the parabolic

$$f(k) = \begin{cases} k^2, & k \geq 0 \\ 0, & k < 0 \end{cases}$$

- 2.11 Solve the following difference equations.
- $y(k+1) - 0.8 y(k) = 0, y(0) = 1$
 - $y(k+1) - 0.8 y(k) = 1(k), y(0) = 0$
 - $y(k+1) - 0.8 y(k) = 1(k), y(0) = 1$
 - $y(k+2) + 0.7 y(k+1) + 0.06 y(k) = \delta(k), y(0) = 0, y(1) = 2$
- 2.12 Find the transfer functions corresponding to the difference equations of Problem 2.2 with input $u(k)$ and output $y(k)$. If no transfer function is defined, explain why.
- 2.13 The following difference equation describes the evolution of the expected price of a commodity⁵

$$p_e(k+1) = (1 - \gamma) p_e(k) + \gamma p(k)$$

where $p_e(k)$ is the expected price after k quarters, $p(k)$ is the actual price after k quarters, and γ is a constant in the range $0 < \gamma < 1$.

- Obtain the transfer function of the system with input $p(k)$ and output $p_e(k)$
- Assuming a zero initial estimate, obtain the price estimate using the transfer function (i) for a fixed actual price of one unit, (ii) for an exponentially decaying price $p(k) = 0.1 + (b)^k$.

⁵ Gujarate, D.N., 1988. Basic Econometrics. McGraw-Hill, NY, p. 547.

- 2.14 The Fibonacci sequence is a set of numbers generated by the difference equation

$$f(k+2) = f(k+1) + f(k), f(0) = 0, f(1) = 1$$

The sequence {0,1,1,2,3,...} describes many phenomena in nature. Show that the sequence is given by

$$f(k) = \frac{1}{\sqrt{5}} \left\{ \left(\frac{1 + \sqrt{5}}{2} \right)^k - \left(\frac{1 - \sqrt{5}}{2} \right)^k \right\}$$

The number $\phi = (1 + \sqrt{5})/2 = 1.618$ is known as the golden ratio. Show that it is the positive solution of the equation

$$\frac{\phi + 1}{\phi} = \frac{\phi}{1}$$

- 2.15 Test the linearity with respect to the input of the systems for which you found transfer functions in Problem 2.12.
- 2.16 If the rational functions of Problems 2.8(c) and (d) are transfer functions of LTI systems, find the difference equation governing each system.
- 2.17 We can use z -transforms to find the sum of integers raised to various powers. This is accomplished by first recognizing that the sum is the solution of the difference equation

$$f(k) = f(k-1) + a(k)$$

where $a(k)$ is the k^{th} term in the summation.

Show that the z -transform of the sum is given by

$$F(z) = \frac{z}{z-1} A(z)$$

then evaluate the following summations

a. $\sum_{k=1}^n k$

b. $\sum_{k=1}^n k^2$

- 2.18 Given the discrete-time system

$$y(k+2) - y(k) = 2u(k)$$

find the impulse response of the system $g(k)$:

- From the difference equation
- Using z -transformation

- 2.19 The following identity provides a recursion for the cosine function:

$$\cos(kx) = 2 \cos((k-1)x) \cos(x) - \cos((k-2)x), \quad k \text{ integer}$$

To verify its validity, let $f(k) = \cos(kx)$ and rewrite the expression as a difference equation. Show that the solution of the difference equation is indeed

$$f(k) = \cos(kx)$$

- 2.20 Repeat Problem 2.19 for the identity

$$\sin(kx) = 2 \sin((k-1)x) \cos(x) - \sin((k-2)x), \quad k \text{ integer}$$

- 2.21 Find the impulse response functions for the systems governed by the following difference equations.

- a. $y(k+1) - 0.5 y(k) = u(k)$
- b. $y(k+2) - 0.1 y(k+1) + 0.8 y(k) = u(k)$

- 2.22 Find the final value for the functions if it exists.

- a. $F(z) = \frac{z}{z^2 - 1.2z + 0.2}$
- b. $F(z) = \frac{z}{z^2 - 0.3z + 2}$

- 2.23 Find the steady-state response of the systems resulting from the sinusoidal input $u(k) = 0.5 \sin(0.4k)$.

- a. $H(z) = \frac{z}{z-0.4}$
- b. $H(z) = \frac{z}{z^2 - 0.4z + 0.03}$

- 2.24 Find the frequency response of a noncausal system whose impulse response sequence is given by

$$\{u(k), u(k) = u(k+K), k = -\infty, \dots, \infty\}$$

Hint: The impulse response sequence is periodic with period K and can be expressed as

$$u^*(t) = \sum_{l=0}^{K-1} \sum_{m=-\infty}^{\infty} u(l+mK) \delta(t-l-mK)$$

- 2.25 The well-known Shannon reconstruction theorem states that any band-limited signal $u(t)$ with bandwidth $\omega_s/2$ can be exactly reconstructed from its samples at a rate $\omega_s = 2\pi/T$. The reconstruction is given by

$$u(t) = \sum_{k=-\infty}^{\infty} u(k) \frac{\sin\left[\frac{\omega_s}{2}(t-kT)\right]}{\frac{\omega_s}{2}(t-kT)}$$

Use the convolution theorem to justify the preceding expression.

- 2.26 Obtain the convolution of the two sequences {1, 1, 1} and {1, 2, 3}.
- Directly
 - Using z -transformation
- 2.27 Obtain the modified z -transforms for the functions of Problems (2.6) and (2.7).
- 2.28 Using the modified z -transform, examine the intersample behavior of the functions $h(k)$ of Problem 2.21. Use delays of (1) $0.3T$, (2) $0.5T$, and (3) $0.8T$. Attempt to obtain the modified z -transform for Problem 2.22 and explain why it is not defined.
- 2.29 The following open-loop systems are to be digitally feedback-controlled. Select a suitable sampling period for each if the closed-loop system is to be designed for the given specifications.
- $G_{ol}(s) = \frac{1}{s+3}$ Time constant = 0.1 s
 - $G_{ol}(s) = \frac{1}{s^2+4s+3}$ Undamped natural frequency = 5 rad/s, damping ratio = 0.7
- 2.30 Repeat Problem 2.29 if the systems have the following sensor delays.
- 0.025 s
 - 0.03 s

Computer exercises

- 2.31 Consider the closed-loop system of Problem 2.29(a).
- Find the impulse response of the **closed-loop** transfer function, and obtain the impulse response sequence for a sampled system output.
 - Obtain the z -transfer function by z -transforming the impulse response sequence.
 - Using MATLAB, obtain the frequency response plots for sampling frequencies $\omega_s = k\omega_b$, $k = 5, 35, 70$.
 - Comment on the choices of sampling periods of part (c).
- 2.32 Repeat Problem 2.31 for the second-order closed-loop system of Problem 2.29(b) with plots for sampling frequencies $\omega_s = k\omega_d$, $k = 5, 35, 70$.
- 2.33 Use MATLAB with a sampling period of 1 s to verify the results of Problem 2.23. Simulate the system for 300 s then change the axes to display the last 50 s only.
- 2.34 Consider the model of evolution of the expected price of a commodity of Problem 2.13

$$p_e(k+1) = (1 - \gamma)p_e(k) + \gamma p(k)$$

where $p_e(k)$ is the expected price after k quarters, $p(k)$ is the actual price after k quarters, and γ is a constant in the range $0 < \gamma < 1$.

- Simulate the system with $\gamma = 0.5$ and a fixed actual price of one unit, and plot the actual and expected prices. Discuss the accuracy of the model prediction.

- b. Repeat part (a) for an exponentially decaying price $p(k) = (0.4)k$.
 - c. Discuss the predictions of the model referring to your simulation results.
- 2.35 Write a MATLAB program that calculates the step response of the system of Example 2.20 by using the step response coefficients and [Eq. \(2.28\)](#).

Modeling of digital control systems

Objectives

After completing this chapter, the reader will be able to do the following:

1. Obtain the transfer function of an analog system with analog-to-digital and digital-to-analog converters, including systems with a time delay.
2. Find the closed-loop transfer function for a digital control system.
3. Find the steady-state tracking error for a closed-loop control system.
4. Find the steady-state error caused by a disturbance input for a closed-loop control system.
5. Simulate a digital control system using SIMULINK.
6. Perform a sensitivity analysis to assess the effect of parameter variations on a system.

As in the case of analog control, mathematical models are needed for the analysis and design of digital control systems. A common configuration for digital control systems is shown in Fig. 3.1. The configuration includes a digital-to-analog converter (DAC), an analog subsystem, and an analog-to-digital converter (ADC). The DAC converts numbers calculated by a microprocessor or computer into analog electrical signals that can be amplified and used to control an analog plant. The analog subsystem includes the plant as well as the amplifiers and actuators necessary to drive it. The output of the plant is periodically measured and converted to a number that can be fed back to the computer using an ADC. In this chapter, we develop models for the various components of this digital control configuration. Many other configurations that include the same components can be similarly analyzed. We begin by developing models for the ADC and DAC, then for the combination of DAC, analog subsystem, and ADC.

Chapter Outline

- 3.1 Analog-to-digital converter (ADC) model 62
- 3.2 Digital-to-analog converter (DAC) model 63
- 3.3 The transfer function of the zero-order hold (ZOH) 63
- 3.4 Effect of the sampler on the transfer function of a cascade 65
- 3.5 DAC, analog subsystem, and analog-to-digital converter (ADC) combination transfer function 68

3.6 Systems with transport lag	76
3.7 The closed-loop transfer function	78
3.8 Analog disturbances in a digital system	80
3.9 Steady-state error and error constants	82
3.9.1 Sampled step input	84
3.9.2 Sampled ramp input	84
3.10 MATLAB commands	86
3.10.1 MATLAB	87
3.10.2 Simulink	87
3.11 Sensitivity analysis	93
3.11.1 Pole sensitivity	95
Further reading	97
Problems	97
Computer exercises	101

3.1 Analog-to-digital converter (ADC) model

Assume that

- ADC outputs are exactly equal in magnitude to their inputs (i.e., quantization errors are negligible).
- The ADC yields a digital output instantaneously.
- Sampling is perfectly uniform (i.e., occurs at a fixed rate).

Then the ADC can be modeled as an ideal sampler with sampling period T as shown in Fig. 3.2.

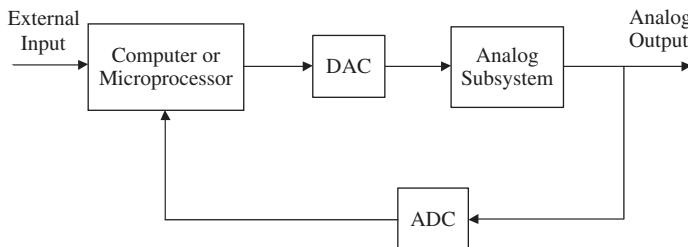


Figure 3.1
Common digital control system configuration.

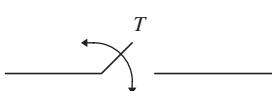


Figure 3.2
Ideal sampler model of an ADC.

Clearly, the preceding assumptions are idealizations that can only be approximately true in practice. Quantization errors are typically small but nonzero; variations in sampling rate occur but are negligible, and physical ADCs have a finite conversion time. Nevertheless, the ideal sampler model is acceptable for most engineering applications.

3.2 Digital-to-analog converter (DAC) model

Assume that

- DAC outputs are exactly equal in magnitude to their inputs.
- The DAC yields an analog output instantaneously.
- DAC outputs are constant over each sampling period.

Then the input-output relationship of the DAC is given by

$$\{u(k)\} \xrightarrow{\text{ZOH}} u(t) = u(k), \quad kT \leq t < (k+1)T, \quad k = 0, 1, 2, \dots \quad (3.1)$$

where $\{u(k)\}$ is the input sequence. This equation describes a **zero-order hold** (ZOH), shown in Fig. 3.3. Other functions may also be used to construct an analog signal from a sequence of numbers. For example, a **first-order hold** constructs analog signals in terms of straight lines, whereas a **second-order hold** constructs them in terms of parabolas.

In practice, the DAC requires a short but nonzero interval to yield an output; its output is not exactly equal in magnitude to its input and may vary slightly over a sampling period. But the model of (3.1) is sufficiently accurate for most engineering applications. The zero-order hold is the most commonly used DAC model and is adopted in most digital control texts. Analyses involving other hold circuits are similar, as seen from Problem 3.2.

3.3 The transfer function of the zero-order hold (ZOH)

To obtain the transfer function of the ZOH, we replace the number or discrete impulse shown in Fig. 3.3 by an impulse $\delta(t)$. The transfer function can then be obtained by

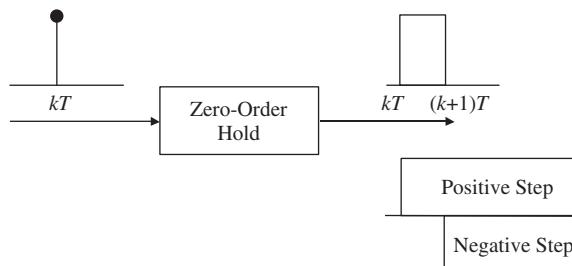


Figure 3.3
Model of a DAC as a zero-order hold.

Laplace transformation of the impulse response. As shown in the figure, the impulse response is a unit pulse of width T . A pulse can be represented as a positive step at time zero followed by a negative step at time T . Using the Laplace transform of a unit step and the time delay theorem for Laplace transforms,

$$\begin{aligned}\mathcal{L}\{\mathbf{1}(t)\} &= \frac{1}{s} \\ \mathcal{L}\{\mathbf{1}(t - T)\} &= \frac{e^{-sT}}{s}\end{aligned}\tag{3.2}$$

where $\mathbf{1}(t)$ denotes a unit step.

Thus, the transfer function of the ZOH is

$$G_{\text{ZOH}}(s) = \frac{1 - e^{-sT}}{s}\tag{3.3}$$

Next, we consider the frequency response of the ZOH:

$$G_{\text{ZOH}}(j\omega) = \frac{1 - e^{-j\omega T}}{j\omega}\tag{3.4}$$

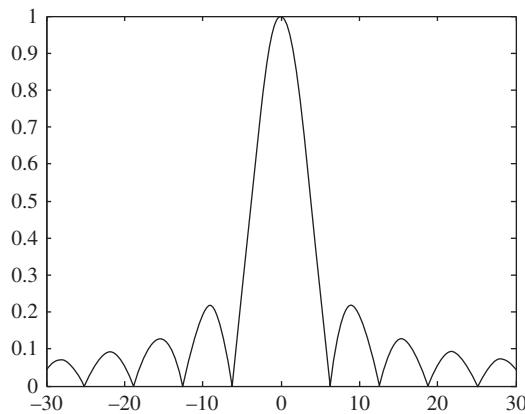
We rewrite the frequency response in the form

$$\begin{aligned}G_{\text{ZOH}}(j\omega) &= \frac{e^{-j\omega T}}{\omega} \left(\frac{e^{j\omega T} - e^{-j\omega T}}{j} \right) \\ &= \frac{e^{-j\omega T}}{\omega} \left(2 \sin\left(\omega \frac{T}{2}\right) \right) = T e^{-j\omega T} \frac{\sin\left(\omega \frac{T}{2}\right)}{\omega \frac{T}{2}}\end{aligned}$$

We now have

$$|G_{\text{ZOH}}(j\omega)| \angle G_{\text{ZOH}}(j\omega) = T \left| \text{sinc}\left(\frac{\omega T}{2}\right) \right| \angle -\frac{\omega T}{2}, -\frac{2\pi}{T} < \omega < \frac{2\pi}{T}\tag{3.5}$$

In the frequency range of interest where the sinc function is positive, the angle of frequency response of the ZOH hold is seen to decrease linearly with frequency, whereas the magnitude is proportional to the **sinc** function. As shown in Fig. 3.4, the magnitude is oscillatory with its peak magnitude equal to the sampling period and occurring at the zero frequency.

**Figure 3.4**

Magnitude of the frequency response of the zero-order hold with $T = 1$ s.

3.4 Effect of the sampler on the transfer function of a cascade

In a discrete-time system including several analog subsystems in cascade and several samplers, the location of the sampler plays an important role in determining the overall transfer function. Assuming that interconnection does not change the mathematical models of the subsystems, the Laplace transform of the output of the system of Fig. 3.5 is given by

$$\begin{aligned} Y(s) &= H_2(s)X(s) \\ &= H_2(s)H_1(s)U(s) \end{aligned} \quad (3.6)$$

Inverse Laplace transforming gives the time response

$$\begin{aligned} y(t) &= \int_0^t h_2(t-\tau)x(\tau)d\tau \\ &= \int_0^t h_2(t-\tau) \left[\int_0^\tau h_1(\tau-\lambda)u(\lambda)d\lambda \right] d\tau \end{aligned} \quad (3.7)$$

Changing the order and variables of integration, we obtain

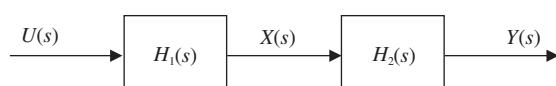


Figure 3.5
Cascade of two analog systems.

$$\begin{aligned}
y(t) &= \int_0^t u(t-\tau) \left[\int_0^\tau h_1(\tau-\lambda)h_2(\lambda)d\lambda \right] d\tau \\
&= \int_0^t u(t-\tau)h_{eq}(\tau)d\tau
\end{aligned} \tag{3.8}$$

where $h_{eq}(t) = \int_0^t h_1(t-\tau)h_2(\lambda)d\lambda$.

Thus, the equivalent impulse response for the cascade is given by the convolution of the cascaded impulse responses. The same conclusion can be reached by inverse-transforming the product of the s -domain transfer functions. The time domain expression shows more clearly that cascading results in a new form for the impulse response. So if the output of the system is sampled to obtain

$$y(iT) = \int_0^{iT} u(iT-\tau)h_{eq}(\tau)d\tau, \quad i = 1, 2, \dots \tag{3.9}$$

it is not possible to separate the three time functions that are convolved to produce it.

By contrast, convolving an impulse-sampled function $u^*(t)$ with a continuous-time signal as shown in Fig. 3.6 results in repetitions of the continuous-time function, each of which is displaced to the location of an impulse in the train. Unlike the earlier situation, the resultant time function is not entirely new, and there is hope of separating the functions that produced it. For a linear time-invariant (LTI) system with impulse-sampled input, the output is given by

$$\begin{aligned}
y(t) &= \int_0^t h(t-\tau)u^*(\tau)d\tau \\
&= \int_0^t h(t-\tau) \left[\sum_{k=0}^{\infty} u(kT)\delta(\tau-kT) \right] d\tau
\end{aligned} \tag{3.10}$$

Changing the order of summation and integration gives

$$\begin{aligned}
y(t) &= \sum_{k=0}^{\infty} u(kT) \int_0^t h(t-\tau)\delta(\tau-kT)d\tau \\
&= \sum_{k=0}^{\infty} u(kT)h(t-kT)
\end{aligned} \tag{3.11}$$

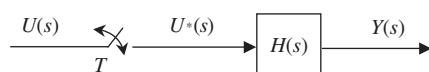


Figure 3.6
Analog system with sampled input.

Sampling the output yields the convolution summation

$$y(iT) = \sum_{k=0}^{\infty} u(kT)h(iT - kT), \quad i = 0, 1, 2, 3, \dots \quad (3.12)$$

As discussed earlier, the convolution summation has the z -transform

$$Y(z) = H(z)U(z) \quad (3.13)$$

or in s -domain notation

$$Y^*(s) = H^*(s)U^*(s) \quad (3.14)$$

If a single block is an equivalent transfer function for the cascade of Fig. 3.5, then its components cannot be separated after sampling. However, if the cascade is separated by samplers, then each block has a sampled output and input as well as a z -domain transfer function. For n blocks not separated by samplers, we use the notation

$$\begin{aligned} Y(z) &= H(z)U(z) \\ &= (H_1 H_2 \dots H_n)(z)U(z) \end{aligned} \quad (3.15)$$

as opposed to n blocks separated by samplers where

$$\begin{aligned} Y(z) &= H(z)U(z) \\ &= H_1(z)H_2(z)\dots H_n(z)U(z) \end{aligned} \quad (3.16)$$

Example 3.1

Find the equivalent sampled impulse response sequence and the equivalent z -transfer function for the cascade of the two analog systems with sampled input

$$H_1(s) = \frac{1}{s+2} \quad H_2(s) = \frac{2}{s+4}$$

1. If the systems are directly connected
2. If the systems are separated by a sampler

Solution

1. In the absence of samplers between the systems, the overall transfer function is

$$\begin{aligned} H(s) &= \frac{2}{(s+2)(s+4)} \\ &= \frac{1}{s+2} - \frac{1}{s+4} \end{aligned}$$

The impulse response of the cascade is

$$h(t) = e^{-2t} - e^{-4t}$$

and the sampled impulse response is

$$h(kT) = e^{-2kT} - e^{-4kT}, \quad k = 0, 1, 2, \dots$$

Example 3.1—cont'd

Thus, the z-domain transfer function is

$$H(z) = \frac{z}{z - e^{-2T}} - \frac{z}{z - e^{-4T}} = \frac{(e^{-2T} - e^{-4T})z}{(z - e^{-2T})(z - e^{-4T})}$$

2. If the analog systems are separated by a sampler, then each has a z-domain transfer function, and the transfer functions are given by

$$H_1(z) = \frac{z}{z - e^{-2T}} \quad H_2(z) = \frac{2z}{z - e^{-4T}}$$

The overall transfer function for the cascade is

$$H(z) = \frac{2z^2}{(z - e^{-2T})(z - e^{-4T})}$$

The partial fraction expansion of the transfer function is

$$H(z) = \frac{2}{e^{-2T} - e^{-4T}} \left[\frac{e^{-2T}z}{z - e^{-2T}} - \frac{e^{-4T}z}{z - e^{-4T}} \right]$$

Inverse z-transforming gives the impulse response sequence

$$\begin{aligned} h(kT) &= \frac{2}{e^{-2T} - e^{-4T}} [e^{-2T}e^{-2kT} - e^{-4T}e^{-4kT}] \\ &= \frac{2}{e^{-2T} - e^{-4T}} [e^{-2(k+1)T} - e^{-4(k+1)T}], \quad k = 0, 1, 2, \dots \end{aligned}$$

Example 3.1 clearly shows the effect of placing a sampler between analog blocks on the impulse responses and the corresponding z-domain transfer function.

3.5 DAC, analog subsystem, and analog-to-digital converter (ADC) combination transfer function

The cascade of a DAC, analog subsystem, and ADC, shown in Fig. 3.7, appears frequently in digital control systems (see Fig. 3.1, for example). Because both the input and the output of the cascade are sampled, it is possible to obtain its z-domain transfer function in terms of the transfer functions of the individual subsystems. The transfer function is derived using the discussion of cascades given in Section 3.4.

Using the DAC model of Section 3.3, and assuming that the transfer function of the analog subsystem is $G(s)$, the transfer function of the DAC and analog subsystem cascade is

$$\begin{aligned} G_{ZA}(s) &= G(s)G_{ZOH}(s) \\ &= (1 - e^{-sT}) \frac{G(s)}{s} \end{aligned} \tag{3.17}$$

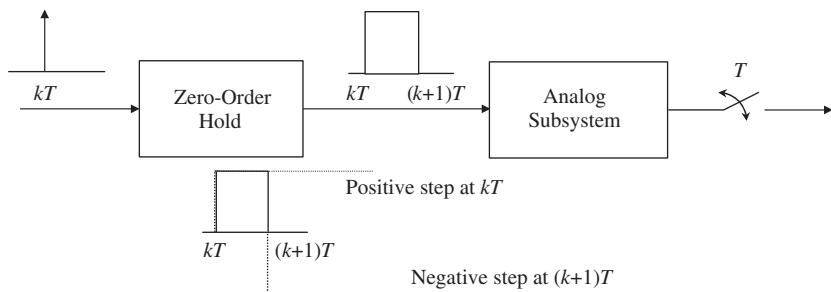


Figure 3.7
Cascade of a DAC, analog subsystem, and ADC.

The corresponding impulse response is

$$\begin{aligned} g_{ZA}(t) &= g(t)^* g_{ZOH}(t) \\ &= g_s(t) - g_s(t - T) \\ g_s(t) &= \mathcal{L}^{-1} \left\{ \frac{G(s)}{s} \right\} \end{aligned} \quad (3.18)$$

The impulse response of (3.18) is the analog system step response minus a second step response delayed by one sampling period. This response is shown in Fig. 3.8 for a second-order underdamped analog subsystem. The analog response of (3.18) is sampled to give the sampled impulse response

$$g_{ZA}(kT) = g_s(kT) - g_s(kT - T) \quad (3.19)$$

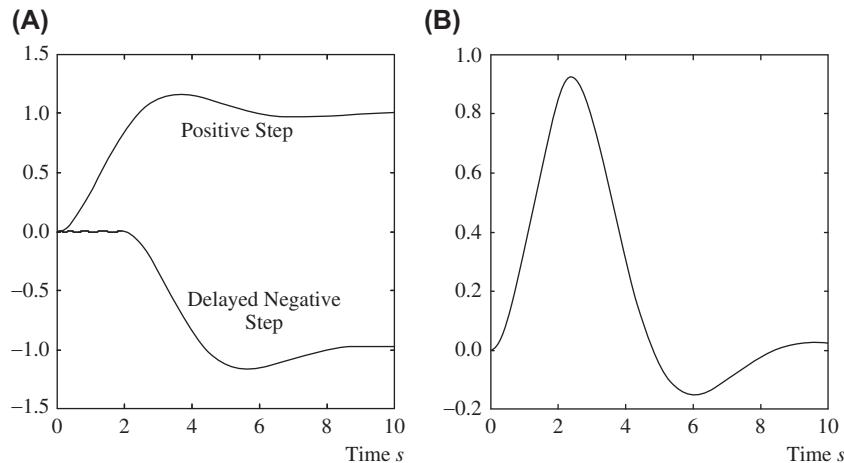


Figure 3.8
Impulse response of a DAC and analog subsystem. (A) Response of an analog system to step inputs. (B) Response of an analog system to a unit pulse input.

By z -transforming, we obtain the z -transfer function of the DAC (zero-order hold), analog subsystem, and ADC (ideal sampler) cascade

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z}\{g_s^*(t)\} \\ &= (1 - z^{-1}) \mathcal{Z}\left\{\mathcal{L}^{-1}\left[\frac{G(s)}{s}\right]^*\right\} \end{aligned} \quad (3.20)$$

The cumbersome notation in (3.20) is used to emphasize that sampling of a time function is necessary before z -transformation. Having made this point, the equation can be rewritten more concisely as

$$G_{ZAS}(z) = (1 - z^{-1}) \mathcal{Z}\left\{\frac{G(s)}{s}\right\} \quad (3.21)$$

Example 3.2

Find $G_{ZAS}(z)$ for the cruise control system for the vehicle shown in Fig. 3.9, where u is the input force, v is the velocity of the car, and b is the viscous friction coefficient.

Solution

We first draw a schematic to represent the cruise control system as shown in Fig. 3.10. Using Newton's law, we obtain the following model:

$$Mv(t) + bv(t) = u(t)$$

which corresponds to the following transfer function:

$$G(s) = \frac{V(s)}{U(s)} = \frac{1}{Ms + b}$$

We rewrite the transfer function in the form

$$G(s) = \frac{K}{\tau s + 1} = \frac{K/\tau}{s + 1/\tau}$$

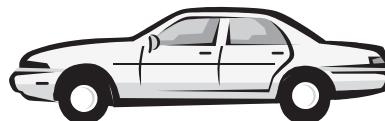


Figure 3.9
Automotive vehicle.

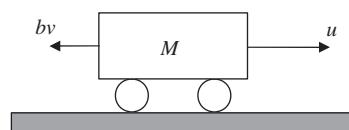


Figure 3.10
Schematic representation of a cruise control system for an automotive vehicle.

Example 3.2—cont'd

where $K = 1/b$ and $\tau = M/b$. The corresponding partial fraction expansion is

$$\frac{G(s)}{s} = \left(\frac{K}{\tau}\right) \left(\frac{A_1}{s} + \frac{A_2}{s + 1/\tau} \right)$$

where

$$A_1 = \frac{1}{s + 1/\tau} \Big|_{s=0} = \tau \quad A_2 = \frac{1}{s} \Big|_{s=-1/\tau} = -\tau$$

Using (3.21) and the z-transform table (see Appendix I), the desired z-domain transfer function is

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{L} \left\{ \left(\frac{K}{\tau} \right) \left[\frac{\tau}{s} - \frac{\tau}{s + 1/\tau} \right] \right\} \\ &= K \left[1 + \frac{z - 1}{z - e^{-T/\tau}} \right] \end{aligned}$$

We simplify to obtain the transfer function

$$G_{ZAS}(z) = K \left[\frac{2z - (1 + e^{-T/\tau})}{z - e^{-T/\tau}} \right]$$

Example 3.3

Find $G_{ZAS}(z)$ for the vehicle position control system shown in Fig. 3.10, where u is the input force, y is the position of the car, and b is the viscous friction coefficient.

Solution

As with the previous example, we obtain the following equation of motion:

$$My''(t) + by'(t) = u(t)$$

and the corresponding transfer function

$$G(s) = \frac{Y(s)}{U(s)} = \frac{1}{s(Ms + b)}$$

We rewrite the transfer function in terms of the system time constant

$$G(s) = \frac{K}{s(\tau s + 1)} = \frac{K/\tau}{s(s + 1/\tau)}$$

where $K = 1/b$ and $\tau = M/b$. The corresponding partial fraction expansion is

Example 3.3—cont'd

$$\frac{G(s)}{s} = \left(\frac{K}{\tau}\right) \left(\frac{A_{11}}{s^2} + \frac{A_{12}}{s} + \frac{A_2}{s + 1/\tau} \right)$$

where

$$A_{11} = \frac{1}{s + 1/\tau} \Big|_{s=0} = \tau \quad A_{12} = \frac{d}{ds} \left[\frac{1}{s + 1/\tau} \right] \Big|_{s=0} = -\tau^2$$

$$A_2 = \frac{1}{s^2} \Big|_{s=-1/\tau} = \tau^2$$

Using (3.21), the desired z-domain transfer function is

$$G_{ZAS}(z) = (1 - z^{-1}) \mathcal{Z} \left\{ K \left[\frac{A_{11}}{s^2} + \frac{A_{12}}{s} + \frac{A_2}{s + 1/\tau} \right] \right\}$$

$$= K \left[\frac{\frac{T}{z-1} - \tau + \frac{\tau(z-1)}{z - e^{-T/\tau}}}{z - e^{-T/\tau}} \right]$$

which can be simplified to

$$G_{ZAS}(z) = K \left[\frac{(T - \tau + \tau e^{-T/\tau})z + [\tau - e^{-T/\tau}(\tau + T)]}{(z-1)(z - e^{-T/\tau})} \right]$$

Example 3.4

Find $G_{ZAS}(z)$ for series R-L circuit shown in Fig. 3.11 with the inductor voltage as output.

Solution

Using the voltage divider rule gives

$$\frac{V_o}{V_{in}} = \frac{Ls}{R + Ls} = \frac{(L/R)s}{1 + (L/R)s} = \frac{\tau s}{1 + \tau s} \quad \tau = \frac{L}{R}$$

Hence, using (3.21), we obtain

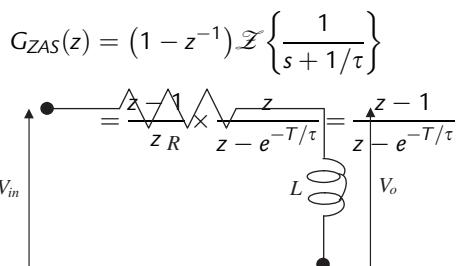


Figure 3.11
Series R-L circuit.

Example 3.5

Find the z-domain transfer function of the furnace sketched in Fig. 3.12, where the inside temperature T_i is the controlled variable, T_w is the wall temperature, and T_o is the outside temperature. Assume perfect insulation so that there is no heat transfer between the wall and the environment with heating provided through a resistor. The control variable u has the dimension of temperature scaled by an amplifier with gain K . The sampling period $T = 1$ s.

Solution

The system can be modeled with the following differential equations:

$$\begin{aligned}\dot{T}_w(t) &= g_{rw}(Ku(t) - T_w(t)) + g_{iw}(T_i(t) - T_w(t)) \\ \dot{T}_i(t) &= g_{iw}(T_w(t) - T_i(t))\end{aligned}$$

where g_{rw} and g_{iw} are the heat transfer coefficients. Laplace transforming and simplifying, we obtain the following transfer function

$$G(s) = \frac{T_i(s)}{U(s)} = \frac{Kg_{rw}g_{iw}}{s^2 + (2g_{iw} + g_{rw})s + g_{rw}g_{iw}}$$

Note that the two poles are real and the transfer function can be rewritten in terms of the poles p_1 and p_2 as

$$G(s) = \frac{Y(s)}{U(s)} = \frac{K}{(s + p_1)(s + p_2)} \quad (3.22)$$

The corresponding partial fraction expansion is

$$\frac{G(s)}{s} = K \left(\frac{A_1}{s} + \frac{A_2}{s + p_1} + \frac{A_3}{s + p_2} \right)$$

where

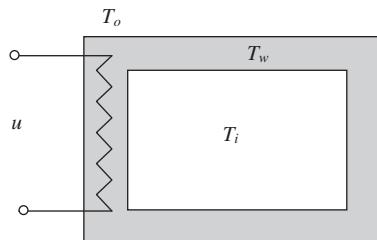


Figure 3.12
Schematic of a furnace.

Example 3.5—cont'd

$$\begin{aligned} A_1 &= \frac{1}{(s + p_1)(s + p_2)} \Big|_{s=0} = \frac{1}{p_1 p_2} \\ A_2 &= \frac{1}{s(s + p_2)} \Big|_{s=-p_1} = -\frac{1}{p_1(p_2 - p_1)} \\ A_3 &= \frac{1}{s(s + p_1)} \Big|_{s=-p_2} = \frac{1}{p_2(p_2 - p_1)} \end{aligned}$$

Using (3.21), the desired z-domain transfer function is

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z} \left\{ K \left[\frac{1}{p_1 p_2 s} - \frac{1}{p_1(p_2 - p_1)(s + p_1)} + \frac{1}{p_2(p_2 - p_1)(s + p_2)} \right] \right\} \\ &= K \left[\frac{1}{p_1 p_2} - \frac{z - 1}{p_1(p_2 - p_1)(z - e^{p_1 T})} + \frac{z - 1}{p_2(p_2 - p_1)(z - e^{p_2 T})} \right] \end{aligned}$$

which can be simplified to

$$G_{ZAS}(z) = K \left[\frac{(p_1 e^{-p_2 T} - p_2 e^{-p_1 T} + p_2 - p_1)z + \begin{pmatrix} p_1 e^{-p_1 T} - p_2 e^{-p_2 T} + p_2 e^{-(p_1+p_2)T} \\ -p_1 e^{-(p_1+p_2)T} \end{pmatrix}}{p_1 p_2 (p_2 - p_1)(z - e^{-p_1 T})(z - e^{-p_2 T})} \right]$$

Example 3.6

Find the z-domain transfer function of an armature-controlled DC motor.

Solution

The system can be modeled by means of the following differential equations:

$$J\ddot{\theta}(t) + b\dot{\theta}(t) = K_t i(t)$$

$$L \frac{di(t)}{dt} + Ri(t) = u(t) - K_e \dot{\theta}(t)$$

$$y(t) = \theta(t)$$

where θ is the position of the shaft (i.e., the output y of the system), i is the armature current, u is the source voltage (i.e., the input of the system), J is the moment of inertia of the motor, b is the viscous friction coefficient, K_t is the torque constant, K_e is the back electromotive force (e.m.f.) constant, R is the electric resistance, and L is the electric inductance. Laplace transforming and simplifying gives the following transfer function:

Example 3.6—cont'd

$$G(s) = \frac{Y(s)}{U(s)} = \frac{K_t}{s[(Js + b)(Ls + R) + K_t K_e]}$$

which can be rewritten as

$$G(s) = \frac{Y(s)}{U(s)} = \frac{K}{s(s + p_1)(s + p_2)}$$

with appropriate values of K , p_1 , and p_2 . The corresponding partial fraction expansion is

$$\frac{G(s)}{s} = K \left(\frac{A_{11}}{s^2} + \frac{A_{12}}{s} + \frac{A_2}{s + p_1} + \frac{A_3}{s + p_2} \right)$$

where

$$\begin{aligned} A_{11} &= \left. \frac{1}{(s + p_1)(s + p_2)} \right|_{s=0} = \frac{1}{p_1 p_2} \\ A_{12} &= \left. \frac{d}{ds} \left[\frac{1}{(s + p_1)(s + p_2)} \right] \right|_{s=0} = -\frac{p_1 + p_2}{p_1^2 p_2^2} \\ A_2 &= \left. \frac{1}{s^2(s + p_2)} \right|_{s=-p_1} = \frac{1}{p_1^2(p_2 - p_1)} \\ A_3 &= \left. \frac{1}{s^2(s + p_1)} \right|_{s=-p_2} = -\frac{1}{p_2^2(p_2 - p_1)} \end{aligned}$$

Using (3.21), the desired z-domain transfer function is

$$\begin{aligned} G_{ZAS}(1 - z^{-1}) \mathcal{Z} \left\{ K \left[\frac{1}{p_1 p_2 s^2} - \frac{p_1 + p_2}{p_1^2 p_2^2 s} + \frac{1}{p_1^2(p_2 - p_1)(s + p_1)} - \frac{1}{p_2^2(p_2 - p_1)(s + p_2)} \right] \right\} \\ = K \left[\frac{T}{p_1 p_2(z - 1)} - \frac{p_1 + p_2}{p_1^2 p_2^2 s} + \frac{1}{p_1(p_2 - p_1)(z - 1)(z - e^{p_1 T})} - \frac{1}{p_2(p_2 - p_1)(z - 1)(z - e^{p_2 T})} \right] \end{aligned}$$

Note that if the velocity of the motor is considered as output (i.e. $y(t) = \dot{\theta}(t)$), we have the transfer function

$$G(s) = \frac{Y(s)}{U(s)} = \frac{K_t}{(Js + b)(Ls + R) + K_t K_e}$$

and the calculations of Example 3.5 can be repeated to obtain the z-domain transfer function (see (3.22)).

In the preceding examples, we observe that if the analog system has a pole at p_s , then $G_{ZAS}(z)$ has a pole at $p_z = e^{p_z T}$. The division by s in (3.21) results in a pole at $z = 1$ that cancels, leaving the same poles as those obtained when sampling and z -transforming the impulse response of the analog subsystem. However, the zeros of the transfer function are different in the presence of a DAC.

3.6 Systems with transport lag

Many physical system models include a transport lag or delay in their transfer functions. These include chemical processes, automotive engines, sensors, digital systems, and so on. In this section, we obtain the z -transfer function $G_{ZAS}(z)$ of (3.21) for a system with transport delay.

The transfer function for systems with a transport delay is of the form

$$G(s) = G_a(s)e^{-T_d s} \quad (3.23)$$

where T_d is the transport delay. As in Section 2.7, the transport delay can be rewritten as

$$\begin{aligned} T_d &= lT - mT, \quad 0 \leq m < 1 \\ l &= \left\lceil \frac{T_d}{T} \right\rceil \\ m &= l - T_d/T \end{aligned} \quad (3.24)$$

where l is a positive integer and $\lceil \cdot \rceil$ denotes the ceiling, i.e., rounding up to the nearest integer. For example, a time delay of 3.1 s with a sampling period T of 1 s corresponds to $l = 4$ and $m = 0.9$. A delay by an integer multiple of the sampling period does not affect the form of the impulse response of the system. Therefore, using the delay theorem and (3.20), the z -transfer function for the system of (3.23) can be rewritten as

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z} \left\{ \mathcal{L}^{-1} \left[\frac{G_a(s)e^{-T(l-m)s}}{s} \right]^* \right\} \\ &= z^{-l} (1 - z^{-1}) \mathcal{Z} \left\{ \mathcal{L}^{-1} \left[\frac{G_a(s)e^{mTs}}{s} \right]^* \right\} \end{aligned} \quad (3.25)$$

From (3.25), we observe that the inverse Laplace transform of the function

$$G_s(s) = \frac{G_a(s)}{s} \quad (3.26)$$

is sampled and z -transformed to obtain the desired z -transfer function. We rewrite (3.26) in terms of $G_s(s)$ as

$$G_{ZAS}(z) = z^{-l} (1 - z^{-1}) \mathcal{Z} \left\{ \mathcal{L}^{-1} [G_s(s)e^{mTs}]^* \right\} \quad (3.27)$$

Using the time advance theorem of Laplace transforms gives

$$G_{ZAS}(z) = z^{-l} (1 - z^{-1}) \mathcal{Z}\{g_s^*(t + mT)\} \quad (3.28)$$

where the impulse-sampled waveform must be used to allow Laplace transformation. The remainder of the derivation does not require impulse sampling, and we can replace the impulse-sampled waveform with the corresponding sequence

$$G_{ZAS}(z) = z^{-l} (1 - z^{-1}) \mathcal{Z}\{g_s(kT + mT)\} \quad (3.29)$$

The preceding result hinges on our ability to obtain the effect of a time advance mT on a sampled waveform. In [Section 2.7](#), we discussed the z -transform of a signal delayed by T and advanced by mT , which is known as the modified z -transform. To express [\(3.29\)](#) in terms of the modified z -transform, we divide the time delay lT into a delay $(l-1)T$ and a delay T and rewrite the transfer function as

$$G_{ZAS}(z) = z^{-(l-1)} (1 - z^{-1}) z^{-1} \mathcal{Z}\{g_s(kT + mT)\} \quad (3.30)$$

Finally, we express the z -transfer function in terms of the modified z -transform

$$G_{ZAS}(z) = \left(\frac{z-1}{z^l}\right) \mathcal{Z}_m\{g_s(kT)\} \quad (3.31)$$

We recall two important modified transforms that are given in [Section 2.7](#):

$$\mathcal{Z}_m\{1(kT)\} = \frac{1}{z-1} \quad (3.32)$$

$$\mathcal{Z}_m\{e^{-pkT}\} = \frac{e^{-mpT}}{z - e^{-pT}} \quad (3.33)$$

Returning to our earlier numerical example, a delay of 3.1 sampling periods gives a delay of four sampling periods and the corresponding z^{-4} term and the modified z -transform of $g(kT)$ with $m = 0.9$.

Example 3.7

If the sampling period is 0.1 s, determine the z -transfer function $G_{ZAS}(z)$ for the system

$$G(s) = \frac{3e^{-0.31s}}{s+3}$$

Solution

First, write the delay in terms of the sampling period as $0.31 = 3.1 \times 0.1 = (4-0.9) \times 0.1$. Thus, $l = 4$ and $m = 0.9$. Next, obtain the partial fraction expansion

$$G_s(s) = \frac{3}{s(s+3)} = \frac{1}{s} - \frac{1}{s+3}$$

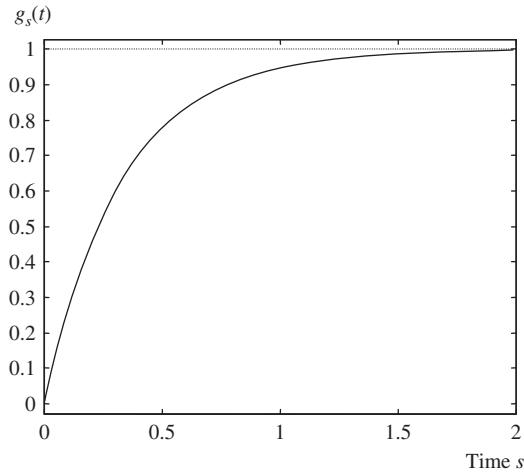
Example 3.7—cont'd

Figure 3.13
Continuous time function $g_s(t)$.

This is the transform of the continuous-time function shown in Fig. 3.13, which must be sampled, shifted, and z-transformed to obtain the desired transfer function. Using the modified z-transforms obtained in Section 2.7, the desired transfer function is

$$\begin{aligned} G_{ZAS}(z) &= \left(\frac{z-1}{z^4} \right) \left\{ \frac{1}{z-1} - \frac{e^{-0.3 \times 0.9}}{z - e^{-0.3}} \right\} \\ &= z^{-4} \left\{ \frac{z - 0.741 - 0.763(z-1)}{z - 0.741} \right\} = \frac{0.237z + 0.022}{z^4(z - 0.741)} \end{aligned}$$

3.7 The closed-loop transfer function

Using the results of Section 3.5, the digital control system of Fig. 3.1 yields the closed-loop block diagram of Fig. 3.14. The block diagram includes a comparator, a digital controller with transfer function $C(z)$, and the ADC-analog subsystem-DAC transfer function $G_{ZAS}(z)$. The controller and comparator are actually computer programs and

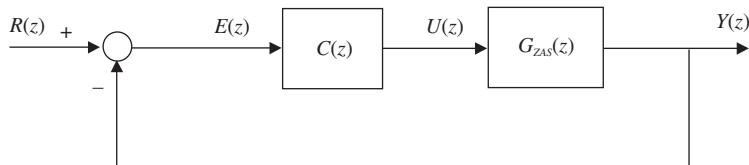


Figure 3.14
Block diagram of a single-loop digital control system.

replace the computer block in Fig. 3.1. The block diagram is identical to those commonly encountered in s -domain analysis of analog systems, with the variable s replaced by z . Hence, the closed-loop transfer function for the system is given by

$$G_{cl}(z) = \frac{C(z)G_{ZAS}(z)}{1 + C(z)G_{ZAS}(z)} \quad (3.34)$$

and the closed-loop characteristic equation is

$$1 + C(z)G_{ZAS}(z) = 0 \quad (3.35)$$

The roots of the equation are the closed-loop system poles, which can be selected for desired time response specifications as in s -domain design. Before we discuss this in some detail, we first examine alternative system configurations and their transfer functions.

When deriving closed-loop transfer functions of other configurations, the results of Section 3.4 must be considered carefully, as seen from Example 3.8.

Example 3.8

Find the Laplace transform of the analog and sampled output for the block diagram of Fig. 3.15.

Solution

The analog variable $x(t)$ has the Laplace transform

$$X(s) = H(s)G(s)D(s)E(s)$$

which involves three multiplications in the s -domain. In the time domain, $x(t)$ is obtained after three convolutions.

From the block diagram

$$E(s) = R(s) - X^*(s)$$

Substituting in the $X(s)$ expression, sampling then gives

$$X(s) = H(s)G(s)D(s)[R(s) - X^*(s)]$$

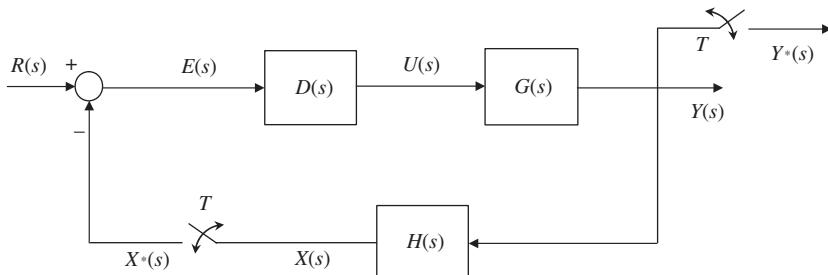


Figure 3.15

Block diagram of a system with sampling in the feedback path.

Example 3.8—cont'd

Thus, the impulse-sampled variable $x^*(t)$ has the Laplace transform

$$X^*(s) = (HGDR)^*(s) - (HGD)^*(s)X^*(s)$$

where, as in the first part of Example 3.1, several components are no longer separable. These terms are obtained as shown in Example 3.1 by inverse Laplace transforming, impulse sampling, and then Laplace transforming the impulse-sampled waveform.

Next, we solve for $X^*(s)$

$$X^*(s) = \frac{(HGDR)^*(s)}{1 + (HGD)^*(s)}$$

and then $E(s)$

$$E(s) = R(s) - \frac{(HGDR)^*(s)}{1 + (HGD)^*(s)}$$

With some experience, the last two expressions can be obtained from the block diagram directly. The combined terms are clearly the ones not separated by samplers in the block diagram.

From the block diagram, the Laplace transform of the output is $Y(s) = G(s)D(s)E(s)$. Substituting for $E(s)$ gives

$$Y(s) = G(s)D(s) \left[R(s) - \frac{(HGDR)^*(s)}{1 + (HGD)^*(s)} \right]$$

Thus, the sampled output is

$$Y^*(s) = (GDR)^*(s) - (GD)^*(s) \frac{(HGDR)^*(s)}{1 + (HGD)^*(s)}$$

With the transformation $z = e^{st}$, we can rewrite the sampled output as

$$Y(z) = (GDR)(z) - (GD)(z) \frac{(HGDR)(z)}{1 + (HGD)(z)}$$

The last equation demonstrates how for some digital systems, no expression is available for the transfer function excluding the input. Nevertheless, the preceding system has a closed-loop characteristic equation similar to (3.35) given by $1 + (HGD)(z) = 0$. This equation can be used in design, as in cases where a closed-loop transfer function is defined.

3.8 Analog disturbances in a digital system

Disturbances are variables that are not included in the system model but affect its response. They can be deterministic, such as load torque in a position control system, or stochastic, such as sensor or actuator noise. However, almost all disturbances are analog and are inputs to the analog subsystem in a digital control loop. We use the results of Section 3.7 to obtain a transfer function with a disturbance input.

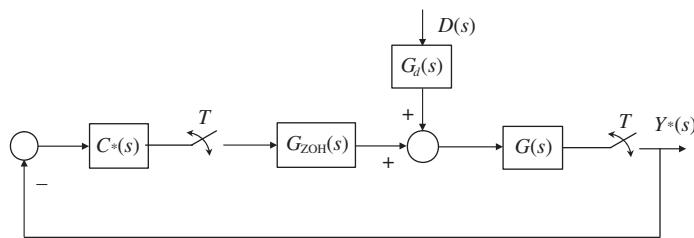


Figure 3.16
Block diagram of a digital system with an analog disturbance.

Consider the system with disturbance input shown in Fig. 3.16. Because the system is linear, the reference input can be treated separately and is assumed to be zero.

The Laplace transform of the impulse-sampled output is

$$Y^*(s) = (GG_dD)^*(s) - (GG_{ZOH})^*(s)C^*(s)Y^*(s) \quad (3.36)$$

Solving for $Y^*(s)$, we obtain

$$Y^*(s) = \frac{(GG_dD)^*(s)}{1 + (GG_{ZOH})^*(s)C^*(s)} \quad (3.37)$$

The denominator involves the transfer function for the zero-order hold, analog subsystem, and sampler. We can therefore rewrite (3.37) using the notation of (3.21) as

$$Y^*(s) = \frac{(GG_dD)^*(s)}{1 + (G_{ZAS})^*(s)C^*(s)} \quad (3.38)$$

or in terms of z as

$$Y(z) = \frac{(GG_dD)(z)}{1 + G_{ZAS}(z)C(z)} \quad (3.39)$$

Example 3.9

Consider the block diagram of Fig. 3.16 with the transfer functions

$$G(s) = \frac{K_p}{s+1}, \quad G_d(s) = \frac{1}{s}, \quad C(z) = K_c$$

Find the steady-state response of the system to an impulse disturbance of strength A .

Solution

We first evaluate

$$G(s)G_d(s)D(s) = \frac{K_p A}{s(s+1)} = K_p A \left[\frac{1}{s} - \frac{1}{s+1} \right]$$

The z-transform of the corresponding impulse response sequence is

Example 3.9—cont'd

$$(GG_dD)(z) = K_p A \left[\frac{z}{z-1} - \frac{z}{z-e^{-T}} \right]$$

Using (3.21), we obtain the transfer function

$$G_{ZAS}(z) = K_p \frac{1 - e^{-T}}{z - e^{-T}}$$

From (3.39), we obtain the sampled output

$$Y(z) = \frac{K_p A \left[\frac{z}{z-1} - \frac{z}{z-e^{-T}} \right]}{1 + K_c \left[K_p \frac{1 - e^{-T}}{z - e^{-T}} \right]}$$

To obtain the steady-state response, we use the final value theorem

$$\begin{aligned} y(\infty) &= (z-1)Y(z)|_{z=1} \\ &= \frac{K_p A}{1 + K_c K_p} \end{aligned}$$

Thus, as with analog systems, increasing the controller gain reduces the error due to the disturbance. Equivalently, an analog amplifier before the point of disturbance injection can increase the gain and reduce the output due to the disturbance and is less likely to saturate the DAC. Note that it is simpler to apply the final value theorem without simplification because terms not involving $(z-1)$ drop out. Because the disturbance is ideally zero, the steady state error due to a disturbance is $e_D(\infty) = 0 - y(\infty) = -y(\infty)$.

3.9 Steady-state error and error constants

The tracking error is the error in tracking the reference input

$$e(t) = r(t) - y(t) \quad (3.40)$$

For a closed-loop system, we can write the transform of the error as

$$E(z) = R(z) - Y(z) = [1 - G_{cl}(z)]R(z) \quad (3.41)$$

where $G_{cl}(z)$ is the closed-loop transfer function. With the final value theorem we obtain the steady-state tracking error

$$e(\infty) = (z-1)E(z)|_{z=1} = (z-1)[1 - G_{cl}(z)]R(z)|_{z=1} \quad (3.42)$$

Next, we consider the unity feedback block diagram shown in Fig. 3.14 subject to standard inputs and determine the associated **tracking error** in each case. The standard inputs considered are the **sampled step**, the **sampled ramp**, and the **sampled parabolic**. As with analog systems, an error constant is associated with each input, and a type number can be defined for any system from which the nature of the error constant can be inferred. All results are obtained by direct application of the final value theorem.

From Fig. 3.14, the tracking error is given by

$$\begin{aligned} E(z) &= \frac{R(z)}{1 + G_{ZAS}(z)C(z)} \\ &= \frac{R(z)}{1 + L(z)} \end{aligned} \quad (3.43)$$

where $L(z)$ denotes the loop gain of the system.

Applying the final value theorem yields the steady-state error

$$\begin{aligned} e(\infty) &= (1 - z^{-1})E(z)|_{z=1} \\ &= \frac{(z - 1)R(z)}{z(1 + L(z))}|_{z=1} \end{aligned} \quad (3.44)$$

The limit exists if all $(z - 1)$ terms in the denominator cancel. This depends on the reference input as well as on the loop gain.

To examine the effect of the loop gain on the limit, rewrite it in the form

$$L(z) = \frac{N(z)}{(z - 1)^n D(z)}, \quad n \geq 0 \quad (3.45)$$

where $N(z)$ and $D(z)$ are numerator and denominator polynomials, respectively, with no unity roots. The following definition plays an important role in determining the steady-state error of unity feedback systems.

Definition 3.1: Type Number

The type number of the system is the number of unity poles in the system z -transfer function.

The loop gain of (3.45) has n poles at unity and is therefore type n . These poles play the same role as poles at the origin for an s -domain transfer function in determining the steady-state response of the system. Note that s -domain poles at zero play the same role as z -domain poles at e^0 .

Substituting from (3.45) in the error expression (3.44) gives

$$\begin{aligned} e(\infty) &= \frac{(z - 1)^{n+1}D(z)R(z)}{z(N(z) + (z - 1)^n D(z))}|_{z=1} \\ &= \frac{(z - 1)^{n+1}D(1)R(z)}{N(1) + (z - 1)^n D(1)}|_{z=1} \end{aligned} \quad (3.46)$$

Next, we examine the effect of the reference input on the steady-state error.

3.9.1 Sampled step input

The z -transform of a sampled unit step input is

$$R(z) = \frac{z}{z - 1}$$

Substituting in (3.44) gives the steady-state error

$$e(\infty) = \left. \frac{1}{1 + L(z)} \right|_{z=1} \quad (3.47)$$

The steady-state error can also be written as

$$e(\infty) = \frac{1}{1 + K_p} \quad (3.48)$$

where K_p is the position error constant given by

$$K_p = L(1) \quad (3.49)$$

Examining (3.45) shows that K_p is finite for type 0 systems and infinite for systems of type 1 or higher. Therefore, the steady-state error for a sampled unit step input is

$$e(\infty) = \begin{cases} \frac{1}{1 + L(1)}, & n = 0 \\ 0, & n \geq 1 \end{cases} \quad (3.50)$$

3.9.2 Sampled ramp input

The z -transform of a sampled unit ramp input is

$$R(z) = \frac{Tz}{(z - 1)^2}$$

Substituting in (3.44) gives the steady-state error

$$\begin{aligned} e(\infty) &= \left. \frac{T}{[z - 1][1 + L(z)]} \right|_{z=1} \\ &= \frac{1}{K_v} \end{aligned} \quad (3.51)$$

where K_v is the velocity error constant. The velocity error constant is thus given by

$$K_v = \frac{1}{T} (z - 1)L(z)|_{z=1} \quad (3.52)$$

From (3.52), the velocity error constant is zero for type 0 systems, finite for type 1 systems, and infinite for type 2 or higher systems. The corresponding steady-state error is

$$e(\infty) = \begin{cases} \infty, & n = 0 \\ \frac{T}{(z-1)L(z)|_{z=1}}, & n = 1 \\ 0, & n \geq 2 \end{cases} \quad (3.53)$$

Similarly, it can be shown that for a sampled parabolic input, an acceleration error constant given by

$$K_a = \frac{1}{T^2}(z-1)^2 L(z)|_{z=1} \quad (3.54)$$

can be defined, and the associated steady-state error is

$$e(\infty) = \begin{cases} \infty, & n \leq 1 \\ \frac{T^2}{(z-1)^2 L(z)|_{z=1}}, & n = 2 \\ 0, & n \geq 3 \end{cases} \quad (3.55)$$

The following examples show how to evaluate the tracking error. In the presence of a disturbance input, the error due to the disturbance $e_D(\infty)$ can be added to the tracking error to obtain the total steady error.

Example 3.10

Find the steady-state position error for the digital position control system with unity feedback and with the transfer functions

$$G_{ZAS}(z) = \frac{K(z+a)}{(z-1)(z-b)}, \quad C(z) = \frac{K_c(z-b)}{z-c}, \quad 0 < a, b, c < 1$$

1. For a sampled unit step input
2. For a sampled unit ramp input

Solution

The loop gain of the system is given by

$$L(z) = C(z)G_{ZAS}(z) = \frac{KK_c(z+a)}{(z-1)(z-c)}$$

Example 3.10—cont'd

The system is type 1. Therefore, it has zero steady-state error for a sampled step input and a finite steady-state error for a sampled ramp input given by

$$e(\infty) = \frac{T}{(z - 1)L(z)|_{z=1}} = \frac{T}{KK_c} \left(\frac{1 - c}{1 + a} \right)$$

Clearly, the steady-state error is reduced by increasing the controller gain and is also affected by the choice of controller pole and zero.

Example 3.11

Find the steady-state error for the analog system

$$G(s) = \frac{K}{s + a} \quad a > 0$$

1. For proportional analog control with a unit step input
2. For proportional digital control with a sampled unit step input

Solution

The transfer function of the system can be written as

$$G(s) = \frac{K/a}{s/a + 1} \quad a > 0$$

Thus, the position error constant for analog control is K/a , and the steady-state error is

$$e(\infty) = \frac{1}{1 + K_p} = \frac{a}{K + a}$$

For digital control, it can be shown that for sampling period T , the DAC-plant-ADC z-transfer function is

$$G_{ZAS}(z) = \frac{K}{a} \left(\frac{1 - e^{-aT}}{z - e^{-aT}} \right)$$

Thus, the position error constant for digital control is

$$K_p = G_{ZAS}(z)|_{z=1} = K/a$$

and the associated steady-state error is the same as that of the analog system with proportional control. In general, it can be shown that the steady-state error for the same control strategy is identical for digital or analog implementation.

3.10 MATLAB commands

The transfer function for the ADC, analog subsystem, and DAC combination can be easily obtained using MATLAB. Assume that the sampling period is equal to 0.1 s and that the transfer function of the analog subsystem is G .

3.10.1 MATLAB

The MATLAB command to obtain a digital transfer function from an analog transfer function is

```
>> g = tf(num, den)
>> gd = c2d(g, 0.1, 'method')
```

where **num** is a vector containing the numerator coefficients of the analog transfer function in descending order, and **den** is a similarly defined vector of denominator coefficients. For example, the numerator polynomial ($2s^2+4s + 3$) is entered as

```
>> num = [2, 4, 3]
```

The term '**method**' specifies the method used to obtain the digital transfer function. For a system with a zero-order hold and sampler (DAC and ADC), we use

```
>> gd = c2d(g, 0.1, 'zoh')
```

For a first-order hold, we use

```
>> gd = c2d(g, 0.1, 'foh')
```

Other options of MATLAB commands are available but are not relevant to the material presented in this chapter.

For a system with a time delay, the discrete transfer function can be obtained using the commands

```
gdelay = tf(num, den, 'inputdelay', Td)%Delay = Td
gdelay_d = c2d(gdelay, 0.1, 'method')
```

3.10.2 Simulink

Simulink is a MATLAB toolbox that provides a graphical language for simulating dynamic systems. The vocabulary of the language is a set of block operations that the user can combine to create a dynamic system. The block operations include continuous and discrete-time, linear and nonlinear, as well as user-defined blocks containing MATLAB functions. Additional blocks called sources provide a wide variety of inputs. Sink blocks provide a means of recording the system response, and the response can then be sent to MATLAB for processing or plotting.

This section provides a brief introduction to Simulink and its simulation of discrete-time systems that we hope will motivate the reader to experiment with this powerful simulation tool. In our view, this is the only way to become proficient in the use of Simulink.

To start Simulink, start MATLAB then click on the Simulink icon of the Home tab or type

```
>> simulink
```

This opens a window from which the user can select “Blank Model” to create a new Simulink model. The model has an icon in the form of four squares called the **Library Browser**, which offers the user many blocks to choose from for simulation. To create a model,

File > New > Model

blocks are selected from the Simulink Library Browser as needed.

We begin with a simple simulation of a source and sink. Click on Sources and find a Step block, and then drag it to the model window. Then click on Sinks and find a Scope and drag it to the model area. To connect two blocks, click on the first block while holding down the Control button, and then click on the second block. Simulink will automatically draw a line from the first block to the second. The simulation diagram for this simple system is shown in [Fig. 3.17](#).

To set the parameters of a block, double-click it, and a window displaying the parameters of the block will appear. Each of the parameters can be selected and changed to the appropriate value for a particular simulation. For example, by clicking on the Step block, we can select the time when the step is applied, the input value before and after the step, and the sampling period. Double-clicking on Scope shows a plot of the step after running the simulation. To run the simulation, click on the arrow above the simulation diagram. The scope will show a plot of the step input.

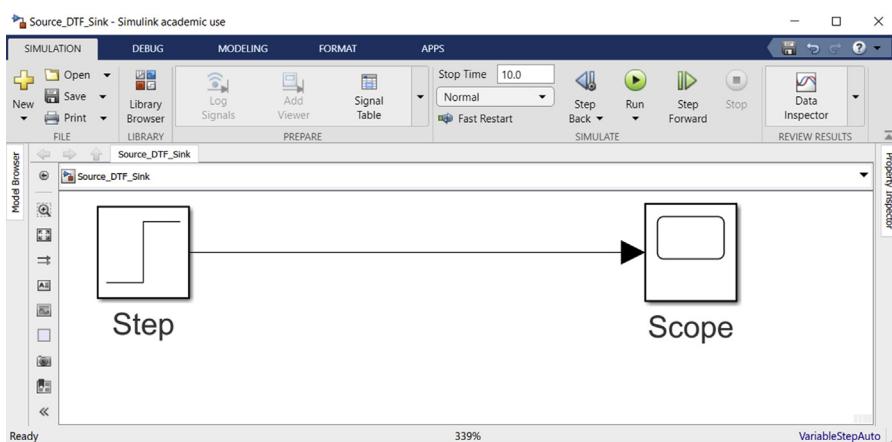


Figure 3.17
Simulation diagram for a step input and scope.

To model a discrete-time system, we select Discrete in the Simulink Library Browser and drag a Discrete Transfer Function block to our model window (see Fig. 3.18). The numerator and denominator of the transfer function are set by clicking on the block, and they are entered as MATLAB vectors, as shown in Fig. 3.19. The sampling period can also be selected in the same window, and we select a value of 0.05 s. Starting with the Source-Sink model, we can select the line joining the source to the sink and then click on Delete

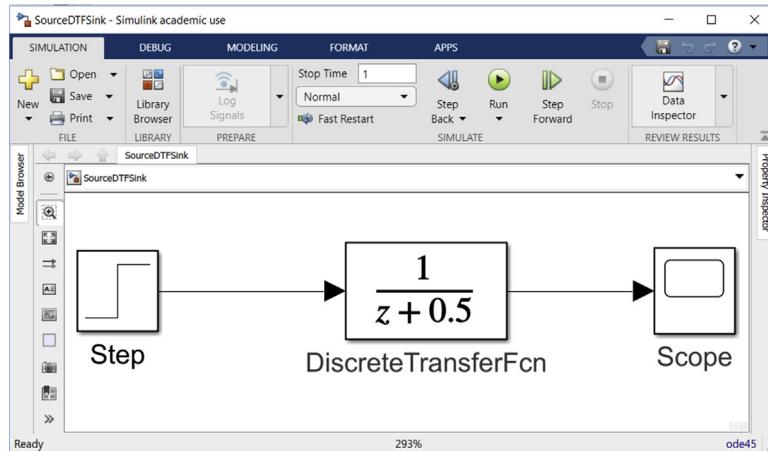


Figure 3.18
Simulation diagram for a step input, discrete transfer function, and scope.

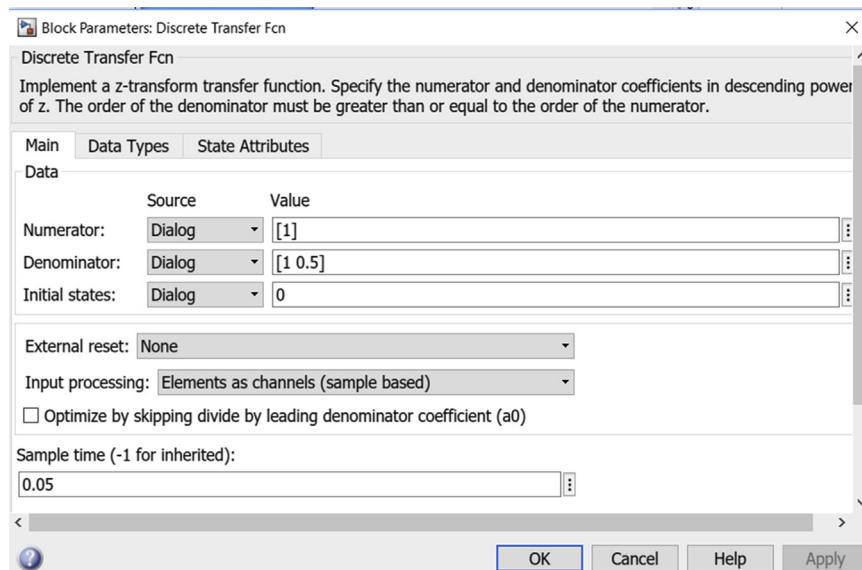


Figure 3.19
Parameters of the discrete transfer function.

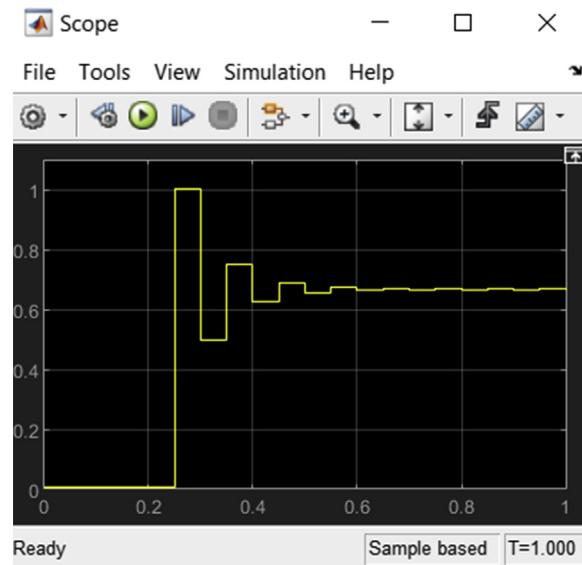
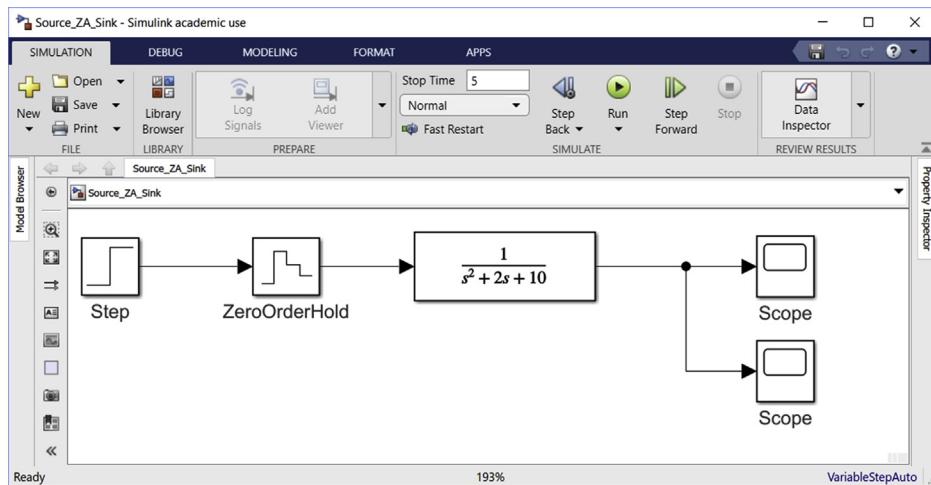


Figure 3.20
Step response of a discrete transfer function.

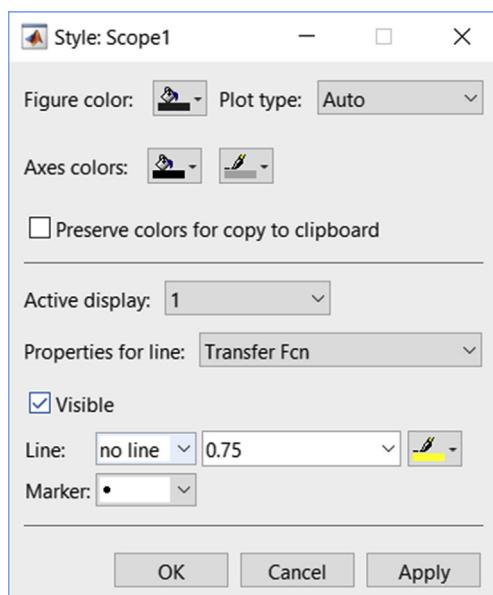
to eliminate it. Then we can join the three blocks using the Control button as before and save the new model. Running the simulation gives the results in Fig. 3.20. The simulation duration is selected in the model window next to the **run** icon. The scale of the plot is selected by right-clicking on the scope plot.

We now discuss simulating an analog system with digital control. The analog transfer function is obtained by clicking on Continuous in the Simulink Library Browser and then dragging the block to the model window. Although it is not required to implement the zero-order hold, the zero-order hold block is available under Discrete, but MATLAB does not have a sampler. The block does not change the simulation results. However, the block allows the user to obtain a plot of the input to the analog system.

The input to any discrete block is sampled automatically at the sampling rate of the discrete block, and the sampling rate can also be chosen for sources and sinks. Fig. 3.21 shows the simulation diagram for the digital control system with two scopes. We set the sampling period for the first zero-order hold to 0.5 and all other blocks inherit the sampling period when set to -1. The parameters of the scope plots are selected by clicking on View (third from the left in Fig. 3.20). For the sampled plot of Fig. 3.23, we select “Style”, then choose “no line” and the “.” Marker as shown in Fig. 3.22. The default setting will give the plot of Fig. 3.24.

**Figure 3.21**

Simulation diagram for a step input, zero-order hold, analog transfer function, and two scopes.

**Figure 3.22**

Selecting the scope parameters.

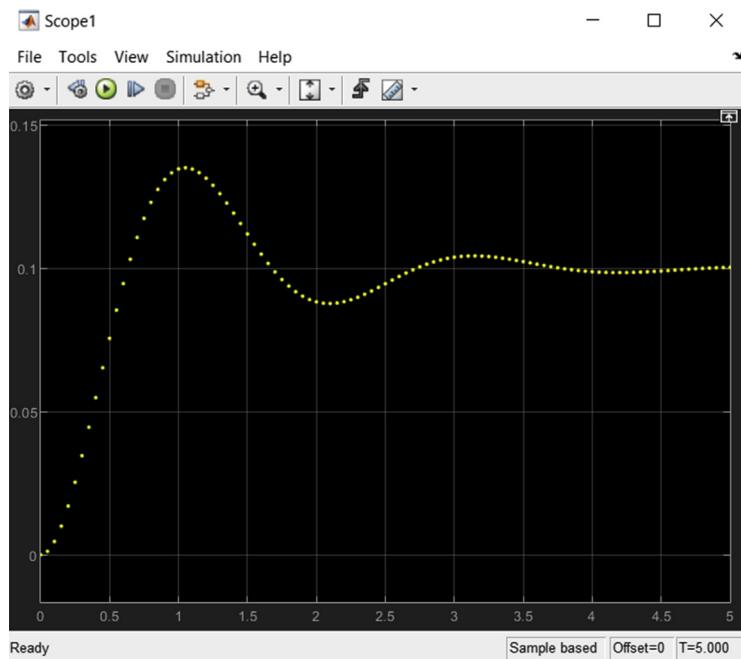


Figure 3.23
Step response of analog system with digital control.

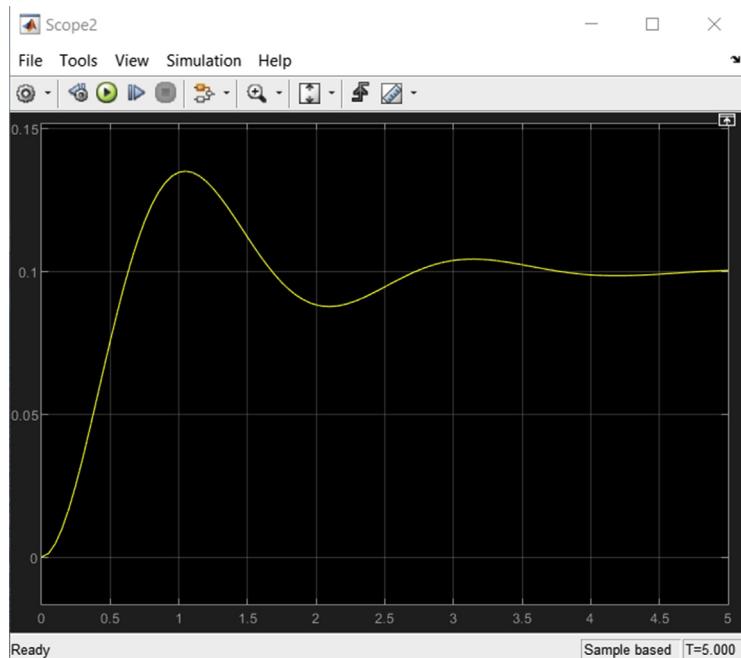


Figure 3.24
Step response of analog system with zero-order hold, without output sampling.

3.11 Sensitivity analysis

Sensitivity is a measure of the effect of a change in a system parameter on a property or characteristic function of a system. Because implementation always results in an approximation of their design, designers seek systems whose characteristics do not change appreciably with parameter changes, i.e., systems with low sensitivity. Such systems are said to be robust.

Bode sensitivity is one of the most commonly used measures of sensitivity in control system analysis and design. Bode sensitivity is the ratio of per unit changes in a system function F to per unit changes in the value of a system parameter p for infinitesimal parameter changes

$$S_p^F = \lim_{\Delta p \rightarrow 0} \frac{\text{per unit change in } F}{\text{per unit change in } p} \quad (3.56)$$

Substituting for the per unit changes gives

$$S_p^F = \lim_{\Delta p \rightarrow 0} \frac{\frac{\Delta F}{F}}{\frac{\Delta p}{p}} = \frac{p}{F} \lim_{\Delta p \rightarrow 0} \frac{\Delta F}{\Delta p} \quad (3.57)$$

From the definition of the derivative, we have the sensitivity function

$$S_p^F = \frac{p}{F} \frac{\partial F}{\partial p} \quad (3.58)$$

Although the sensitivity definition is for infinitesimal parameter changes, approximate values can be obtained for small perturbations by substituting a ratio of changes for the derivative

$$S_p^F = \frac{p}{F} \frac{\partial F}{\partial p} \approx \frac{p}{F} \left[\frac{\Delta F}{\Delta p} \right]_{\text{small } \Delta p} \quad (3.59)$$

The perturbation is approximately given by

$$\Delta F \approx \frac{\Delta p}{p} F S_p^F, \text{ small } \Delta p \quad (3.60)$$

Example 3.12

- (i) Determine the sensitivity of the closed-loop transfer function

$$F = \frac{K/a}{1 + K/a} = \frac{K}{K + a}$$

to changes in the parameter a . What is the effect of increasing K on the sensitivity?

Example 3.12—cont'd

- (ii) For the nominal values $K = 10$, $a = 1$, calculate the sensitivity and the approximate percent change in F when $\Delta a = 0.01a$.
- (iii) Compare the results of the closed-loop system to those of the open-loop transfer function $F = K/a$

Solution

For the closed-loop transfer function, the sensitivity is

$$S_a^F = \frac{a}{F} \frac{\partial F}{\partial a} = \frac{a(K+a)}{K} \times \frac{-K}{(K+a)^2}$$

$$S_a^F = -\frac{a}{K+a}$$

Because the term K appears in the denominator, increasing K decreases the magnitude of the sensitivity. Negative sensitivity values indicate that increasing the parameter values reduces the value of the transfer function.

For nominal values $K = 10$, $a = 1$, the sensitivity is

$$S_a^F \Big|_{\begin{array}{l} K=10 \\ a=1 \end{array}} = -\frac{a}{K+a} \Big|_{\begin{array}{l} K=10 \\ a=1 \end{array}} = -\frac{1}{11} = -0.0909$$

The small perturbation $\Delta a = 0.01a$, results in the perturbation

$$\frac{\Delta F}{F} \times 100 \approx \frac{\Delta a}{a} S_a^F \times 100 = 0.01 \left(\frac{-1}{11} \right) \times 100$$

$$\approx -0.1\%$$

This shows that the closed-loop transfer function F is robust with respect to changes in the parameter a .

For the open-loop transfer function, the sensitivity is

$$S_a^F = \frac{a}{F} \frac{\partial F}{\partial a} = \frac{a^2}{K} \times \frac{-K}{a^2} = -1$$

This negative sensitivity value is independent of the numerical values of K and a .

The small perturbation $\Delta a = 0.01a$, results in the percent perturbation

$$\frac{\Delta F}{F} \times 100 \approx \frac{\Delta a}{a} S_a^F \times 100 = 0.01(-1) \times 100 = -1$$

As one would expect for unit sensitivity, a 1% change in the parameter a results in a 1% decrease in F and the system is very sensitive to parameter changes. This compares with a 0.1% for the closed-loop transfer function. This example demonstrates the reduction in sensitivity that feedback provides to a closed-loop system.

3.11.1 Pole sensitivity

In control system analysis and design, we are particularly interested in the effect of changes in a parameter p on the locations of closed-loop pole z_{cl}

$$S_p^{z_{cl}} = \frac{p}{z_{cl}} \frac{\partial z_{cl}}{\partial p} \quad (3.61)$$

While it may be obvious that z_{cl} depends on a particular parameter p , it is often difficult to solve for the pole as a function of the parameter. While this may seem as a formidable obstacle to calculating pole sensitivity, a simple solution can be found by considering the closed-loop characteristic equation. As long as the parameter of interest enters linearly in the equation, then simple algebraic manipulation allows us to rewrite the equation in the form

$$0 = D_{np}(z_{cl}) + pD_p(z_{cl})$$

where $D_{np}(z_{cl})$ includes all terms free of the parameter p and $D_p(z_{cl})$ includes all terms including the parameter p .

Differentiate partially w.r.t. p

$$\begin{aligned} 0 &= \left\{ p \frac{\partial D_p(z_{cl})}{\partial z_{cl}} + \frac{\partial D_p(z_{cl})}{\partial z_{cl}} \right\} \frac{\partial z_{cl}}{\partial p} + D_p(z_{cl}) \\ S_p^{z_{cl}} &= \frac{p}{z_{cl}} \frac{\partial z_{cl}}{\partial p} = -\frac{p}{z_{cl}} \times \frac{D_p(z_{cl})}{p \frac{\partial D_p(z_{cl})}{\partial z_{cl}} + D_p(z_{cl})} \end{aligned} \quad (3.62)$$

Example 3.13

Find the sensitivity of the closed loop pole of the system

$$F = \frac{K/(z+a)}{1+K/(z+a)} = \frac{K}{z+a+K}$$

to changes in the magnitude of the open-loop pole a .

Solution

The closed-loop characteristic equation is

$$z + K + a = 0$$

The system has only one closed pole at $z_{cl} = -(K + a)$. The polynomial of parameter-free terms is $D_{np}(z_{cl}) = z_{cl} + K$ and the remaining term is simply the parameter a . Thus, D_p is unity and the sensitivity is

$$\begin{aligned} S_a^{z_{cl}} &= -\frac{a}{z_{cl}} \times \frac{D_p(z_{cl})}{\frac{\partial D_{np}(z_{cl})}{\partial z_{cl}} + a \frac{\partial D_p(z_{cl})}{\partial z_{cl}}} = \frac{(-a)}{-(K+a)} \times \frac{1}{1+a} \\ &= \frac{1}{K+a} \times \frac{a}{1+a} \end{aligned}$$

As in Example 3.12, increasing the value of K decreases the magnitude of the sensitivity.

If the parameter of interest is the gain K , then the closed-loop characteristic equation is in the form

$$1 + G(z_{cl})H(z_{cl}) = 1 + \frac{KN(z_{cl})}{D(z_{cl})} = 0$$

where $D(z)$ is the denominator of the open-loop gain and $N(z)$ is its numerator.

Multiplying by the denominator gives

$$D(z_{cl}) + KN(z_{cl}) = 0$$

In this case, $D_{np}(z_{cl}) = D(z_{cl})$, $D_p(z_{cl}) = N(z_{cl})$, and the sensitivity is

$$S_K^{z_{cl}} = -\frac{K}{z_{cl}} \times \frac{\frac{N(z_{cl})}{\partial z_{cl}}}{K \frac{\partial N(z_{cl})}{\partial z_{cl}} + \frac{\partial D(z_{cl})}{\partial z_{cl}}} \quad (3.63)$$

Example 3.14

For the closed-loop system with transfer function

$$G_{cl}(z) = \frac{K}{z^2 - 2az + K}$$

- (i) Obtain the sensitivity in terms of the system parameters and show that it has a discontinuity at $K = a^2$.
- (ii) For $a = 0.05$, plot the magnitude of the sensitivity for gain values in the range $[0, a^2/2] \cup [2a^2, 2a^2+a]$ and discuss the effect of changing the gain on the pole sensitivity function.

Solution

We first solve for the poles in terms of the gain K

$$z_{cl} = a \pm \sqrt{a^2 - K}$$

Using Eq. (3.63), we have the sensitivity function

$$\begin{aligned} S_K^{z_{cl}} &= -\frac{K}{2z_{cl}(z_{cl} - a)} \\ &= -\frac{K}{2(a \pm \sqrt{a^2 - K})(\pm \sqrt{a^2 - K})} \\ &= -\frac{K/2}{a^2 - K \pm 2a\sqrt{a^2 - K}} \end{aligned}$$

The sensitivity function has a singularity (infinite value) at $a^2 = K$ (zero denominator).

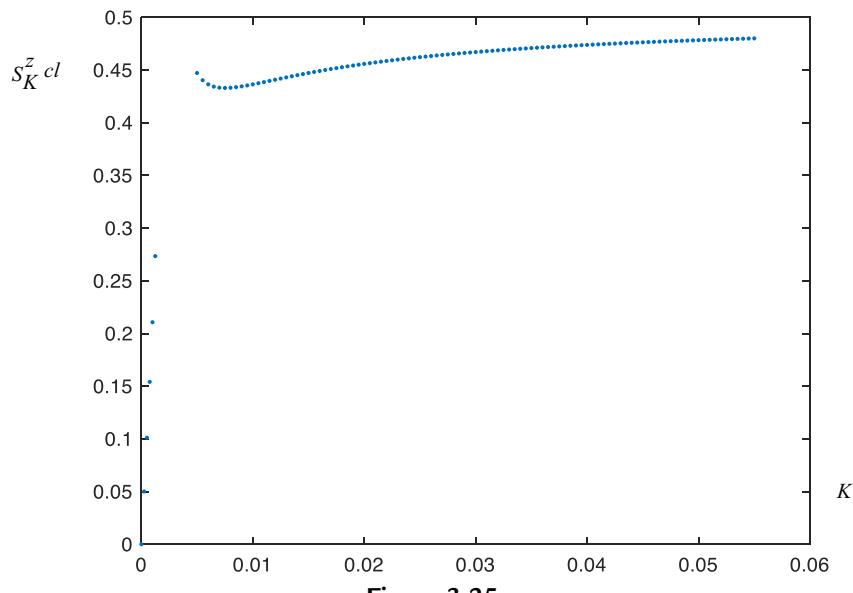


Figure 3.25
Closed-loop pole sensitivity to changes in the gain.

Fig. 3.25 shows that the sensitivity increases with the gain for low gain values. The discontinuity corresponds to the critical gain case with two equal real poles. For complex poles associated with gains above the critical gain, the magnitude of the sensitivity dips slightly then increases steadily with the gain. One would think that increasing the gain should decrease rather than increase sensitivity. However, in this case increasing the gain changes the locations of the closed-loop poles and consequently increases the sensitivity.

Further reading

- Franklin, G.F., Powell, J.D., Workman, M.L., 1998. Digital Control of Dynamic Systems. Addison-Wesley, Boston, MA.
- Jacquot, R.G., 1981. Modern Digital Control Systems. Marcel Dekker, New York.
- Katz, P., 1981. Digital Control Using Microprocessors. Prentice Hall International, Englewood Cliffs, NJ.
- Kuo, B.C., 1992. Digital Control Systems. Saunders, Ft. Worth, TX.
- Ogata, K., 1987. Discrete-Time Control Systems. Prentice Hall, Englewood Cliffs, NJ.

Problems

- 3.1 Find the magnitude and phase at frequency $\omega = 1$ of a zero-order hold with sampling period $T = 0.1$ s.

- 3.2 The first-order hold uses the last two numbers in a sequence to generate its output. The output of both the zero-order hold and the first-order hold is given by

$$u(t) = u(kT) + a \frac{u(kT) - u[(k-1)T]}{T} (t - kT), \quad kT \leq t \leq (k+1)T \quad k = 0, 1, 2, \dots$$

with $a = 0, 1$, respectively.

- For a discrete impulse input, obtain and sketch the impulse response for the preceding equation with $a = 0$ and $a = 1$.
- Write an impulse-sampled version of the preceding impulse response, and show that the transfer function is given by

$$G_H(s) = \frac{1}{s} (1 - e^{-sT}) \left[1 - ae^{-sT} + \frac{a}{sT} (1 - e^{-sT}) \right]$$

- Obtain the transfer functions for the zero-order hold and for the first-order hold from the preceding transfer function. Verify that the transfer function of the zero-order hold is the same as that obtained in [Section 3.3](#).

- 3.3 A beer filtration process can be described by the following transfer function¹:

$$G(s) = \frac{0.0216s - 0.0031}{s^2 + 0.4576s + 0.0868}$$

Obtain the transfer function of the system with zero-order hold and sampler for a sampling period of $T = 0.1$.

- 3.4 Many chemical processes can be modeled by the following transfer function:

$$G(s) = \frac{K}{\tau s + 1} e^{-T_d s}$$

where K is the gain, τ is the time constant, and T_d is the time delay. Obtain the transfer function $G_{ZAS}(z)$ for the system in terms of the system parameters. Assume that the time delay T_d is a multiple of the sampling period T .

- Obtain the transfer function of a point mass (m) with force as input and displacement as output, neglecting actuator dynamics; then find $G_{ZAS}(z)$ for the system.
- For an internal combustion engine, the transfer function with injected fuel flow rate as input and fuel flow rate into the cylinder as output is given by²

¹ M. Lees, L. Wang, PID controller design for industrial beer filtration, Proceedings fifth Australian Control Conference, pp. 306–311, 2005.

² Moskwa, J., 1988. Automotive Engine Modeling and Real time Control, MIT doctoral thesis.

$$G(s) = \frac{\varepsilon\tau s + 1}{\tau s + 1}$$

where τ is a time constant and ε is known as the fuel split parameter. Obtain the transfer function $G_{ZAS}(z)$ for the system in terms of the system parameters.

- 3.7 Repeat Problem 3.6 including a delay of 25 ms in the transfer function with a sampling period of 10 ms.
- 3.8 Find the equivalent sampled impulse response sequence and the equivalent z -transfer function for the cascade of the two analog systems with sampled input

$$H_1(s) = \frac{1}{s+6} \quad H_2(s) = \frac{10}{s+1}$$

- a. If the systems are directly connected
 b. If the systems are separated by a sampler
- 3.9 Obtain expressions for the analog and sampled outputs from the block diagrams shown in Fig. P3.9
- 3.10 For the system with inner-loop feedback, find an expression for the sampled output $Y^*(s)$ (Fig. P3.10).

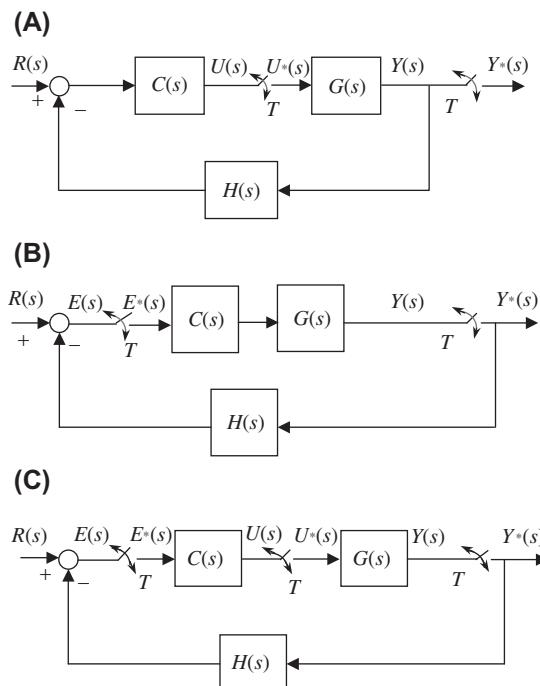


Figure P3.9
Block diagrams for systems with multiple samplers.

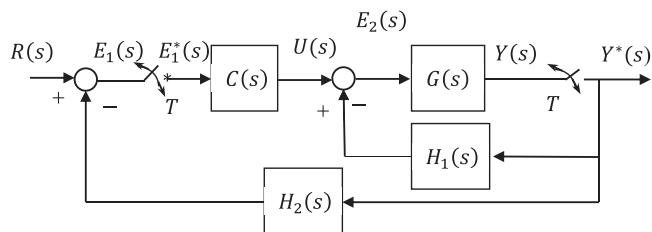


Figure P3.10
System with inner-loop feedback.

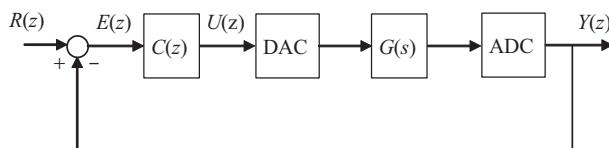


Figure P3.11
Block diagram for a closed-loop system with digital control.

- 3.11 For the unity feedback system shown in Fig. P3.11, we are given the analog subsystem

$$G(s) = \frac{s + 8}{s + 5}$$

The system is digitally controlled with a sampling period of 0.02 s. The controller transfer function was selected as

$$C(z) = \frac{0.35z}{z - 1}$$

- a. Find the z -transfer function for the analog subsystem with DAC and ADC.
- b. Find the closed-loop transfer function and characteristic equation.
- c. Find the steady-state error for a sampled unit step and a sampled unit ramp.
Comment on the effect of the controller on steady-state error.

- 3.12 Find the steady-state error for a unit step disturbance input for the systems shown in Fig. P3.12 with a sampling period of 0.03 s and the transfer functions

$$G_d(s) = \frac{2}{s + 1} \quad G(s) = \frac{4(s + 2)}{s(s + 3)} \quad C^*(s) = \frac{e^{sT} - 0.95}{e^{sT} - 1}$$

- 3.13 For the following systems with unity feedback, find
- a. The position error constant
 - b. The velocity error constants
 - c. The steady-state error for a unit step input

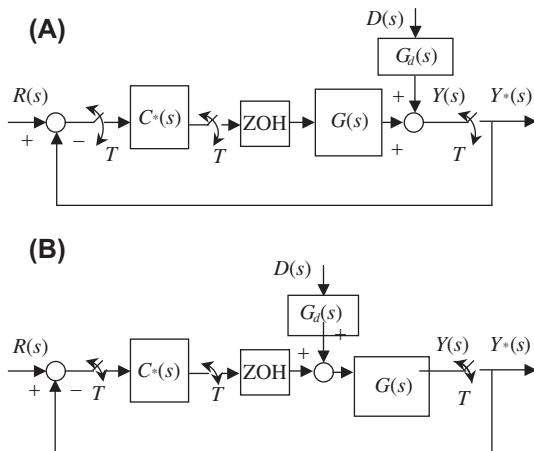


Figure P3.12
Block diagrams for systems with disturbance inputs.

- d. The steady-state error for a unit ramp input
 - (i) $G(z) = \frac{0.4(z+0.2)}{(z-1)(z-0.1)}$
 - (ii) $G(z) = \frac{0.5(z+0.2)}{(z-0.1)(z-0.8)}$

Computer exercises

3.14 For the analog system with a sampling period of 0.05 s

$$G(s) = \frac{10(s+2)}{s(s+5)}$$

- a. Obtain the transfer function for the system with sampled input and output.
 - b. Obtain the transfer function for the system with DAC and ADC.
 - c. Obtain the unit step response of the system with sampled output and analog input.
 - d. Obtain the poles of the systems in (a) and (b), and the output of (c), and comment on the differences between them.
- 3.15 For the system of Problem 3.11
- a. Obtain the transfer function for the analog subsystem with DAC and ADC.
 - b. Obtain the step response of the open-loop analog system and the closed-loop digital control system and comment on the effect of the controller on the time response.
 - c. Obtain the frequency response of the digital control system, and verify that 0.02 s is an acceptable choice of sampling period. Explain briefly why the sampling period is chosen based on the closed-loop rather than the open-loop dynamics.

- 3.16 Consider the internal combustion engine model of Problem 3.6. Assume that for the operational conditions of interest, the time constant τ is approximately 1.2 s, whereas the parameter ε can vary in the range 0.4–0.6. The digital cascade controller

$$C(z) = \frac{0.02z}{z - 1}$$

was selected to improve the time response of the system with unity feedback. Simulate the digital control system with $\varepsilon = 0.4, 0.5$, and 0.6 , and discuss the behavior of the controller in each case.

- 3.17 Simulate the continuous-discrete system discussed in Problem 3.11 and examine the behavior of both the continuous output and the sampled output. Repeat the simulation with a 10% error in the plant gain. Discuss the simulation results, and comment on the effect of the parameter error on disturbance rejection.

Stability of digital control systems

Objectives

After completing this chapter, the reader will be able to do the following:

1. Determine the input–output stability of a z -transfer function.
2. Determine the asymptotic stability of a z -transfer function.
3. Determine the internal stability of a digital feedback control system.
4. Determine the stability of a z -polynomial using the Routh–Hurwitz criterion.
5. Determine the stability of a z -polynomial using the Jury criterion.
6. Determine the stable range of a parameter for a z -polynomial.
7. Determine the closed-loop stability of a digital system using the Nyquist criterion.
8. Determine the gain margin and phase margin of a digital system.

Stability is a basic requirement for digital and analog control systems. Digital control is based on samples and is updated every sampling period, and there is a possibility that the system will become unstable between updates. This obviously makes stability analysis different in the digital case. We examine different definitions and tests of the stability of linear time-invariant (LTI) digital systems based on transfer function models. In particular, we consider input–output stability and internal stability. We provide several tests for stability: the Routh–Hurwitz criterion, the Jury criterion, and the Nyquist criterion. We also define the gain margin and phase margin for digital systems.

Chapter Outline

- 4.1 Definitions of stability 104**
- 4.2 Stable z -domain pole locations 105**
- 4.3 Stability conditions 106**
 - 4.3.1 Asymptotic stability 106
 - 4.3.2 BIBO stability 107
 - 4.3.3 Internal stability 110
- 4.4 Stability determination 114**
 - 4.4.1 MATLAB 114
 - 4.4.2 Routh–Hurwitz criterion 115

4.5 Jury test 117**4.6 Nyquist criterion** 122

4.6.1 Phase margin and gain margin 129

Resources 137**Problems** 138**Computer exercises** 139

4.1 Definitions of stability

The most commonly used definitions of stability are based on the magnitude of the system response in the steady state. If the steady-state response is unbounded, the system is said to be unstable. In this chapter, we discuss two stability definitions that concern the boundedness or exponential decay of the system output. The first stability definition considers the system output due to its initial conditions. To apply it to transfer function models, we need the assumption that no pole-zero cancellation occurs in the transfer function. Reasons for this assumption are given in Chapter 8 and are discussed further in the context of state-space models.

Definition 4.1: Asymptotic stability

A system is said to be asymptotically stable if its response to any initial conditions decays to zero asymptotically in the steady state—that is, the response due to the initial conditions satisfies

$$\lim_{k \rightarrow \infty} y(k) = 0 \quad (4.1)$$

If the response due to the initial conditions remains bounded but does not decay to zero, the system is said to be **marginally stable**.

The second definition of stability concerns the forced response of the system for a bounded input. A bounded input satisfies the condition

$$\begin{aligned} |u(k)| &< b_u, \quad k = 0, 1, 2, \dots \\ 0 &< b_u < \infty \end{aligned} \quad (4.2)$$

For example, a bounded sequence satisfying the constraint $|u(k)| < 3$ is shown in Fig. 4.1.

Definition 4.2: Bounded-input–bounded-output stability

A system is said to be bounded-input–bounded-output (BIBO) stable if its response to any bounded input remains bounded—that is, for any input satisfying (4.2), the output satisfies

$$\begin{aligned} |y(k)| &< b_y, \quad k = 0, 1, 2, \dots \\ 0 &< b_y < \infty \end{aligned} \quad (4.3)$$

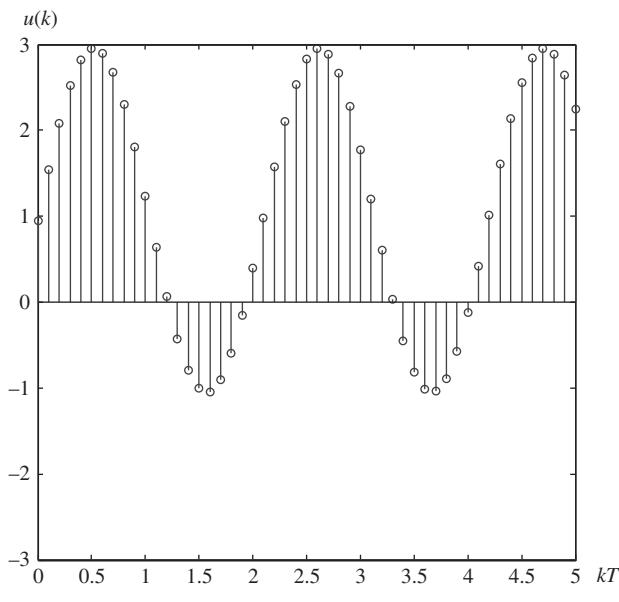


Figure 4.1
Bounded sequence with bound $b_u = 3$.

4.2 Stable z-domain pole locations

The examples provided in Chapter 2 show that the locations of the poles of the system z -transfer function determine the time response. The implications of this fact for system stability are now examined more closely.

Consider the sampled exponential and its z -transform

$$p^k, \quad k = 0, 1, 2, \dots \leftrightarrow \frac{z}{z - p} \quad (4.4)$$

where p is real or complex. Then the time sequence for large k is given by

$$|p|^k \rightarrow \begin{cases} 0, & |p| < 1 \\ 1, & |p| = 1 \\ \infty, & |p| > 1 \end{cases} \quad (4.5)$$

Any time sequence can be described by

$$f(k) = \sum_{i=1}^n A_i p_i^k, \quad k = 0, 1, 2, \dots \leftrightarrow F(z) = \sum_{i=1}^n A_i \frac{z}{z - p_i} \quad (4.6)$$

where A_i are partial fraction coefficients and p_i are z -domain poles. Hence, we conclude that the sequence is bounded if its poles lie in the closed unit disc (i.e., on or inside the unit circle) and decays exponentially if its poles lie in the open unit disc (i.e., inside the unit circle). This conclusion allows us to derive stability conditions based on the locations of the system poles. Note that the case of repeated poles on the unit circle

corresponds to an unbounded time sequence (see, for example, the transform of the sampled ramp).

Although the preceding conclusion is valid for complex as well as real-time sequences, we will generally restrict our discussions to real-time sequences. For real-time sequences, the poles and partial fraction coefficients in Eq. (4.6) are either real or complex conjugate pairs.

4.3 Stability conditions

The analysis of Section 4.2 allows the derivation of conditions for asymptotic and BIBO stability based on transfer function models. It is shown that, in the absence of pole-zero cancellation, conditions for BIBO stability and asymptotic stability are identical.

4.3.1 Asymptotic stability

Theorem 4.1 gives conditions for asymptotic stability.

Theorem 4.1: Asymptotic stability

In the absence of pole-zero cancellation, an LTI digital system is asymptotically stable if its transfer function poles are in the open unit disc and marginally stable if the poles are in the closed unit disc with no repeated poles on the unit circle.

Proof

Consider the LTI system governed by the constant coefficient difference equation

$$\begin{aligned} y(k+n) + a_{n-1}y(k+n-1) + \cdots + a_1y(k+1) + a_0y(k) \\ = b_mu(k+m) + b_{m-1}u(k+m-1) + \cdots + b_1u(k+1) + b_0u(k), \quad k = 0, 1, 2, \dots \end{aligned}$$

with initial conditions $y(0), y(1), y(2), \dots, y(n-1)$. Using the z-transform of the output, we observe that the response of the system due to the initial conditions with the input zero is of the form

$$Y(z) = \frac{N(z)}{z^n + a_{n-1}z^{n-1} + \cdots + a_1z + a_0}$$

where $N(z)$ is a polynomial dependent on the initial conditions. Because transfer function zeros arise from transforming the input terms, they have no influence on the response due to the initial conditions. The denominator of the output z-transform is the same as the denominator of the z-transfer function in the absence of pole-zero cancellation. Hence, the poles of the function $Y(z)$ are the poles of the system transfer function.

$Y(z)$ can be expanded as partial fractions of the form (4.6). Thus, the output due to the initial conditions is bounded for system poles in the closed unit disc with no repeated poles on the unit circle. It decays exponentially for system poles in the open unit disc (i.e., inside the unit circle).

Theorem 4.1 applies even if pole-zero cancellation occurs, provided that the poles that cancel are stable. This follows from the fact that stability testing is essentially a search for unstable poles, with the system being declared stable if none are found. Invisible but stable poles would not lead us to a wrong conclusion. However, invisible but unstable poles would. The next example shows how to determine asymptotic stability using Theorem 4.1.

Example 4.1

Determine the asymptotic stability of the following systems:

- (a) $H(z) = \frac{4(z-2)}{(z-2)(z-0.1)}$
- (b) $H(z) = \frac{4(z-0.2)}{(z-0.2)(z-0.1)}$
- (c) $H(z) = \frac{5(z-0.3)}{(z-0.2)(z-0.1)}$
- (d) $H(z) = \frac{8(z-0.2)}{(z-0.1)(z-1)}$

Solution

Theorem 4.1 can only be used for transfer functions (a) and (b) if their poles and zeros are not canceled. Ignoring the zeros, which do not affect the response to the initial conditions, (a) has a pole outside the unit circle and the poles of (b) are inside the unit circle. Hence, (a) is unstable, whereas (b) is asymptotically stable.

Theorem 4.1 can be applied to the transfer functions (c) and (d). The poles of (c) are all inside the unit circle, and the system is therefore asymptotically stable. However, (d) has one pole on the unit circle and is only marginally stable.

4.3.2 BIBO stability

BIBO stability concerns the response of a system to a bounded input. The response of the system to any input is given by the convolution summation

$$y(k) = \sum_{i=0}^k h(k-i)u(i), \quad k = 0, 1, 2, \dots \quad (4.7)$$

where $h(k)$ is the impulse response sequence.

It may seem that a system should be BIBO stable if its impulse response is bounded. To show that this is generally false, let the impulse response of a linear system be bounded and strictly positive with lower bound b_{h1} and upper bound b_{h2} —that is,

$$0 < b_{h1} < h(k) < b_{h2} < \infty, \quad k = 0, 1, 2, \dots \quad (4.8)$$

Then using the bound (4.8) in (4.7) gives the inequality

$$|y(k)| = \sum_{i=0}^k h(k-i)u(i) > b_{h1} \sum_{i=0}^k u(i), \quad k = 0, 1, 2, \dots \quad (4.9)$$

which is unbounded as $k \rightarrow \infty$ for the bounded input $u(k) = 1, k = 0, 1, 2, \dots$

Theorem 4.2 establishes necessary and sufficient conditions for BIBO stability of a discrete-time linear system.

Theorem 4.2

A discrete-time linear system is BIBO stable if and only if its impulse response sequence is absolutely summable—that is,

$$\sum_{i=0}^{\infty} |h(i)| < \infty \quad (4.10)$$

Proof

- Necessity (only if)** To prove necessity by contradiction, assume that the system is BIBO stable but does not satisfy (4.10). Then consider the input sequence given by

$$u(k-i) = \begin{cases} 1, & h(i) \geq 0 \\ -1, & h(i) < 0 \end{cases}$$

The corresponding output is

$$y(k) = \sum_{i=0}^k |h(i)|$$

which is unbounded as $k \rightarrow \infty$. This contradicts the assumption of BIBO stability.

- Sufficiency (if)** To prove sufficiency, we assume that (4.10) is satisfied and then show that the system is BIBO stable. Using the bound (4.2) in the convolution summation (4.7) gives the inequality

$$|y(k)| \leq \sum_{i=0}^k |h(i)| |u(k-i)| < b_u \sum_{i=0}^k |h(i)|, \quad k = 0, 1, 2, \dots$$

which remains bounded as $k \rightarrow \infty$ if (4.10) is satisfied.

Because the z -transform of the impulse response is the transfer function, BIBO stability can be related to pole locations as shown in Theorem 4.3.

Theorem 4.3

A discrete-time linear system is BIBO stable if and only if the poles of its transfer function lie inside the unit circle.

Proof

Applying (4.6) to the impulse response and transfer function shows that the impulse response is bounded if the poles of the transfer function are in the closed unit disc and decays exponentially if the poles are in the open unit disc. It has already been established that systems with a bounded impulse response that does not decay exponentially are not BIBO stable. Thus, for BIBO stability, the system poles must lie inside the unit circle.

To prove sufficiency, assume an exponentially decaying impulse response (i.e., poles inside the unit circle). Let A_r be the coefficient of largest magnitude and $|p_s| < 1$ be the system pole of largest magnitude in (4.6). Then the impulse response (assuming no repeated poles for simplicity) is bounded by

$$|h(k)| = \left| \sum_{i=1}^n A_i p_i^k \right| \leq \sum_{i=1}^n |A_i| |p_i|^k \leq n |A_r| |p_s|^k, \quad k = 0, 1, 2, \dots$$

Hence, the impulse response decays exponentially at a rate determined by the largest system pole. Substituting in (4.10) gives

$$\sum_{i=0}^{\infty} |h(i)| \leq n |A_r| \sum_{i=0}^{\infty} |p_s|^i = n |A_r| \frac{1}{1 - |p_s|} < \infty$$

Thus, the condition is sufficient by Theorem 4.2.

Example 4.2

Investigate the BIBO stability of the class of systems with the impulse response

$$h(k) = \begin{cases} K, & 0 \leq k \leq m < \infty \\ 0, & \text{elsewhere} \end{cases}$$

where K is a finite constant.

Solution

The impulse response satisfies

$$\sum_{i=0}^{\infty} |h(i)| = \sum_{i=0}^m |h(i)| = (m+1)|K| < \infty$$

Example 4.2—cont'd

Using condition (4.10), the systems are all BIBO stable. This is the class of **finite impulse response** (FIR) systems (i.e., systems whose impulse response is nonzero over a finite interval). Thus, we conclude that all FIR systems are BIBO stable.

Example 4.3

Investigate the BIBO stability of the systems discussed in Example 4.1.

Solution

After pole-zero cancellation, the transfer functions (a) and (b) have all poles inside the unit circle and are therefore BIBO stable. The transfer function (c) has all poles inside the unit circle and is stable; (d) has a pole on the unit circle and is not BIBO stable.

The preceding analysis and examples show that for LTI systems, with no pole-zero cancellation, BIBO and asymptotic stability are equivalent and can be investigated using the same tests. Hence, the term **stability** is used in the sequel to denote either BIBO or asymptotic stability with the assumption of no unstable pole-zero cancellation. Pole locations for a stable system (inside the unit circle) are shown in Fig. 4.2.

4.3.3 Internal stability

So far, we have only considered stability as applied to an open-loop system. For closed-loop systems, these results are applicable to the closed-loop transfer function. However, the stability of the closed-loop transfer function is not always sufficient for proper system operation because some of the internal variables may be unbounded. In a feedback control system, it is essential that all the signals in the loop be bounded when bounded exogenous inputs are applied to the system.

Consider the unity feedback digital control system of Fig. 4.3 where, for simplicity, a disturbance input is added to the controller output before the ADC. We consider that

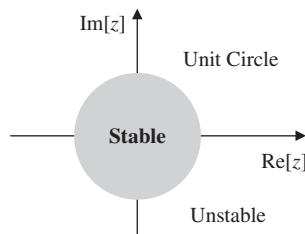


Figure 4.2

Stable pole locations in the z -plane. $\text{Im}[z]$ denotes the imaginary part and $\text{Re}[z]$ denotes the real part of z .

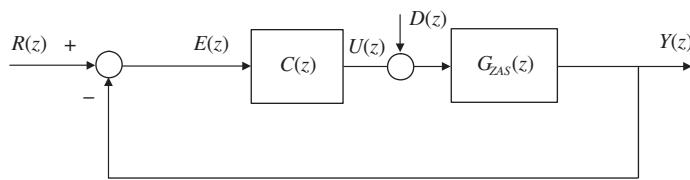


Figure 4.3
Digital control system with disturbance $D(z)$.

system as having two outputs, Y and U , and two inputs, R and D . Thus, the transfer functions associated with the system are given by

$$\begin{bmatrix} Y(z) \\ U(z) \end{bmatrix} = \begin{bmatrix} \frac{C(z)G_{ZAS}(z)}{1 + C(z)G_{ZAS}(z)} & \frac{G_{ZAS}(z)}{1 + C(z)G_{ZAS}(z)} \\ \frac{C(z)}{1 + C(z)G_{ZAS}(z)} & \frac{C(z)G_{ZAS}(z)}{1 + C(z)G_{ZAS}(z)} \end{bmatrix} \begin{bmatrix} R(z) \\ D(z) \end{bmatrix} \quad (4.11)$$

Clearly, it is not sufficient to prove that the output of the controlled system Y is bounded for bounded reference input R because the controller output U can be unbounded. In addition, the system output must be bounded when a different input is applied to the system—namely, in the presence of a disturbance. This suggests the following definition of stability.

Definition 4.3: Internal stability

If all the transfer functions that relate system inputs (R and D) to the possible system outputs (Y and U) are BIBO stable, then the system is said to be internally stable.

Because internal stability guarantees the stability of the transfer function from R to Y , among others, it is obvious that an internally stable system is also externally stable (i.e., the system output Y is bounded when the reference input R is bounded). However, external stability does not, in general, imply internal stability.

We now provide some results that allow us to test internal stability.

Theorem 4.4

The system shown in Fig. 4.3 is internally stable if and only if all its closed-loop poles are in the open unit disc.

Proof

1. **Necessity (only if)** To prove necessity, we write $C(z)$ and $G_{ZAS}(z)$ as ratios of coprime polynomials (i.e., polynomials with no common factors):

$$C(z) = \frac{N_C(z)}{D_C(z)} \quad G_{ZAS}(z) = \frac{N_G(z)}{D_G(z)} \quad (4.12)$$

Substituting in (4.11), we rewrite it as

$$\begin{bmatrix} Y \\ U \end{bmatrix} = \frac{1}{D_C D_G + N_C N_G} \begin{bmatrix} N_C N_G & D_C N_G \\ N_C D_G & -N_C N_G \end{bmatrix} \begin{bmatrix} R \\ D \end{bmatrix} \quad (4.13)$$

where we have dropped the argument z for brevity. If the system is internally stable, then the four transfer functions in (4.11) have no poles on or outside the unit circle. Because of the factorization of (4.12) is coprime, the polynomial $D_C D_G + N_C N_G$ has no zeros that cancel with all four numerators. Thus, the polynomial has no zeros on or outside the unit circle.

2. **Sufficiency (if)** From (4.13), it is evident that if the characteristic polynomial $D_C D_G + N_C N_G$ has no zeros on or outside the unit circle, then all the transfer functions are asymptotically stable and the system is internally stable.

Theorem 4.5

The system in Fig. 4.3 is internally stable if and only if the following two conditions hold:

1. The characteristic polynomial $1 + C(z)G_{ZAS}(z)$ has no zeros on or outside the unit circle.
2. The loop gain $C(z)G_{ZAS}(z)$ has no pole-zero cancellation on or outside the unit circle.

Proof

1. **Necessity (only if)** Condition 1 is clearly necessary by Theorem 4.4. To prove the necessity of condition 2, we first factor $C(z)$ and $G_{ZAS}(z)$ as in (4.12) to write the characteristic polynomial in the form $D_C D_G + N_C N_G$. We also have that $C(z) G_{ZAS}(z)$ is equal to $N_C N_G / D_C D_G$. Assume that condition 2 is violated and that there exists Z_0 , $|Z_0| \geq 1$, which is a zero of $D_C D_G$ as well as a zero of $N_C N_G$. Then clearly Z_0 is also a zero of the characteristic polynomial $D_C D_G + N_C N_G$, and the system is unstable. This establishes the necessity of condition 2.
2. **Sufficiency (if)** By Theorem 4.4, condition 1 implies internal stability unless unstable pole-zero cancellation occurs in the characteristic polynomial $1 + C(z)G_{ZAS}(z)$. We therefore have internal stability if condition 2 implies the absence of unstable pole-zero cancellation. If the loop gain $C(z)G_{ZAS}(z) = N_C N_G / D_C D_G$ has no unstable pole-zero cancellation, then

Proof—cont'd

$1 + C(z)G_{ZAS}(z) = [D_c D_G + N_c N_G]/D_c D_G$ does not have unstable pole-zero cancellation, and the system is internally stable.

Example 4.4

An isothermal chemical reactor where the product concentration is controlled by manipulating the feed flow rate is modeled by the following transfer function¹:

$$G(s) = \frac{0.5848(-0.3549s + 1)}{0.1828s^2 + 0.8627 + 1}$$

Determine $G_{ZAS}(Z)$ with a sampling period $T = 0.1$, and then verify that the closed-loop system with the feedback controller

$$C(z) = \frac{-10(z - 0.8149)(z - 0.7655)}{(z - 1)(z - 1.334)}$$

is not internally stable.

Solution

The discretized process transfer function is

$$G_{ZAS}(z) = (1 - z^{-1})\mathcal{Z}\left\{\frac{G(s)}{s}\right\} = \frac{-0.075997(z - 1.334)}{(z - 0.8149)(z - 0.7655)}$$

The transfer function from the reference input to the system output is given by

$$\frac{Y(z)}{R(z)} = \frac{C(z)G_{ZAS}(z)}{1 + C(z)G_{ZAS}(z)} = \frac{0.75997}{z - 0.24}$$

The system appears to be asymptotically stable with all its poles inside the unit circle. However, the system is not internally stable as seen by examining the transfer function

$$\frac{U(z)}{R(z)} = \frac{C(z)}{1 + C(z)G_{ZAS}(z)} = \frac{-10(z - 0.8149)(z - 0.7655)}{(z - 0.24)(z - 1.334)}$$

which has a pole at 1.334 outside the unit circle. The control variable is unbounded even when the reference input is bounded. In fact, the system violates condition 2 of Theorem 4.5 because the pole at 1.334 cancels in the loop gain

$$C(z)G_{ZAS}(z) = \frac{-10(z - 0.8149)(z - 0.7655)}{(z - 1)(z - 1.334)} \times \frac{-0.075997(z - 1.334)}{(z - 0.8149)(z - 0.7655)}$$

¹ Bequette, B.W., 2003. Process Control: Modeling, Design, and Simulation, Prentice Hall, Upper Saddle River, NJ.

4.4 Stability determination

The simplest method for determining the stability of a discrete-time system given its z-transfer function is by finding the system poles. This can be accomplished using a suitable numerical algorithm based on Newton's method.

4.4.1 MATLAB

The roots of a polynomial are obtained using one of the MATLAB commands

```
>> roots(den)  
>> zpk(g)
```

where **den** is a vector of denominator polynomial coefficients. The command **zpk** factorizes the numerator and denominator of the transfer function **g** and displays it. The poles of the transfer function can be obtained with the command **pole** and then sorted with the command **dsort** in order of decreasing magnitude.

Alternatively, one may use the command **ddamp**, which yields the pole locations (eigenvalues), the damping ratio, and the undamped natural frequency. For example, given a sampling period of 0.1 s and the denominator polynomial with coefficients

```
>> den = [1.0, 0.2, 0.0, 0.4]
```

the command is

```
>> ddamp(den, 0.1)
```

The command yields the output:

Eigenvalue	Magnitude	Equiv. damping	Equiv. freq. (rad/sec)
0.2306 + 0.7428i	0.7778	0.1941	12.9441
0.2306 - 0.7428i	0.7778	0.1941	12.9441
-0.6612	0.6612	0.1306	31.6871

The MATLAB command

```
>> T = feedback(g, gf, ±1)
```

calculates the closed-loop transfer function **T** using the forward transfer function **g** and the feedback transfer function **gf**. For negative feedback, the third argument is **-1** or is omitted. For unity feedback, we replace the argument **gf** with **1**. We can solve for the poles of the closed-loop transfer function as before using **zpk** or **ddamp**.

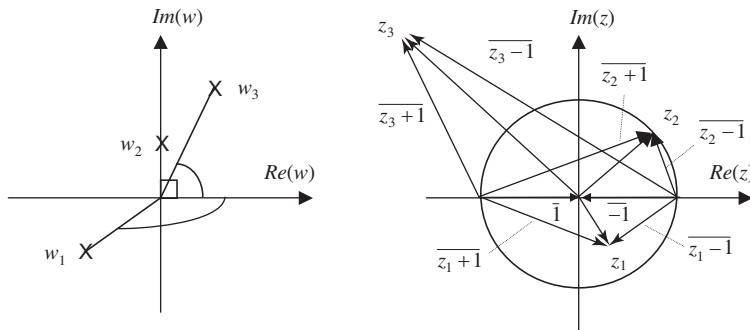


Figure 4.4
Angles associated with bilinear transformation.

4.4.2 Routh–Hurwitz criterion

The Routh–Hurwitz criterion determines conditions for left half plane (LHP) polynomial roots and cannot be directly used to investigate the stability of discrete-time systems. The **bilinear** transformation

$$z = \frac{1+w}{1-w} \Leftrightarrow w = \frac{z-1}{z+1} \quad (4.14)$$

transforms the inside of the unit circle to the LHP.

To verify this property, consider the three cases shown in Fig. 4.4. They represent the mapping of a point in the LHP, a point in the right half plane (RHP), and a point on the jw axis. The angle of w after bilinear transformation is

$$\angle w = \angle(z-1) - \angle(z+1) \quad (4.15)$$

For a point inside the unit circle, the angle of w is of a magnitude greater than 90 degrees, which corresponds to points in the LHP. For a point on the unit circle, the angle is ± 90 degrees, which corresponds to points on the imaginary axis, and for points outside the unit circle, the magnitude of the angle is less than 90 degrees, which corresponds to points in the RHP.

The bilinear transformation allows the use of the Routh–Hurwitz criterion for the investigation of discrete-time system stability. For the general z -polynomial, we have the transform pairs

$$F(z) = a_n z^n + a_{n-1} z^{n-1} + \cdots + a_0 \xrightarrow{z = \frac{1+w}{1-w}} a_n \left(\frac{1+w}{1-w} \right)^n + a_{n-1} \left(\frac{1+w}{1-w} \right)^{n-1} + \cdots + a_0 \quad (4.16)$$

The Routh–Hurwitz approach becomes progressively more difficult as the order of the z -polynomial increases. But for low-order polynomials, it easily gives stability conditions. For high-order polynomials, a symbolic manipulation package can be used to perform the necessary algebraic manipulations. The Routh–Hurwitz approach is demonstrated in Example 4.5.

Example 4.5

Find stability conditions for

1. The first-order polynomial $a_1z + a_0$, $a_1 > 0$
2. The second-order polynomial $a_2z^2 + a_1z + a_0$, $a_2 > 0$

Solution

1. The stability of the first-order polynomial can be easily determined by solving for its root. Hence, the stability condition is

$$\left| \frac{a_0}{a_1} \right| < 1 \quad (4.17)$$

2. The roots of the second-order polynomial are in general given by

$$z_{1,2} = \frac{-a_1 \pm \sqrt{a_1^2 - 4a_0a_2}}{2a_2} \quad (4.18)$$

Thus, it is not easy to determine the stability of the second-order polynomial by solving for its roots. For a monic polynomial (coefficient of z^2 is unity), the constant term is equal to the product of the poles. Hence, for pole magnitudes less than unity, we obtain the necessary stability condition

$$\left| \frac{a_0}{a_2} \right| < 1 \quad (4.19)$$

or equivalently

$$-a_0 < a_2 \text{ and } a_0 < a_2$$

This condition is also sufficient in the case of complex conjugate poles where the two poles are of equal magnitude. The condition is only necessary for real poles because the product of a number greater than unity and a number less than unity can be less than unity. For example, for poles at 0.01 and 10, the product of the two poles has magnitude 0.1, which satisfies (4.19), but the system is clearly unstable.

Substituting the bilinear transformation in the second-order polynomial gives

$$a_2 \left(\frac{1+w}{1-w} \right)^2 + a_1 \left(\frac{1+w}{1-w} \right) + a_0$$

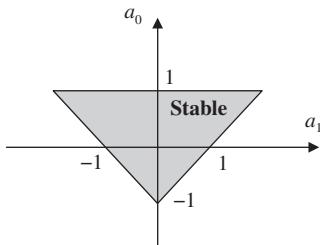
Example 4.5—cont'd

Figure 4.5
Stable parameter range for a second-order z -polynomial.

which reduces to

$$(a_2 - a_1 + a_0)w^2 + 2(a_2 - a_0)w + (a_2 + a_1 + a_0)$$

By the Routh–Hurwitz criterion, it can be shown that the poles of the second-order w -polynomial remain in the LHP if and only if its coefficients are all positive. Hence, the stability conditions are given by

$$\begin{aligned} a_2 - a_1 + a_0 &> 0 \\ a_2 - a_0 &> 0 \\ a_2 + a_1 + a_0 &> 0 \end{aligned} \tag{4.20}$$

Adding the first and third conditions gives

$$a_2 + a_0 > 0 \Rightarrow -a_0 < a_2$$

This condition, obtained earlier in (4.19), is therefore satisfied if the three conditions of (4.20) are satisfied. The reader can verify through numerical examples that if real roots satisfying conditions (4.20) are substituted in (4.18), we obtain roots between -1 and $+1$.

Without loss of generality, the coefficient a_2 can be assumed to be unity, and the stable parameter range can be depicted in the a_0 versus a_1 parameter plane as shown in Fig. 4.5.

4.5 Jury test

It is possible to investigate the stability of z -domain polynomials directly using the **Jury test** for real coefficients or the **Schur–Cohn test** for complex coefficients. These tests involve determinant evaluations as in the Routh–Hurwitz test for s -domain polynomials but are more time-consuming. The Jury test is given next.

Theorem 4.6

For the polynomial

$$F(z) = a_n z^n + a_{n-1} z^{n-1} + \cdots + a_1 z + a_0 = 0, \quad a_n > 0 \quad (4.21)$$

the roots of the polynomial are inside the unit circle if and only if

- $$\begin{aligned} (1) \quad & F(1) > 0 \\ (2) \quad & (-1)^n F(-1) > 0 \\ (3) \quad & |a_0| < a_n \\ (4) \quad & |b_0| > |b_{n-1}| \\ (5) \quad & |c_0| > |c_{n-2}| \\ & \vdots \\ (n+1) \quad & |r_0| > |r_2| \end{aligned} \quad (4.22)$$

where the terms in the $n + 1$ conditions are calculated from [Table 4.1](#).

The entries of the table are calculated as follows:

$$\begin{aligned} b_k &= \begin{vmatrix} a_0 & a_{n-k} \\ a_n & a_k \end{vmatrix}, \quad k = 0, 1, \dots, n-1 \\ c_k &= \begin{vmatrix} b_0 & b_{n-k-1} \\ b_{n-1} & b_k \end{vmatrix}, \quad k = 0, 1, \dots, n-2 \\ & \vdots \\ r_0 &= \begin{vmatrix} s_0 & s_3 \\ s_3 & s_0 \end{vmatrix}, \quad r_1 = \begin{vmatrix} s_0 & s_2 \\ s_3 & s_1 \end{vmatrix}, \quad r_2 = \begin{vmatrix} s_0 & s_1 \\ s_3 & s_2 \end{vmatrix} \end{aligned} \quad (4.23)$$

Table 4.1: Jury table.

Row	z^0	z^1	z^2	...	z^{n-k}	...	z^{n-1}	z^n
1	a_0	a_1	a_2	...	a_{n-k}	...	a_{n-1}	a_n
2	a_n	a_{n-1}	a_{n-2}	...	a_k	...	a_1	a_0
3	b_0	b_1	b_2	...	b_{n-k}	...	b_{n-1}	
4	b_{n-1}	b_{n-2}	b_{n-3}	...	b_k	...	b_0	
5	c_0	c_1	c_2	c_{n-2}	
6	c_{n-2}	c_{n-3}	c_{n-4}	c_0	
.		
.		
.		
2 $n-5$	s_0	s_1	s_2	s_3		
2 $n-4$	s_3	s_2	s_1	s_0		
2 $n-3$	r_0	r_1	r_2					

Based on the Jury table and the Jury stability conditions, we make the following observations:

1. The first row of the Jury table is a listing of the coefficients of the polynomial $F(z)$ in order of increasing power of z .
2. The number of rows of table 2 $n - 3$ is always odd, and the coefficients of each even row are the same as the odd row directly above it with the order of the coefficients reversed.
3. There are $n + 1$ conditions in (4.22) that correspond to the $n + 1$ coefficients of $F(z)$.
4. Condition 3 through $n + 1$ of (4.22) are calculated using the coefficient of the first column of the Jury table together with the last coefficient of the preceding row. The middle coefficient of the last row is never used and need not be calculated.
5. Conditions 1 and 2 of (4.22) are calculated from $F(z)$ directly. If one of the first two conditions is violated, we conclude that $F(z)$ has roots on or outside the unit circle without the need to construct the Jury table or test the remaining conditions.
6. Condition 3 of (4.22), with $a_n = 1$, requires the constant term of the polynomial to be less than unity in magnitude. The constant term is simply the product of the roots and must be smaller than unity for all the roots to be inside the unit circle.
7. Condition (4.22) reduce to conditions (4.19) and (4.20) for first and second-order systems, respectively, where the Jury table is simply one row.
8. For higher-order systems, applying the Jury test by hand is laborious, and it is preferable to test the stability of a polynomial $F(z)$ using a computer-aided design (CAD) package.
9. If the coefficients of the polynomial are functions of system parameters, the Jury test can be used to obtain the stable ranges of the system parameters.

Example 4.6

Test the stability of the polynomial.

$$F(z) = z^5 + 2.6z^4 - 0.56z^3 - 2.05z^2 + 0.0775z + 0.35 = 0$$

We compute the entries of the Jury table using the coefficients of the polynomial (Table 4.2).

Table 4.2: Jury table for Example 4.6.

Row	z^0	z^1	z^2	z^3	z^4	z^5
1	0.35	0.0775	-2.05	-0.56	2.6	1
2	1	2.6	-0.56	-2.05	0.0775	0.35
3	-0.8775	-2.5729	-0.1575	1.854	0.8325	
4	0.8325	1.854	-0.1575	-2.5729	-0.8775	
5	0.0770	0.7143	0.2693	0.5151		
6	0.5151	0.2693	0.7143	0.0770		
7	-0.2593	-0.0837	-0.3472			

Example 4.6—cont'd

The first two conditions require the evaluation of $F(z)$ at $z = \pm 1$.

1. $F(1) = 1 + 2.6 - 0.56 - 2.05 + 0.0775 - 0.35 = 1.4175 > 0$
2. $(-1)^5 F(-1) = (-1)(-1 + 2.6 + 0.56 - 2.05 - 0.0775 + 0.35) = -0.3825 < 0$
Conditions 3 through 6 can be checked quickly using the entries of the first column of the Jury table.
3. $|0.35| < 1$
4. $|-0.8775| > |0.8325|$
5. $|0.0770| < |0.5151|$
6. $|-0.2593| < |-0.3472|$

Conditions 2, 5, and 6 are violated, and the polynomial has roots on or outside the unit circle. In fact, the polynomial can be factored as

$$F(z) = (z - 0.7)(z - 0.5)(z + 0.5)(z + 0.8)(z + 2.5) = 0$$

and has a root at -2.5 outside the unit circle. Note that the number of conditions violated is not equal to the number of roots outside the unit circle and that condition 2 is sufficient to conclude the instability of $F(z)$.

Example 4.7

Find the stable range of the gain K for the unity feedback digital cruise control system of Example 3.2 with the analog plant transfer function

$$G(s) = \frac{K}{s+3}$$

and with digital-to-analog converter (DAC) and analog-to-digital converter (ADC) if the sampling period is 0.02s .

Solution

The transfer function for analog subsystem ADC and DAC is

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z} \left\{ \mathcal{L}^{-1} \left[\frac{G(s)}{s} \right] \right\} \\ &= (1 - z^{-1}) \mathcal{Z} \left\{ \mathcal{L}^{-1} \left[\frac{K}{s(s+3)} \right] \right\} \end{aligned}$$

Using the partial fraction expansion

$$\frac{K}{s(s+3)} = \frac{K}{3} \left[\frac{1}{s} - \frac{1}{s+3} \right]$$

Example 4.7—cont'd

we obtain the transfer function

$$G_{ZAS}(z) = \frac{1.9412 \times 10^{-2}K}{z - 0.9418}$$

For unity feedback, the closed-loop characteristic equation is

$$1 + G_{ZAS}(z) = 0$$

which can be simplified to

$$z - 0.9418 + 1.9412 \times 10^{-2}K = 0$$

The stability conditions are

$$\begin{aligned} 0.9418 - 1.9412 \times 10^{-2}K &< 1 \\ -0.9418 + 1.9412 \times 10^{-2}K &< 1 \end{aligned}$$

Thus, the stable range of K is

$$-3 < K < 100.03$$

Example 4.8

Find the stable range of the gain K for the vehicle position control system (see Example 3.3) with the analog plant transfer function

$$G(s) = \frac{10K}{s(s + 10)}$$

and with DAC and ADC if the sampling period is 0.05s.

Solution

The transfer function for analog subsystem ADC and DAC is

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z} \left\{ \mathcal{L}^{-1} \left[\frac{G(s)}{s} \right] \right\} \\ &= (1 - z^{-1}) \mathcal{Z} \left\{ \mathcal{L}^{-1} \frac{10K}{s^2(s + 10)} \right\} \end{aligned}$$

Using the partial fraction expansion

$$\frac{10K}{s^2(s + 10)} = 0.1K \left[\frac{10}{s^2} - \frac{1}{s} + \frac{1}{s + 10} \right]$$

we obtain the transfer function

$$G_{ZAS}(z) = \frac{1.0653 \times 10^{-2}K(z + 0.8467)}{(z - 1)(z - 0.6065)}$$

Example 4.8—cont'd

For unity feedback, the closed-loop characteristic equation is

$$1 + G_{ZAS}(z) = 0$$

which can be simplified to

$$\begin{aligned} & (z - 1)(z - 0.6065) + 1.0653 \times 10^{-2}K(z + 0.8467) \\ &= z^2 + (1.0653 \times 10^{-2}K - 1.6065)z + 0.6065 + 9.02 \times 10^{-3}K = 0 \end{aligned}$$

The stability conditions are

1. $F(1) = 1 + (1.0653 \times 10^{-2}K - 1.6065) + 0.6065 + 9.02 \times 10^{-3}K > 0 \Leftrightarrow K > 0$
2. $F(-1) = 1 - (1.0653 \times 10^{-2}K - 1.6065) + 0.6065 + 9.02 \times 10^{-3}K > 0 \Leftrightarrow K < 1967.582$
3. $|0.6065 + 0.0902K| < 1 \Leftrightarrow + (0.6065 + 0.0902K) < 1 - (0.6065 + 0.0902K) < 1$
 $\Leftrightarrow -178.104 < K < 43.6199$

The three conditions yield the stable range

$$0 < K < 43.6199$$

4.6 Nyquist criterion

The Nyquist criterion allows us to answer two questions:

1. Does the system have closed-loop poles outside the unit circle?
2. If the answer to the first question is yes, how many closed-loop poles are outside the unit circle?

The Nyquist criterion is particularly useful in situations where the frequency response can be obtained experimentally and used to obtain a polar or Nyquist plot.

We begin by considering the closed-loop characteristic polynomial

$$p_{cl}(z) = 1 + C(z)G(z) = 1 + L(z) = 0 \quad (4.24)$$

where $L(z)$ denotes the loop gain. We rewrite the characteristic polynomial in terms of the numerator of the loop gain N_L and its denominator D_L in the form

$$p_{cl}(z) = 1 + \frac{N_L(z)}{D_L(z)} = \frac{N_L(z) + D_L(z)}{D_L(z)} \quad (4.25)$$

We observe that the zeros of the rational function are the closed-loop poles, whereas its poles are the open-loop poles of the system. We assume that we are given the number of open-loop poles outside the unit circle, and we denote this number by P . The number of closed-loop poles outside the unit circle is denoted by Z and is unknown.

To determine Z , we use some simple results from complex analysis. We first give the following definition.

Definition 4.4: Contour

A contour is a closed directed simple (does not cross itself) curve.

An example of a contour is shown in Fig. 4.6. In the figure, shaded text denotes a vector. Recall that in the complex plane the vector connecting any point a to a point z is the vector $(z - a)$. We can calculate the net angle change for the term $(z - a)$ as the point z traverses the contour in the shown (counterclockwise) direction by determining the net number of rotations of the corresponding vector. From Fig. 4.6, we observe that the net rotation is one full turn or 360 degrees for the point a_1 , which is inside the contour. The net rotation is zero for the point a_2 , which is outside the contour. If the point in question corresponds to a zero, then the rotation gives a numerator angle; if it is a pole, we have a denominator angle. The net angle change for a rational function is the change in the angle of the numerator minus the change in the angle of the denominator. So for Fig. 4.6, we have one counterclockwise rotation because of a_1 and no rotation as a result of a_2 for a net angle change of one counterclockwise rotation. Angles are typically measured in the counterclockwise direction and we count clockwise rotations as negative.

The preceding discussion shows how to determine the number of zeros of a rational function in a specific region; given the number of poles, we perform the following steps:

1. Select a closed contour surrounding the region.
2. Compute the net angle change for the function as we traverse the contour once.
3. The net angle change or number of counterclockwise rotations N is equal to the number of closed-loop poles inside the contour Z_{in} minus the angle of open-loop poles inside the contour P_{in} , that is, the number of counterclockwise encirclements of the origin is given by

$$N = Z_{in} - P_{in} \quad (4.26)$$

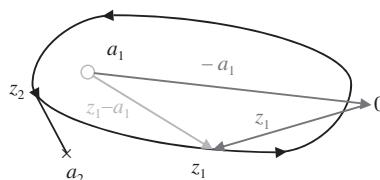


Figure 4.6
Closed contours (shaded letters denote vectors).

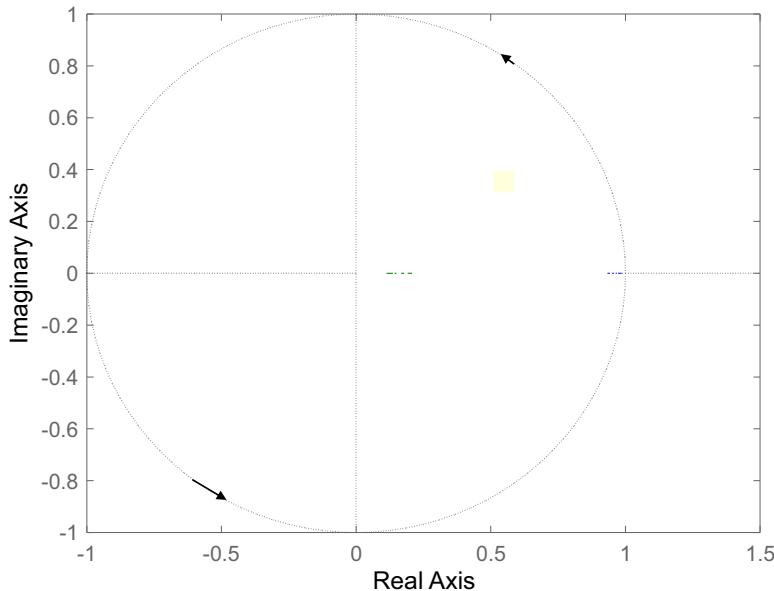


Figure 4.7
Contour for stability determination.

To use this to determine closed-loop stability, we need to count the closed-loop poles outside the unit circle. Franklin et al. observed that for an n th order system

$$N = Z_{in} - P_{in} = (n - Z) - (n - P) = P - Z$$

We can therefore determine the number of closed-loop poles as

$$Z = P - N$$

The stability contour is the unit circle shown in Fig. 4.7 and is traversed in the counterclockwise direction.

The value of the loop gain on the unit circle is $L(e^{j\omega T})$, which is the frequency response of the discrete-time system for angles ωT in the interval $[-\pi, \pi]$. The values obtained for negative frequencies $L(e^{-j\omega T})$ are simply the complex conjugate of the values $L(e^{j\omega T})$ and need not be separately calculated. Because the order of the numerator is equal to that of the denominator or less, points on the large circle map to zero or to a single constant value. **Note that traversing the contour counterclockwise implies increasing ω** and following the direction typically shown on Nyquist plots.

We can simplify the test by plotting $L(e^{j\omega T})$ as we traverse the contour and then counting its encirclements of the point $-1 + j0$. As Fig. 4.8 shows, this is equivalent to plotting $p_{cl}(z)$ and counting encirclements of the origin.

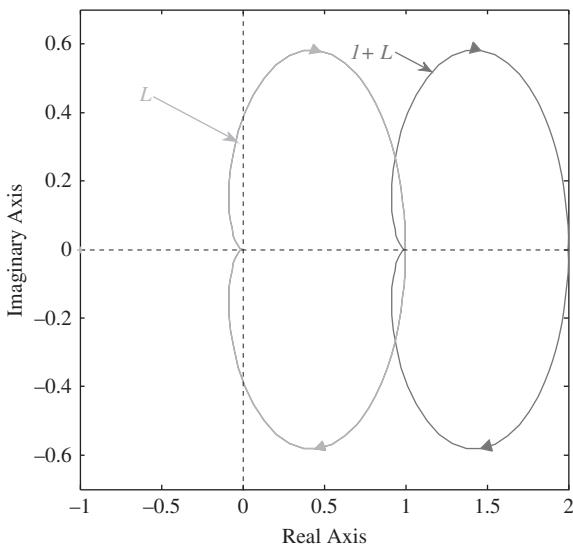


Figure 4.8
Nyquist plots of L and $1 + L$.

If the system has open-loop poles on the unit circle, the contour passes through poles and the test fails. To avoid this, we modify the contour to avoid these open-loop poles. The most common case is a pole at unity for which the modified contour is shown in Fig. 4.9. The contour includes an additional circular arc of infinitesimal radius. The contour reduces to the one shown in Fig. 4.10. For m poles at unity, the loop gain is given by

$$L(z) = \frac{N_L(s)}{(z-1)^m D(s)} \quad (4.27)$$

where N_L and D have no unity roots. The value of the transfer function on the circular arc is approximately given by

$$L(z)|_{z \rightarrow 1+\varepsilon e^{j\theta}} = \frac{K}{\varepsilon e^{jm\theta}}, \quad \theta \in \left[-\frac{\pi}{2}, \frac{\pi}{2}\right] \quad (4.28)$$

where K is equal to $N_L(1)/D(1)$ and $(z-1) = \varepsilon e^{j\theta}$, with ε the radius of the small circular arc. Therefore, the small circle maps to a large circle, and traversing the small circle once causes

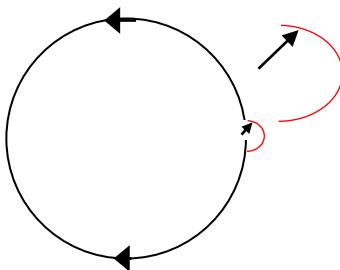


Figure 4.9
Modified contour for stability determination.

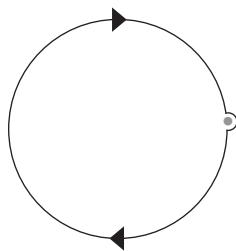


Figure 4.10

Simplification of the modified contour for stability determination with $(-1, 0)$ shown in gray.

a net denominator angle change of $-m\pi$ radians (counterclockwise) (i.e., m half circles). The net angle change for the quotient on traversing the small circular arc is thus $m\pi$ radians (clockwise). We conclude that for a type m system, the Nyquist contour will include m large clockwise semicircles. We now summarize the results obtained in Theorem 4.7.

Theorem 4.7: Nyquist criterion

Let the number of counterclockwise encirclements of the point $(-1, 0)$ for a loop gain $L(z)$ when traversing the stability contour be N (i.e., $-N$ for clockwise encirclements), where $L(z)$ has P open-loop poles inside the contour. Then the system has Z closed-loop poles outside the unit circle with Z given by

$$Z = (-N) + P \quad (4.29)$$

Corollary

An open-loop stable system is closed-loop stable if and only if its Nyquist plot does not encircle the point $(-1, 0)$ (i.e., if $N=0$).

Although counting encirclements appears complicated, it is actually quite simple using the following recipe:

1. Starting at a distant point, move toward the point $(-1, 0)$.
2. Count all lines of the stability contour crossed. Count each line with an arrow pointing from your left to your right as negative and every line with an arrow pointing from your right to your left as positive.
3. The net number of lines counted is equal to the number of clockwise encirclements of the point $(-1, 0)$.

The recipe is demonstrated in the following two examples.

Example 4.9

Consider a linearized model of a furnace:

$$G(s) = \frac{T_i(s)}{U(s)} = \frac{K g_{rw} g_{iw}}{s^2 + (2g_{iw} + g_{rw})s + g_{rw}g_{iw}}$$

During the heating phase, we have the model²

$$G(s) = \frac{1}{s^2 + 3s + 1}$$

Determine the closed-loop stability of the system with digital control and a sampling period of 0.01.

Solution

In Example 3.5, we determined the transfer function of the system to be

$$G_{ZAS}(z) = K \left[\frac{(p_1 e^{-p_2 T} - p_2 e^{-p_1 T} + p_2 - p_1)z + (p_1 e^{-p_1} - p_2 e^{-p_2} - p_2 e^{-(p_1+p_2)} - p_1 e^{-(p_1+p_2)})}{p_1 p_2 (p_2 - p_1)(z - e^{-p_1 T})(z - e^{-p_2 T})} \right]$$

Substituting the values of the parameters in the general expression gives the z-transfer function

$$G_{ZAS}(z) = 10^{-5} \frac{4.95z + 4.901}{z^2 - 1.97z + 0.9704}$$

The Nyquist plot of the system is shown in Fig. 4.11. The plot shows that the Nyquist plot does not encircle the point $(-1, 0)$ —that is, $N=0$. Because the open-loop transfer function

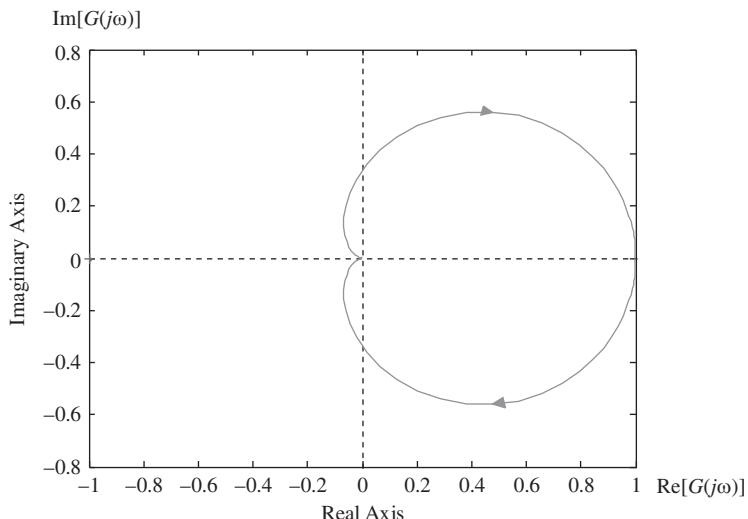


Figure 4.11
Nyquist plot of the furnace transfer function.

Example 4.9—cont'd

has no unstable poles ($P = 0$), the system is closed-loop stable ($Z = 0$). Note that there is a free path to the point $(-1, 0)$ without crossing any lines of the plot because there are no encirclements.

² Hagglund, T., Tengall, A., 1995. An automatic tuning procedure for unsymmetrical processes. In: Proceedings of the European Control Conference, Rome, 1995.

Example 4.10

Use the Nyquist criterion to investigate the stability of the type-1 transfer function

$$G(z) = \frac{10}{(z-1)(z-0.1)}$$

Solution

The Nyquist plot of the system is shown in Fig. 4.12. For increasing ω , we move on the plot in the direction of the arrow, which corresponds to the Nyquist contour. If we include one large clockwise semicircle, we obtain two clockwise encirclements, and we have

$$Z = -N = 2 \text{ unstable poles}$$

Alternatively, starting at a far point and proceeding to the point $(-1, 0)$, we cross two lines that go from our right to our left (in the direction of the arrows). This indicates two clockwise encirclements and two unstable poles.

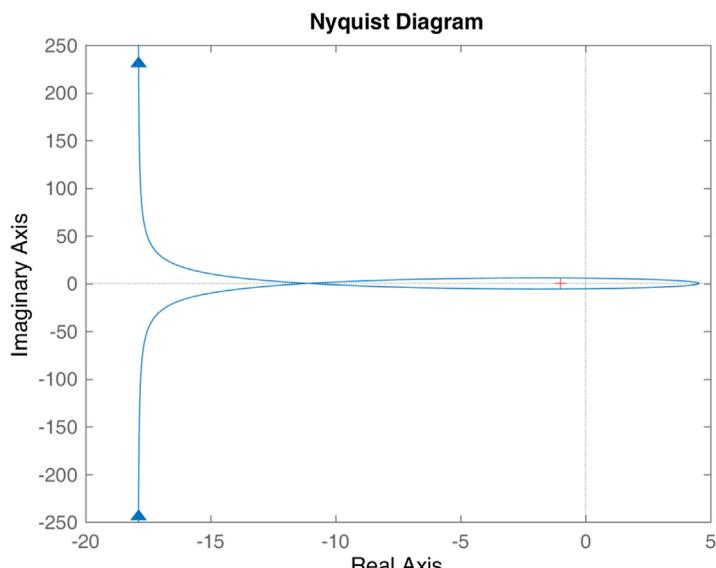


Figure 4.12
Nyquist plot of the system of Example 4.10.

4.6.1 Phase margin and gain margin

In practice, the stability of a mathematical model is not sufficient to guarantee acceptable system performance or even to guarantee the stability of the physical system that the model represents. This is because of the approximate nature of mathematical models and our imperfect knowledge of the system parameters. We therefore need to determine how far the system is from instability. This degree of stability is known as **relative stability**. To keep our discussion simple, we restrict it to open-loop stable systems where zero encirclements guarantee stability. For open-loop stable systems that are nominally closed-loop stable, the distance from instability can be measured by the distance between the set of points of the Nyquist plot and the point $(-1, 0)$.

Typically, the distance between a set of points and the single point $(-1, 0)$ is defined as the minimum distance over the set of points. However, it is more convenient to define relative stability in terms of two distances: a magnitude distance and an angular distance. The two distances are given in the following definitions.

Definition 4.5: Gain margin

The gain margin is the gain perturbation that makes the system marginally stable.

Definition 4.6: Phase margin

The phase margin is the negative phase perturbation that makes the system marginally stable.

The two stability margins become clearer by examining the block diagram of Fig. 4.13. If the system has a multiplicative gain perturbation $\Delta G(z) = \Delta K$, then the gain margin is the magnitude of ΔK that makes the system on the verge of instability. If the system has a multiplicative gain perturbation $\Delta G(z) = e^{-j\Delta\theta}$, then the gain margin is the lag angle $\Delta\theta$ that makes the system on the verge of instability. Clearly, the perturbations corresponding to the gain margin and phase margin are limiting values, and satisfactory behavior would require smaller model perturbations.

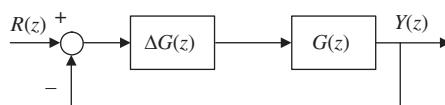


Figure 4.13
Model perturbation $\Delta G(z)$.

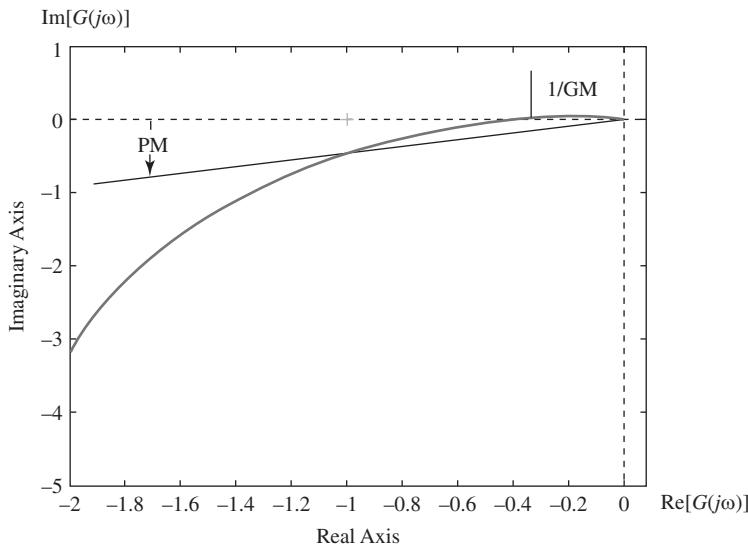
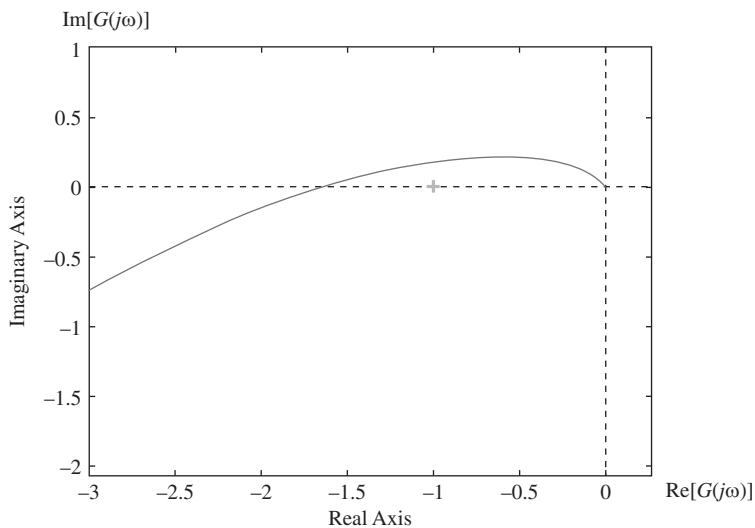


Figure 4.14
Nyquist plot with phase margin and gain margin.

The Nyquist plot of Fig. 4.14 shows the gain margin and phase margin for a given polar plot (the positive frequency portion of the Nyquist plot). Recall that each point on the plot represents a complex number, which is represented by a vector from the origin. Scaling the plot with a gain ΔK results in scaled vectors without rotation. Thus, the vector on the negative real axis is the one that reaches the point $(-1, 0)$ if appropriately scaled, and the magnitude of that vector is the reciprocal of the gain margin. On the other hand, multiplication by $e^{-j\Delta\theta}$ rotates the plot clockwise without changing the magnitudes of the vectors, and it is the vector of magnitude unity that can reach the point $(-1, 0)$ if rotated by the phase margin.

For an unstable system, a counterclockwise rotation or a reduction in gain is needed to make the system on the verge of instability. The system will have a negative phase margin and a gain margin less than unity, which is also negative if it is expressed in decibels—that is, in units of $20 \log\{|G(j\omega)|\}$. The polar plot of a system with negative gain margin and phase margin is shown in Fig. 4.15.

The gain margin can be obtained analytically by equating the imaginary part of the frequency response to zero and solving for the real part. The phase margin can be obtained by equating the magnitude of the frequency response to unity and solving for angle, and then adding 180 degrees. However, because only approximate values are needed in practice, it is easier to use MATLAB to obtain both margins. In some cases, the intercept with the real axis can be obtained as the value where $z = -1$ provided that the system has no pole at -1 (i.e., the frequency response has no discontinuity at the folding frequency $\omega_s/2$).

**Figure 4.15**

Nyquist plot with negative gain margin (dBs) and phase margin.

It is sometimes convenient to plot the frequency response using the Bode plot, but the availability of frequency response plotting commands in MATLAB reduces the necessity for such plots. The MATLAB commands for obtaining frequency response plots (which work for both continuous-time and discrete-time systems) are

```
>> nyquist(gd) %Nyquist plot
>> bode(gd) %Bode plot
```

It is also possible to find the gain and phase margins with the command

```
>> [gm, pm] = margin(gd)
```

An alternative form of the command is

```
>> margin(gd)
```

The latter form shows the gain margin and phase margin on the Bode plot of the system. We can also obtain the phase margin and gain margin using the Nyquist plot by clicking on the plot and selecting

Characteristics
All stability margins

The concepts of the gain margin and phase margin and their evaluation using MATLAB are illustrated in the following example.

Example 4.11

Determine the closed-loop stability of the digital control system for the furnace model of Example 4.9 with a discrete-time first-order actuator of the form

$$G_a(z) = \frac{0.9516}{z - 0.9048}$$

and a sampling period of 0.01. If an amplifier of gain $K=5$ is added to the actuator, how does the value of the gain affect closed-loop stability?

Solution

We use MATLAB to obtain the z-transfer function of the plant and actuator:

$$G_a(z)G_{ZAS}(z) = 10^{-5} \frac{4.711z + 4.644}{z^3 - 2.875z^2 + 2.753z - 0.8781}$$

The Nyquist plot for the system, Fig. 4.16, is obtained with no additional gain and then for a gain $K=5$. We also show the plot in the vicinity of the point $(-1, 0)$ in Fig. 4.17, from which we see that the system with $K=5$ encircles the point twice clockwise.

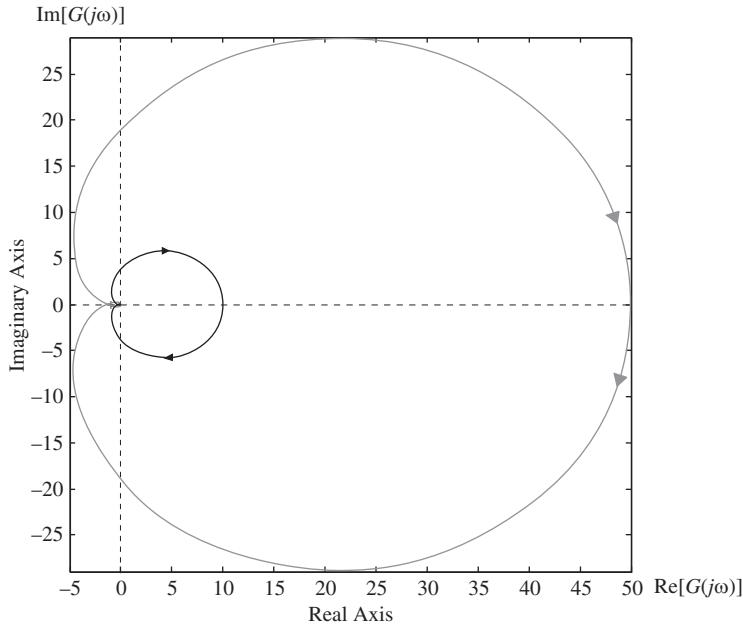
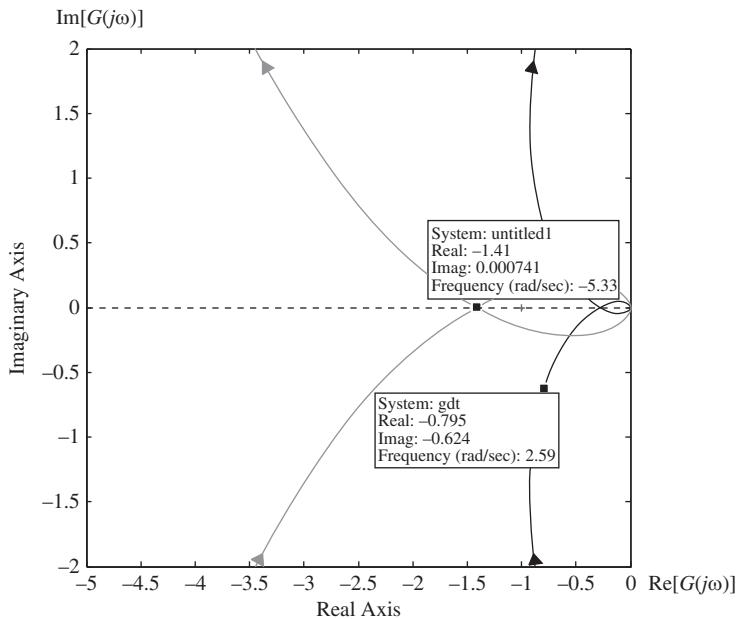


Figure 4.16

Nyquist plot for the furnace and actuator ($K=1$, black, $K=5$, gray).

Example 4.11—cont'd**Figure 4.17**

Nyquist plot for the furnace and actuator in the vicinity of the point $(-1, 0)$ ($K=1$, black, $K=5$, gray).

We count the encirclements by starting away from the point $(-1, 0)$ and counting the lines crossed as we approach it. We cross the gray curve twice, and at each crossing the arrow for increasing ω is moving from our right to our left (i.e., two counterclockwise encirclements). The system is unstable, and the number of closed-loop poles outside the unit circle is given by

$$\begin{aligned} Z &= (-N) + P \\ &= 2 + 0 \end{aligned}$$

For the original gain of unity, the intercept with the real axis is at a magnitude of approximately 0.28 and can be increased by a factor of about 3.5 before the system becomes unstable.

Example 4.11—cont'd

At a magnitude of unity, the phase is about 38 degrees less negative than the instability value of -180 degrees. We therefore have a gain margin of about 3.5 and a phase margin of about 38 degrees. Using MATLAB, we find approximately the same values for the margins:

$\gg [gm, pm] = \text{margin}(\text{gtd})$

$gm = 3.4817$

$pm = 37.5426$

Thus, an additional gain of over three or an additional phase lag of over 37 degrees can be tolerated without causing instability. However, such perturbations may cause a significant deterioration in the time response of the system. Perturbations in gain and phase may actually occur upon implementing the control, and the margins are needed for successful implementation. In fact, the phase margin of the system is rather low, and a controller may be needed to improve the response of the system.

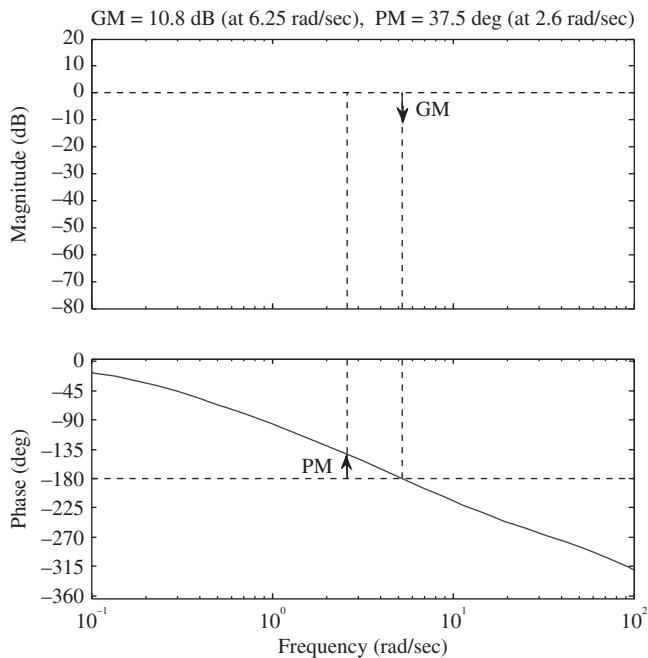


Figure 4.18

Phase margin and gain margin for the oven control system shown on the Bode plot.

To obtain the Bode plot showing the phase margin and gain margin of Fig. 4.18, we use the following command:

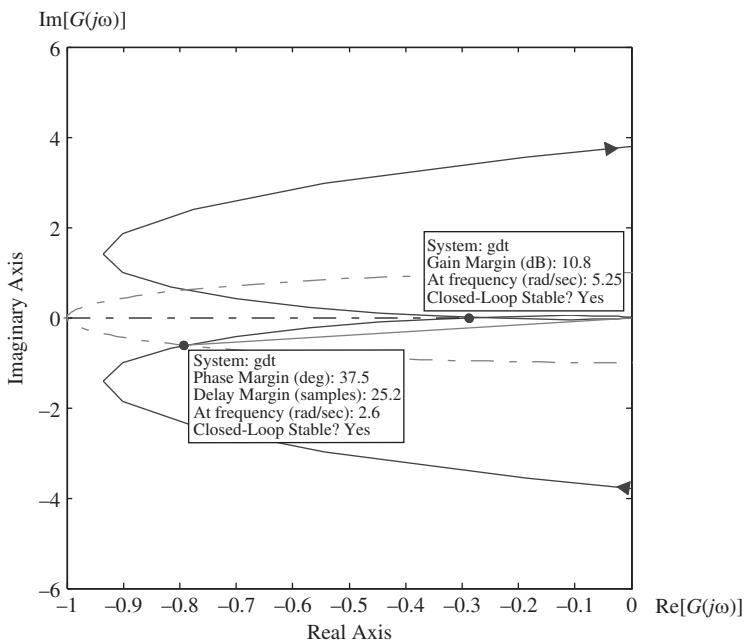
$\gg \text{margin}(\text{gtd})$

Example 4.11—cont'd

1. The phase margin for unity gain shown on the plot is as obtained with the first form of the command **margin**, but the gain margin is in dBs. The values are nevertheless identical as verified with the MATLAB command:

```
>> 20 * log10(gm)
ans = 10.8359
```

2. The gain margin and phase margin can also be obtained using the Nyquist command as shown in Fig. 4.19.

**Figure 4.19**

Phase margin and gain margin for the oven control system shown on the Nyquist plot.

Example 4.12

Determine the closed-loop stability of the digital control system for the position control system with analog transfer function

$$G(s) = \frac{10}{s(s+1)}$$

and with a sampling period of 0.01. If the system is stable, determine the gain margin and the phase margin.

Example 4.12—cont'd**Solution**

We first obtain the transfer function for the analog plant with ADC and DAC. The transfer function is given by

$$G_{ZAS}(z) = 4.983 \times 10^{-4} \frac{z + 0.9967}{(z - 1)(z - 0.99)}$$

Note that the transfer function has a pole at unity because the analog transfer function has a pole at the origin or is type I. Although such systems require the use of the modified Nyquist contour, this has no significant impact on the steps required for stability testing using the Nyquist criterion. The Nyquist plot obtained using the MATLAB command **nyquist** is shown in Fig. 4.20.

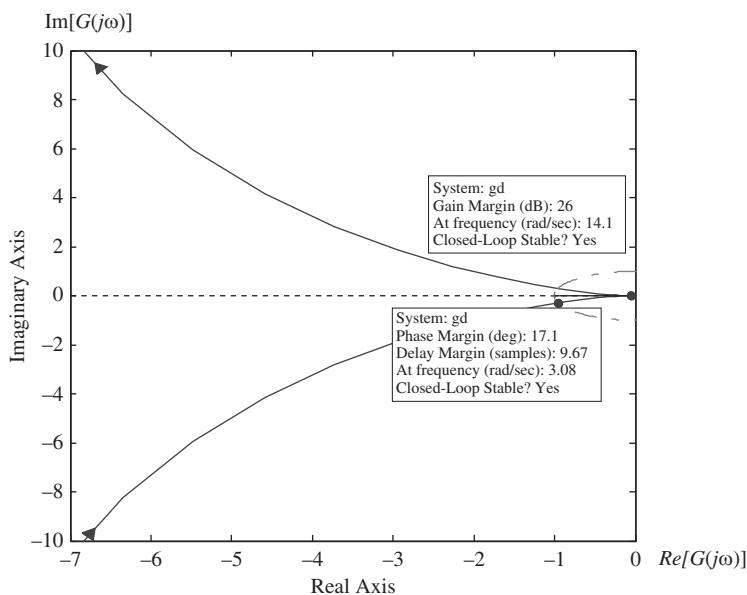
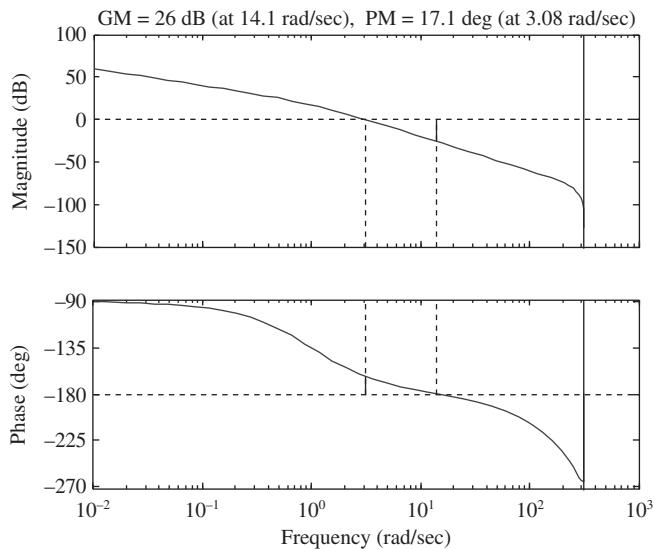


Figure 4.20
Nyquist plot for the position control system of Example 4.12.

The plot does not include the large semicircle corresponding to the small semicircle on the modified contour in Fig. 4.9. However, this does not prevent us from investigating stability. It is obvious that the contour does not encircle the point $(-1, 0)$ because the point is to the left of the observer moving along the polar plot (lower half). In addition, we can reach the $(-1, 0)$ point without crossing any of the lines of the Nyquist plot. The system is stable because the number of closed-loop poles outside the unit circle is given by

$$\begin{aligned} Z &= (-N) + P \\ &= 0 + 0 = 0 \end{aligned}$$

The gain margin is 17.1 degree, and the phase margin is 26dB. The gain margin and phase margin can also be obtained using the margin command as shown in Fig. 4.21.

Example 4.12—cont'd**Figure 4.21**

Bode diagram with phase margin and gain margin for the position control system of Example 4.12.

Resources

- Franklin, G.F., Powell, J.D., Workman, M.L., 1990. Digital Control of Dynamic Systems. Addison-Wesley, Reading, MA.
- Gupta, S.C., Hasdorff, L., 1970. Fundamentals of Automatic Control. Wiley, New York.
- Jury, E.I., 1973. Theory and Applications of the z-Transform Method. Krieger, Huntington, NY.
- Kuo, B.C., 1992. Digital Control Systems. Saunders, Fort Worth, TX.
- Ogata, K., 1987. Digital Control Engineering. Prentice Hall, Englewood Cliffs, NJ.
- Oppenheim, A.V., Willsky, A.S., Young, I.T., 1983. Signals and Systems. Prentice Hall, Englewood Cliffs, NJ.
- Ragazzini, J.R., Franklin, G.F., 1958. Sampled-Data Control Systems. McGraw-Hill, New York.

Problems

- 4.1 Determine the asymptotic stability and the BIBO stability of the following systems:
- $y(k+2) + 0.8y(k+1) + 0.07y(k) = 2u(k+1) + 0.2u(k)$ $k = 0, 1, 2, \dots$
 - $y(k+2) - 0.8y(k+1) + 0.07y(k) = 2u(k+1) + 0.2u(k)$ $k = 0, 1, 2, \dots$
 - $y(k+2) + 0.1y(k+1) + 0.9y(k) = 3.0u(k)$ $k = 0, 1, 2, \dots$
- 4.2 Biochemical reactors are used in different processes such as waste treatment and alcohol fermentation. By considering the dilution rate as the manipulated variable and the biomass concentration as the measured output, the biochemical reactor can be modeled by the following transfer function in the vicinity of an unstable steady-state operating point³:

$$G(s) = \frac{5.8644}{-5.888s + 1}$$

Determine $G_{ZAS}(z)$ with a sampling rate $T = 0.1$, and then consider the feedback controller

$$C(z) = -\frac{z - 1.017}{z - 1}$$

Verify that the resulting feedback system is not internally stable.

- 4.3 The transfer function of a furnace is in the form (see Example 3.5)

$$1 + G_{ZAS}(z) = 1 + K \frac{z - a}{(z - e^{p_1 T})(z - e^{p_2 T})} = 0$$

$$z^2 + [K - (e^{p_1 T} + e^{p_2 T})]z + e^{(p_1 + p_2)T} - Ka = 0$$

Show that the conditions for the stability of the system with unity feedback are

$$(1 + e^{p_1 T})(1 + e^{p_2 T}) - K(1 + a) > 0$$

$$(1 - e^{p_1 T})(1 - e^{p_2 T}) + K(1 - a) > 0$$

$$e^{(p_1 + p_2)T} < 1 + Ka$$

- 4.4 Use the Routh–Hurwitz criterion to investigate the stability of the following systems:
- $G(z) = \frac{5(z-2)}{(z-0.1)(z-0.8)}$
 - $G(z) = \frac{10(z+0.1)}{(z-0.7)(z-0.9)}$
- 4.5 Repeat Problem 4.4 using the Jury criterion.

³ Bequette, B.W., 2003. Process Control: Modeling, Design, and Simulation. Prentice Hall, Upper Saddle River, NJ.

- 4.6 Obtain the impulse response for the systems shown in Problem 4.4, and verify the results obtained using the Routh–Hurwitz criterion. Also determine the exponential rate of decay for each impulse response sequence.
- 4.7 Use the Routh–Hurwitz criterion to find the stable range of K for the closed-loop unity feedback systems with loop gain
- $G(z) = \frac{K(z-1)}{(z-0.1)(z-0.8)}$
 - $G(z) = \frac{K(z+0.1)}{(z-0.7)(z-0.9)}$
- 4.8 Repeat Problem 4.7 using the Jury criterion.
- 4.9 Use the Jury criterion to determine the stability of the following polynomials:
- $z^5 + 0.2z^4 + z^2 + 0.3z - 0.1 = 0$
 - $z^5 - 0.25z^4 + 0.1z^3 + 0.4z^2 + 0.3z - 0.1 = 0$
- 4.10 Determine the stable range of the parameter a for the closed-loop unity feedback systems with loop gain
- $G(z) = \frac{1.1(z-1)}{(z-a)(z-0.8)}$
 - $G(z) = \frac{1.2(z+0.1)}{(z-a)(z-0.9)}$
- 4.11 For a gain of 0.5, derive the gain margin and phase margin of the systems of Problem 4.7 analytically. Let $T=1$ with no loss of generality because the value of ωT in radians is all that is needed for the solution. Explain why the phase margin is not defined for the system shown in Problem 4.7(a). Hint: The gain margin is obtained by finding the point where the imaginary part of the frequency response is zero. The phase margin is obtained by finding the point where the magnitude of the frequency response is unity.

Computer exercises

- 4.12 Write a computer program to perform the Routh-Hurwitz test using a suitable CAD tool.
- 4.13 Write a computer program to perform the Jury test using a suitable CAD tool.
- 4.14 Write a computer program that uses the Jury test program of Exercise 4.13 to determine the stability of a system with an uncertain gain K in a given range $[K_{min}, K_{max}]$. Verify the answers obtained for Problem 4.7 using your program.
- 4.15 Show how the program written for Exercise 4.14 can be used to test the stability of a system with uncertain zero location. Use the program to test the effect of a $\pm 20\%$ variation in the location of the zero for the systems shown in Problem 4.7, with a fixed gain equal to half the critical value.
- 4.16 Show how the program written for Exercise 4.14 can be used to test the stability of a system with uncertain pole location. Use the program to test the effect of a $\pm 20\%$

variation in the location of the first pole for the systems shown in Problem 4.7, with a fixed gain equal to half the critical value.

- 4.17 Simulate the closed-loop systems shown in Problem 4.7 with a unit step input and
 (a) gain K equal to half the critical gain and (b) gain K equal to the critical gain.
 Discuss their stability using your simulation results.
- 4.18 For unity gain, obtain the Nyquist plots of the systems shown in Problem 4.7 using MATLAB and determine the following:
- The intersection with the real axis using the Nyquist plot and then using the Bode plot
 - The stable range of positive gains K for the closed-loop unity feedback systems
 - The gain margin and phase margin for a gain $K=0.5$
- 4.19 For twice the nominal gain, use MATLAB to obtain the Nyquist and Bode plots of the systems of the furnace control system of Example 4.11 with a sampling period of 0.01 and determine the following:
- The intersection with the real axis using the Nyquist plot and then using the Bode plot
 - The stable range of additional positive gains K for the closed-loop unity feedback systems
 - The gain margin and phase margin for twice the nominal gain
- 4.20 In many applications, there is a need for accurate position control at the nanometer scale. This is known as **nanopositioning** and is now feasible because of advances in nanotechnology. The following transfer function represents a single-axis nano-positioning system⁴:

$$G(s) = \frac{4.29 \times 10^{10} (s^2 + 631.2s + 9.4 \times 10^6)}{(s^2 + 178.2s + 6 \times 10^6)(s^2 + 412.3s + 16 \times 10^6)} \\ \times \frac{(s^2 + 638.8s + 45 \times 10^6)}{(s^2 + 209.7s + 56 \times 10^6)(s + 5818)}$$

- Obtain the DAC—analog system—ADC transfer function for a sampling period of 100ms, and determine its stability using the Nyquist criterion.
- Obtain the DAC—analog system—ADC transfer function for a sampling period of 1ms, and determine its stability using the Nyquist criterion.
- Plot the closed-loop step response of the system of (b), and explain the stability results of (a) and (b) based on your plot.

⁴ Sebastian, A., Salapaka, S.M., 2005. Design methodologies of robust nano-positioning. IEEE Trans. Control Syst. Technol. 13 (6), 868–876.

Analog control system design

Objectives

After completing this chapter, the reader will be able to do the following:

1. Obtain root locus plots for analog systems.
2. Characterize a system's step response based on its root locus plot.
3. Design proportional (P), proportional-derivative (PD), proportional-integral (PI), and proportional-integral-derivative (PID) controllers in the s -domain.
4. Tune PID controllers using the Ziegler-Nichols approach.

Analog controllers can be implemented using analog components or approximated with digital controllers using standard analog-to-digital transformations. In addition, direct digital control system design in the z -domain is very similar to the s -domain design of analog systems. Thus, a review of classical control design is the first step toward understanding the design of digital control systems. This chapter reviews the design of analog controllers in the s -domain and prepares the reader for the digital controller design methods presented in Chapter 6. The reader is assumed to have had some exposure to the s -domain and its use in control system design.

Chapter Outline

- 5.1 Root locus 142**
- 5.2 Root locus using MATLAB 146**
- 5.3 Design specifications and the effect of gain variation 147**
- 5.4 Root locus design 149**
 - 5.4.1 Proportional control 151
 - 5.4.2 Proportional-derivative (PD) control 152
 - 5.4.3 Proportional-integral (PI) control 162
 - 5.4.4 Proportional-integral-derivative (PID) control 168
- 5.5 Empirical tuning of PID controllers 171**
- References 176**
- Further reading 176**
- Problems 176**
- Computer exercises 178**

5.1 Root locus

The **root locus** method provides a quick means of predicting the closed-loop behavior of a system based on its open-loop **poles** and **zeros**. The method is based on the properties of the closed-loop **characteristic equation**

$$1 + KL(s) = 0 \quad (5.1)$$

where the gain K is a design parameter and $L(s)$ is the loop gain of the system. We assume a loop gain of the form

$$L(s) = \frac{\prod_{i=1}^{n_z} (s - z_i)}{\prod_{j=1}^{n_p} (s - p_j)} \quad (5.2)$$

where $z_i, i = 1, 2, \dots, n_z$, are the open-loop system zeros and $p_j, j = 1, 2, \dots, n_p$ are the open-loop system poles. It is required to determine the loci of the closed-loop poles of the system (root loci) as K varies between zero and infinity.¹ Because of the relationship between pole locations and the time response, this gives a preview of the closed-loop system behavior for different K .

The complex equality (5.1) is equivalent to the two real equalities:

- *Magnitude condition* $K|L(s)| = 1$
- *Angle condition* $\angle L(s) = \pm(2m + 1)180^\circ, m = 0, 1, 2, \dots$

Using (5.1) or the preceding conditions, the following rules for sketching root loci can be derived:

1. The number of root locus branches is equal to the number of open-loop poles of $L(s)$.
2. The root locus branches start at the open-loop poles and end at the open-loop zeros or at infinity.
3. The real axis root loci have an odd number of poles plus zeros to their right.
4. The branches going to infinity asymptotically approach the straight lines defined by the angle

$$\theta_a = \pm \frac{(2m + 1)180^\circ}{n_p - n_z}, \quad m = 0, 1, 2, \dots \quad (5.3)$$

and the intercept

¹ In rare cases, negative gain values are allowed, and the corresponding root loci are obtained. Negative root loci are not addressed in this text.

$$\sigma_a = \frac{\sum_{i=1}^{n_p} p_i - \sum_{j=1}^{n_z} z_j}{n_p - n_z} \quad (5.4)$$

5. **Breakaway points** (points of departure from the real axis) correspond to local maxima of K , whereas **break-in points** (points of arrival at the real axis) correspond to local minima of K .
6. The **angle of departure** from a complex pole p_n is given by

$$180^\circ - \sum_{i=1}^{n_p-1} \angle(p_n - p_i) + \sum_{j=1}^{n_z} \angle(p_n - z_j) \quad (5.5)$$

The **angle of arrival** at a complex zero is similarly defined.

Example 5.1

Sketch the root locus plots for the loop gains

1. $L(s) = \frac{1}{(s+1)(s+3)}$
2. $L(s) = \frac{1}{(s+1)(s+3)(s+5)}$
3. $L(s) = \frac{s+5}{(s+1)(s+3)}$

Comment on the effect of adding a pole or a zero to the loop gain.

Solution

The root loci for the three loop gains as obtained using MATLAB are shown in Fig. 5.1. We now discuss how these plots can be sketched using root locus sketching rules.

1. Using rule 1, the function has two root locus branches. By rule 2, the branches start at -1 and -3 and go to infinity. By rule 3, the real axis locus is between (-1) and (-3) . Rule 4 gives the asymptote angles

$$\begin{aligned} \theta_a &= \pm \frac{(2m+1)180^\circ}{2}, \quad m = 0, 1, 2, \dots \\ &= \pm 90^\circ, \quad \pm 270^\circ, \dots \end{aligned}$$

and the intercept

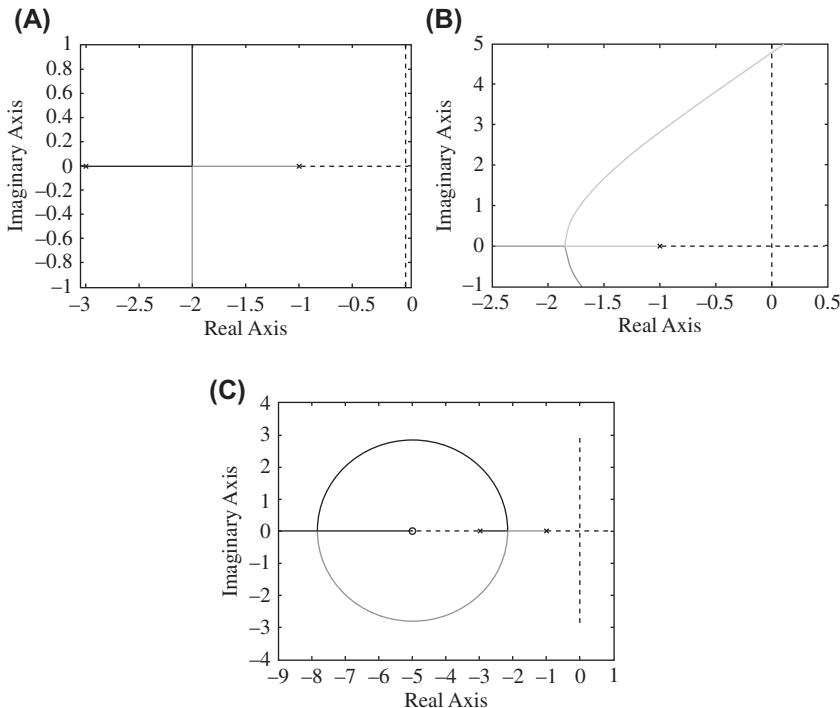
$$\sigma_a = \frac{-1 - 3}{2} = -2$$

To find the breakaway point using Rule 5, we express real K using the characteristic equation as

$$K = -(\sigma + 1)(\sigma + 3) = -(\sigma^2 + 4\sigma + 3)$$

We then differentiate with respect to σ and equate to zero for a maximum to obtain

$$\frac{dK}{d\sigma} = 2\sigma + 4 = 0$$

Example 5.1—cont'd**Figure 5.1**

Root loci of second- and third-order systems. (A) Root locus of a second-order system.
(B) Root locus of a third-order system. (C) Root locus of a second-order system with zero.

Hence, the breakaway point is at $\sigma_b = -2$. This corresponds to a maximum of K because the second derivative is equal to -2 (negative). It can be easily shown that for any system with only two real axis poles, the breakaway point is midway between the two poles.

2. The root locus has three branches, with each branch starting at one of the open-loop poles $(-1, -3, -5)$. The real axis loci are between -1 and -3 and to the left of -5 . The branches all go to infinity, with one branch remaining on the negative real axis and the other two breaking away. The breakaway point is given by the maximum of the real gain K

$$K = -(\sigma + 1)(\sigma + 3)(\sigma + 5)$$

Differentiating gives

$$\begin{aligned} -\frac{dK}{d\sigma} &= (\sigma + 1)(\sigma + 3) + (\sigma + 3)(\sigma + 5) + (\sigma + 1)(\sigma + 5) \\ &= 3\sigma^2 + 18\sigma + 23 \\ &= 0 \end{aligned}$$

which yields $\sigma_b = -1.845$ or -4.155 . The first value is the desired breakaway point because it lies on the real axis locus between the poles and -1 and -3 . The second value

Example 5.1—cont'd

corresponds to a negative gain value and is therefore inadmissible. The gain at the breakaway point can be evaluated from the magnitude condition and is given by

$$K = -(-1.845 + 1)(-1.845 + 3)(-1.845 + 5) = 3.079$$

The asymptotes are defined by the angles

$$\begin{aligned}\theta_a &= \pm \frac{(2m+1)180^\circ}{3}, \quad m = 0, 1, 2, \dots \\ &= \pm 60^\circ, \quad \pm 180^\circ, \dots\end{aligned}$$

and the intercept by

$$\sigma_a = \frac{-1 - 3 - 5}{3} = -3$$

The closed-loop characteristic equation

$$s^3 + 9s^2 + 25s + 15 + K = 0$$

corresponds to the **Routh table**

s^3	1	23
s^2	9	$15 + K$
s^1	$\frac{192 - K}{9}$	
s^0	$15 + K$	

Thus, at $K = 192$, a zero row results. This value defines the auxiliary equation

$$9s^2 + 207 = 0$$

Thus, the intersection with the $j\omega$ -axis is $\pm j4.796$ rad/s. The intersection can also be obtained by factorizing the characteristic polynomial at the critical gain

$$s^2 + 9s^2 + 25s + 15 + 192 = (s + 9)(s + j4.796)(s - j4.796)$$

3. The root locus has two branches as in (1), but now one of the branches ends at the zero. From the characteristic equation, the gain is given by

$$K = -\frac{(\sigma + 1)(\sigma + 3)}{\sigma + 5}$$

Differentiating gives

$$\begin{aligned}\frac{dK}{d\sigma} &= -\frac{(\sigma + 1 + \sigma + 3)(\sigma + 5) - (\sigma + 1)(\sigma + 3)}{(\sigma + 5)^2} \\ &= -\frac{\sigma^2 + 10\sigma + 17}{(\sigma + 5)^2} \\ &= 0\end{aligned}$$

Example 5.1—cont'd

which yields $\sigma_b = -2.172$ or -7.828 . The first value is the breakaway point because it lies between the poles, whereas the second value is to the left of the zero and corresponds to the break-in point. The second derivative

$$\begin{aligned}\frac{d^2K}{d\sigma^2} &= \frac{(2\sigma + 10)(\sigma + 5) - 2(\sigma^2 + 10\sigma + 17)}{(\sigma + 5)^3} \\ &= -16/(\sigma + 5)^3\end{aligned}$$

is negative for the first value and positive for the second value. Hence, K has a maximum at the first value and a minimum at the second. It can be shown that the root locus is a circle centered at the zero with radius given by the geometric mean of the distances between the zero and the two real poles.

Clearly, adding a pole pushes the root locus branches toward the right hand plane (RHP), whereas adding a zero pulls them back into the left hand plane (LHP). Thus, adding a zero allows the use of higher gain values without destabilizing the system. In practice, the allowable increase in gain is limited by the cost of the associated increase in control effort and by the possibility of driving the system outside the linear range of operation.

5.2 Root locus using MATLAB

While the above rules together with (5.1) allow the sketching of root loci for any loop gain of the form (5.2), it is often sufficient to use a subset of these rules to obtain the root loci. For higher-order or more complex situations, it is easier to use a CAD tool like MATLAB. These packages do not actually use root locus sketching rules. Instead, they numerically solve for the roots of the characteristic equation as K is varied in a given range and then display the root loci.

The MATLAB command to obtain root locus plots is “**rlocus**.” To obtain the root locus of the system

$$G(s) = \frac{s + 5}{s^2 + 2s + 10}$$

using MATLAB enter

```
>> g = tf([1,5],[1,2,10]);
>> rlocus(g);
```

To obtain specific points on the root locus and the corresponding data, we simply click on the root locus. Dragging the mouse allows us to change the referenced point to obtain more data.

5.3 Design specifications and the effect of gain variation

The objective of control system design is to construct a system that has a desirable response to standard inputs. A desirable transient response is one that is sufficiently fast without excessive oscillations. A desirable steady-state response is one that follows the desired output with sufficient accuracy. In terms of the response to a unit step input, the transient response is characterized by the following criteria:

1. *Time constant τ .* Time required to reach about 63% of the final value.
2. *Rise time T_r .* Time to go from 10% to 90% of the final value.
3. *Percentage overshoot (PO).*

$$PO = \frac{\text{Peak value} - \text{Final value}}{\text{Final value}} \times 100\%$$

4. *Peak time T_p .* Time to first peak of an oscillatory response.
5. *Settling time T_s .* Time after which the oscillatory response remains within a specified percentage (usually 2%) of the final value.

Clearly, the percentage overshoot and the peak time are intended for use with an oscillatory response (i.e., for a system with at least one pair of complex conjugate poles). For a single complex conjugate pair, these criteria can be expressed in terms of the pole locations.

Consider the second-order system

$$G(s) = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2} \quad (5.6)$$

where ζ is the damping ratio and ω_n is the undamped natural frequency. Then criteria 3 through 5 are given by

$$PO = e^{-\frac{\pi\zeta}{\sqrt{1-\zeta^2}}} \times 100\% \quad (5.7)$$

$$T_p = \frac{\pi}{\omega_d} = \frac{\pi}{\omega_n \sqrt{1 - \zeta^2}} \quad (5.8)$$

$$T_s \approx \frac{4}{\zeta\omega_n} \quad (5.9)$$

From (5.7) and (5.9), the damping ratio ζ is an indicator of the oscillatory nature of the response, with excessive oscillations occurring at low ζ values. Hence, ζ is used as a measure of the relative stability of the system. From (5.8), the time to first peak drops as the undamped natural frequency ω_n increases. Hence, ω_n is used as a measure of speed of response. For higher-order systems, these measures and Eqs. (5.7) through (5.9) can provide approximate answers if the time response is dominated by a single pair of complex conjugate poles. This occurs if additional poles and zeros are far in the left half plane or almost cancel. For systems with zeros, the percentage overshoot is higher than predicted by (5.7) unless the zero is located far in the LHP or almost cancels with a pole. However, Eq. (5.7) through (5.9) are always used in design because of their simplicity.

Thus, the design process reduces to the selection of pole locations and the corresponding behavior in the time domain. The root locus summarizes information on the time response of a closed-loop system as dictated by the pole locations in a single plot. Together with the previously stated criteria, it provides a powerful design tool, as demonstrated by the next example.

Example 5.2

Discuss the effect of gain variation on the time response of the position control system described in Example 3.3 with

$$L(s) = \frac{1}{s(s+p)}$$

Solution

The root locus of the system is similar to that of Example 5.1(1), and it is shown in Fig. 5.2A for $p = 4$. As the gain K is increased, the closed-loop system poles become complex conjugate; then the damping ratio ζ decreases progressively. Thus, the relative stability of the system deteriorates for high gain values. However, large gain values are required to reduce the steady-state error of the system due to a unit ramp, which is given by

$$e(\infty)\% = \frac{100}{K_v} = \frac{100p}{K}$$

In addition, increasing K increases the undamped natural frequency ω_n (i.e., the magnitude of the pole), and hence the speed of response of the system increases. Thus, the chosen gain must be a compromise value that is large enough for a low steady-state error and an acceptable speed of response but small enough to avoid excessive oscillations. The time response for a gain of 10 is shown in Fig. 5.2B.

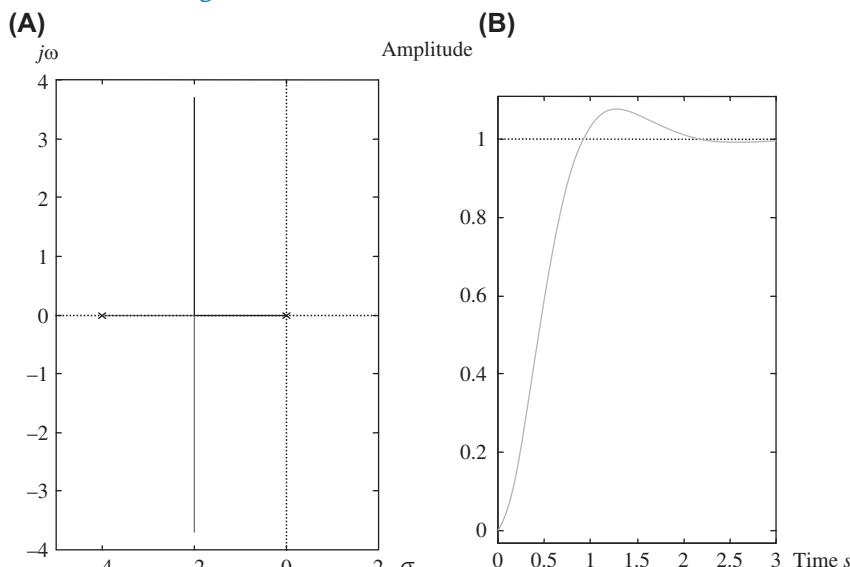


Figure 5.2

Use of the root locus in the design of a second-order system. (A) Root locus for $p = 4$. (B) Step response for $K = 10$.

Example 5.2 illustrates an important feature of design—namely, that it typically involves a compromise between conflicting requirements. The designer must always remember this when selecting design specifications so as to avoid overoptimizing some design criteria at the expense of others.

Note that the settling time of the system does not change in this case when the gain is increased. For a higher-order system or a system with a zero, this is usually not the case. Yet, for simplicity, the second-order Eq. (5.7) through (5.9) are still used in design. The designer must always be alert to errors that this simplification may cause. In practice, design is an iterative process where the approximate results from the second-order approximation are checked and, if necessary, the design is repeated until satisfactory results are obtained.

5.4 Root locus design

Laplace transformation of a time function yields a function of the complex variable s that contains information about the transformed time function. We can therefore use the poles of the s -domain function to characterize the behavior of the time function without inverse transformation. Fig. 5.3 shows pole locations in the s -domain and the associated time functions. Real poles are associated with an exponential time response that decays for LHP poles and increases for RHP poles. The magnitude of the pole determines the rate of exponential change. A pole at the origin is associated with a unit step. Complex conjugate poles are associated with an oscillatory response that decays exponentially for LHP poles and increases exponentially for RHP poles. The real part of the pole determines the rate of exponential change, and the imaginary part determines the frequency of oscillations.

Imaginary axis poles are associated with sustained oscillations.

The objective of control system design in the s -domain is to indirectly select a desirable time response for the system through the selection of closed-loop pole locations. The simplest means of shifting the system poles is through the use of an amplifier or proportional controller. If this fails, then the pole locations can be more drastically altered by adding a dynamic controller with its own open-loop poles and zeros.

As Examples 5.1 and 5.2 illustrate, adding a zero to the system allows the improvement of its time response because it pulls the root locus into the LHP. Adding a pole at the origin increases the type number of the system and reduces its steady-state error but may adversely affect the transient response. If an improvement of both transient and steady-state performance is required, then it may be necessary to add two zeros as well as a pole at the origin. At times, more complex controllers may be needed to achieve the desired design objectives.

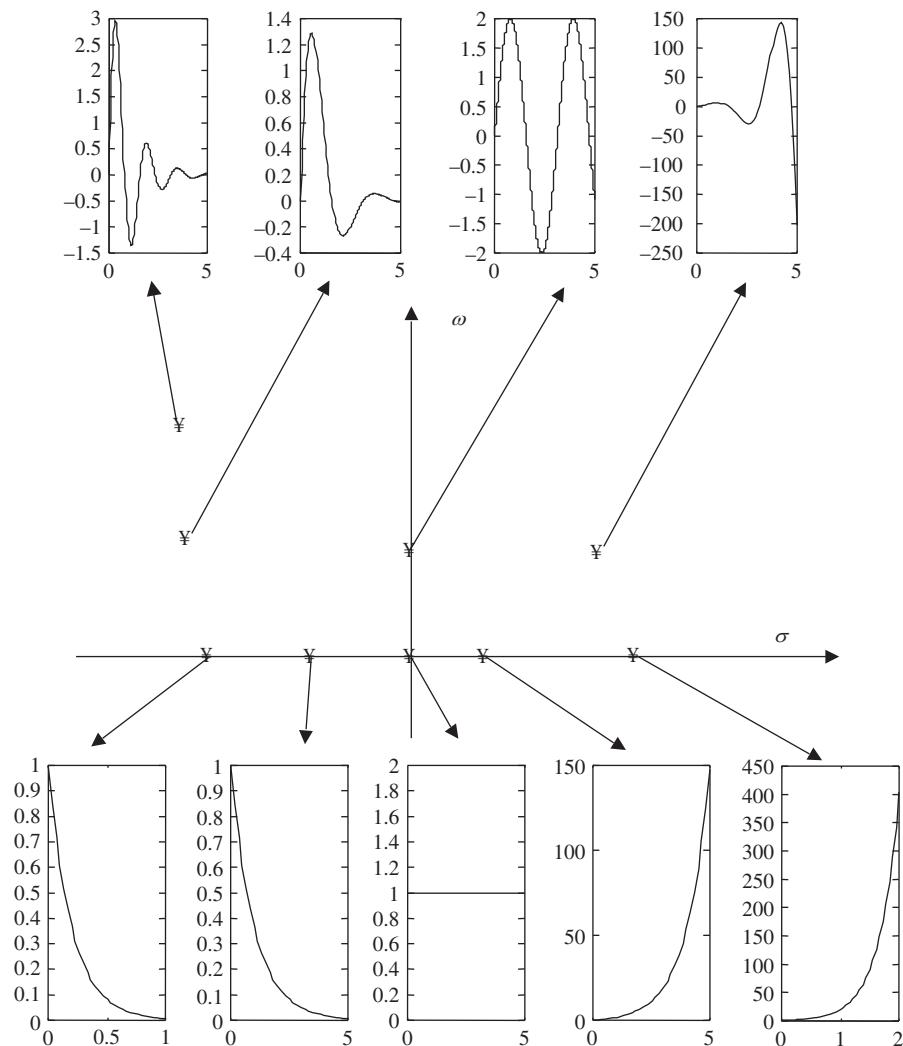


Figure 5.3
Pole locations and the associated time responses.

The controller could be added in the forward path, in the feedback path, or in an inner loop. A prefilter could also be added before the control loop to allow more freedom in design. Several controllers could be used simultaneously, if necessary, to meet all the design specifications. Examples of these control configurations are shown in Fig. 5.4.

In this section, we review the design of analog controllers. We restrict the discussion to proportional (P), proportional-derivative (PD), proportional-integral (PI), and proportional-integral-derivative (PID) control. Similar procedures can be developed for the design of lead, lag, and lag-lead controllers.

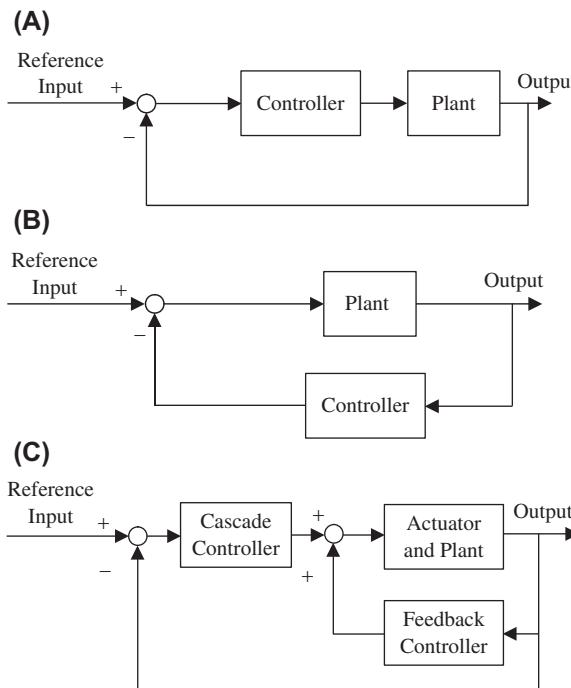


Figure 5.4

Control configurations. (A) Cascade compensation. (B) Feedback compensation. (C) Inner-loop feedback compensation.

5.4.1 Proportional control

Gain adjustment or proportional control allows the selection of closed-loop pole locations from among the poles given by the root locus plot of the system loop gain. For lower-order systems, it is possible to design proportional control systems analytically, but a sketch of the root locus is still helpful in the design process, as seen from Example 5.3.

Example 5.3

A position control system (see Example 3.3) with load angular position as output and motor armature voltage as input consists of an armature-controlled DC motor driven by a power amplifier together with a gear train. The overall transfer function of the system is

$$G(s) = \frac{1}{s(s+p)}$$

Design a proportional controller for the system to obtain

1. A specified damping ratio ζ
2. A specified undamped natural frequency ω_n

Example 5.3—cont'd**Solution**

The root locus of the system was discussed in Example 5.2 and shown in Fig. 5.2A. The root locus remains in the LHP for all positive gain values. The closed-loop characteristic equation of the system is given by

$$s(s+p) + K = s^2 + 2\zeta\omega_n s + \omega_n^2 = 0$$

Equating coefficients gives

$$p = 2\zeta\omega_n \quad K = \omega_n^2$$

which can be solved to yield

$$\omega_n = \sqrt{K} \quad \zeta = \frac{p}{2\sqrt{K}}$$

Clearly, with one free parameter either ζ or ω_n can be selected, but not both. We now select a gain value that satisfies the design specifications.

1. If ζ is given and p is known, then the gain of the system and its undamped natural frequency are obtained from the equations

$$K = \left(\frac{p}{2\zeta}\right)^2 \quad \omega_n = \frac{p}{2\zeta}$$

2. If ω_n is given and p is known, then the gain of the system and its damping ratio are obtained from the equations

$$K = \omega_n^2 \quad \zeta = \frac{p}{2\omega_n}$$

Example 5.3 reveals some of the advantages and disadvantages of proportional control. The design is simple, and this simplicity carries over to higher-order systems if a CAD tool is used to assist in selecting the pole locations. Using cursor commands, MATLAB allows the designer to select desirable pole locations from the root locus plot directly. The step response of the system can then be examined using the MATLAB command **step**. But the single free parameter available limits the designer's choice to one design criterion. If more than one aspect of the system time response must be improved, a dynamic controller is needed.

5.4.2 Proportional-derivative (PD) control

As seen from Example 5.1, adding a zero to the loop gain improves the time response in the system. Adding a zero is accomplished using a cascade or feedback controller of the form

$$C(s) = K_p + K_d s = K_d(s+a) \quad (5.10)$$

$$a = K_p/K_d$$

This is known as a **proportional-derivative**, or **PD, controller**. The derivative term is only approximately realizable and is also undesirable because differentiating a noisy input results in large errors. However, if the derivative of the output is measured, an equivalent controller is obtained without differentiation. Thus, PD compensation is often feasible in practice.

The design of PD controllers depends on the specifications given for the closed-loop system and on whether a feedback or cascade controller is used. For a cascade controller, the system block diagram is shown in Fig. 5.4A and the closed-loop transfer function is of the form

$$\begin{aligned} G_{cl}(s) &= \frac{G(s)C(s)}{1 + G(s)C(s)} \\ &= \frac{K_d(s + a)N(s)}{D(s) + K_d(s + a)N(s)} \end{aligned} \quad (5.11)$$

where $N(s)$ and $D(s)$ are the numerator and denominator of the open-loop gain, respectively. Pole-zero cancellation occurs if the loop gain has a pole at $(-a)$. In the absence of pole-zero cancellation, the closed-loop system has a zero at $(-a)$, which may drastically alter the time response of the system. In general, the zero results in greater percentage overshoot than is predicted using (5.7).

Fig. 5.5 shows feedback compensation including a preamplifier in cascade with the feedback loop and an amplifier in the forward path. We show that both amplifiers are often needed. The closed-loop transfer function is

$$\begin{aligned} G_{cl}(s) &= \frac{K_p K_a G(s)}{1 + K_a G(s) C(s)} \\ &= \frac{K_p K_a N(s)}{D(s) + K_a K_d(s + a) N(s)} \end{aligned} \quad (5.12)$$

where K_a is the feedforward amplifier gain and K_p is the preamplifier gain. Note that although the loop gain is the same for both cascade and feedback compensation, the closed-loop system does not have a zero at $(-a)$ in the feedback case. If $D(s)$ has a pole at $(-a)$, the feedback-compensated system has a closed-loop pole at $(-a)$ that appears to cancel with a zero in the root locus when in reality it does not.

In both feedback and cascade compensation, two free design parameters are available and two design criteria can be selected. For example, if the settling time and percentage

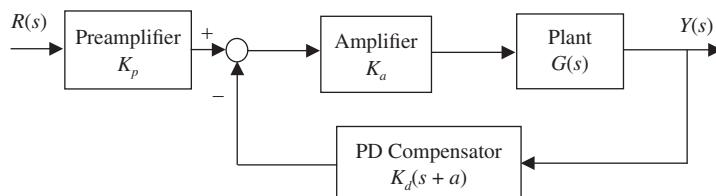


Figure 5.5
Block diagram of a PD-feedback-compensated system.

overshoot are specified, the steady-state error can only be checked after the design is completed and cannot be independently chosen.

Example 5.4

Design a PD controller for the type 1 system described in Example 5.3 to meet the following specifications:

1. Specified ζ and ω_n
2. Specified ζ and steady-state error $e(\infty)\%$ due to a ramp input

Consider both cascade and feedback compensation and compare them using a numerical example.

Solution

The root locus of the PD-compensated system is of the form of Fig. 5.1C. This shows that the system gain can be increased with no fear of instability. Even though this example is solved analytically, a root locus plot is needed to give the designer a feel for the variation in pole locations with gain.

With a PD controller the closed-loop characteristic equation is of the form

$$\begin{aligned}s^2 + ps + K(s + a) &= s^2 + (p + K)s + Ka \\ &= s^2 + 2\zeta\omega_n s + \omega_n^2\end{aligned}$$

where $K = K_d$ for cascade compensation and $K = K_a K_d$ for feedback compensation.

Equating coefficients gives the equations

$$Ka = \omega_n^2 \quad p + K = 2\zeta\omega_n$$

1. In this case, there is no difference between (K, a) in cascade and in feedback compensation, but the feedback case requires a preamplifier with the correct gain to yield zero steady-state error due to unit step. We examine the steady-state error in part 2. In either case, solving for K and a gives

$$K = 2\zeta\omega_n - p \quad a = \frac{\omega_n^2}{2\zeta\omega_n - p}$$

2. For cascade compensation, the velocity error constant of the system is

$$K_v = \frac{Ka}{p} = \frac{100}{e(\infty)\%}$$

The undamped natural frequency is fixed at

$$\omega_n = \sqrt{Ka} = \sqrt{pK_v}$$

Solving for K and a gives

$$K = 2\zeta\sqrt{pK_v} - p \quad a = \frac{pK_v}{2\zeta\sqrt{pK_v} - p}$$

For feedback compensation, the steady-state error in tracking a unit step is given by

$$\lim_{s \rightarrow 0} sE(s) = \lim_{s \rightarrow 0} s[1 - G_{cl}(s)]R(s) = [1 - G_{cl}(0)]$$

Example 5.4—cont'd

Thus, the steady-state error can only be set to zero if $G_c(0) = 1$. With preamplifier gain K_p and cascade amplifier gain K_a , as in Fig. 5.5, the closed-loop transfer function is given by

$$G_{cl}(s) = \frac{K_p K_a}{s^2 + (p + K)s + K_a}$$

The condition for zero steady-state error due to step is $K_p K_a = K_a$. To ensure that this condition also makes the error due a unit ramp a finite value, we consider the error expression

$$\begin{aligned} R(s) - Y(s) &= R(s) \left[1 - \frac{K_p K_a}{s^2 + (p + K)s + K_a} \right] \\ &= R(s) \frac{s^2 + (p + K)s + K_a - K_p K_a}{s^2 + (p + K)s + K_a} \end{aligned}$$

Using the final value theorem gives the steady-state error due to a unit ramp input as

$$e(\infty)\% = \lim_{s \rightarrow 0} s \left[\frac{1}{s^2} \frac{s^2 + (p + K)s + K_a - K_p K_a}{s^2 + (p + K)s + K_a} \right] \times 100\%$$

This error is infinite unless the amplifier gain is selected such that $K_p K_a = K_a$. The steady-state error is then given by

$$e(\infty)\% = \frac{p + K}{K_a} \times 100\%$$

The steady-state error $e(\infty)$ is simply the percentage error divided by 100. Hence, using the equations governing the closed-loop characteristic equation

$$K_a = \frac{p + K}{e(\infty)} = \frac{2\zeta\omega_n}{e(\infty)} = \omega_n^2$$

the undamped natural frequency is fixed at

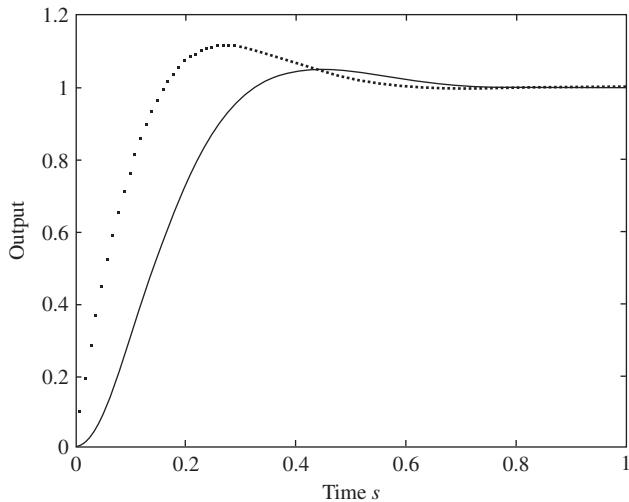
$$\omega_n = \frac{2\zeta}{e(\infty)}$$

Then solving for K and a we obtain

$$\begin{aligned} K &= \frac{4\zeta^2}{e(\infty)} - p \\ a &= \frac{4\zeta^2}{e(\infty)(4\zeta^2 - pe(\infty))} \end{aligned}$$

Note that, unlike cascade compensation, ω_n can be freely selected if the steady-state error is specified and ζ is free. To further compare cascade and feedback compensation, we consider the system with the pole $p = 4$, and require $\zeta = 0.7$ and $\omega_n = 10$ rad/s for part 1. These values give $K = K_d = 10$ and $a = 10$. In cascade compensation, (5.11) gives the closed-loop transfer function

$$G_{cl}(s) = \frac{10(s + 10)}{s^2 + 14s + 100}$$

Example 5.4—cont'd**Figure 5.6**

Step response of PD cascade (dotted) and feedback (solid) compensated systems with a given damping ratio and undamped natural frequency.

For feedback compensation, amplifier gains must be appropriately selected for zero steady-state error due to step. As shown in Part (a), this requires the closed-loop transfer function to satisfy $G_{cl}(0) = 1$ and yields the condition $K_p K_a = K_a$. For example, one may select

$$\begin{aligned}K_p &= 10, \quad K_a = 10 \\K_d &= 1, \quad a = 10\end{aligned}$$

Substituting in (5.12) gives the closed-loop transfer function

$$G_{cl}(s) = \frac{100}{s^2 + 14s + 100}$$

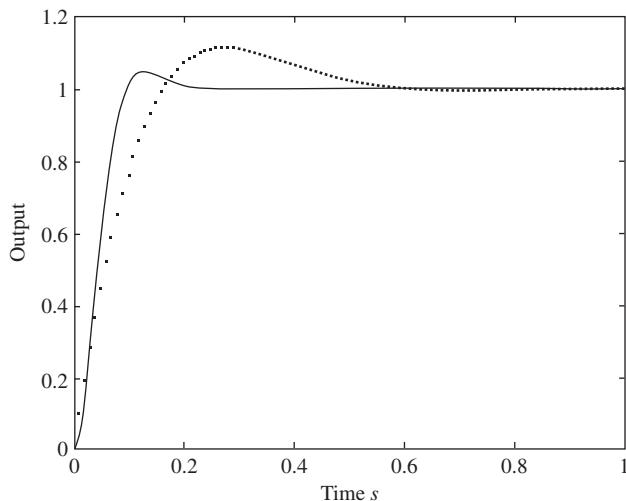
The responses of the cascade- and feedback-compensated systems are shown together in Fig. 5.6. The PO for the feedback case can be predicted exactly using (5.7) and is equal to about 4.6%. For cascade compensation, the PO is higher due to the presence of the zero. The zero is at a distance from the imaginary axis less than one and a half times the negative real part of the complex conjugate poles. Therefore, its effect is significant, and the PO increases to over 10% with a faster response.

For part 2 with $p = 4$, we specify $\zeta = 0.7$ and a steady-state error of 4%. Cascade compensation requires $K = 10$, $a = 10$. These are identical to the values of part 1 and correspond to an undamped natural frequency $\omega_n = 10$ rad/s.

For feedback compensation we obtain $K = 45$, $a = 27.222$. Using (5.12) gives the closed-loop transfer function

$$G_{cl}(s) = \frac{1225}{s^2 + 49s + 1225}$$

with $\omega_n = 35$ rad/s.

Example 5.4—cont'd**Figure 5.7**

Step response of PD cascade (dotted) and feedback (solid) compensated systems with a given damping ratio and steady-state error.

The responses of cascade- and feedback-compensated systems are shown in Fig. 5.7. The PO for the feedback-compensated case is still 4.6% as calculated using (5.7). For cascade compensation, the PO is higher due to the presence of the zero that is close to the complex conjugate poles.

Although the response of the feedback-compensated system is superior, several amplifiers are required for its implementation with high gains. The high gains may cause nonlinear behavior such as saturation in some situations.

Having demonstrated the differences between cascade and feedback compensation, we restrict our discussion in the sequel to cascade compensation. Similar procedures can be developed for feedback compensation.

Example 5.4 is easily solved analytically because the plant is only second order. For higher-order systems, the design is more complex, and a solution using CAD tools is preferable. We develop design procedures using CAD tools based on the classical graphical solution methods. These procedures combine the convenience of CAD tools and the insights that have made graphical tools successful. The procedures find a controller transfer function such that the angle of its product with the loop gain function at the desired pole location is an odd multiple of 180° . From the angle condition, this ensures

that the desired location is on the root locus of the compensated system. The angle contribution required from the controller for a desired closed-loop pole location s_{cl} is

$$\theta_C = \pm 180^\circ - \angle L(s_{cl}) \quad (5.13)$$

where $L(s)$ is the open-loop gain with numerator $N(s)$ and denominator $D(s)$. For a PD controller, the controller angle is simply the angle of the zero at the desired pole location. Applying the angle condition at the desired closed-loop location, it can be shown that the zero location is given by

$$a = \frac{\omega_d}{\tan(\theta_C)} + \zeta \omega_n \quad (5.14)$$

The proof of (5.13) and (5.14) is straightforward and is left as an exercise (see Problem 5.5).

In some special cases, a satisfactory design is obtained by cancellation, or near cancellation, of a system pole with the controller zero. The desired specifications are then satisfied by tuning the reduced transfer function's gain. In practice, exact pole-zero cancellation is impossible. However, with near cancellation the effect of the pole-zero pair on the time response is usually negligible.

A key to the use of powerful calculators and CAD tools in place of graphical methods is the ease with which transfer functions can be evaluated for any complex argument using direct calculation or cursor commands. The following CAD procedure exploits the MATLAB command **evalfr** to obtain the design parameters for a specified damping ratio and undamped natural frequency. The command **evalfr** evaluates a transfer function **g** for any complex argument **s** as follows:

>> evalfr(g, s)

The **angle** command gives the angle of any complex number. The complex value and its angle can also be evaluated using any hand calculator.

Procedure 5.1: Given ζ and ω_n

MATLAB or Calculator

1. Calculate the angle **theta** of the loop gain function evaluated at the desired location s_{cl} , and subtract the angle from π using a hand calculator or the MATLAB commands **evalfr** and **angle**.
2. Calculate the zero location using Eq. (5.14) using **tan(theta)**, where **theta** is in radians.
3. Calculate the magnitude of the numerator of the new loop gain function, including the controller zero, using the commands **abs** and **evalfr**, then calculate the gain using the magnitude condition.
4. Check the time response of the PD-compensated system, and modify the design to meet the desired specifications if necessary. Most currently available calculators cannot perform this step.

The following MATLAB function calculates the gain and zero location for PD control.

Procedure 5.1: Given ζ and ω_n —cont'd

```
% L is the open loop gain.
% zeta and wn specify the desired closed-loop pole.
% scl is the closed-loop pole, theta is the controller angle at scl
% k (a) are the corresponding gain (zero location).
function [k, a, scl] = pdcon(zeta, wn, L).
scl = wn*exp(j*(pi-acos(zeta))); % Find the desired closed-loop% pole location.
theta = pi - angle(evalfr(L, scl)); % Calculate the controller
% angle.
a = wn * sqrt(1-zeta^2)/tan(theta) + zeta*wn; % Calculate the
% controller zero.
Lcomp = L*tf([1, a],1); % Include the controller zero.
k = 1/abs(evalfr(Lcomp, scl)); % Calculate the gain that yields the
% desired pole.
```

For a specified steady-state error, the system gain is fixed and the zero location is varied. Other design specifications require varying parameters other than the gain K . Root locus design with a free parameter other than the gain is performed using Procedure 5.2.

Procedure 5.2: Given steady-state error and ζ

1. Obtain the error constant from the steady-state error, and determine a system parameter that remains free after the error constant is fixed for the system with PD control.
2. Rewrite the closed-loop characteristic equation of the PD-controlled system in the form

$$1 + K_f G_f(s) = 0 \quad (5.15)$$

where K_f is a gain dependent on the free system parameter and $G_f(s)$ is a function of s .

3. Obtain the value of the free parameter K_f corresponding to the desired closed-loop pole location using the MATLAB command **rlocus**. As in Procedure 5.1, K_f can be obtained by applying the magnitude condition using MATLAB or a calculator.
4. Calculate the free parameter from the gain K_f
5. Check the time response of the PD-compensated system, and modify the design to meet the desired specifications if necessary.

Example 5.5

Using a CAD package, design a PD controller for the type 1 position control system of Example 3.3 with transfer function

$$G(s) = \frac{1}{s(s+4)}$$

Example 5.5—cont'd

to meet the following specifications:

1. $\zeta = 0.7$ and $\omega_n = 10$ rad/s
2. $\zeta = 0.7$ and 4% steady-state error due to a unit ramp input.

Solution

1. We solve the problem using Procedure 5.1 and the MATLAB function **pdcon**. Fig. 5.8 shows the root locus of the uncompensated system with the desired pole location at the intersection of the radial line for a damping ratio of 0.7 and the circular arc for an undamped natural frequency of 10. A compensator angle of 67.2° is obtained using (5.13) with a hand calculator or MATLAB. The MATLAB function **pdcon** gives

```
>> [k, a, scl] = pdcon(0.7, 10, tf(1, [1, 4, 0]))
```

$$k = 10.0000$$

$$a = 10.0000$$

$$scl = -7.0000 + 7.1414i$$

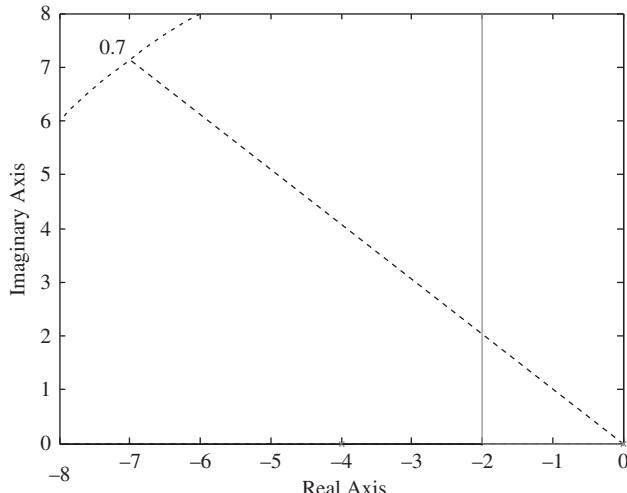


Figure 5.8

Root locus plot of uncompensated systems.

Fig. 5.9 shows the root locus of the compensated plot with the cursor at the desired pole location and a corresponding gain of 10. The results are approximately equal to those obtained analytically in Example 5.4.

2. The specified steady-state error gives

$$K_v = \frac{100}{e(\infty)\%} = \frac{100}{4\%} = 25 = \frac{Ka}{4} \Rightarrow Ka = 100$$

The closed-loop characteristic equation of the PD-compensated system is given by

$$1 + K \frac{s + a}{s(s + 4)} = 0$$

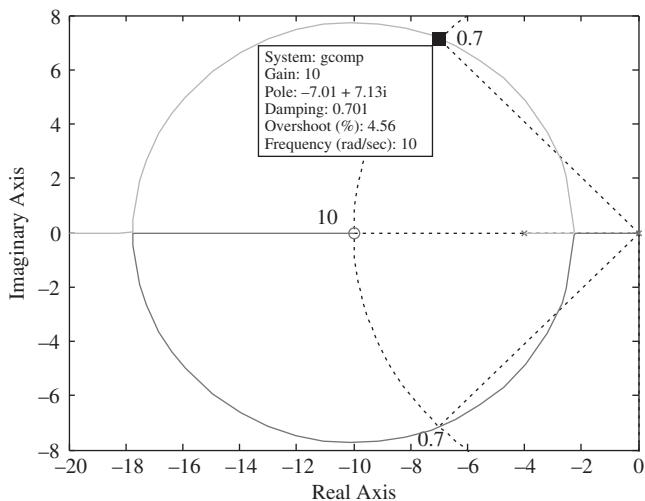
Example 5.5—cont'd

Figure 5.9
Root locus plot of PD-compensated systems.

Let K vary with a so that their product Ka remains equal to 100, then Procedure 5.2 requires that the characteristic equation be rewritten as

$$1 + K \frac{s}{s^2 + 4s + 100} = 0$$

The corresponding root locus is a circle centered at the origin as shown in Fig. 5.10 with the cursor at the location corresponding to the desired damping ratio. The desired location is at

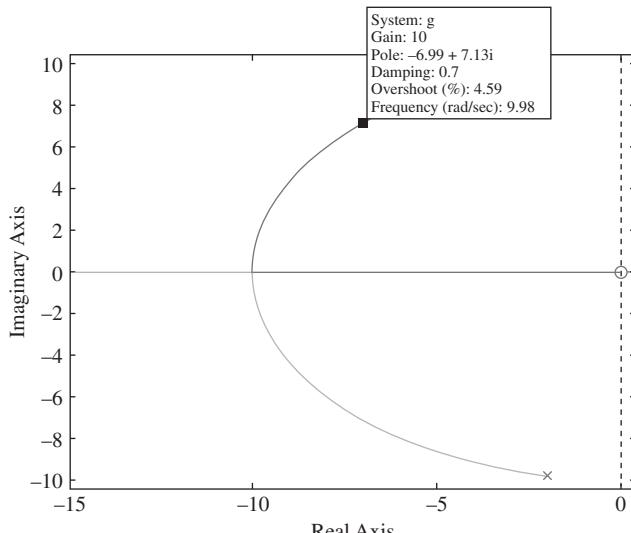


Figure 5.10
Root locus plot of PD-compensated systems with Ka fixed.

Example 5.5—cont'd

the intersection of the root locus with the $\zeta = 0.7$ radial line. The corresponding gain value is $K = 10$, which yields $a = 10$ —that is, the same values as in Example 5.3. We obtain the value of K using the MATLAB commands

```
>> g = tf([1, 0], [1, 4, 100]); rlocus(g)
```

(Click on the root locus and drag the mouse until the desired gain is obtained.)

The time responses of the two designs are identical and were obtained earlier as the cascade-compensated responses of Figs. 5.6 and 5.7, respectively.

5.4.3 Proportional-integral (PI) control

Increasing the type number of the system drastically improves its steady-state response. If an **integral controller** is added to the system, its type number is increased by one, but its transient response deteriorates or the system becomes unstable. If a proportional control term is added to the integral control, the controller has a pole and a zero. The transfer function of the **proportional-integral (PI) controller** is

$$C(s) = K_p + \frac{K_i}{s} = K_p \frac{s + a}{s} \quad (5.16)$$

$$a = K_i/K_p$$

and is used in cascade compensation. An integral term in the feedback path is equivalent to a differentiator in the forward path and is therefore undesirable (see Problem 5.7).

PI design for a plant transfer function $G(s)$ can be viewed as PD design for the plant $G(s)/s$. Thus, Procedure 5.1 or 5.2 can be used for PI design. However, a better design is often possible by placing the controller zero close to the pole at the origin so that the controller pole and zero “almost cancel.” An almost canceling pole-zero pair has a negligible effect on the time response. Thus, the PI controller results in a small deterioration in the transient response with a significant improvement in the steady-state error. The following procedure can be used for PI controller design.

Procedure 5.3

1. Design a proportional controller for the system to meet the transient response specifications (i.e., place the dominant closed-loop system poles at a desired location $s_{cl} = -\zeta\omega_n \pm j\omega_d$).
2. Add a PI controller with the zero location specified by

$$a = \frac{\omega_n}{\zeta + \sqrt{1 - \zeta^2 / \tan(\phi)}} \quad (5.17)$$

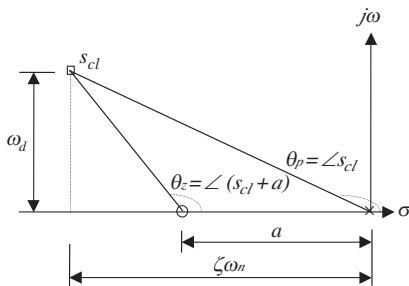


Figure 5.11
Pole-zero diagram of a PI controller.

Procedure 5.3—cont'd

or

$$a = \frac{\zeta \omega_n}{10} \quad (5.18)$$

where ϕ is a small angle ($3 \rightarrow 5^\circ$).

3. Tune the gain of the system to move the closed-loop pole closer to s_{cl} .
4. Check the system time response.

If a PI controller is used, it is implicitly assumed that proportional control meets the transient response but not the steady-state error specifications. Failure of the first step in Procedure 5.3 indicates that a different controller (one that improves both transient and steady-state behavior) must be used.

To prove (5.17) and justify (5.18), we use the pole-zero diagram of Fig. 5.11. The figure shows the angle contribution of the controller at the closed-loop pole location s_{cl} . The contribution of the open-loop gain $L(s)$ is not needed and is not shown.

Proof

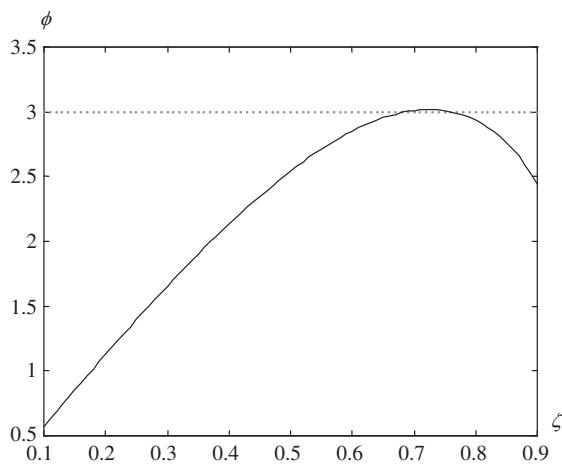
The controller angle at s_{cl} is

$$-\phi = \theta_z - \theta_p = (180^\circ - \theta_p) - (180^\circ - \theta_z)$$

From Fig. 5.11, the tangents of the two angles in the preceding equation are

$$\tan(180^\circ - \theta_p) = \frac{\omega_d}{\zeta \omega_n} = \frac{\sqrt{1 - \zeta^2}}{\zeta}$$

$$\tan(180^\circ - \theta_z) = \frac{\omega_d}{\zeta \omega_n - a} = \frac{\sqrt{1 - \zeta^2}}{\zeta - x}$$

Proof—cont'd**Figure 5.12**

Plot of the controller angle ϕ at the underdamped closed-loop pole versus ζ , where x is the ratio (a/ω_n) . Next, we use the trigonometric identity

$$\tan(A - B) = \frac{\tan(A) - \tan(B)}{1 + \tan(A)\tan(B)}$$

to obtain

$$\tan(\phi) = \frac{\frac{\sqrt{1 - \zeta^2}}{\zeta - x} - \frac{\sqrt{1 - \zeta^2}}{\zeta}}{1 + \frac{1 - \zeta^2}{\zeta(\zeta - x)}} = \frac{x\sqrt{1 - \zeta^2}}{1 - \zeta x}$$

Solving for x , we have

$$x = \frac{1}{\zeta + \sqrt{1 - \zeta^2}/\tan(\phi)}$$

Multiplying by ω_n gives (5.17).

If the controller zero is chosen using (5.18), then $x = \zeta/10$. Solving for ϕ we obtain

$$\phi = \tan^{-1} \left(\frac{\zeta\sqrt{1 - \zeta^2}}{10 - \zeta^2} \right)$$

This yields the plot of Fig. 5.12, which clearly shows an angle ϕ of 3° or less.

The use of Procedure 5.1 or Procedure 5.3 to design PI controllers is demonstrated in Example 5.6.

Example 5.6

Design a controller for the position control system

$$G(s) = \frac{1}{s(s+10)}$$

to perfectly track a ramp input and have a dominant pair with a damping ratio of 0.7 and an undamped natural frequency of 4 rad/s.

Solution**Design #1**

Apply Procedure 5.1 to the modified plant

$$G_i(s) = \frac{1}{s^2(s+10)}$$

This plant is unstable for all gains, as seen from its root locus plot in Fig. 5.13. The controller must provide an angle of about 111° at the desired closed-loop pole location. Substituting in (5.14) gives a zero at -1.732 . Then moving the cursor to the desired pole location on the root locus of the compensated system (Fig. 5.14) gives a gain of about 40.6.

The design can also be obtained analytically by writing the closed-loop characteristic polynomial as

$$\begin{aligned} s^3 + 10s^2 + Ks + Ka &= (s + \alpha)(s^2 + 2\zeta\omega_n s + \omega_n^2) \\ &= s^3 + (\alpha + 2\zeta\omega_n)s^2 + (2\zeta\omega_n\alpha + \omega_n^2)s + \alpha\omega_n^2 \end{aligned}$$

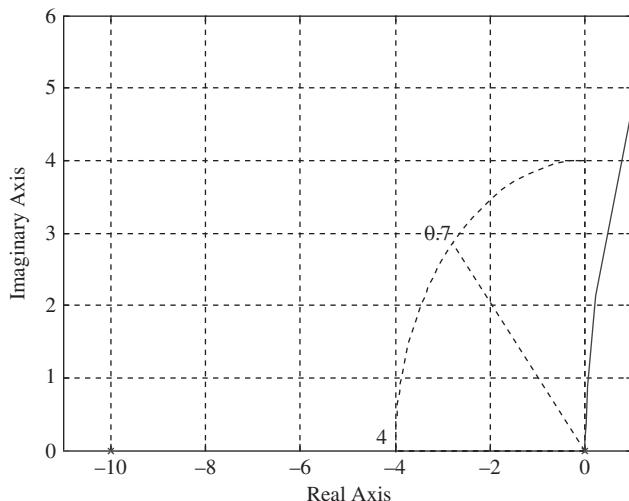
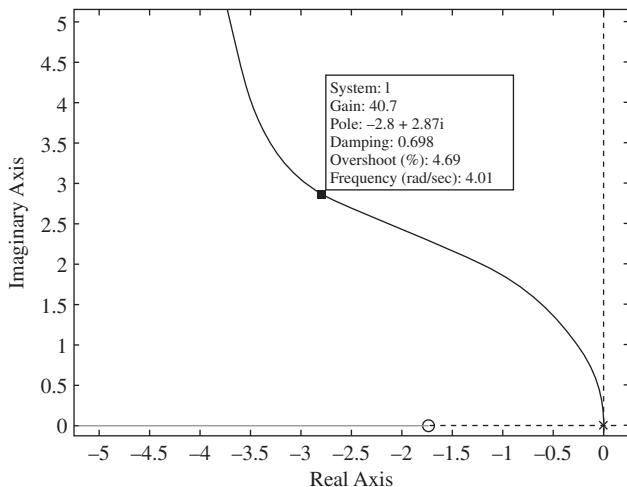


Figure 5.13
Root locus of a system with an integrator.

Example 5.6—cont'd**Figure 5.14**

Root locus of a PI-compensated system.

Then equating coefficients gives

$$\alpha = 10 - 2\zeta\omega_n = 10 - 2(0.7)(4) = 4.4$$

$$K = \omega_n(2\zeta\alpha + \omega_n) = 4 \times [2(0.7)(4.4) + 4] = 40.64$$

$$a = \frac{\alpha\omega_n^2}{K} = \frac{4.4 \times 4^2}{40.64} = 1.732$$

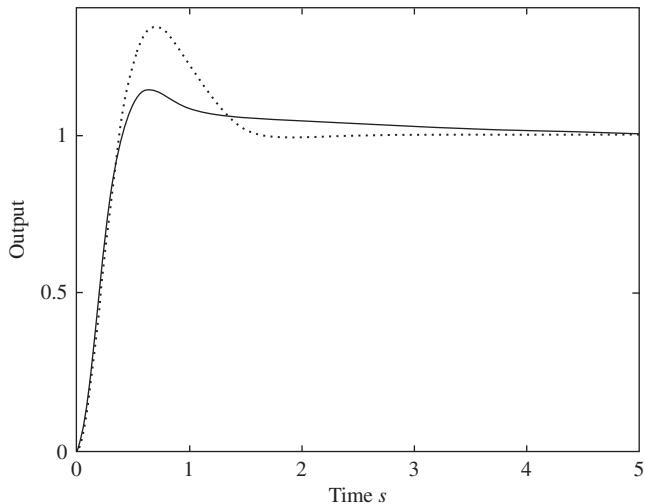
which are approximately the same as the values obtained earlier.

The MATLAB commands to obtain the zero location and the corresponding gain are

```
>> g = tf([1, 10, 0, 0]); scl = 4*(-0.7+j*sqrt(1-0.7^2))
scl = -2.8000 + 2.8566i
>> theta = pi-angle(polyval([1, 10, 0, 0], scl))
theta = 1.9285
>> a = imag(scl)/tan(theta) - real(scl)
a = 1.7323
>> k = 1/abs(evalfr(g*tf([1, a], 1), scl))
k = 40.6400
```

The closed-loop transfer function for the preceding design (Design #1) is

$$\begin{aligned} G_{cl}(s) &= \frac{40.64(s + 1.732)}{s^3 + 10s^2 + 40.64s + 69.28} \\ &= \frac{40.64(s + 1.732)}{(s + 4.4)(s^2 + 5.6s + 16)} \end{aligned}$$

Example 5.6—cont'd**Figure 5.15**

Step response of a PI-compensated system: Design #1 (dotted) and Design #2 (solid).

The system has a zero close to the closed-loop poles, which results in excessive overshoot in the time response of Fig. 5.15 (shown together with the response for an alternative design).

Design #2

Next, we apply Procedure 5.3 to the same problem. The proportional control gain for a damping ratio of 0.7 is approximately 51.02 and yields an undamped natural frequency of 7.143 rad/s. This is a faster design than required and is therefore acceptable. Then we use (5.17) and obtain the zero location

$$\begin{aligned} \alpha &= \frac{\omega_n}{\zeta + \sqrt{1 - \zeta^2}/\tan(-\phi)} \\ &= \frac{7.143}{0.7 + \sqrt{1 - 0.49}/\tan(3^\circ)} \cong 0.5 \end{aligned}$$

If (5.18) is used, we have

$$\alpha = \frac{\zeta \omega_n}{10} = \frac{7.143 \times 0.7}{10} \cong 0.5$$

That is, the same zero value is obtained.

The closed-loop transfer function for this design is

$$\begin{aligned} G_{cl}(s) &= \frac{50(s + 0.5)}{s^3 + 10s^2 + 50s + 25} \\ &= \frac{50(s + 0.5)}{(s + 0.559)(s^2 + 9.441s + 44.7225)} \end{aligned}$$

where the gain value has been slightly reduced to bring the damping ratio closer to 0.7. This actually does give a damping ratio of about 0.7 and an undamped natural frequency of 6.6875 rad/s for the dominant closed-loop poles.

Example 5.6—cont'd

The time response for this design (Design #2) is shown, together with that of Design #1, in Fig. 5.15. The percentage overshoot for Design #2 is much smaller because its zero almost cancels with a closed-loop pole. Design #2 is clearly superior to Design #1.

5.4.4 Proportional-integral-derivative (PID) control

If both the transient and steady-state response of the system must be improved, then neither a PI nor a PD controller may meet the desired specifications. Adding a zero (PD) may improve the transient response but does not increase the type number of the system. Adding a pole at the origin increases the type number but may yield an unsatisfactory time response even if one zero is also added. With a **proportional-integral-derivative (PID) controller**, two zeros and a pole at the origin are added. This both increases the type number and allows satisfactory reshaping of the root locus.

The transfer function of a PID controller is given by

$$C(s) = K_p + \frac{K_i}{s} + K_d s = K_d \frac{s^2 + 2\zeta\omega_n s + \omega_n^2}{s}$$

$$2\zeta\omega_n = K_p/K_d, \quad \omega_n^2 = K_i/K_d \quad (5.19)$$

where K_p , K_i , and K_d are the proportional, integral, and derivative gain, respectively.

The zeros of the controller can be real or complex conjugate, allowing the cancellation of real or complex conjugate LHP poles if necessary. In some cases, good design can be obtained by canceling the pole closest to the imaginary axis. The design then reduces to a PI controller design that can be completed by applying Procedure 5.3. Alternatively, one could apply Procedure 5.1 or 5.2 to the reduced transfer function with an added pole at the origin. A third approach to PID design is to follow Procedure 5.3 with the proportional control design step modified to PD design. The PD design is completed using Procedure 5.1 or 5.2 to meet the transient response specifications. PI control is then added to improve the steady-state response. Examples 5.7 and 5.8 illustrate these design procedures.

Example 5.7

Design a PID controller for an armature-controlled direct current (DC) motor with transfer function

$$G(s) = \frac{1}{s(s+1)(s+10)}$$

to obtain zero steady-state error due to ramp, a damping ratio of 0.7, and an undamped natural frequency of 4 rad/s.

Example 5.7—cont'd**Solution**

Cancelling the pole at -1 with a zero yields the transfer function

$$G(s)C_{PD}(s) = \frac{1}{s(s+10)}$$

where $C_{PD}(s) = s + 1$ is the transfer function of the PD controller. This is identical to the transfer function of Example 5.6 where we designed a PI controller using two different approaches. Because the second design gave much better results, we use it to complete our PID design. Hence, the overall PID controller is given by

$$C(s) = 50 \frac{(s+1)(s+0.5)}{s}$$

This design is henceforth referred to as Design #1.

A second design (Design #2) is obtained by first selecting a PD controller to meet the transient response specifications. We seek an undamped natural frequency of 5 rad/s in anticipation of the effect of adding PI control. The PD controller is designed using the MATLAB commands (using the function **pdcon**)

```
>> [k, a, scl] = pdcon(0.7, 5, tf(1, [1, 11, 10, 0]))
```

$k = 43.0000$

$a = 2.3256$

$s_{cl} = -3.5000 + 3.5707i$

The PI zero is obtained using (5.17) and the command

```
>> b = 5/(0.7+sqrt(1-.49)/tan(3*pi/180))
```

$b = 0.3490$

Note that this result is almost identical to the zero obtained location obtained more easily using (5.18). For a better transient response, the gain is reduced to 40 and the controller transfer function for Design #2 is

$$C(s) = 40 \frac{(s+0.349)(s+2.326)}{s}$$

The step responses for Designs #1 and #2 are shown in Fig. 5.16. Clearly, Design #1 is superior because the zeros in Design #2 result in excessive overshoot. The plant transfer function favors pole cancellation in this case because the remaining real axis pole is far in the LHP. If the remaining pole is at, say, -3 , the second design procedure would give better results. The lesson to be learned is that there are no easy solutions in design. There are recipes with which the designer should experiment until satisfactory results are obtained.

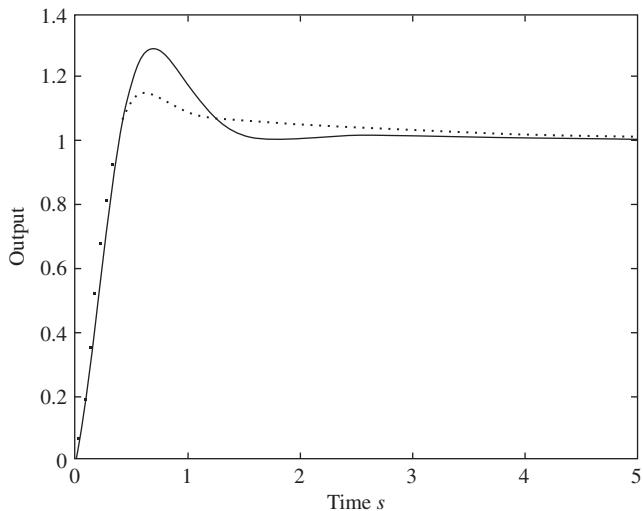
Example 5.7—cont'd

Figure 5.16
Time response for Design #1 (dotted) and Design #2 (solid).

Example 5.8

Design a PID controller to obtain zero steady-state error due to step, a damping ratio of 0.7, and an undamped natural frequency of at least 4 rad/s for the transfer function

$$G(s) = \frac{1}{(s+10)(s^2 + 2s + 10)}$$

Solution

The system has a pair of complex conjugate poles that slow down its time response and a third pole that is far in the LHP. Canceling the complex conjugate poles with zeros and adding the integrator yields the transfer function

$$G(s) = \frac{1}{s(s+10)}$$

The root locus of the system is similar to Fig. 5.2A, and we can increase the gain without fear of instability. The closed-loop characteristic equation of the compensated system with gain K is

$$s^2 + 10s + K = 0$$

Equating coefficients as in Example 5.3, we observe that for a damping ratio of 0.7 the undamped natural frequency is

$$\omega_n = \frac{10}{2\zeta} = \frac{5}{0.7} = 7.143 \text{ rad/s}$$

Example 5.8—cont'd

This meets the design specifications. The corresponding gain is 51.02, and the PID controller is given by

$$C(s) = 51.02 \frac{s^2 + 2s + 10}{s}$$

In practice, pole-zero cancellation may not occur, but near cancellation is sufficient to obtain a satisfactory time response, as shown in Fig. 5.17.

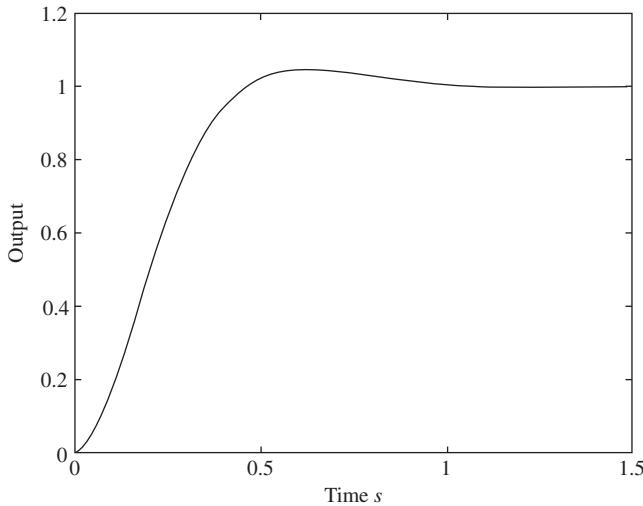


Figure 5.17
Step response of the PID-compensated system of Example 5.8.

5.5 Empirical tuning of PID controllers

In industrial applications, PID controllers are often tuned empirically. Typically, the controller parameters are selected based on a simple process model using a suitable tuning rule. This allows us to address (1) load disturbance rejection specifications (which are often a major concern in process control) and (2) the presence of a time delay in the process. We first write the PID controller transfer function (5.19) in the form

$$C(s) = K_p \left(1 + \frac{1}{T_i s} + T_d s \right) \quad (5.20)$$

$$T_i = K_p/K_i, \quad T_d = K_d/K_p$$

where T_i denotes the **integral time constant** and T_d denotes the **derivative time constant**. The three controller parameters K_p , T_i , and T_d have a clear physical meaning. Increasing K_p (i.e., increasing the proportional action) provides a faster but more oscillatory response.

The same behavior results from increasing the integral action by decreasing the value of T_i . Finally, increasing the value of T_d leads to a slower but more stable response.

These considerations allow tuning the controller by a trial-and-error procedure. However, this can be time consuming, and the achieved performance depends on the skill of the designer. Fortunately, tuning procedures are available to simplify the PID design.

Typically, the parameters of the process model are determined assuming a first-order-plus-dead-time model. That is,

$$G(s) = \frac{K}{\tau s + 1} e^{-Ls} \quad (5.21)$$

where K is the process gain, τ is the process (dominant) time constant, and L is the (apparent) dead time of the process. The response for the approximate model of (5.21) due to a step input of amplitude A is given by the expression

$$KA \left\{ 1 - \exp \left[-\frac{t-L}{\tau} \right] \right\} u(t-L) \quad (5.22)$$

The rising exponential of (5.22) has a steady-state level of KA and a tangent to it at the initiation of its rise and with the initial slope reaches the steady-state level after one time constant τ . In addition, the rising exponential reaches 63% of the final value after one time constant. Hence, we can estimate the parameters of the approximate model based on the step response of the process through the **tangent method**.² The method consists of the following steps.

Tangent method

1. Obtain the step response of the process experimentally.
2. Draw a tangent to the step response at the inflection point as shown in Fig. 5.18.
3. Compute the process gain as the ratio of the steady-state change in the process output y to the amplitude of the input step A .
4. Compute the apparent dead time L as the time interval between the application of the step input and the intersection of the tangent line with the time axis.
5. Determine the sum $\tau+L$ (from which the value of τ can be easily computed) as the time interval between the application of the step input and the intersection of the tangent line with the straight representing the final steady-state value of the process output. Alternatively, the value of $\tau+L$ can be determined as the time interval between the application of the step input and the time when the process output attains 63.2% of its final value. Note that if the dynamics of the process can be perfectly described by a first-order-plus-dead-time model, the values of τ obtained in the two methods are identical.

² See Åström and Hägglund (2006) and Visioli (2006) for a detailed description and analysis of different methods for estimating the parameters of a first-order-plus-dead-time system.

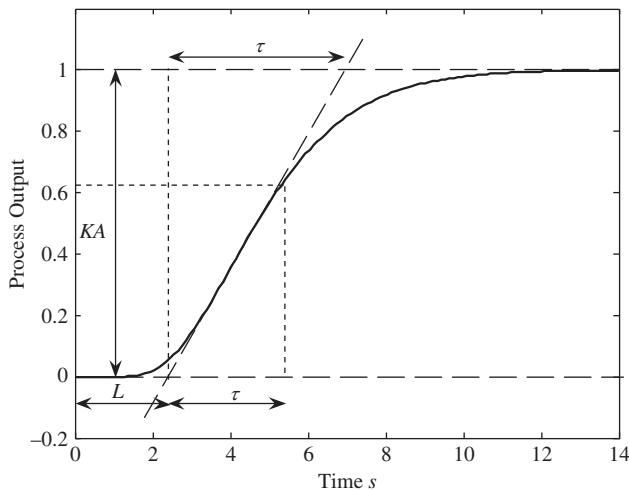
Tangent method—cont'd

Figure 5.18
Application of the tangent method.

Given the process model parameters, several tuning rules are available for determining the PID controller parameter values, but different rules address different design specifications. The most popular tuning rules are those attributed to **Ziegler-Nichols**. Their aim is to provide satisfactory load disturbance rejection. [Table 5.1](#) shows the Ziegler-Nichols rules for P, PI, and PID controllers. Although the rules are empirical, they are consistent with the physical meaning of the parameters. For example, consider the effect of the derivative action in PID control governed by the third row of [Table 5.1](#). The derivative action provides added damping to the system, which increases its relative stability. This allows us to increase both the proportional and integral action while maintaining an acceptable time response. We demonstrate the Ziegler-Nichols procedure using Example 5.9.

Table 5.1: Ziegler-Nichols tuning rules for a first-order-plus-dead-time model of the process.

Controller type	K_p	T_i	T_d
P	$\frac{\tau}{KL}$	—	—
PI	$0.9 \frac{\tau}{KL}$	$3 L$	—
PID	$1.2 \frac{\tau}{KL}$	$2 L$	$0.5 L$

Example 5.9

Consider the control system shown in Fig. 5.19, where the process has the following transfer function:

$$G(s) = \frac{1}{(s+1)^4} e^{-0.2s}$$

Estimate a first-order-plus-dead-time model of the process, and design a PID controller by applying the Ziegler-Nichols tuning rules.

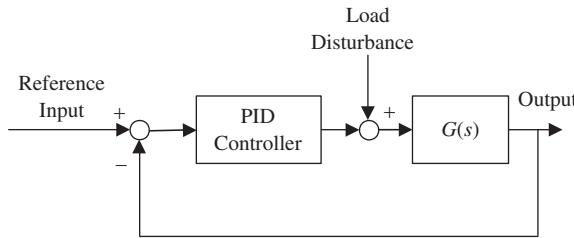


Figure 5.19
Block diagram of the process of Example 5.9.

Solution

The process step response is shown in Fig. 5.20. By applying the tangent method, a first-order-plus-dead-time model with gain $K = 1$, $L = 1.55$ and a delay $\tau = 3$ is estimated. The Ziegler-Nichols rules of Table 5.1 provide the following PID parameters: $K_p = 2.32$, $T_i = 3.1$, and $T_d = 0.775$. Fig. 5.21 shows the response of the closed-loop system due to a step input at $t = 0$ followed by a step load disturbance input at $t = 50$.

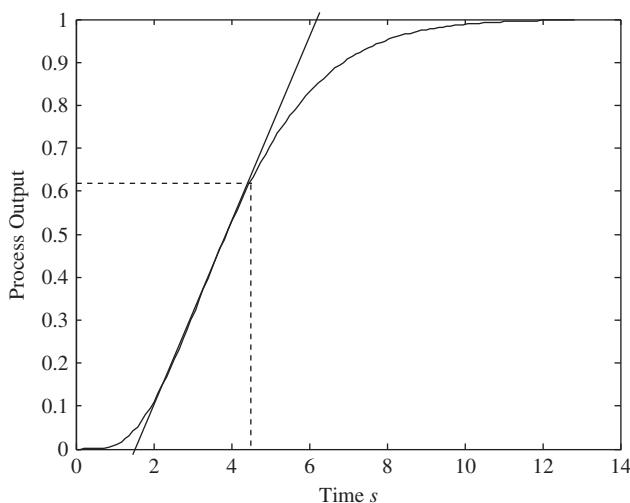
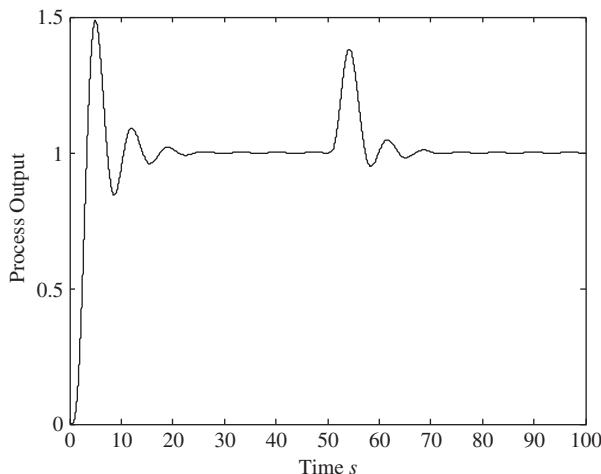


Figure 5.20
Application of the tangent method in Example 5.9.

Example 5.9—cont'd**Figure 5.21**

Process output with the PID controller tuned with the Ziegler-Nichols method.

The system response is oscillatory, which is a typical feature of the Ziegler-Nichols method, whose main objective is to satisfy load disturbance response specifications. However, the effect of the load step disturbance is minimized.

Ziegler and Nichols also devised tuning rules based on using the critical gain K_c (that is, the gain margin) and the period of the corresponding oscillations T_c . The data needed for this procedure can be obtained experimentally using proportional feedback control by progressively increasing the controller gain until the system becomes marginally stable, and then recording the controller gain and the period of the resulting oscillations. The controller parameters are then computed using the rules of Table 5.2. The main drawback of this approach is that it requires the system to operate on the verge of instability, which is often detrimental. However, this problem can be overcome using a relay in the feedback loop to cause oscillatory behavior, known as a limit cycle.³ Controller tuning using a

Table 5.2: Ziegler-nichols tuning rules for the closed-loop method.

Controller type	K_p	T_i	T_d
P	$0.5 K_c$	—	—
PI	$0.4 K_c$	$0.8 T_c$	—
PID	$0.6 K_c$	$0.5 T_c$	$0.125 T_c$

³ See Åström and Hägglund (2006) for details on the relay-feedback methodology.

closed-loop system with proportional control is explored through a computer exercise (see Problem 5.17).

References

- Åström, K.J., Hägglund, T., 2006. Advanced PID Controllers. ISA Press, Research Triangle Park, NJ.
Visioli, A., 2006. Practical PID Control. Springer, UK.

Further reading

- D'Azzo, J.J., Houpis, C.H., 1988. Linear Control System Analysis and Design. McGraw-Hill, New York.
Kuo, B.C., 1995. Automatic Control Systems. Prentice Hall, Englewood Cliffs, NJ.
Nise, N.S., 1995. Control Systems Engineering. Benjamin Cummings, Hoboken, NJ.
O'Dwyer, A., 2006. Handbook of PI and PID Tuning Rules. Imperial College Press, London, UK.
Ogata, K., 1990. Modern Control Engineering. Prentice Hall, Englewood Cliffs, NJ.
Van de Vegte, J., 1986. Feedback Control Systems. Prentice Hall, Englewood Cliffs, NJ.

Problems

- 5.1 Prove that for a system with two real poles and a real zero

$$L(s) = \frac{s + a}{(s + p_1)(s + p_2)}, \quad a < p_1 < p_2 \text{ or } p_1 < p_2 < a$$

the breakaway point is at a distance of $\sqrt{(a - p_1)(a - p_2)}$ from the zero.

- 5.2 Use the result of Problem 5.1 to draw the root locus of the system

$$KL(s) = \frac{K(s + 4)}{s(s + 2)}$$

- 5.3 Sketch the root loci of the following systems:

a. $KL(s) = \frac{K}{s(s+2)(s+5)}$

b. $KL(s) = \frac{K(s+2)}{s(s+3)(s+5)}$

- 5.4 Consider the speed control of an armature controlled DC motor with the transfer function

$$G(s) = \frac{\Omega(s)}{V_a(s)} = \frac{K_m}{Ts + 1}$$

where $\Omega(s)$ is the angular velocity and $V_a(s)$ is the armature voltage.

Design a PI controller $C(s) = K_p \frac{s+a}{s}$ to obtain a the closed-loop transfer function with a desired time constant τ .

- 5.5 Consider the system in 5.3(b) with a required steady-state error of 20%, and an adjustable PI controller zero location (not fixed at -2). Show that the corresponding closed-loop characteristic equation is given by

$$1 + K \left(\frac{s+a}{s} \right) \frac{1}{(s+3)(s+5)} = 0$$

Next, rewrite the equation as

$$1 + K_f G_f(s) = 0$$

where $K_f = K$, $K.a$ is constant, and $G_f(s)$ is a function of s , and examine the effect of shifting the zero on the closed-loop poles.

- a. Design the system for a dominant second-order pair with a damping ratio of 0.5. What is ω_n for this design?
 - b. Obtain the time response using a CAD program. How does the time response compare with that of a second-order system with the same ω_n and ζ as the dominant pair? Give reasons for the differences.
 - c. Discuss briefly the trade-off between error, speed of response, and relative stability in this problem.
- 5.6 Prove Eqs. (5.13) and (5.14), and justify the design Procedures 5.1 and 5.2.
- 5.7 Show that a PI feedback controller is undesirable because it results in a differentiator in the forward path. Discuss the step response of the closed-loop system.
- 5.8 Design a controller for the transfer function

$$G(s) = \frac{1}{(s+1)(s+5)}$$

to obtain (a) zero steady-state error due to step, (b) a settling time of less than 2 s, and (c) an undamped natural frequency of 5 rad/s. Obtain the response due to a unit step, and find the percentage overshoot, the time to the first peak, and the steady-state error percentage due to a ramp input.

- 5.9 Repeat Problem 5.8 with a required settling time less than 0.5 s and an undamped natural frequency of 10 rad/s.
- 5.10 Consider the oven temperature control system of Example 3.5 with transfer function

$$G(s) = \frac{K}{s^2 + 3s + 1}$$

- a. Design a proportional controller for the system to obtain a percentage overshoot less than 5%.

- b. Design a controller for the system to reduce the steady-state error due to step to zero without significant deterioration in the transient response.

5.11 For the inertial system governed by the differential equation

$$\ddot{\theta} = \tau$$

design a feedback controller to stabilize the system and reduce the percentage overshoot below 10% with a settling time of less than 4 s.

Computer exercises

5.12 Consider the oven temperature control system described in Example 3.5 with transfer function

$$G(s) = \frac{K}{s^2 + 3s + 10}$$

- a. Obtain the step response of the system with a PD cascade controller with gain 80 and a zero at -5 .
- b. Obtain the step response of the system with PD feedback controller with a zero at -5 and unity gain and forward gain of 80.
- c. Why are the root loci identical for both systems?
- d. Why are the time responses different although the systems have the same loop gains?
- e. Complete a comparison table using the responses of (a) and (b), including the percentage overshoot, the time to first peak, the settling time, and the steady-state error. Comment on the results, and explain the reason for the differences in the response.

5.13 Use Simulink to examine a practical implementation of the cascade controller described in Exercise 5.12. The compensator transfer function includes a pole because PD control is only approximately realizable. The controller transfer is of the form

$$C(s) = 80 \frac{0.2s + 1}{0.02s + 1}$$

- a. Simulate the system with a step reference input both with and without a saturation block with saturation limits ± 5 between the controller and plant. Export the output to MATLAB for plotting (you can use a Scope block and select “Save data to workspace”).
- b. Plot the output of the system with and without saturation together, and comment on the difference between the two step responses.

- 5.14 Consider the system

$$G(s) = \frac{1}{(s+1)^4}$$

and apply the Ziegler-Nichols procedure based on [Table 5.1](#) to design a PID controller. Obtain the response due to a unit step input as well as a unit step disturbance signal.

- 5.15 Write a computer program that implements the estimation of a first-order-plus-dead-time transfer function with the tangent method and then determine the PID parameters using the Ziegler-Nichols formula. Apply the program to the system

$$G(s) = \frac{1}{(s+1)^8}$$

and simulate the response of the control system when a set-point step change and a load disturbance step are applied. Discuss the choice of the time constant value based on the results.

- 5.16 Apply the script of Exercise 5.15 to the system

$$G(s) = \frac{1}{(s+1)^2}$$

and simulate the response of the control system when a set-point step change and a load disturbance step are applied. Compare the results obtained with those of Problem 5.15.

- 5.17 Use the Ziegler-Nichols closed-loop method to design a PID controller for the system

$$G(s) = \frac{1}{(s+1)^4}$$

based on [Table 5.2](#). Obtain the response due to a unit step input together with a unit step disturbance signal. Compare the results with those of Problem 5.14.

Digital control system design

Objectives

After completing this chapter, the reader will be able to do the following:

1. Sketch the z -domain root locus for a digital control system, or obtain it using MATLAB.
2. Obtain and tune a digital controller from an analog design.
3. Design a digital controller in the z -domain directly using the root locus approach.
4. Design a digital controller using frequency domain techniques.
5. Design a digital controller directly using the synthesis approach of Ragazzini.
6. Design a digital control system with finite settling time.

To design a digital control system, we seek a z -domain transfer function or difference equation model of the controller that meets given design specifications. The controller model can be obtained from the model of an analog controller that meets the same design specifications. Alternatively, the digital controller can be designed in the z -domain using procedures that are almost identical to s -domain analog controller design. We discuss both approaches in this chapter. We begin by introducing the z -domain root locus.

Chapter Outline

6.1 z -domain root locus 182

6.2 z -domain digital control system design 184

 Observation 186

 Remarks 186

 6.2.1 z -domain contours 187

 6.2.2 Proportional control design in the z -domain 190

6.3 Digital implementation of analog controller design 195

 6.3.1 Differencing methods 196

Backward differencing 198

 6.3.2 Pole-zero matching 199

 6.3.3 Bilinear transformation 201

 6.3.4 Empirical digital PID controller tuning 214

6.4 Direct z -domain digital controller design 216

6.5 Frequency response design 221

6.6 Direct control design 229

6.7 Finite settling time design 234

6.7.1 Eliminating intersample oscillation 239

Further reading 248**Problems 248****Computer exercises 251****6.1 z-domain root locus**

In Chapter 3, we showed that the closed-loop characteristic equation of a digital control system is of the form

$$1 + C(z)G_{ZAS}(z) = 0 \quad (6.1)$$

where $C(z)$ is the controller transfer function and $G_{ZAS}(z)$ is the transfer function of the DAC, analog subsystem, and ADC combination. If the controller is assumed to include a constant gain multiplied by a rational z -transfer function, then Eq. (6.1) is equivalent to

$$1 + KL(z) = 0 \quad (6.2)$$

where $L(z)$ is the open-loop gain.

Eq. (6.2) is identical in form to the s -domain characteristic Eq. (5.1) with the variable s replaced by z . Thus, all the rules derived for Eq. (5.1) are applicable to Eq. (6.2) and can be used to obtain z -domain root locus plots. The plots can also be obtained using the root locus plots of most computer-aided design (CAD) programs. Thus, we can use the MATLAB command **rlocus** for z -domain root loci.

Example 6.1

Obtain the root locus plot and the critical gain for the first-order type 1 system with loop gain

$$L(z) = \frac{1}{z - 1}$$

Solution

Using root locus rules gives the root locus plot in Fig. 6.1, which can be obtained using the MATLAB command **rlocus**. The root locus lies entirely on the real axis between the open-loop pole and the open-loop zero. For a stable discrete system, real axis z -plane poles must lie between the point $(-1, 0)$ and the point $(1, 0)$. The critical gain for the system corresponds to the point $(-1, 0)$. The closed-loop characteristic equation of the system is

$$z - 1 + K = 0$$

Substituting $z = -1$ gives the critical gain $K_{cr} = 2$, as shown on the root locus plot.

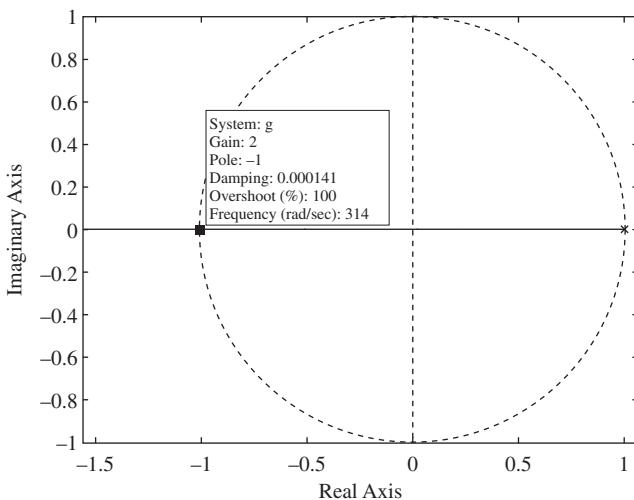
Example 6.1—cont'd

Figure 6.1
 Root locus of a type 1 first-order system.

Example 6.2

Obtain the root locus plot and the critical gain for the second-order type 1 system with loop gain

$$L(z) = \frac{1}{(z-1)(z-0.5)}$$

Solution

Using root locus rules gives the root locus plot in Fig. 6.2, which has the same form as the root locus of Example 5.1 (1) but is entirely in the right-hand plane (RHP). The breakaway point is midway between the two open-loop poles at $z_b = 0.75$. The critical gain now occurs at the intersection of the root locus with the unit circle. To obtain the critical gain value, first write the closed-loop characteristic equation

$$(z-1)(z-0.5) + K = z^2 - 1.5z + K + 0.5 = 0$$

On the unit circle, the closed-loop poles are complex conjugate and of magnitude unity. Hence, the magnitude of the poles satisfies the equation

$$|z_{1,2}|^2 = K_{cr} + 0.5 = 1$$

where K_{cr} is the critical gain. The critical gain is equal to 0.5, which, from the closed-loop characteristic equation, corresponds to unit circle poles at

$$z_{1,2} = 0.75 \pm j0.661$$

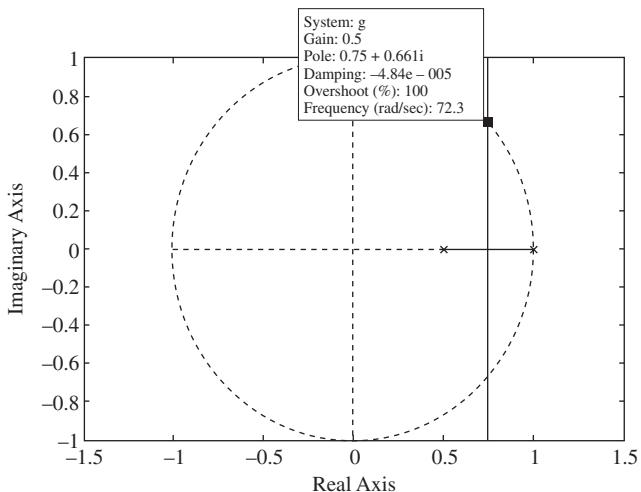
Example 6.2—cont'd

Figure 6.2
z-root locus of a type 1 s-order system.

6.2 z-domain digital control system design

In Chapter 5, we were able to design analog control systems by selecting their poles and zeros in the s -domain to correspond to the desired time response. This approach was based on the relation between any time function and its s -domain poles and zeros. If the time function is sampled and the resulting sequence is z -transformed, the z -transform contains information about the transformed time sequence and the original time function. The poles of the z -domain function can therefore be used to characterize the sequence, and possibly the sampled continuous time function, without inverse z -transformation. However, this latter characterization is generally more difficult than characterization based on the s -domain functions described by Fig. 5.3.

Fig. 6.3 shows z -domain pole locations and the associated temporal sequences. As in the continuous case, positive real poles are associated with exponentials. Unlike the continuous case, the exponentials decay for poles inside the unit circle and increase for poles outside it. In addition, negative real poles are associated with sequences of alternating signs. Poles on the unit circle are associated with a response of constant magnitude. For complex conjugate poles, the response is oscillatory, with the rate of decay determined by the pole distance from the origin and the frequency of oscillations determined by the magnitude of the pole angle. Complex conjugate poles on the unit circle are associated with sustained oscillations.

Because it is easier to characterize the time functions using s -domain poles, it may be helpful to reexamine Fig. 5.3 and compare it to Fig. 6.3. Comparing the figures suggests a

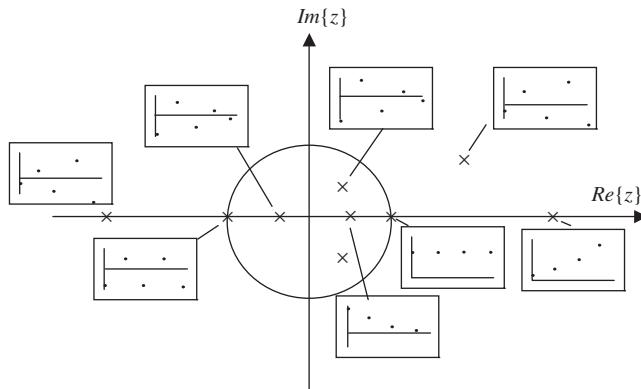


Figure 6.3
z-domain pole locations and the associated temporal sequences.

relationship between s -domain and z -domain poles that greatly simplifies z -domain pole characterization. To obtain the desired relationship, we examine two key cases, both in the s -domain and in the z -domain. One case yields a real pole and the other a complex conjugate pair of poles. More complex time functions can be reduced to these cases by partial fraction expansion of the transforms. The two cases are summarized in [Tables 6.1](#) and [6.2](#).

Using the two tables, it appears that if $F(s)$ has a pole at $-\alpha$, $F(z)$ has a pole at $e^{-\alpha T}$, and if $F(s)$ has poles at $-\zeta\omega_n + j\omega_d$, $F(z)$ has poles at $e^{(-\zeta\omega_n + j\omega_d)T}$. We therefore make the following observation.

Table 6.1: Time functions and real poles.

Continuous	Laplace transform	Sampled	z -Transform
$f(t) = \begin{cases} e^{-\alpha t}, & t \geq 0 \\ 0, & t < 0 \end{cases}$	$F(s) = \frac{1}{s + \alpha}$	$f(kT) = \begin{cases} e^{-\alpha kT}, & k \geq 0 \\ 0, & k < 0 \end{cases}$	$F(z) = \frac{z}{z - e^{-\alpha T}}$

Table 6.2: Time functions and complex conjugate poles.

Continuous	Laplace transform	Sampled	z -Transform
$f(t) = \begin{cases} e^{-\alpha t} \sin(\omega_d t), & t \geq 0 \\ 0, & t < 0 \end{cases}$	$F(s) = \frac{\omega_d}{(s + \alpha)^2 + \omega_d^2}$	$f(kT) = \begin{cases} e^{-\alpha kT} \sin(\omega_d kT), & k \geq 0 \\ 0, & k < 0 \end{cases}$	$F(z) = \frac{\sin(\omega_d T) e^{-\alpha T_z}}{z^2 - 2\cos(\omega_d T) e^{-\alpha T} + e^{-2\alpha T}}$

Observation

If the Laplace transform $F(s)$ of a continuous-time function $f(t)$ has a pole p_s , then the z -transform $F(z)$ of its sampled counterpart $f(kT)$ has a pole at

$$p_z = e^{p_s T} \quad (6.3)$$

where T is the sampling period.

Remarks

1. The preceding observation is valid for a unit step with its s -domain pole at the origin because the z -transform of a sampled step has a pole at $1 = e^0$.
2. There is no general mapping of s -domain zeros to z -domain zeros.
3. From Eq. (6.3), the z -domain poles in the complex conjugate case are given by

$$\begin{aligned} p_z &= e^{\sigma T} e^{j\omega_d T} \\ &= e^{\sigma T} e^{j(\omega_d T + k2\pi)}, \quad k = 0, 1, 2, \dots \end{aligned}$$

Thus, pole locations are a periodic function of the damped natural frequency ω_d with period $(2\pi/T)$ (i.e., the sampling angular frequency ω_s). The mapping of distinct s -domain poles to the same z -domain location is clearly undesirable in situations where a sampled waveform is used to represent its continuous counterpart. The strip of width ω_s over which no such ambiguity occurs (frequencies in the range $[(-\omega_s/2), \omega_s/2]$ rad/s) is known as the *primary strip* (Fig. 6.4). The width of this strip can clearly be increased by faster sampling, with the choice of suitable sampling rate dictated by the nature of the continuous time function. We observe that the minimum sampling frequency for good correlation between the analog and digital signals is twice the frequency ω_d as expected based on the sampling theorem of Section 2.9.

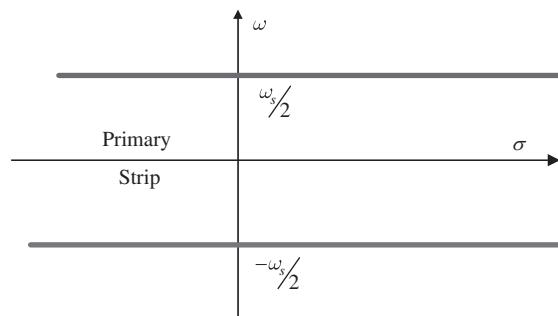


Figure 6.4
Primary strip in the s -plane.

6.2.1 z-domain contours

Using the observation (6.3), one can use s -domain contours over which certain characteristics of the function poles are fixed to obtain the shapes of similar z -domain contours. In particular, the important case of a second-order underdamped system yields [Table 6.3](#).

The information in the table is shown in [Figs. 6.5 and 6.6](#) and can be used to predict the time responses of [Fig. 6.3](#). From [Fig. 6.5](#), we see that negative values of σ correspond to the inside of the unit circle in the z -plane, whereas positive values correspond to the outside of the unit circle. The unit circle is the $\sigma = 0$ contour. Both [Table 6.3](#) and [Fig. 6.5](#) also show that large positive σ values correspond to circles of large radii, whereas large negative σ values correspond to circles of small radii. In particular, a positive infinite σ corresponds to the point at ∞ and a negative infinite σ corresponds to the origin of the z -plane.

From [Fig. 6.6](#), we see that larger ω_d values correspond to larger angles, with $\omega_d = \pm\omega_s/2$ corresponding to $\pm\pi$. As observed earlier using a different argument, a system with poles outside this range (outside the primary strip) does not have a one-to-one correspondence between s -domain and z -domain poles.

Table 6.3: Pole Contours in the s -Domain and the z -Domain.

Contour	s -Domain poles	Contour	z -Domain poles
$\sigma = \text{constant}$	Vertical line	$ z = e^{\sigma T} = \text{constant}$	Circle
$\omega_d = \text{constant}$	Horizontal line	$\angle z = \omega_d T = \text{constant}$	Radial line

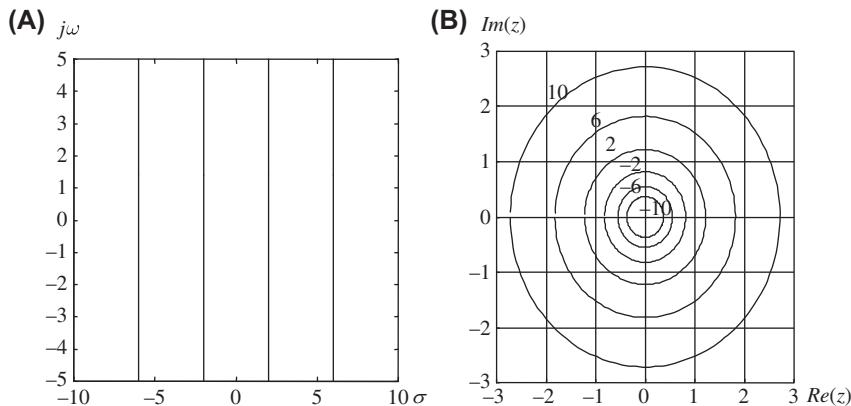


Figure 6.5

Constant σ contours in the s -plane and in the z -plane. (A) Constant σ contours in the s -plane.
(B) Constant σ contours in the z -plane.

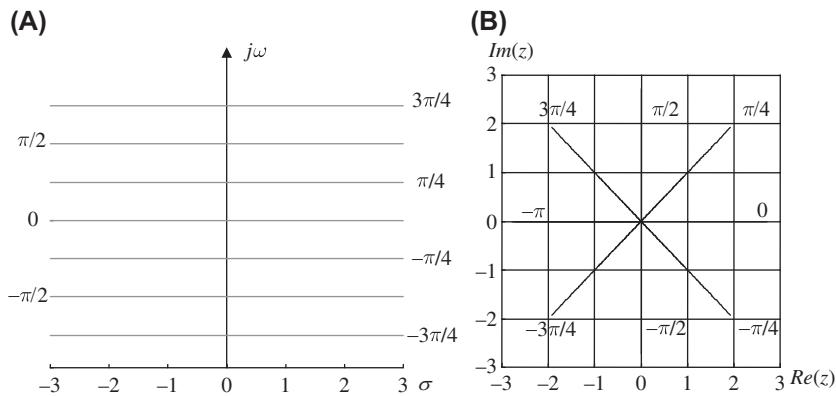


Figure 6.6

Constant ω_d contours in the s -plane and the z -plane. (A) Constant ω_d contours in the s -plane.
(B) Constant ω_d contours in the z -plane.

The z -domain characteristic polynomial for a second-order underdamped system is

$$(z - e^{(-\zeta\omega_n + j\omega_d)T})(z - e^{(-\zeta\omega_n - j\omega_d)T}) = z^2 - 2 \cos(\omega_d T) e^{-\zeta\omega_n T} z + e^{-2\zeta\omega_n T}$$

Hence, the poles of the system are given by

$$z_{1,2} = e^{-\zeta\omega_n T} \angle \pm \omega_d T \quad (6.4)$$

This confirms that constant $\zeta\omega_n$ contours are circles, whereas constant ω_d contours are radial lines.

Constant ζ lines are logarithmic spirals that get smaller for larger values of ζ . The spirals are defined by the equation

$$|Z| = e^{\frac{-\zeta\theta}{\sqrt{1-\zeta^2}}} = e^{\frac{-\zeta(\pi\theta^\circ/180^\circ)}{\sqrt{1-\zeta^2}}} \quad (6.5)$$

where $|z|$ is the magnitude of the pole and θ is its angle. Constant ω_n contours are defined by the equation

$$|Z| = e^{-\sqrt{(\omega_n T)^2 - \theta^2}} \quad (6.6)$$

Fig. 6.7 shows constant ζ contours and constant ω_n contours.

To prove Eq. (6.5), rewrite the complex conjugate poles as

$$z_{1,2} = e^{-\zeta\omega_n T} \angle \pm \omega_d T = |z| \angle \pm \theta \quad (6.7)$$

or equivalently,

$$\theta = \omega_d T = \omega_n T \sqrt{1 - \zeta^2} \quad (6.8)$$

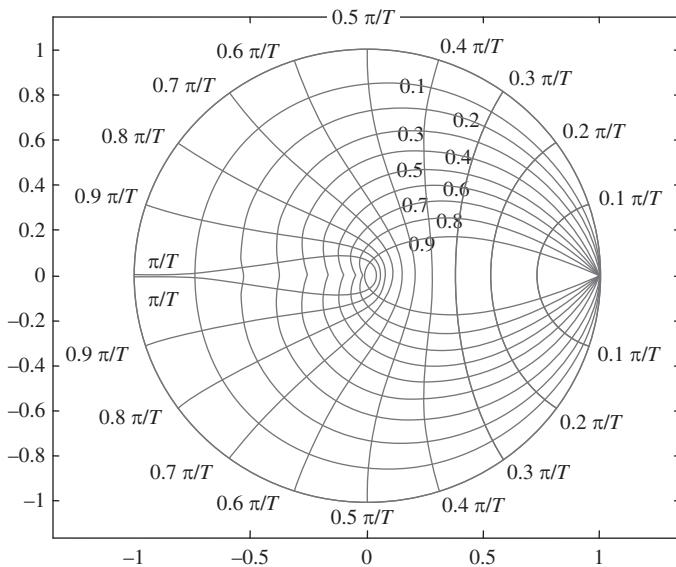


Figure 6.7
Constant ζ and ω_n contours in the z -plane.

$$|z| = e^{-\zeta \omega_n T} \quad (6.9)$$

Eliminating $\omega_n T$ from Eq. (6.9) using (6.8), we obtain the spiral Eq. (6.5). The proof of Eq. (6.6) is similar and is left as an exercise (Problem 6.2).

The following observations can be made by examining the spiral Eq. (6.5):

1. For every ζ value, there are two spirals corresponding to negative and positive angles θ . The negative θ spiral is below the real axis and is the mirror image of the positive θ spiral.
2. For a given spiral, the magnitude of the pole drops logarithmically with its angle.
3. At the same angle θ , increasing the damping ratio gives smaller pole magnitudes. Hence, the spirals are smaller for larger ζ values.
4. All spirals start at $\theta = 0$, $|z| = 1$ but end at different points.
5. For a given damping ratio and angle θ , the pole magnitude can be obtained by substituting in Eq. (6.5). For a given damping ratio and pole magnitude, the pole angle can be obtained by substituting in the equation

$$\theta = \frac{\sqrt{1 - \zeta^2}}{\zeta} |\ln(|z|)| \quad (6.10)$$

The standard contours can all be plotted using a handheld calculator. Their intersection with the root locus can then be used to obtain the pole locations for desired closed-loop characteristics. However, it is more convenient to obtain the root loci and contours using a

CAD tool, especially for higher-order systems. The unit circle and constant ζ and ω_n contours can be added to root locus plots obtained with CAD packages to provide useful information on the significance of z -domain pole locations. In MATLAB, this is accomplished using the command

>> zgrid(zeta, wn)

where **zeta** is a vector of damping ratios and **wn** is a vector of the undamped natural frequencies for the contours.

Clearly, the significance of pole and zero locations in the z -domain is completely different from identical locations in the s -domain. For example, the stability boundary in the z -domain is the unit circle, not the imaginary axis. The characterization of the time response of a discrete-time system based on z -domain information is more complex than the analogous process for continuous-time systems based on the s -domain information discussed in Chapter 5. These are factors that slightly complicate z -domain design, although the associated difficulties are not insurmountable.

The specifications for z -domain design are similar to those for s -domain design. Typical design specifications are as follows:

Time constant. This is the time constant of exponential decay for the continuous envelope of the sampled waveform. The sampled signal is therefore not necessarily equal to a specified portion of the final value after one time constant. The time constant is defined as

$$\tau = \frac{1}{\zeta \omega_n} \quad (6.11)$$

Settling time. The settling time is defined as the period after which the envelope of the sampled waveform stays within a specified percentage (usually 2%) of the final value. It is a multiple of the time constant depending on the specified percentage. For a 2% specification, the settling time is given by

$$T_s = \frac{4}{\zeta \omega_n} \quad (6.12)$$

Frequency of oscillations ω_d . This frequency is equal to the angle of the dominant complex conjugate poles divided by the sampling period.

Other design criteria such as the percentage overshoot, the damping ratio, and the undamped natural frequency can also be defined analogously to the continuous case.

6.2.2 Proportional control design in the z -domain

Proportional control involves the selection of a DC gain value that corresponds to a time response satisfying design specifications. As in s -domain design, a satisfactory time

response is obtained by tuning the gain to select a dominant closed-loop pair in the appropriate region of the complex plane. Analytical design is possible for low-order systems but is more difficult than its analog counterpart. Example 6.3 illustrates the design of proportional digital controllers.

Example 6.3

Design a proportional controller for the digital system described in Example 6.2 with a sampling period $T = 0.1$ s to obtain

1. A damped natural frequency of 5 rad/s
2. A time constant of 0.5 s
3. A damping ratio of 0.7

Solution

After some preliminary calculations, the design results can be easily obtained using the **rlocus** command of MATLAB. The following calculations, together with the information provided by a cursor command, allow us to determine the desired closed-loop pole locations:

1. The angle of the pole is $\omega_d T = 5 \times 0.1 = 0.5$ rad or 28.65 degrees.
2. The reciprocal of the time constant is $\zeta \omega_n = 1/0.5 = 2$ rad/s. This yields a pole magnitude of $e^{-\zeta \omega_n T} = 0.82$.
3. The damping ratio given can be used directly to locate the desired pole.

Using MATLAB, we obtain the results shown in [Table 6.4](#). The corresponding sampled step response plots obtained using the command **step** (MATLAB) are shown in [Fig. 6.8](#). As expected, the higher gain designs are associated with a low damping ratio and a more oscillatory response.

[Table 6.4](#) results can also be obtained analytically using the characteristic equation for the complex conjugate poles of [Eq. \(6.4\)](#). The system's closed-loop characteristic equation is

$$z^2 - 1.5z + K + 0.5 = z^2 - 2 \cos(\omega_d T) e^{-\zeta \omega_n T} z + e^{-2\zeta \omega_n T}$$

Equating coefficients gives the two equations

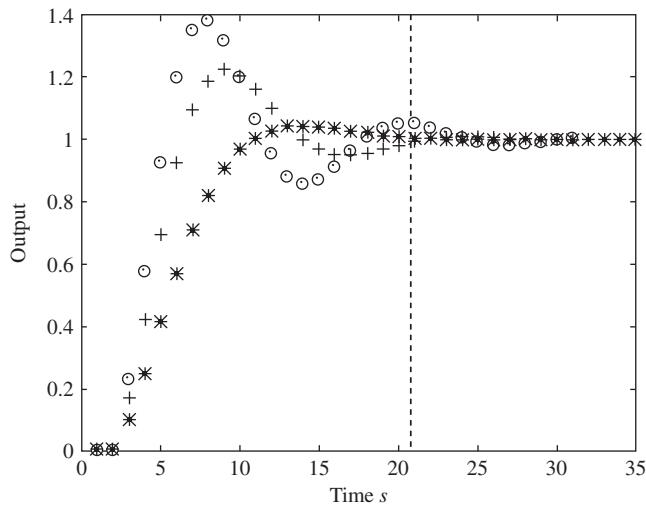
$$\begin{aligned} z^1: \quad 1.5 &= 2 \cos(\omega_d T) e^{-\zeta \omega_n T} \\ z^0: \quad K + 0.5 &= e^{-2\zeta \omega_n T} \end{aligned}$$

1. From the z^1 equation,

$$\zeta \omega_n = \frac{1}{T} \ln \left(\frac{1.5}{2 \cos(\omega_d T)} \right) = 10 \left| \ln \left(\frac{1.5}{2 \cos(0.5)} \right) \right| = 1.571$$

Table 6.4: Proportional control design results.

Design	Gain	ζ	ω_n rad/s
(a)	0.23	0.3	5.24
(b)	0.17	0.4	4.60
(c)	0.10	0.7	3.63

Example 6.3—cont'd**Figure 6.8**

Time response for the designs of Table 6.4: (a) \odot , (b) $+$, (c) $*$.

In addition,

$$\omega_d^2 = \omega_n^2(1 - \zeta^2) = 25$$

Hence, we obtain the ratio

$$\frac{\omega_d^2}{(\zeta\omega_n)^2} = \frac{1 - \zeta^2}{\zeta^2} = \frac{25}{(1.571)^2}$$

This gives a damping ratio $\zeta = 0.3$ and an undamped natural frequency $\omega_n = 5.24 \text{ rad/s}$. Finally, the z^0 equation gives a gain

$$K = e^{-2\zeta\omega_n T} - 0.5 = e^{-2 \times 1.571 \times 0.1} - 0.5 = 0.23$$

2. From Eq. (6.11) and the z^1 equation, we obtain

$$\zeta\omega_n = \frac{1}{\tau} = \frac{1}{0.5} = 2 \text{ rad/s}$$

$$\omega_d = \frac{1}{T} \cos^{-1} \left(\frac{1.5e^{\zeta\omega_n T}}{2} \right) = 10 \cos^{-1} (0.75e^{0.2}) = 4.127 \text{ rad/s}$$

Solving for ζ in the equality

$$\frac{\omega_d^2}{(\zeta\omega_n)^2} = \frac{1 - \zeta^2}{\zeta^2} = \frac{(4.127)^2}{2^2}$$

gives a damping ratio $\zeta = 0.436$ and an undamped natural frequency $\omega_n = 4.586 \text{ rad/s}$. The gain for this design is

$$K = e^{-2\zeta\omega_n T} - 0.5 = e^{-2 \times 2 \times 0.1} - 0.5 = 0.17$$

Example 6.3—cont'd

3. For a damping ratio of 0.7, the z^1 equation obtained by equating coefficients remains nonlinear and is difficult to solve analytically. The equation now becomes

$$1.5 = 2 \cos(0.0714\omega_n) e^{-0.07\omega_n}$$

The equation can be solved numerically by trial and error with a calculator to obtain the undamped natural frequency $\omega_n = 3.63$ rad/s. The gain for this design is

$$K = e^{-0.14 \times 3.63} - 0.5 = 0.10$$

This controller can also be designed graphically by drawing the root locus and a segment of the constant ζ spiral and finding their intersection. But the results obtained graphically are often very approximate, and the solution is difficult for all but a few simple root loci.

Example 6.4

Consider the vehicle position control system of Example 3.3 with the transfer function

$$G(s) = \frac{1}{s(s+5)}$$

Design a proportional controller for the unity feedback digital control system with analog process and a sampling period $T = 0.04$ s to obtain

1. A steady-state error of 10% due to a ramp input
2. A damping ratio of 0.7

Solution

The analog transfer function together with a DAC and ADC has the z -transfer function

$$G_{ZAS}(z) = \frac{7.4923 \times 10^{-4}(z + 0.9355)}{(z - 1)(z - 0.8187)}$$

and the closed-loop characteristic equation is

$$\begin{aligned} 1 + KG_{ZAS}(z) &= z^2 - (1.8187 - 7.4923 \times 10^{-4}K)z + 0.8187 - 7.009 \times 10^{-4}K \\ &= z^2 - 2 \cos(\omega_d T) e^{-\zeta \omega_n T} z + e^{-2\zeta \omega_n T} \end{aligned}$$

The equation involves three parameters: ζ , ω_n , and K . As in Example 6.3, equating coefficients yields two equations that we can use to evaluate two unknowns. The third parameter must be obtained from a design specification.

1. The system is type 1, and the velocity error constant is

$$\begin{aligned} K_v &= \frac{1}{T} \left. \frac{z-1}{z} KG(z) \right|_{z=1} \\ &= K \frac{7.4923 \times 10^{-4}(1 + 0.9355)}{(0.04)(1 - 0.8187)} \\ &= \frac{K}{5} \end{aligned}$$

Example 6.4—cont'd

This is identical to the velocity error constant for the analog proportional control system. In both cases, a steady-state error due to ramp of 10% is achieved with

$$\begin{aligned}\frac{K}{5} &= K_v = \frac{100}{e(\infty)\%} \\ &= \frac{100}{10} = 10\end{aligned}$$

Hence, the required gain is $K = 50$.

2. As in Example 6.3, the analytical solution for constant ζ is difficult. But the design results are easily obtained using the root locus cursor command of any CAD program. As shown in Fig. 6.9, moving the cursor to the $\zeta = 0.7$ contour yields a gain value of approximately 11.7.

The root locus of Fig. 6.10 shows that the critical gain for the system K_{cr} is approximately 268, and the system is therefore stable at the desired gain and can meet the design specifications for both (1) and (2). However, other design criteria, such as the damping ratio and the undamped natural frequency, should be checked. Their values can be obtained using a CAD cursor command or by equating characteristic equation coefficients as in Example 6.3. For the gain of 50 selected in (1), the root locus of Fig. 6.10 and the cursor give a damping ratio of 0.28. This corresponds to the highly oscillatory response of Fig. 6.11, which is likely to be unacceptable in practice. For the gain of 11.7 selected in (2), the steady-state error is 42.7% due to a unit ramp. It is therefore clear that to obtain the steady-state error specified together with an acceptable transient response, proportional control is inadequate.

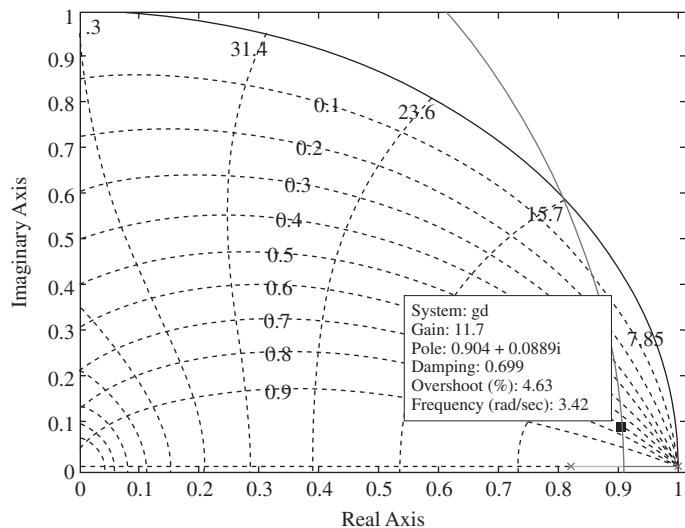


Figure 6.9
Root locus for the constant ζ design.

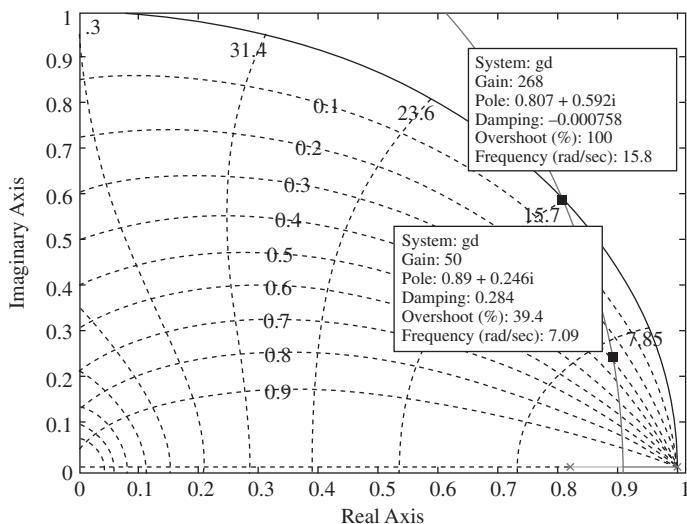
Example 6.4—cont'd

Figure 6.10
Root locus for $K = 50$.

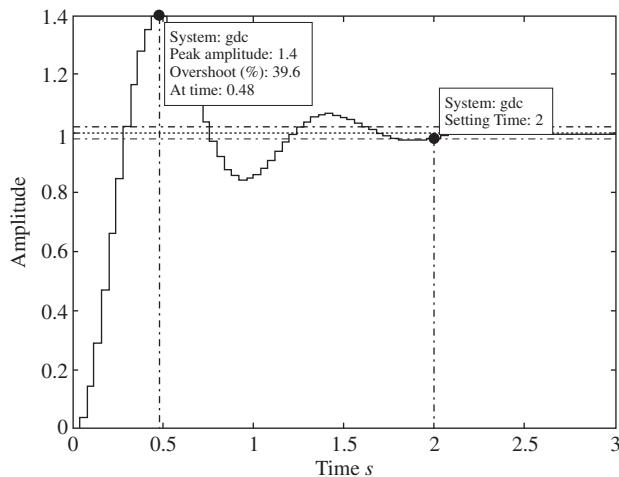


Figure 6.11
Time response for $K = 50$.

6.3 Digital implementation of analog controller design

This section introduces an indirect approach to digital controller design. The approach is based on designing an analog controller for the analog subsystem and then obtaining an equivalent digital controller and using it to digitally implement the desired control. The

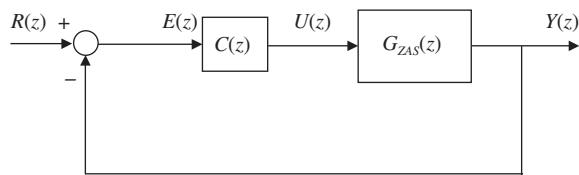


Figure 6.12

Block diagram of a single-loop digital control system.

digital controller can be obtained using a number of recipes that are well known in the field of signal processing, where they are used in the design of digital filters. In fact, a controller can be viewed as a filter that attenuates some dynamics and accentuates others so as to obtain the desired time response. We limit our discussion of digital filters and the comparison of various recipes for obtaining them from analog filters to differencing methods, pole-zero matching, and bilinear transformation. The system configuration we consider is shown in Fig. 6.12. The system includes (1) a z -transfer function model of a DAC, analog subsystem, and ADC and (2) a cascade controller. We begin with a general procedure to obtain a digital controller using analog design.

Procedure 6.1

1. Design a controller $C_a(s)$ for the analog subsystem to meet the desired design specifications.
2. Map the analog controller to a digital controller $C(z)$ using a suitable transformation.
3. Tune the gain of the transfer function $C(z)G_{ZAS}(z)$ using proportional z -domain design to meet the design specifications.
4. Check the sampled time response of the digital control system and repeat steps 1 to 3, if necessary, until the design specifications are met.

Step 2 of Procedure 6.1—that is, the transformation from an analog to a digital filter—must satisfy the following requirements:

1. A stable analog filter (poles in the left half plane (LHP)) must transform to a stable digital filter.
2. The frequency response of the digital filter must closely resemble the frequency response of the analog filter in the frequency range $0 \rightarrow \omega_s/2$ where ω_s is the sampling frequency.

Most filter transformations satisfy these two requirements to varying degrees. However, this is not true of all analog-to-digital transformations, as illustrated by the following section.

6.3.1 Differencing methods

An analog filter can be represented by a transfer function or differential equation.

Numerical analysis provides standard approximations of the derivative so as to obtain the solution to a differential equation. The approximations reduce a differential equation to a difference equation and could thus be used to obtain the difference equation of a digital

filter from the differential equation of an analog filter. We examine two approximations of the derivative: forward differencing and backward differencing.

Forward differencing

The forward differencing approximation of the derivative is

$$\dot{y}(k) \cong \frac{1}{T} [y(k+1) - y(k)] \quad (6.13)$$

The approximation of the second derivative can be obtained by applying Eq. (6.13) twice—that is,

$$\begin{aligned}\ddot{y}(k) &\cong \frac{1}{T} [\dot{y}(k+1) - \dot{y}(k)] \\ &\cong \frac{1}{T} \left\{ \frac{1}{T} [y(k+2) - y(k+1)] - \frac{1}{T} [y(k+1) - y(k)] \right\} \\ &= \frac{1}{T^2} \{y(k+2) - 2y(k+1) + y(k)\}\end{aligned} \quad (6.14)$$

Approximations of higher-order derivatives can be similarly obtained. Alternatively, one may consider the Laplace transform of the derivative and the z -transform of the difference in Eq. (6.13). This yields the mapping

$$sY(s) \rightarrow \frac{1}{T} [z - 1]Y(z) \quad (6.15)$$

Therefore, the direct transformation of an s -transfer function to a z -transfer function is possible using the substitution

$$s \rightarrow \frac{z - 1}{T} \quad (6.16)$$

Example 6.5: Forward difference

Apply the forward difference approximation of the derivative to the second-order analog filter

$$C_a(s) = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2}$$

and examine the stability of the resulting digital filter for a stable analog filter.

Solution

The filter is equivalent to the differential equation

$$\dot{y}(t) + 2\zeta\omega_n\dot{y}(t) + \omega_n^2 y(t) = \omega_n^2 u(t)$$

where $y(t)$ is the filter output and $u(t)$ is the filter input. The approximation of the first derivative by Eq. (6.13) and the second derivative by Eq. (6.14) gives the difference equation

$$\frac{1}{T^2} \{y(k+2) - 2y(k+1) + y(k)\} + 2\zeta\omega_n \frac{1}{T} [y(k+1) - y(k)] + \omega_n^2 y(k) = \omega_n^2 u(k)$$

Example 6.5: Forward difference—cont'd

Multiplying by T^2 and rearranging terms, we obtain the digital filter

$$y(k+2) + 2[\zeta\omega_n T - 1]y(k+1) + \left[(\omega_n T)^2 - 2\zeta\omega_n T + 1\right]y(k) = (\omega_n T)^2 u(k)$$

Equivalently, we obtain the transfer function of the filter using the simpler transformation Eq. (6.16)

$$\begin{aligned} C(z) &= \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2} \Big|_s = \frac{z-1}{T} \\ &= \frac{(\omega_n T)^2}{z^2 + 2[\zeta\omega_n T - 1]_z + \left[(\omega_n T)^2 - 2\zeta\omega_n T + 1\right]} \end{aligned}$$

For a stable analog filter, we have $\zeta > 0$ and $\omega_n > 0$ (positive denominator coefficients are sufficient for a second-order polynomial). However, the digital filter is unstable if the magnitude of the constant term in its denominator polynomial is greater than unity. This gives the instability condition

$$\begin{aligned} (\omega_n T)^2 - 2\zeta\omega_n T + 1 &> 1 \\ \text{i.e., } \zeta &< \omega_n T / 2 \end{aligned}$$

For example, a sampling period of 0.2 s and an undamped natural frequency of 10 rad/s yield unstable filters for any underdamped analog filter.

Backward differencing

The backward differencing approximation of the derivative is

$$\dot{y}(k) \cong \frac{1}{T} [y(k) - y(k-1)] \quad (6.17)$$

The approximation of the second derivative can be obtained by applying Eq. (6.17) twice—that is,

$$\begin{aligned} \ddot{y}(k) &\cong \frac{1}{T} [\dot{y}(k) - \dot{y}(k-1)] \\ &\cong \frac{1}{T} \left\{ \frac{1}{T} [y(k) - y(k-1)] - \frac{1}{T} [y(k-1) - y(k-2)] \right\} \\ &= \frac{1}{T^2} \{y(k) - 2y(k-1) + y(k-2)\} \end{aligned} \quad (6.18)$$

Approximations of higher-order derivatives can be similarly obtained. One may also consider the Laplace transform of the derivative and the z -transform of the difference in Eq. (6.17). This yields the substitution

$$s \rightarrow \frac{z-1}{zT} \quad (6.19)$$

Example 6.6: Backward difference

Apply the backward difference approximation of the derivative to the second-order analog filter

$$C_a(s) = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2}$$

and examine the stability of the resulting digital filter for a stable analog filter.

Solution

We obtain the transfer function of the filter using Eq. (6.19)

$$\begin{aligned} C(z) &= \left. \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2} \right|_{s=\frac{z-1}{zT}} \\ &= \frac{(\omega_n T z)^2}{[(\omega_n T)^2 - 2\zeta\omega_n T + 1]z^2 - 2[\zeta\omega_n T + 1]z + 1} \end{aligned}$$

The stability conditions for the digital filter are (see Chapter 4)

$$\begin{aligned} &[(\omega_n T)^2 + 2\zeta\omega_n T + 1] + 2[\zeta\omega_n T + 1] + 1 > 0 \\ &[(\omega_n T)^2 + 2\zeta\omega_n T + 1] - 1 > 0 \\ &[(\omega_n T)^2 + 2\zeta\omega_n T + 1] - 2[\zeta\omega_n T + 1] + 1 > 0 \end{aligned}$$

The conditions are all satisfied for $\zeta > 0$ and $\omega_n > 0$ —that is, for all stable analog filters.

6.3.2 Pole-zero matching

We know from Eq. (6.3) that discretization maps an s -plane pole at p_s to a z -plane pole at $e^{p_s T}$ but that no rule exists for mapping zeros. In pole-zero matching, a discrete approximation is obtained from an analog filter by mapping both poles and zeros using Eq. (6.3). If the analog filter has n poles and m zeros, then we say that the filter has $n-m$ zeros at infinity. For $n-m$ zeros at infinity, we add $n-m$ or $n-m-1$ digital filter zeros at unity. If the zeros are not added, it can be shown that the resulting system will include a time delay (see Problem 6.6). The second choice gives a strictly proper filter where the computation of the output is easier, since it only requires values of the input at past sampling points. Finally, we adjust the gain of the digital filter so that it is equal to that of the analog filter at a critical frequency dependent on the filter. For a low-pass filter, α is selected so that the gains are equal at DC; for a bandpass filter, they are set equal at the center of the pass band.

For an analog filter with transfer function

$$G_a(s) = K \frac{\prod_{i=1}^m (s - a_i)}{\prod_{j=1}^n (s - b_j)} \quad (6.20)$$

we have the digital filter

$$G(z) = \alpha K \frac{(z + 1)^{n-m-1} \prod_{i=1}^m (z - e^{a_i T})}{\prod_{j=1}^n (z - e^{b_j T})} \quad (6.21)$$

where α is a constant selected for equal filter gains at a critical frequency. For example, for a low-pass filter, α is selected to match the DC gains using $G(1) = G_a(0)$, while for a high-pass filter, it is selected to match the high-frequency gains using $G(-1) = G_a(\infty)$. Setting $z = e^{j\omega T} = -1$ (i.e., $\omega T = \pi$) is equivalent to selecting the folding frequency $\omega_s/2$, which is the highest frequency allowable without aliasing. Pole-zero matched digital filters can be obtained using the MATLAB command

```
>> g = c2d(ga, T, 'matched')
```

Example 6.7

Find a pole-zero matched digital filter approximation for the analog filter

$$G_a(s) = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2}$$

If the damping ratio is equal to 0.5 and the undamped natural frequency is 5 rad/s, determine the transfer function of the digital filter for a sampling period of 0.1 s. Check your answer using MATLAB and obtain the frequency response of the digital filter.

Solution

The filter has a zero at the origin and complex conjugate poles at $s_{1,2} = -\zeta\omega_n \pm j\omega_d$. We apply the pole-zero matching transformation to obtain

$$G(z) = \frac{\alpha(z + 1)}{z^2 - 2e^{-\zeta\omega_n T} \cos(\omega_d T)z + e^{-2\zeta\omega_n T}}$$

The analog filter has two zeros at infinity, and we choose to add one digital filter zeros at -1 for a strictly proper filter. The difference equation for the filter is

$$\begin{aligned} y(k+2) &= 2e^{-\zeta\omega_n T} \cos(\omega_d T)y(k+1) - e^{-2\zeta\omega_n T}y(k) \\ &\quad + \alpha(u(k+1) - u(k)) \end{aligned}$$

Thus, the computation of the output only requires values of the input at earlier sampling points, and the filter is easily implementable. The gain α is selected in order for the digital filter to have the same DC gain as the analog filter, which is equal to unity.

For the given numerical values, we have the damping ratio $\omega_d = 5\sqrt{1 - 0.5^2} = 4.33$ rad/s and the filter transfer function

$$G(z) = \frac{0.09634(z + 1)}{z^2 - 1.414z + 0.6065}$$

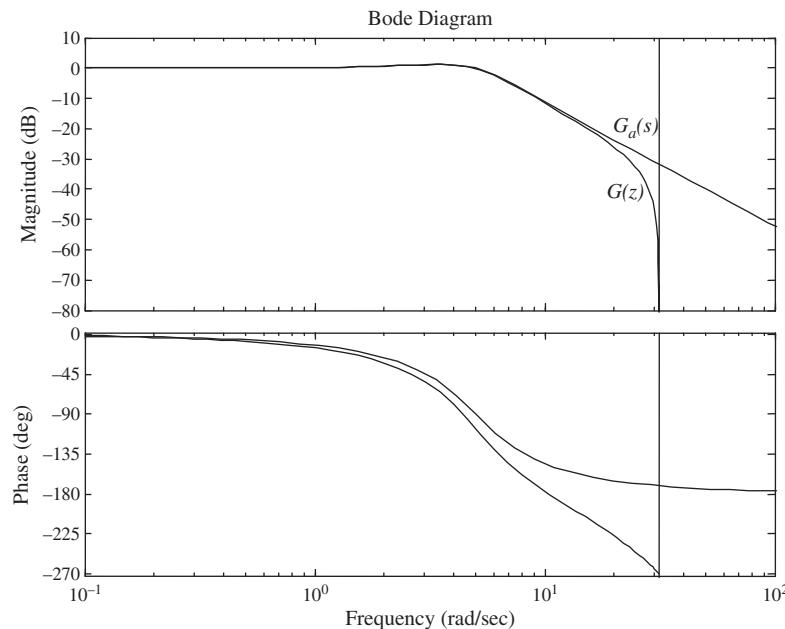


Figure 6.13
Frequency response of digital filter for Example 6.7.

The following MATLAB commands give the transfer function:

```
>>wn = 5; zeta = 0.5; % Undamped natural frequency, damping ratio.
>>ga = tf([wn^2],[1,2*zeta*wn,wn^2]) % Analog transfer function.
Transfer function:
25.
-
s^2 + 5 s + 25.
>> g = c2d(ga,.1,'matched') % Transformation with a sampling period 0.1.
Transfer function:
0.09634 z + 0.09634
-
z^2 - 1.414 z + 0.6065
Sampling time: 0.1
```

The frequency responses of the analog and digital filters obtained using MATLAB are shown in Fig. 6.13. Note that the frequency responses are almost identical in the low frequency range but become different at high frequencies.

6.3.3 Bilinear transformation

The relationship

$$s = c \frac{z - 1}{z + 1} \quad (6.22)$$

with a linear numerator and a linear denominator and a constant scale factor c is known as a bilinear transformation. The relationship can be obtained from the equality $z = e^{sT}$ using the first-order approximation

$$s = \frac{1}{T} \ln(z) \approx \frac{2}{T} \left[\frac{z - 1}{z + 1} \right] \quad (6.23)$$

where the constant $c = 2/T$. A digital filter $C(z)$ is obtained from an analog filter $C_a(s)$ by the substitution

$$C(z) = C_a(s) \Big|_{s=c\left[\frac{z-1}{z+1}\right]} \quad (6.24)$$

The resulting digital filter has the frequency response

$$\begin{aligned} C(e^{j\omega T}) &= C_a(s) \Big|_{s=c\left[\frac{e^{j\omega T}-1}{e^{j\omega T}+1}\right]} \\ &= C_a \left(c \left[\frac{e^{j\omega T/2} - e^{-j\omega T/2}}{e^{j\omega T/2} + e^{-j\omega T/2}} \right] \right) \end{aligned}$$

Thus, the frequency responses of the digital and analog filters are related by

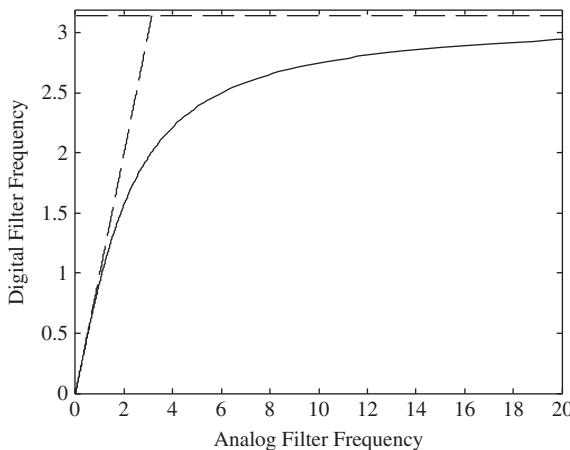
$$C(e^{j\omega T}) = C_a \left(j \tan \left[\frac{\omega T}{2} \right] \right) \quad (6.25)$$

Evaluating the frequency response at the folding frequency $\omega_s/2$ gives

$$\begin{aligned} C(e^{j\omega_s T/2}) &= C_a \left(j \tan \left[\frac{\omega_s T}{4} \right] \right) \\ &= C_a \left(j \tan \left[\frac{2\pi}{4} \right] \right) = C_a(j\infty) \end{aligned}$$

We observe that bilinear mapping squeezes the entire frequency response of the analog filter for a frequency range $0 \rightarrow \infty$ into the frequency range $0 \rightarrow \omega_s/2$. This implies the absence of aliasing (which makes the bilinear transformation a popular method for digital filter design) but also results in distortion or warping of the frequency response. The relationship between the frequency ω_a of the analog filter and the associated frequency ω of the digital filter for the case $c = 2/T$ —namely,

$$\omega_a = \frac{2}{T} \tan \left(\frac{\omega T}{2} \right)$$

**Figure 6.14**

Relationship between analog filter frequencies ω_a and the associated digital filter frequencies with bilinear transformation.

is plotted in Fig. 6.14 for $T = 1$. Note that, in general, if the sampling period is sufficiently small so that $\omega \ll \pi/T$, then

$$\tan\left(\frac{\omega T}{2}\right) \approx \frac{\omega T}{2}$$

and therefore $\omega_a \approx \omega$, so that the effect of the warping is negligible.

In any case, the distortion of the frequency response can be corrected at a single frequency ω_0 using the **prewarping** equality

$$C(e^{j\omega_0 T}) = C_a\left(j \tan\left[\frac{\omega_0 T}{2}\right]\right) = C_a(j\omega_0) \quad (6.26)$$

The equality holds provided that the constant c is chosen as

$$c = \frac{\omega_0}{\tan\left(\frac{\omega_0 T}{2}\right)} \quad (6.27)$$

The choice of the prewarping frequency ω_0 depends on the mapped filter. In control applications, a suitable choice of ω_0 is the 3-dB frequency for a PI or PD controller and the upper 3-dB frequency for a PID controller. This is explored further in design examples.

In MATLAB, the bilinear transformation is accomplished using the following command:

```
>> gd = c2d(g, tc, 'tustin')
```

where \mathbf{g} is the analog system and \mathbf{tc} is the sampling period. If prewarping is requested at a frequency \mathbf{w} , then the command is

```
>> gd = c2d(g, tc, 'prewarp', w)
```

Example 6.8

Design a digital filter by applying the bilinear transformation to the analog filter

$$C_a(s) = \frac{1}{0.1s + 1} \quad (6.28)$$

with $T = 0.1$ s. Examine the warping effect and then apply prewarping at the 3-dB frequency.

Solution

By applying the bilinear transformation Eqs. (6.22) to (6.28), we obtain

$$C(z) = \frac{1}{0.1 \frac{z-1}{z+1} + 1} = \frac{z+1}{3z-1}$$

The Bode plots of $C_a(s)$ (solid line) and $C(z)$ (dash-dot line) are shown in Fig. 6.15, where the warping effect can be evaluated. We select the 3-dB frequency $\omega_0 = 10$ as a prewarping frequency and apply Eq. (6.27) to obtain

$$C(z) = \frac{1}{0.1 \frac{10}{\tan(\frac{10 \cdot 0.01}{2})} \frac{z-1}{z+1} + 1} \cong \frac{0.35z + 0.35}{z - 0.29}$$

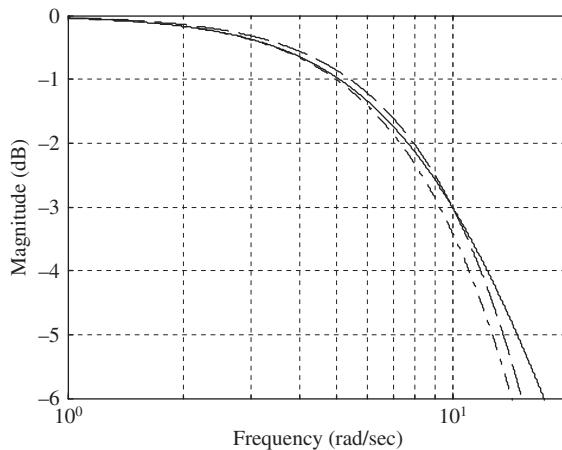


Figure 6.15

Bode plots of the analog filter (*solid*) and the digital filter obtained with (*dashed*) and without prewarping (*dash-dot*).

The corresponding Bode plot is shown again in Fig. 6.15 (dashed line). It coincides with the Bode plot of $C(s)$ at $\omega_0 = 10$. Note that for lower values of the sampling period, the three Bode plots tend to coincide.

Another advantage of bilinear transformation is that it maps points in the LHP to points inside the unit circle and thus guarantees the stability of a digital filter for a stable analog filter. This property was discussed in Section 4.42 and is clearly demonstrated in Figure 4.4.

Bilinear transformation of the analog PI controller gives the following digital PI controller:

$$\begin{aligned} C(z) &= K \frac{(s+a)}{s} \Big|_{s=c} \left[\frac{z-1}{z+1} \right] \\ &= K \left(\frac{a+c}{c} \right) \frac{z + \left(\frac{a-c}{a+c} \right)}{z-1} \end{aligned} \quad (6.29)$$

The digital PI controller increases the type of the system by one and can therefore be used to improve steady-state error. As in the analog case, it has a zero that reduces the deterioration of the transient response due to the increase in system type. The PI controller of Eq. (6.29) has a numerator order equal to its denominator order. Hence, the calculation of its output from its difference equation requires knowledge of the input at the current time. Assuming negligible computational time, the controller is approximately realizable.

Bilinear transformation of the analog PD controller gives the digital PD controller

$$\begin{aligned} C(z) &= K(s+a) \Big|_{s=c} \left[\frac{z-1}{z+1} \right] \\ &= K(a+c) \frac{z + \left(\frac{a-c}{a+c} \right)}{z+1} \end{aligned} \quad (6.30)$$

This includes a zero that can be used to improve the transient response and a pole at $z = -1$ that occurs because the continuous time system is not proper (see Problem 6.9). A pole at $z = -1$ corresponds to an unbounded frequency response at the folding frequency, as $e^{j\omega_s T/2} = e^{j\pi} = -1$, and must therefore be eliminated. However, eliminating the undesirable pole would result in an unrealizable controller. An approximately realizable PD controller is obtained by replacing the pole at $z = -1$ with a pole at the origin to obtain

$$C(z) = K(a+c) \frac{z + \left(\frac{a-c}{a+c} \right)}{z} \quad (6.31)$$

A pole at the origin is associated with a term that decays as rapidly as possible so as to have the least effect on the controller dynamics. However, this variation from direct transformation results in additional distortion of the analog filter and complication of the digital controller design and doubles the DC gain $C(1)$ of the controller. To provide the best approximation of the continuous-time controller, disregarding subsequent gain tuning, the gain K can be halved.

Bilinear transformation of the analog PID controller gives the digital PD controller

$$\begin{aligned} C(z) &= K \frac{(s+a)(s+b)}{s} \Big|_{s=c} \left[\frac{z-1}{z+1} \right] \\ &= K \frac{(a+c)(b+c)}{c} \frac{\left[z + \left(\frac{a-c}{a+c} \right) \right] \left[z + \left(\frac{b-c}{b+c} \right) \right]}{(z+1)(z-1)} \end{aligned}$$

The controller has two zeros that can be used to improve the transient response and a pole at $z = 1$ to improve the steady-state error. As with PD control, transforming an improper transfer function yields a pole at $z = -1$, which must be replaced by a pole at the origin to yield a transfer function with a bounded frequency response at the folding frequency. The resulting transfer function is approximately realizable and is given by

$$C(z) = K \frac{(a+c)(b+c)}{c} \frac{\left[z + \left(\frac{a-c}{a+c} \right) \right] \left[z + \left(\frac{b-c}{b+c} \right) \right]}{z(z-1)} \quad (6.32)$$

As in the case of PD control, the modification of the bilinearly transformed transfer function results in distortion and doubling the DC gain $C(1)$ that can be reduced by halving the gain K .

Using Procedure 6.1 and Eqs. (6.29), (6.31) and (6.32), respectively, digital PI, PD, and PID controllers can be designed to yield satisfactory transient and steady-state performance.

Example 6.9

Design a digital controller for a DC motor speed control system (see Example 3.6) where the (type 0) analog plant has the transfer function

$$G(s) = \frac{1}{(s+1)(s+10)}$$

to obtain zero steady-state error due to a unit step, percentage overshoot less than 5%, and a settling time of about 1 s.

Solution

The design is completed following Procedure 6.1. First, an analog controller is designed for the given plant. For zero steady-state error due to unit step, the system type must be increased by one. A PI controller effects this increase, but the location of its zero must be chosen so as to obtain an acceptable transient response. The simplest possible design is obtained by pole-zero cancellation and is of the form

$$C_a(s) = K \frac{s+1}{s}$$

The corresponding loop gain is

$$C_a(s)G(s) = \frac{K}{s(s+10)}$$

Hence, the closed-loop characteristic equation of the system is

$$s(s+10) + K = s^2 + 2\zeta\omega_n s + \omega_n^2$$

Example 6.9—cont'd

Equating coefficients gives $\zeta\omega_n = 5$ rad/s and the settling time

$$T_s = \frac{4}{\zeta\omega_n} = \frac{4}{5} = 0.8 \text{ s}$$

as required.

For percentage overshoot of 5%, we need a damping ratio

$$\zeta = \frac{|\ln(0.05)|}{\sqrt{|\ln(0.05)|^2 + \pi^2}} \approx 0.69$$

For percentage overshoot less than 5%, we choose a damping ratio of 0.7. The damping ratio of the analog system can be set equal to 0.7 by appropriate choice of the gain K . The gain selected at this stage must often be tuned after filter transformation to obtain the same damping ratio for the digital controller. We solve for the undamped natural frequency

$$\omega_n = 10/(2\zeta) = 10/(2 \times 0.7) = 7.142 \text{ rad/s}$$

The corresponding analog gain is

$$K = \omega_n^2 = 51.02$$

We therefore have the analog filter

$$C_a(s) = 51.02 \frac{s+1}{s}$$

Next, we select a suitable sampling period for an undamped natural frequency of about 7.14 rad/s. We select $T = 0.02 \text{ s} < 2\pi/(40\omega_d)$, which corresponds to a sampling frequency higher than 40 times the damped natural frequency (see Chapter 2). The model of the analog plant together with an ADC and sampler is

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} \\ &= 1.8604 \times 10^{-4} \frac{z + 0.9293}{(z - 0.8187)(z - 0.9802)} \end{aligned}$$

Bilinear transformation of the PI controller, with gain K included as a free parameter, gives

$$C(z) = 1.01K \frac{z - 0.9802}{z - 1}$$

Because the analog controller was obtained using pole-zero cancellation, near pole-zero cancellation occurs when the digital controller $C(z)$ is multiplied by $G_{ZAS}(z)$. The gain can now be tuned for a damping ratio of 0.7 using a CAD package with the root locus of the loop gain $C(z)G_{ZAS}(z)$. From the root locus, shown in Fig. 6.16, at $\zeta = 0.7$, the gain K is about 46.7, excluding the 1.01 gain of $C(z)$ (i.e., a net gain of 47.2). The undamped natural frequency is $\omega_n = 6.85 \text{ rad/s}$. This yields the approximate settling time

$$\begin{aligned} T_s &= \frac{4}{\zeta\omega_n} = \frac{4}{6.85 \times 0.7} \\ &= 0.83 \text{ s} \end{aligned}$$

Example 6.9—cont'd

The settling time is acceptable but is slightly worse than the settling time for the analog controller. The step response of the closed-loop digital control system shown in Fig. 6.17 is also acceptable and confirms the estimated settling time. The net gain value of 47.2, which meets the design specifications, is significantly less than the gain value of 51.02 for the analog design. This demonstrates the need for tuning the controller gain after mapping the analog controller to a digital controller.

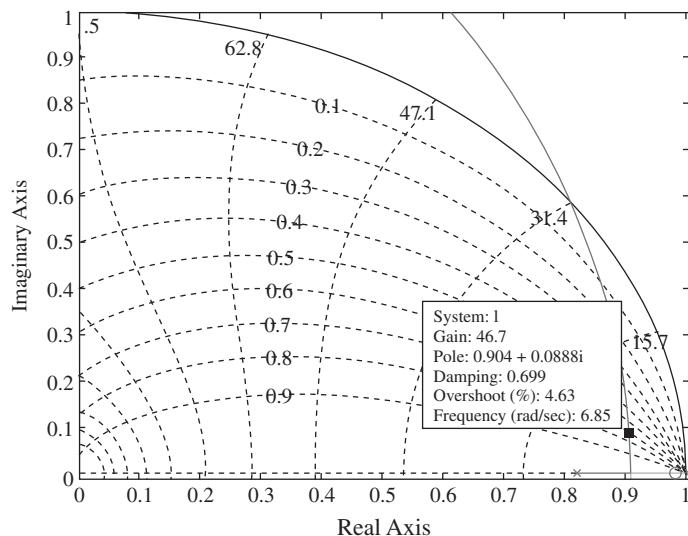


Figure 6.16
Root locus for PI design.

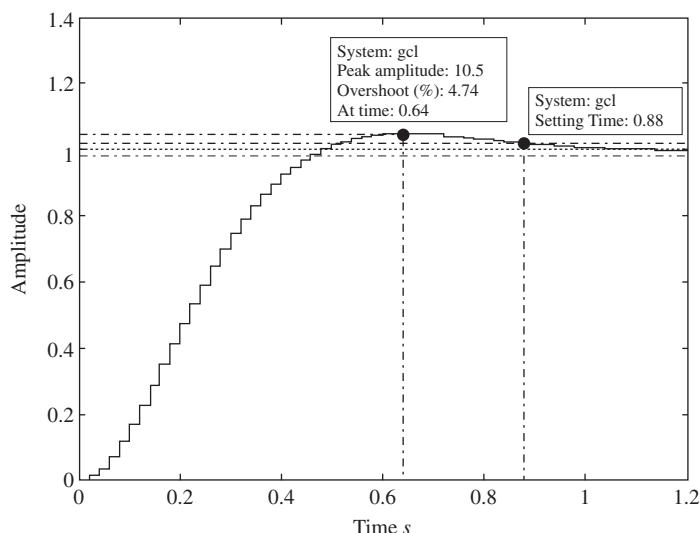


Figure 6.17
Step response for PI design with $K = 47.2$.

Example 6.10

Design a digital controller for the DC motor position control system of Example 3.6, where the (type 1) analog plant has the transfer function

$$G(s) = \frac{1}{s(s+1)(s+10)}$$

to obtain a settling time of about 1 s and percentage overshoot of 5%.

Solution

Using Procedure 6.1, we first observe that an analog PD controller is needed to improve the system transient response. Pole-zero cancellation yields the simple design

$$C_d(s) = K(s+1)$$

As shown in Example 6.9, a damping ratio of 0.7 corresponds to a percentage overshoot of about 5%. We can solve for the undamped natural frequency analytically, or we can use a CAD package to obtain the values $K = 51.02$ and $\omega_n = 7.143$ rad/s for $\zeta = 0.7$.

A sampling period of 0.02 s is appropriate because it is less than $2\pi/(40\omega_d)$. The plant with the ADC and DAC has the z-transfer function

$$G_{ZAS}(z) = (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} = 1.2629 \times 10^{-6} \frac{(z + 0.2534)(z + 3.535)}{(z - 1)(z - 0.8187)(z - 0.9802)}$$

Bilinear transformation of the PD controller gives

$$C(z) = K \frac{z - 0.9802}{z} = K(1 - 0.9802z^{-1})$$

The root locus of the system with PD control (Fig. 6.18) gives a gain K of 2160 and an undamped natural frequency of 6.51 rad/s at a damping ratio $\zeta = 0.7$. The settling time for this design is about

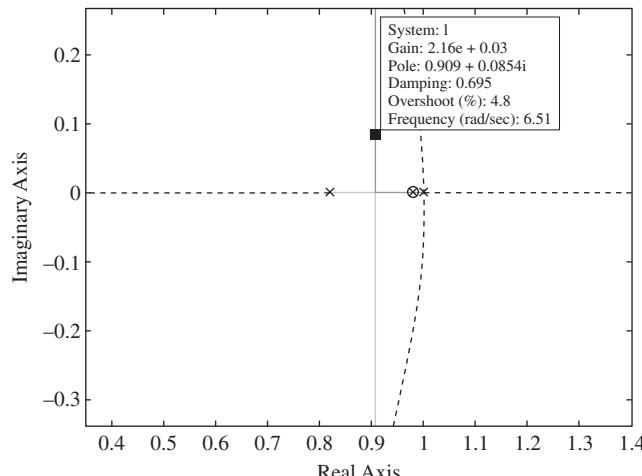


Figure 6.18
Root locus for PD design.

Example 6.10—cont'd

$$T_s = \frac{4}{\zeta \omega_n} = \frac{4}{0.7 \times 6.51} = 0.88 \text{ s}$$

which meets the design specifications.

Checking the step response with MATLAB gives Fig. 6.19 with a settling time of 0.94 s, a peak time of 0.68 s, and a 5% overshoot. The time response shows a slight deterioration from the characteristics of the analog system but meets all the design specifications. In some cases, the deterioration may necessitate repeatedly modifying the digital design or modifying the analog design and then mapping it to the z-domain until the resulting digital filter meets the desired specifications.

Note that for a prewarping frequency $\omega_0 = 1 \text{ rad/s}$, the 3-dB frequency of the PD controller, $\omega_0 T = 0.02 \text{ rad}$ and $\tan(\omega_0 T/2) = \tan(0.01) \approx 0.01$. Hence, Eq. (6.25) is approximately valid without prewarping, and prewarping has a negligible effect on the design.

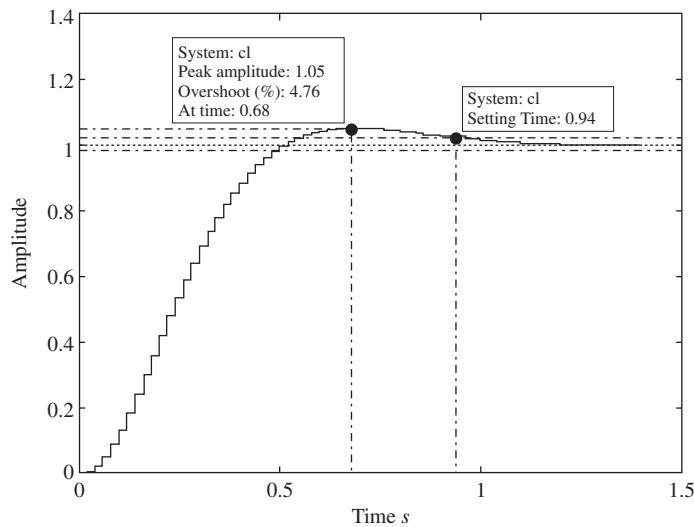


Figure 6.19
Time step response for PD design with $K = 2160$.

Example 6.11

Design a digital controller for a speed control system, where the analog plant has transfer function

$$G(s) = \frac{1}{(s+1)(s+3)}$$

to obtain a time constant of less than 0.3 s, a dominant pole damping ratio of at least 0.7, and zero steady-state error due to a step input.

Example 6.11—cont'd**Solution**

The root locus of the analog system is shown in Fig. 6.20. To obtain zero steady-state error due to a step input, the system type must be increased to one by adding an integrator in the forward path. However, adding an integrator results in significant deterioration of the time response or in instability. If the pole at -1 is canceled, the resulting system is stable but has $\zeta\omega_n = 1.5$ —that is, a time constant of $2/3$ s and not less than 0.3 s as specified. Using a PID controller provides an additional zero that can be used to stabilize the system and satisfy the remaining design requirements.

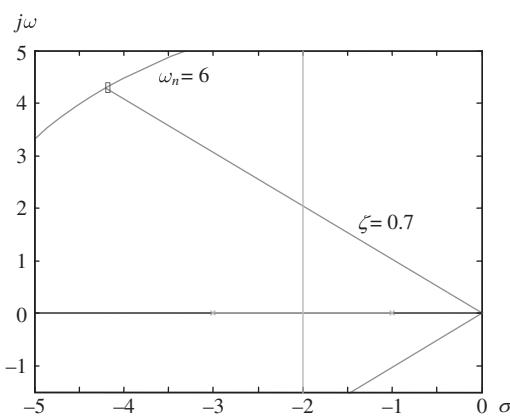


Figure 6.20
Root locus for the analog speed control system.

For time constant τ of 0.3 s, we have $\zeta\omega_n = 1/\tau \geq 3.33$ rad/s. A choice of $\zeta = 0.7$ and ω_n of about 6 rad/s meets the design specifications. The design appears conservative, but we choose a larger undamped natural frequency than the minimum needed in anticipation of the deterioration due to adding PI control. We first design a PD controller to meet these specifications using MATLAB. We obtain the controller angle of about 52.4 degrees using the angle condition. The corresponding zero location is

$$\alpha = \frac{6\sqrt{1 - (0.7)^2}}{\tan(52.4^\circ)} + (0.7)(6) \approx 7.5$$

The root locus for the system with PD control (Fig. 6.21) shows that the system with $\zeta = 0.7$ has ω_n of about 6 rad/s and meets the transient response specifications with a gain of 4.4 and $\zeta\omega_n = 4.2$. Following the PI-design procedure, we place the second zero of the PID controller at one-tenth this distance from the $j\omega$ axis to obtain

$$C_a(s) = K \frac{(s + 0.4)(s + 7.5)}{s}$$

To complete the analog PID design, the gain must be tuned to ensure that $\zeta = 0.7$. Although this step is not needed, we determine the gain $K \approx 5.8$, and $\omega_n = 6.7$ rad/s (Fig. 6.22) for later comparison to the actual gain value used in the digital design. The analog design meets the transient response specification with $\zeta\omega_n = 4.69 > 3.33$, and the dynamics allow us to choose a sampling period of 0.025 s ($\omega_s > 50\omega_d$).

Example 6.11—cont'd

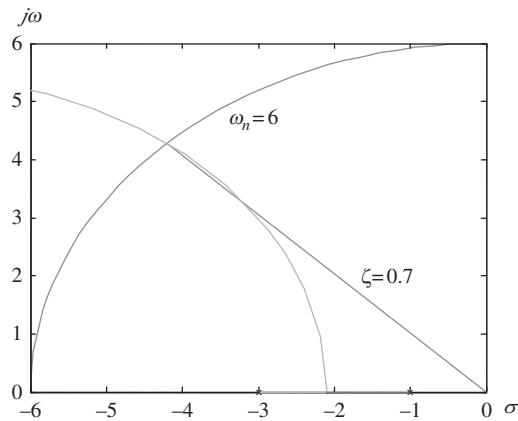


Figure 6.21
Root locus of a PD-controlled system.

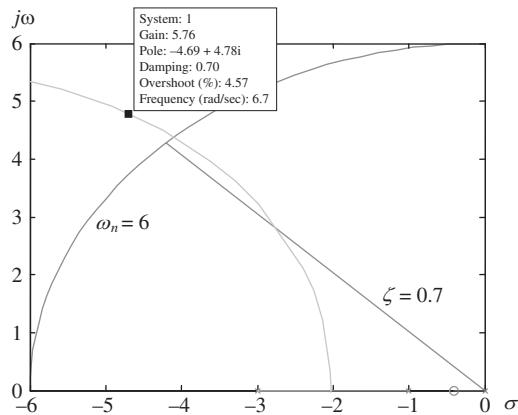


Figure 6.22
Root locus of an analog system with PID control.

The model of the analog plant with DAC and ADC is

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} \\ &= 1.170 \times 10^{-3} \frac{z + 0.936}{(z - 0.861)(z - 0.951)} \end{aligned}$$

Bilinear transformation and elimination of the pole at -1 yields the digital PID controller

$$C(z) = 47.975K \frac{(z - 0.684)(z - 0.980)}{z(z - 1)}$$

The root locus for the system with digital PID control is shown in Fig. 6.23, and the system is seen to be **minimum phase** (i.e., its zeros are inside the unit circle).

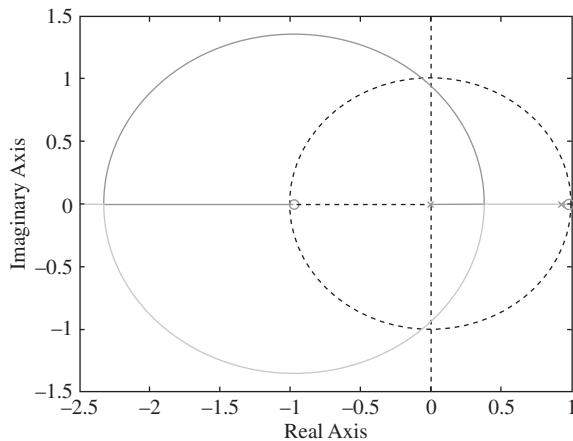
Example 6.11—cont'd

Figure 6.23
Root locus of a system with digital PID control.

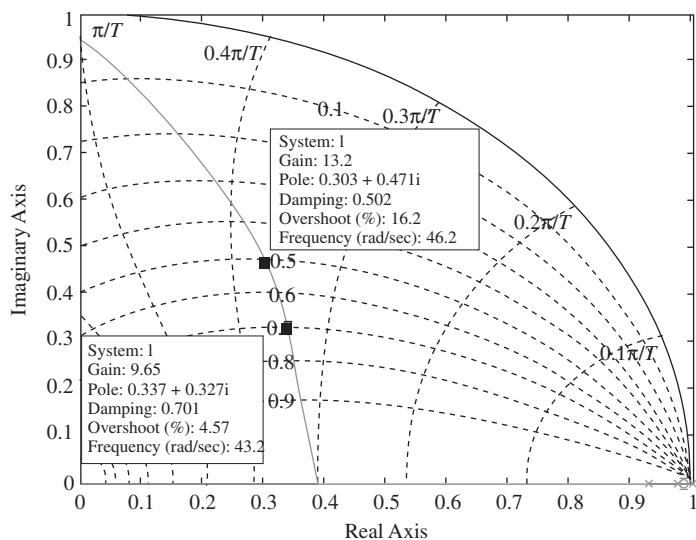
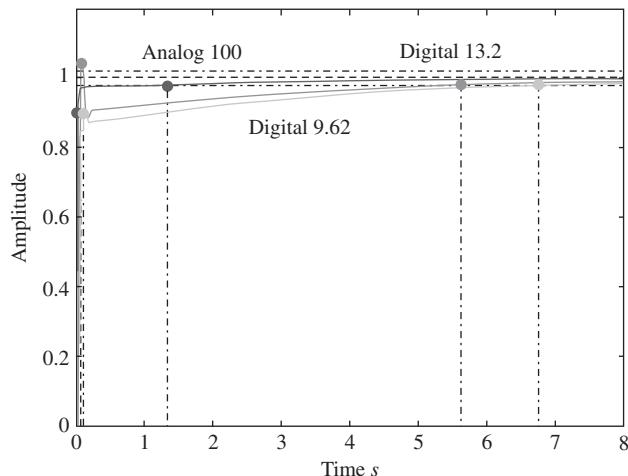


Figure 6.24
Detail of the root locus of a system with digital PID control.

For design purposes, we zoom in on the most significant portion of the root locus and obtain the plot of Fig. 6.24. With $K = 9.62, 13.2$, the system has $\zeta = 0.7, 0.5$, and $\omega_n = 43.2, 46.2$ rad/s, respectively. Both designs have a sufficiently fast time constant, but the second damping ratio is less than the specified value of 0.7. The time response of the two digital systems and for analog control with $K = 100$ are shown in Fig. 6.25. Lower gains give an unacceptably slow analog design. The time response for the high-gain digital design is very fast. However, it has an overshoot of over 4% but has a settling time of 5.63 s. The digital design for $\zeta = 0.7$ has a much slower time response than its analog counterpart.

Example 6.11—cont'd**Figure 6.25**

Time step response for the digital PID design with $K = 9.62$ (light gray), $K = 13.2$ (dark gray), and for analog design (black).

It is possible to improve the design by trial and error, including redesign of the analog controller, but the design with $\zeta = 0.5$ may be acceptable. One must weigh the cost of redesign against that of relaxing the design specifications for the particular application at hand. The final design must be a compromise between speed of response and relative stability.

6.3.4 Empirical digital PID controller tuning

As explained in Section 5.5, the parameters of a PID controller are often selected by means of tuning rules. This concept can also be exploited to design a digital PID controller. The reader can show (Problem 6.8) that bilinear transformation of the PID controller expression (5.20) yields

$$C(z) = K_p \left(1 + \frac{1}{T_i} \frac{T}{2} \frac{z+1}{z-1} + T_d \frac{2}{T} \frac{z-1}{z+1} \right) \quad (6.33)$$

Bilinear transformation of the PID controller results in a pole at $z = -1$ because the derivative part is not proper (see Section 12.4.1). As in Section 6.3.3, we avoid an unbounded frequency response at the folding frequency by replacing the pole at $z = -1$ with a pole at $z = 0$ and dividing the gain by two. The resulting transfer function is

$$C(z) = \frac{K_p}{2} \left(1 + \frac{T}{2T_i} \frac{z+1}{z-1} + \frac{2T_d}{T} \frac{z-1}{z} \right)$$

If parameters K_p , T_i , and T_d are obtained by means of a tuning rule as in the analog case, then the expression of the digital controller is obtained by substituting in the previous expression. The transfer function of a zero-order hold can be approximated by truncating the series expansions as

$$\frac{G_{ZOH}(s)}{T} = \frac{1 - e^{-sT}}{Ts} \cong \frac{1 - 1 + Ts - (Ts)^2/2 + \dots}{Ts} = 1 - \frac{Ts}{2} + \dots \cong e^{-\frac{T}{2}s}$$

Thus, the presence of the ZOH can be considered as an additional time delay equal to half of the sampling period. The tuning rules of Table 5.1 can then be applied to a system with a delay equal to the sum of the process time delay and a delay of $T/2$ due to the zero-order hold.

Example 6.12

Design a digital PID controller with sampling period $T = 0.1$ for the analog plant of Example 5.9

$$G(s) = \frac{1}{(s+1)^4} e^{-0.2s}$$

by applying the Ziegler–Nichols tuning rules of Table 5.1.

Solution

A first-order-plus-dead-time model of the plant was obtained in Example 5.9 using the tangent method with gain $K = 1$, a dominant time constant $\tau = 3$, and an apparent time delay $L = 1.55$. The apparent time delay for digital control is obtained by adding half of the value of the sampling period (0.05). This gives $L = 1.55 + 0.05 = 1.6$. The application of the tuning rules of Table 5.1 yields

$$K_p = 1.2 \frac{\tau}{KL} = 2.25$$

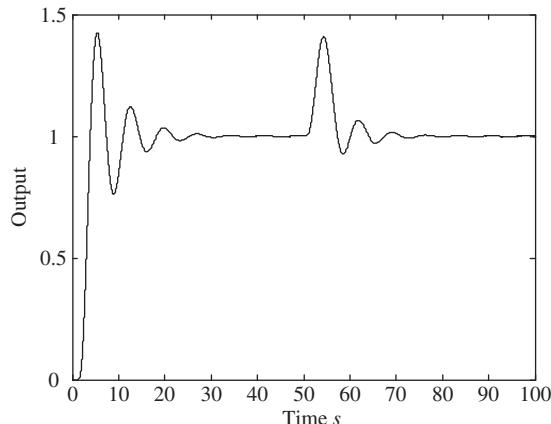
$$T_i = 2L = 3.2$$

$$T_d = 0.5L = 0.8$$

Thus, the digital PID controller has the transfer function

$$C(z) = \frac{19.145z^2 - 35.965z + 16.895}{z(z-1)}$$

The response of the digital control system due to a unit step reference input applied at time $t = 0$ and to a unit step change in the control variable at time $t = 50$ is shown in Fig. 6.26. The response is similar to the result obtained with the analog PID controller in Example 5.9.

Example 6.12—cont'd**Figure 6.26**

Process output with the digital PID controller tuned with the Ziegler-Nichols method.

6.4 Direct z-domain digital controller design

Obtaining digital controllers from analog designs involves approximation that may result in significant controller distortion. In addition, the locations of the controller poles and zeros are often restricted to subsets of the unit circle. For example, bilinear transformation of the term $(s+a)$ gives $[z-(c-a)/(c+a)]$, as seen from Eqs. (6.29), (6.31), (6.32). This yields only RHP zeros because a is almost always smaller than c . The plant poles are governed by $p_z = e^{p_s T}$, where p_s and p_z are the s -domain and z -domain poles, respectively, and can be canceled with RHP zeros. Nevertheless, the restrictions on the poles and zeros in Eqs. (6.29), (6.31), (6.32) limit the designer's ability to reshape the system root locus.

Another complication in digital approximation of analog filters is the need to have a pole at 0 in place of the pole at -1 , as obtained by direct digital transformation, to avoid an unbounded frequency response at the folding frequency. This may result in a significant difference between the digital and analog controllers and may complicate the design process considerably.

Alternatively, it is possible to directly design controllers in the less familiar z -plane. The controllers used are typically of the same form as those discussed in Section 6.3, but the poles of the controllers are no longer restricted as in s -domain-to- z -domain mapping. Thus, the poles can now be in either the LHP or the RHP as needed for design.

Because of the similarity of s -domain and z -domain root loci, Procedures 5.1, 5.2, and 5.3 are applicable with minor changes in the z -domain. However, Equation (5.14) is no longer

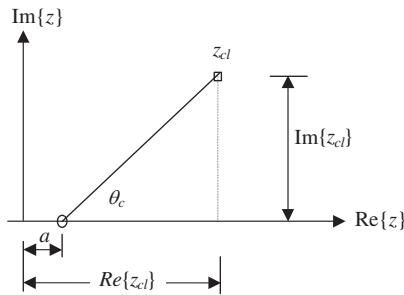


Figure 6.27
PD compensator zero.

valid because the real and imaginary components of the complex conjugate poles are different in the z -domain. In addition, a digital controller with a zero and no poles is not realizable, and a pole must be added to the controller. To minimize the effect of the pole on the time response, it is placed at the origin. This is analogous to placing an s -plane pole far in the LHP to minimize its effect. We must now derive an expression similar to (5.14) for the digital PID controller.

Using the expression for the complex conjugate z -domain poles Eq. (6.4), we obtain (Fig. 6.27)

$$\begin{aligned} a &= \text{Re}\{z_{cl}\} - \frac{\text{Im}\{z_{cl}\}}{\tan(\theta_a)} \\ &= e^{-\zeta\omega_n T} \cos(\omega_d T) - \frac{e^{-\zeta\omega_n T} \sin(\omega_d T)}{\tan(\theta_a)} \end{aligned} \quad (6.34)$$

where θ_a is the angle of the controller zero. θ_a is given by

$$\begin{aligned} \theta_a &= \theta_c + \theta_p \\ &= \theta_c + \theta_{zcl} \end{aligned} \quad (6.35)$$

In Eq. (6.35), θ_{zcl} is the angle of the controller pole and θ_c is the controller angle contribution at the closed-loop pole location. The sign of the second term in Eq. (6.34) is negative, unlike (5.14), because the real part of a stable z -domain pole can be positive. In addition, a digital controller with a zero and no poles is not realizable and a pole, at the origin or other locations inside the unit circle, must first be added to the system before computing the controller angle. The computation of the zero location is simple using the following MATLAB function:

```
% Digital PD controller design with pole at origin and zero to be selected.
function[c,zcl] = dpdcon(zeta,wn,g,T)
% g (L) is the uncompensated (compensated) loop gain.
% zeta and wn specify the desired closed-loop pole, T = sampling period.
% zcl is the closed-loop pole, theta is the angle of the.
```

```
% compensator zero at zcl.
% The corresponding gain is "k" and the zero is "a".
wdT = wn*T*sqrt(1-zeta^2); % Pole angle: T*damped natural frequency.
rzcl = exp(-zeta*wn*T)*cos(wdT); % Real part of the closed-loop pole.
izcl = exp(-zeta*wn*T)*sin(wdT); % Imaginary part of the closed-loop pole.
zcl = rzcl + j*izcl; % Complex closed-loop pole.
% Find the angle of the compensator zero. Include the contribution
% of the pole at the origin.
theta = pi - angle(evalfr(g,zcl)) + angle(zcl);
a = rzcl-izcl/tan(theta); % Calculate the zero location.
c = zpk([a],[0],1,T); % Calculate the compensator transfer function.
L = c*g; % Loop gain.
k = 1/abs(evalfr(l, zcl)); % Calculate the gain.
C = k*c; %Include the correct gain.
```

Although Eq. (6.34) is useful in some situations, in many others, PD or PID controllers can be more conveniently obtained by simply canceling the slow poles of the system with zeros. Design by pole-zero cancellation is the simplest possible and should be explored before more complex designs are attempted.

The main source of difficulty in z -domain design is the fact that the stable region is now the unit circle as opposed to the much larger left half of the s -plane. In addition, the selection of pole locations in the z -domain directly is less intuitive and is generally more difficult than s -domain pole selection. Pole selection and the entire design process are significantly simplified by the use of CAD tools and the availability of constant ζ contours, constant ω_n contours, and cursor commands.

No new theory is needed to introduce z -domain design, and we proceed directly to design examples. We repeat Examples 6.9, 6.10, and 6.11 using direct digital design to demonstrate its strengths and weaknesses compared to the indirect design approach.

Example 6.13

Design a digital controller for a DC motor speed control system where the (type 0) analog plant has the transfer function

$$G(s) = \frac{1}{(s+1)(s+10)}$$

to obtain zero steady-state error due to a unit step, percentage overshoot less than 5%, and a settling time of about 1 s.

Solution

First, selecting $T = 0.1$ s, we obtain the z -transfer function

$$G_{ZAS}(z) = (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} = 3.55 \times 10^{-3} \frac{z + 0.694}{(z - 0.368)(z - 0.905)}$$

Example 6.13—cont'd

To perfectly track a step input, the system type must be increased by one and a PI controller is therefore used. The controller has a pole at $z = 1$ and a zero to be selected to meet the remaining design specifications. Directly canceling the pole at $z = 0.905$ gives a design that is almost identical to that of Example 6.9 and meets the design specifications.

Example 6.14

Design a digital controller for the DC motor position control system of Example 3.6, where the (type 1) analog plant has the transfer function

$$G(s) = \frac{1}{s(s+1)(s+10)}$$

for a settling time of less than 1 s and percentage overshoot of about 5%.

Solution

For a sampling period of 0.01 s, the plant, ADC, and DAC have the z-transfer function

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} \\ &= 1.6217 \times 10^{-7} \frac{(z + 0.2606)(z + 3.632)}{(z - 1)(z - 0.9048)(z - 0.99)} \end{aligned}$$

Using a digital PD controller improves the system transient response as in Example 6.10. Pole-zero cancellation yields the simple design

$$C(z) = K \frac{z - 0.99}{z}$$

which includes a pole at $z = 0$ to make the controller realizable. The design is almost identical to that of Example 6.10 and meets the desired transient response specifications with a gain of 4,580.

Example 6.15

Design a digital controller for the DC motor speed control system where the analog plant has the transfer function

$$G(s) = \frac{1}{(s+1)(s+3)}$$

for time constant of less than 0.3 s, percentage overshoot less than 5%, and zero steady-state error due to a step input.

Example 6.15—cont'd**Solution**

The plant is type 0 and is the same as in Example 6.10 with a sampling period $T = 0.005$ s:

$$G_{ZAS}(z) = (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} = 1.2417 \times 10^{-5} \frac{z + 0.9934}{(z - 0.9851)(z - 0.995)}$$

For zero steady-state error due to step, the system type must be increased by one by adding a pole at $z = 1$. A controller zero can be used to cancel the pole at $z = 0.995$, leaving the loop gain

$$L(z) = 1.2417 \times 10^{-5} \frac{z + 0.9934}{(z - 1)(z - 0.9851)}$$

The system root locus of Fig. 6.28 shows that the closed-loop poles are close to the unit circle at low gains, and the system is unstable at higher gains. Clearly, an additional zero is needed to meet the design specification, and a PID controller is required. To make the controller realizable, a pole must be added at $z = 0$. The simplest design is then to cancel the pole closest to but not on the unit circle, giving the loop gain

$$L(z) = 1.2417 \times 10^{-5} \frac{z + 0.9934}{z(z - 1)}$$

Adding the zero to the transfer function, we obtain the root locus of Fig. 6.29 and select a gain of 20,200. The corresponding time response is shown in Fig. 6.30. The time response shows less than 5% overshoot with a fast time response that meets all design specifications. The design is better than that of Example 6.11, where the digital controller was obtained via analog design.

Although it may be possible to improve the analog design to obtain better results than those of Example 6.11, this requires trial and error as well as considerable experience. By contrast, the digital design is obtained here directly in the z -domain without the need for trial and error. This demonstrates that direct design in the z -domain using CAD tools can be easier than indirect design.

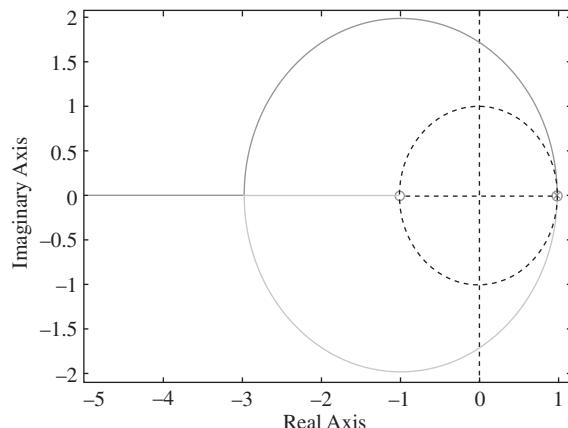


Figure 6.28
Root locus for digital PI control.

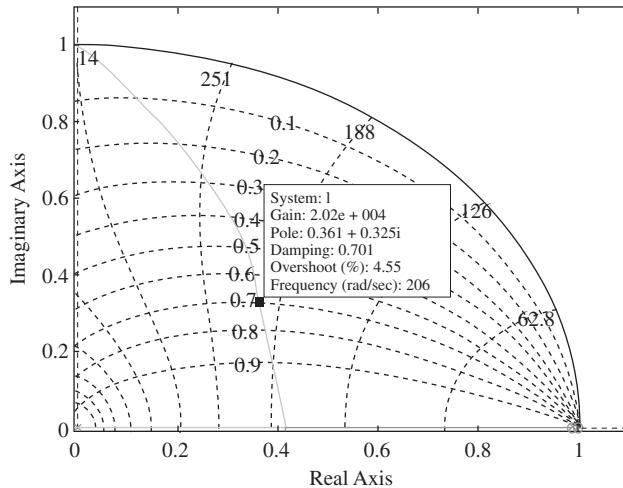
Example 6.15—cont'd

Figure 6.29
Root locus for digital PID control design by pole-zero cancellation.

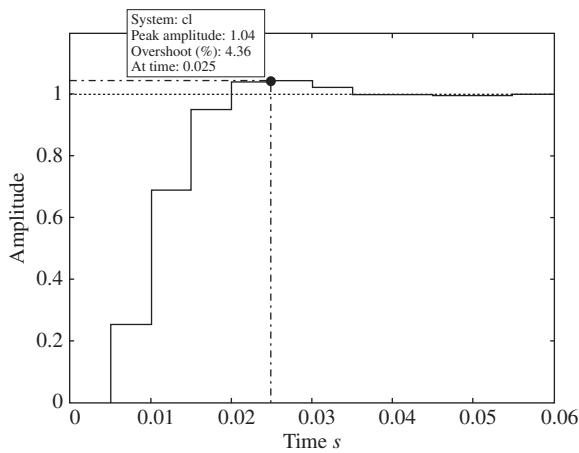


Figure 6.30
Time step response for digital PID control.

6.5 Frequency response design

Frequency response design approaches, especially design based on Bode plots, are very popular in the design of analog control systems. They exploit the fact that the Bode plots of rational loop gains (i.e., ones with no delays) can be approximated with straight lines. Further, if the transfer function is minimum phase, the phase can be determined from the

plot of the magnitude, and this allows us to simplify the design of a compensator that provides specified stability properties. These specifications are typically given in terms of the phase margin and the gain margin. As with other design approaches, the design is greatly simplified by the availability of CAD packages.

Unfortunately, discrete transfer functions are not rational functions in $j\omega$ because the frequency is introduced through the substitution $z = e^{j\omega T}$ in the z -transfer function (see Section 2.8). Hence, the simplification provided by methods based on Bode plots for analog systems is lost. A solution to this problem is to bilinearly transform the z -plane into a new plane, called the w -plane, where the corresponding transfer function is rational and where the Bode approximation is valid. For this purpose, we recall the bilinear transformation

$$w = c \frac{z - 1}{z + 1} \quad \text{where} \quad c = \frac{2}{T} \quad (6.36)$$

from Section 6.3.3, which maps points in the LHP into points inside the unit circle. To transform the inside of the unit circle to the LHP, we use the inverse bilinear transformation

$$z = \frac{1 + \frac{wT}{2}}{1 - \frac{wT}{2}} \quad (6.37)$$

This transforms the transfer function of a system from the z -plane to the w -plane. The w -plane is a complex plane whose imaginary part is denoted by v . To express the relationship between the frequency ω in the s -plane and the frequency v in the w -plane, we let $s = j\omega$ and therefore $z = e^{j\omega T}$. Substituting in Eq. (6.36) and using steps similar to those used to derive Eq. (6.25), we have

$$w = jv = j \frac{2}{T} \tan \frac{\omega T}{2} \quad (6.38)$$

From Eq. (6.38), as ω varies in the interval $[0, \pi/T]$, z moves on the unit circle and v goes from 0 to infinity. This implies that there is a distortion or warping of the frequency scale between v and ω (see Section 6.3.3). This distortion is significant, especially at high frequencies. However, if $\omega \ll \omega_s/2 = \pi/T$, we have from Eq. (6.38) that $\omega \approx v$ and the distortion is negligible.

In addition to the problem of frequency distortion, the transformed transfer function $G(w)$ has two characteristics that can complicate the design: (1) the transfer function will always have a pole-zero deficit of zero (i.e., the same number of poles as zeros), and (2) the bilinear transformation Eq. (6.36) can introduce RHP zeros and result in a nonminimum phase system.

Nonminimum phase systems limit the achievable performance. For example, because some root locus branches start at the open-loop pole and end at zeros, the presence of RHP zeros limits the stable range of gains K . Thus, attempting to reduce the steady-state error

or to speed up the response by increasing the system gain can lead to instability. These limitations clearly make w -plane design challenging. Examples 6.16, 6.17, and 6.18 illustrate w -plane design and how to overcome its limitations.

We summarize the steps for controller design in the w -plane in the following procedure.

Procedure 6.2

1. Select a sampling period and obtain the transfer function $G_{ZAS}(z)$ of the discretized process.
2. Transform $G_{ZAS}(z)$ into $G(w)$ using Eq. (6.37)
3. Draw the Bode plot of $G(j\nu)$, and use analog frequency response methods to design a controller $C(w)$ that satisfies the frequency domain specifications.
4. Transform the controller back into the z -plane by means of Eq. (6.36), thus determining $C(z)$.
5. Verify that the performance obtained is satisfactory.

In controller design problems, the specifications are often given in terms of the step response of the closed-loop system, such as the settling time and the percentage overshoot. As shown in Chapter 5, the percentage overshoot specification yields the damping ratio ζ , which is then used with the settling time to obtain the undamped natural frequency ω_n . Thus, we need to obtain frequency response specifications from ζ and ω_n to follow Procedure 6.2.

The relationship between the time domain criteria and the frequency domain criteria is in general quite complicated. However, the relationship is much simpler if the closed-loop system can be approximated by the second-order underdamped transfer function

$$T(s) = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2} \quad (6.39)$$

The corresponding loop gain with unity feedback is given by

$$L(s) = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s} \quad (6.40)$$

We substitute $s = j\omega$ to obtain the corresponding frequency response and equate the square of its magnitude to unity:

$$|L(j\omega)|^2 = \frac{1}{(\omega/\omega_n)^4 + 4\zeta^2(\omega/\omega_n)^2} = 1 \quad (6.41)$$

The magnitude of the loop gain, as well as its square, is unity at the gain crossover frequency. For the second-order underdamped case, we now have the relation

$$\omega_{gc} = \omega_n \left[\sqrt{4\zeta^4 + 1} - 2\zeta^2 \right]^{1/2} \quad (6.42)$$

Next, we consider the phase margin and derive

$$PM = 180^\circ + \angle G(j\omega_{gc}) = \tan^{-1} \left(\frac{2\zeta}{\left[\sqrt{4\zeta^4 + 1} - 2\zeta^2 \right]^{1/2}} \right)$$

The last expression can be approximated by

$$PM \approx 100\zeta \quad (6.43)$$

Eqs. (6.42) and (6.43) provide the transformations we need to obtain frequency domain specifications from step response specifications. Together with Procedure 6.2, the equations allow us to design digital control systems using the w -plane approach. The following three examples illustrate w -plane design using Procedure 6.2.

Example 6.16

Consider the cruise control system of Example 3.2, where the analog process is

$$G(s) = \frac{1}{s+1}$$

Transform the corresponding $G_{ZAS}(z)$ to the w -plane. By considering both $T = 0.1$ and $T = 0.01$, evaluate the role of the sampling period by analyzing the corresponding Bode plots.

Solution

When $T = 0.1$ we have

$$G_{ZAS}(z) = \frac{0.09516}{z - 0.9048}$$

and by applying Eq. (6.37), we obtain

$$G_1(w) = \frac{-0.05w + 1}{w + 1}$$

The transfer function $G_1(w)$ can be obtained using the MATLAB command

```
>> Gw = d2c(Gzas, 'tustin')
```

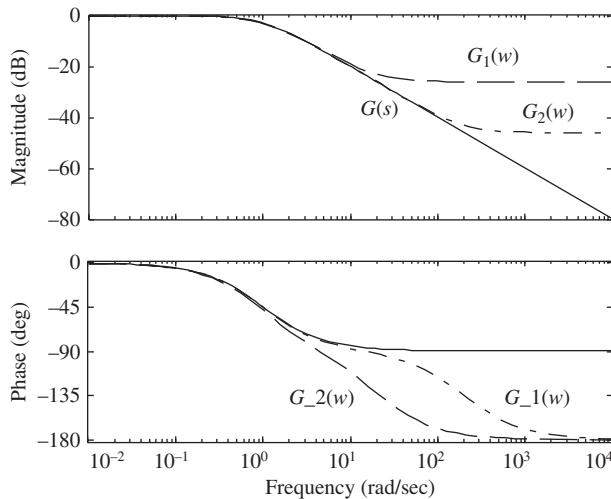
When $T = 0.01$ we have

$$G_{ZAS}(z) = \frac{0.00995}{z - 0.99}$$

and, again by applying Eq. (6.37), we obtain

$$G_2(w) = \frac{-0.005w + 1}{w + 1}$$

The Bode plots of $G(s)$, $G_1(w)$, and $G_2(w)$ are shown in Fig. 6.31. For both sampling periods, the pole in the w -plane is in the same position as the pole in the s -plane. However, both $G_1(w)$ and $G_2(w)$ have a zero, whereas $G(s)$ does not. This results in a big difference between the frequency response of the analog system and that of the digital systems at high frequencies. However, the influence of the zero on the system dynamics is clearly less significant when the sampling period is smaller. Note that for both sampling periods, distortion in the

Example 6.16—cont'd**Figure 6.31**

Bode plots for Example 6.16.

low-frequency range is negligible. For both systems, the gain as w goes to zero is unity, as is the DC gain of the analog system. This is true for the choice $c = 2/T$ in Eq. (6.36). Other possible choices are not considered because they do not yield DC gain equality.

Example 6.17

Consider a DC motor speed control system where the (type 0) analog plant has the transfer function

$$G(s) = \frac{1}{(s+1)(s+10)}$$

Design a digital controller by using frequency response methods to obtain zero steady-state error due to a unit step, an overshoot less than 10%, and a settling time of about 1 s.

Solution

From the given specification, we have that the controller in the w -plane must contain a pole at the origin. For 10% overshoot, we calculate the damping ratio as

$$\zeta = \frac{|\ln(0.1)|}{\sqrt{|\ln(0.1)|^2 + \pi^2}} \approx 0.6$$

Using the approximate expression Eq. (6.43), the phase margin is about 100 times the damping ratio of the closed-loop system, and the required phase margin is about 60 degrees. For a settling time of 1 s, we calculate the undamped natural frequency

$$\omega_n = \frac{4}{\zeta T_s} \approx 6.7 \text{ rad/s}$$

Example 6.17—cont'd

Using Eq. (6.42), we obtain the gain crossover $\omega_{gc} = 4.8$ rad/s.

A suitable sampling period for the selected dynamics is $T = 0.02$ s (see Example 6.9). The discretized process is then determined as

$$G_{ZAS}(z) = (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} = 1.8604 \times 10^{-4} \frac{z + 0.9293}{(z - 0.8187)(z - 0.9802)}$$

Using Eq. (6.37), we obtain the w -plane transfer function

$$G(w) = \frac{-3.6519 \cdot 10^{-6} (w + 2729)(w - 100)}{(w + 9.967)(w + 1)}$$

Note the two additional zeros (with respect to $G(s)$) do not significantly influence the system dynamics in the range of frequencies of interest for the design. The two poles are virtually in the same position as the poles of $G(s)$. The simplest design that meets the desired specifications is to insert a pole at the origin, to cancel the dominant pole at -1 , and to increase the gain until the required gain crossover frequency is attained. Thus, the resulting controller transfer function is

$$C(w) = 54 \frac{w + 1}{w}$$

The Bode plot of the loop transfer function $C(w)G(w)$, together with the Bode plot of $G(w)$, is shown in Fig. 6.32. The figure also shows the phase and gain margins. By transforming the controller back to the z -plane using Eq. (6.36), we obtain

$$C(z) = \frac{54.54z - 53.46}{z - 1}$$

$GM = 25.7$ dB (at 32.2 rad/sec),
 $PM = 61.3$ deg (at 4.86 rad/sec)

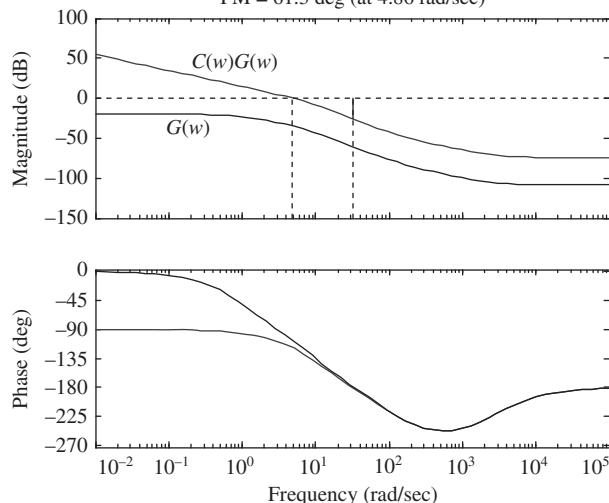


Figure 6.32

Bode plots of $C(w)G(w)$ and $G(w)$ of Example 6.17.

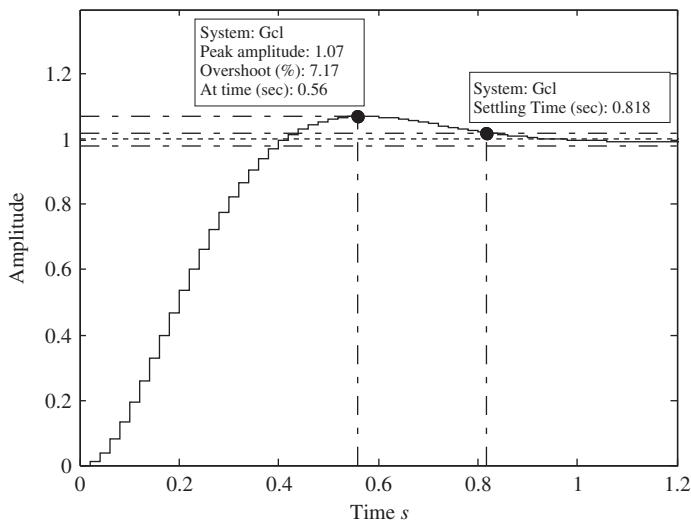
Example 6.17—cont'd

Figure 6.33
Closed-loop step response for Example 6.17.

The transfer function $C(z)$ can be obtained using the MATLAB command

`>> Cz = c2d(Cw, 0.02, 'matched')`

The corresponding discretized closed-loop step response is plotted in Fig. 6.33 and clearly meets the design specifications.

Example 6.18

Consider the DC motor speed control system with transfer function

$$G(s) = \frac{1}{(s+1)(s+3)}$$

Design a digital controller using frequency response methods to obtain zero steady-state error due to a unit step, an overshoot less than 10%, and a settling time of about 1 s. Use a sampling period $T = 0.2$ s.

Solution

As in Example 6.17, (1) the given steady-state error specification requires a controller with a pole at the origin, (2) the percentage overshoot specification requires a phase margin of about 60 degrees, and (3) the settling time yields a gain crossover frequency of about 5 rad/s. The transfer function of the system with DAC and ADC is

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} \\ &= 0.015437 \frac{z + 0.7661}{(z - 0.8187)(z - 0.5488)} \end{aligned}$$

Example 6.18—cont'd

Transforming to the w -plane using Eq. (6.37), we obtain

$$G(w) = \frac{-12.819 \cdot 10^{-4}(w + 75.5)(w - 10)}{(w + 2.913)(w + 0.9967)}$$

Note again that the two poles are virtually in the same locations as the poles of $G(s)$. Here the RHP zero at $w = 10$ must be taken into account in the design because the required gain crossover frequency is about 5 rad/s. To achieve the required phase margin, both poles must be canceled with two controller zeros. For a realizable controller, we need at least two controller poles. In addition to the pole at the origin, we select a high-frequency controller pole so as not to impact the frequency response in the vicinity of the crossover frequency. Next, we adjust the system gain to achieve the desired specifications.

As a first attempt, we select a high-frequency pole in $w = -20$ and increase the gain to 78 for a gain crossover frequency of 4 rad/s with a phase margin of 60° . The controller is therefore

$$C(w) = 78 \frac{(w + 2.913)(w + 0.9967)}{w(w + 20)}$$

The corresponding open-loop Bode plot is shown in Fig. 6.34 together with the Bode plot of $G(w)$. Transforming the controller back to the z -plane using Eq. (6.36), we obtain

$$C(z) = C(w) \Big|_{w=\frac{2}{T} \left[\frac{z-1}{z+1} \right]} = \frac{36.9201z^2 - 50.4902z + 16.5897}{(z-1)(z-0.3333)}$$

The resulting discretized closed-loop step response, which satisfies the given requirements, is plotted in Fig. 6.35.

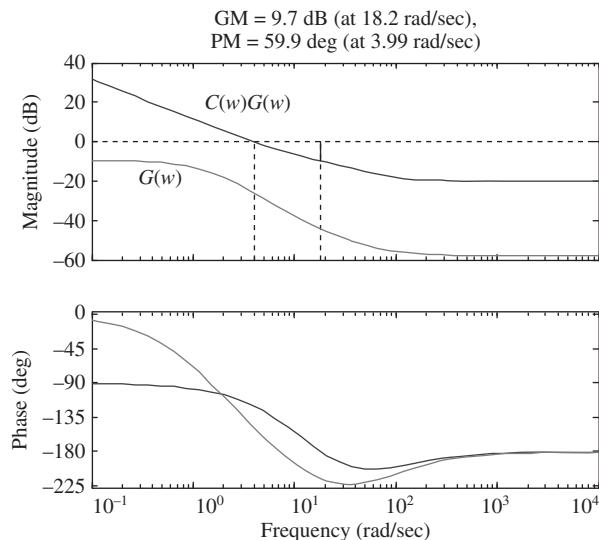


Figure 6.34
Bode plots of $C(w)G(w)$ and $G(w)$ for Example 6.18.

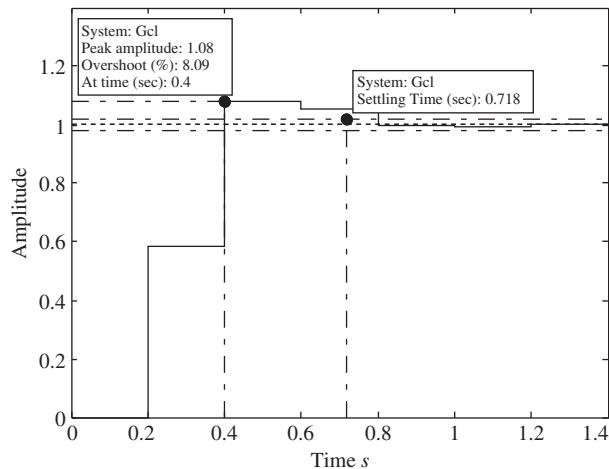
Example 6.18—cont'd

Figure 6.35
Closed-loop step response for Example 6.18.

6.6 Direct control design

In some control applications, the desired transfer function of the closed-loop system is known from the design specification. For a particular system configuration, it is possible to calculate the controller transfer function for a given plant from the desired closed-loop transfer function. This approach to design is known as **synthesis**. Clearly, the resulting controller must be realizable for this approach to yield a useful design.

We consider the block diagram of Fig. 6.12 with known closed-loop transfer function $G_{cl}(z)$. The controller transfer function $C(z)$ can be computed analytically, starting from the expression of the desired closed-loop transfer function $G_{cl}(z)$:

$$G_{cl}(z) = \frac{C(z)G_{ZAS}(z)}{1 + C(z)G_{ZAS}(z)}$$

We solve for the controller transfer function

$$C(z) = \frac{1}{G_{ZAS}(z)} \frac{G_{cl}(z)}{1 - G_{cl}(z)} \quad (6.44)$$

For the control system to be implementable, the controller must be causal and must ensure the asymptotic stability of the closed-loop control system.

For a causal controller, the closed-loop transfer function $G_{cl}(z)$ must have at least the same pole-zero deficit as $G_{ZAS}(z)$. In other words, the delay in $G_{cl}(z)$ must be at least as long as

the delay in $G_{ZAS}(z)$. We see this by examining the degree of the numerator and the degree of the denominator in Eq. (6.44). We first write the transfer functions in the form

$$G_{ZAS}(z) = \frac{N_{ZAS}(z)}{D_{ZAS}(z)} = \frac{N_{ol}(z)}{D_{ZAS}(z)} z^{-l_z}$$

$$G_{cl}(z) = \frac{N_{cl}(z)}{D_{cl}(z)} = \frac{N_{cl1}(z)}{D_{cl}(z)} z^{-l_c}$$

Substituting in Eq. (6.44) gives

$$C(z) = \frac{D_{ZAS}(z)}{N_{ZAS}(z)} \left[\frac{N_{cl}(z)}{D_{cl}(z) - N_{cl}(z)} \right] = \frac{D_{ZAS}(z)}{N_{ol}(z)z^{-l_z}} \left[\frac{N_{cl1}(z)z^{-l_c}}{D_{cl}(z) - N_{cl1}(z)z^{-l_c}} \right]$$

For the controller to be causal, the closed-loop transfer function must have a delay $l_c \geq l_z$, and pole-zero deficit equal or larger than that of $G_{ZAS}(z)$.

Recall from Chapter 4 that if unstable pole-zero cancellation occurs, the system is input-output stable but not asymptotically stable. This is because the response due to the initial conditions is unaffected by the zeros and is affected by the unstable poles, even if they cancel with a zero. Hence, one must be careful when designing closed-loop control systems to avoid unstable pole-zero cancellations. This implies that the set of zeros of $G_{cl}(z)$ must include all the zeros of $G_{ZAS}(z)$ that are outside the unit circle. Suppose that the process has unstable pole $z = \bar{z}, |\bar{z}| > 1$ —namely

$$G_{ZAS}(z) = \frac{G_1(z)}{z - \bar{z}}$$

In view of Eq. (6.44), we avoid unstable pole-zero cancellation by requiring

$$1 - G_{cl}(z) = \frac{1}{1 + C(z) \frac{G_1(z)}{z - \bar{z}}} = \frac{z - \bar{z}}{z - \bar{z} + C(z)G_1(z)}$$

In other words, $z = \bar{z}$ must be a zero of $1 - G_{cl}(z)$.

An additional condition can be imposed to address steady-state accuracy requirements. In particular, if zero steady-state error due to a step input is required, the condition must be

$$G_{cl}(1) = 1$$

This condition is easily obtained by applying the Final Value theorem (see Section 2.3.4). The proof is left as an exercise for the reader.

Summarizing, the conditions required for the choice of $G_{cl}(z)$ are as follows:

- $G_{cl}(z)$ must have at least the same pole-zero deficit as $G_{ZAS}(z)$ (causality).
- $G_{cl}(z)$ must contain as zeros all the zeros of $G_{ZAS}(z)$ that are outside the unit circle (stability).

- The zeros of $1 - G_{cl}(z)$ must include all the poles of $G_{ZAS}(z)$ that are outside the unit circle (stability).
- $G_{cl}(1) = 1$ (zero steady-state error).

The choice of a suitable closed-loop transfer function is clearly the main obstacle in the application of the direct design method. The correct choice of closed-loop poles and zeros to meet the design requirements can be difficult if there are additional constraints on the control variable because of actuator limitations. Further, the performance of the control system relies heavily on an accurate process model.

To partially address the first problem, instead of directly selecting the poles and the zeros of the (first- or second-order) closed-loop system directly in the z -plane, a continuous-time closed-loop system can be specified based on desired properties of the closed-loop system, and then $G_{cl}(z)$ is obtained by applying pole-zero matching. The following ad hoc procedure often yields the desired transfer function.

Procedure 6.3

1. Select the desired settling time T_s and the desired maximum overshoot.
2. Select a suitable continuous-time closed-loop first-order or second-order closed-loop system with unit gain.
3. Obtain $G_{cl}(z)$ by converting the s -plane pole location to the z -plane pole location using pole-zero matching.
4. Verify that $G_{cl}(z)$ meets the conditions for causality, stability, and steady-state error. If not, modify $G_{cl}(z)$ until the conditions are met.

Example 6.19

Design a digital controller for a DC motor speed control system where the (type 0) analog plant has the transfer function

$$G(s) = \frac{1}{(s+1)(s+10)}$$

to obtain zero steady-state error due to a unit step and a settling time of about 4 s. The sampling period is chosen as $T = 0.02$ s.

Solution

As in Example 6.9, the discretized process transfer function is

$$G_{ZAS}(z) = (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} = 1.8604 \times 10^{-4} \frac{z + 0.9293}{(z - 0.8187)(z - 0.9802)}$$

Note that there are no poles or zeros outside the unit circle and the pole-zero deficit is unity.

Example 6.19—cont'd

We choose the desired continuous-time closed-loop transfer function as second-order with damping ratio 0.88, undamped natural frequency 1.15 rad/s, and with unity gain to meet the settling time and steady-state error specifications. The resulting transfer function is

$$G_{cl}(s) = \frac{1.322}{s^2 + 2.024s + 1.322}$$

The desired closed-loop transfer function is obtained using pole-zero matching:

$$G_d(z) = 0.25921 \cdot 10^{-3} \frac{z+1}{z^2 - 1.96z + 0.9603}$$

By applying Eq. (6.44), we have

$$C(z) = \frac{1.3932(z - 0.8187)(z - 0.9802)(z + 1)}{(z - 1)(z + 0.9293)(z - 0.9601)}$$

We observe that stable pole-zero cancellation occurs and that an integral term is present in the controller, as expected, because the zero steady-state error condition is addressed. The closed-loop step response is shown in Fig. 6.36.

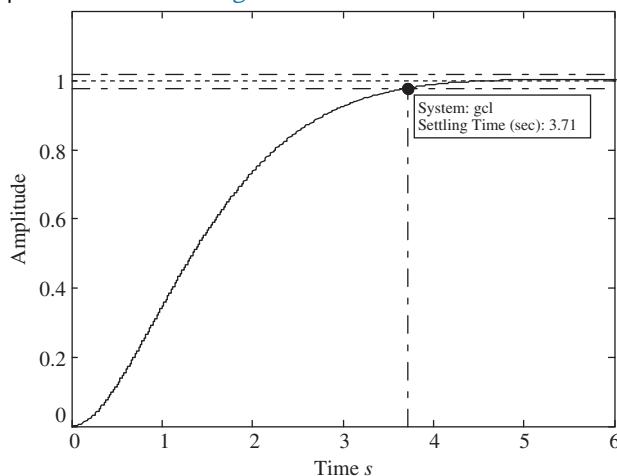


Figure 6.36
Step response for the direct design of Example 6.19.

Example 6.20

Design a digital controller for the type 0 analog plant

$$G(s) = \frac{-0.16738(s - 9.307)(s + 6.933)}{(s^2 + 0.3311s + 9)}$$

to obtain zero steady-state error due to a unit step, a damping ratio of 0.7, and a settling time of about 1 s. The sampling period is chosen as $T = 0.01$ s.

Example 6.20—cont'd**Solution**

The discretized process transfer function is

$$G_{ZAS}(z) = (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} = \frac{-0.16738(z - 1.095)(z - 0.9319)}{z^2 - 1.996z + 0.9967}$$

Note that there is a zero outside the unit circle and this must be included in the desired closed-loop transfer function, which is therefore selected as

$$G_{cl}(z) = \frac{K(z - 1.095)}{z^2 - 1.81z + 0.8269}$$

The closed-loop poles are selected as in Example 6.9, and the value $K = -0.17789$ results from solving the equation $G_{cl}(1) = 1$. By applying Eq. (6.44), we have

$$C(z) = \frac{1.0628(z^2 - 1.81z + 0.8269)}{(z - 1)(z - 0.9319)(z - 0.6321)}$$

Note the stable cancellation of the pole at 0.9319 with a zero and the presence of the pole at $z = 1$. The digital closed-loop step response is shown in Fig. 6.37. The undershoot is due to the unstable zero that cannot be modified by the direct design method. The settling time is less than the required value because the presence of the zero was not considered when selecting the closed-loop poles based on the required damping ratio and settling time.

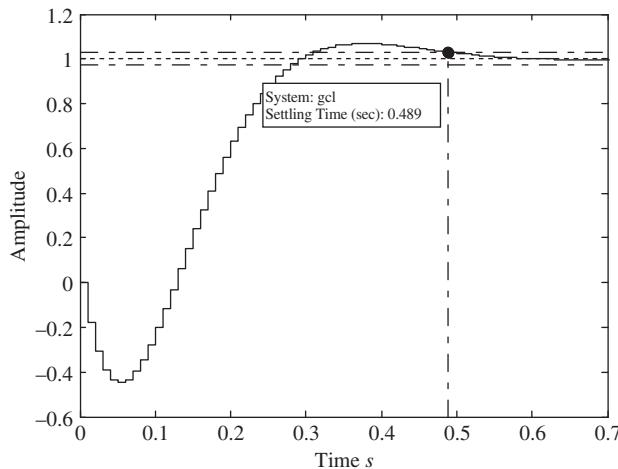


Figure 6.37

Step response for the direct design of Example 6.20.

Example 6.21

Design a digital controller for the type 0 analog plant

$$G(s) = \frac{1}{10s + 1} e^{-5s}$$

to obtain zero steady-state error that results from a unit step and a settling time of about 10 s with no overshoot. The sampling period is chosen as $T = 1$ s.

Solution

The discretized process transfer function is

$$G_{ZAS}(z) = (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} = \frac{0.09516}{z - 0.9048} z^{-5}$$

To meet the causality requirements, a delay of five sampling periods must be included in the desired closed-loop transfer function. Then the settling time of 10 s (including the time delay) is achieved by considering a closed-loop transfer function with a pole at $z = 0.5$ —namely, by selecting

$$G_{cl}(z) = \frac{K}{z - 0.5} z^{-5}$$

Setting $G_{cl}(1) = 1$ yields $K = 0.5$, and then applying Eq. (6.44), we have

$$C(z) = \frac{5.2543(z - 0.9048)z^5}{z^6 - 0.5z^5 - 0.5}$$

The resulting closed-loop step response is shown in Fig. 6.38.

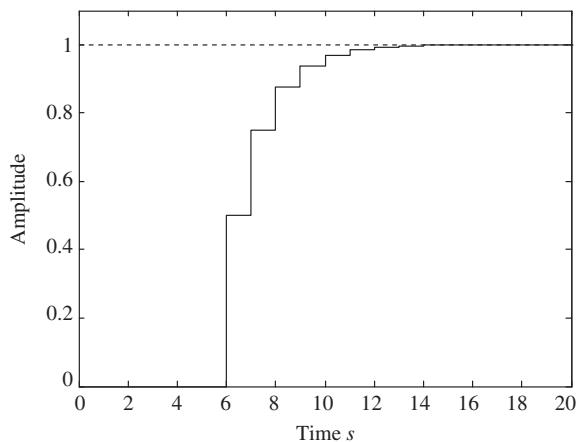


Figure 6.38
Step response for the direct design of Example 6.21.

6.7 Finite settling time design

Continuous-time systems can only reach the desired output asymptotically after an infinite time period. By contrast, digital control systems can be designed to settle at the reference

output after a finite time period and follow it exactly thereafter. By following the direct control design method described in [Section 6.6](#), if all the poles and zeros of the discrete-time process are inside the unit circle, an attractive choice is to select

$$G_{cl}(z) = z^{-k} \quad (6.45)$$

where k must be greater than or equal to the intrinsic delay of the discretized process—namely, the difference between the degree of the denominator and the degree of the numerator of the discrete process transfer function.

Disregarding the time delay, the definition implies that a unit step is tracked perfectly starting at the first sampling point. From [Eq. \(6.44\)](#) we have the **deadbeat controller**

$$C(z) = \frac{1}{G_{ZAS}(z)} \left[\frac{z^{-k}}{1 - z^{-k}} \right] = \frac{1}{G_{ZAS}(z)} \left[\frac{1}{z^k - 1} \right] \quad (6.46)$$

In this case, the only design parameter is the sampling period T , and the overall control system design is very simple. However, finite settling time designs may exhibit undesirable intersample behavior (oscillations) because the control is unchanged between two consecutive sampling points. Further, the control variable can easily assume values that may cause saturation of the DAC or exceed the limits of the actuator in a physical system, resulting in unacceptable system behavior. The behavior of finite settling time designs such as the deadbeat controller must therefore be carefully checked before implementation.

Example 6.22

Design a deadbeat controller for a DC motor speed control system where the (type 0) analog plant has transfer function

$$G(s) = \frac{1}{(s + 1)(s + 10)}$$

and the sampling period is initially chosen as $T = 0.02$ s. Redesign the controller with $T = 0.1$ s.

Solution

Because the discretized process transfer function is

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} \\ &= 1.8604 \times 10^{-4} \frac{z + 0.9293}{(z - 0.8187)(z - 0.9802)} \end{aligned}$$

we have no poles and zeros outside or on the unit circle and a deadbeat controller can therefore be designed by setting

$$G_{cl}(z) = z^{-1}$$

Example 6.22—cont'd

Note that the difference between the order of the denominator and the order of the numerator is one. By applying Eq. (6.46), we have

$$C(z) = \frac{1}{G_{ZAS}(z)} \left[\frac{1}{z-1} \right] = 5375.0533 \frac{(z-0.8187)(z-0.9802)}{(z-1)(z+0.9293)}$$

The resulting sampled and analog closed-loop step response is shown in Fig. 6.39, whereas the corresponding control variable is shown in Fig. 6.40. It appears that, as expected, the sampled process output attains its steady-state value after just one sample—namely, at time $t = T = 0.02$ s, but between samples the output oscillates wildly and the control variable assumes very high values. In other words, the oscillatory behavior of the control variable causes unacceptable intersample oscillations. This can be ascertained analytically by considering the block diagram shown in Fig. 6.41. Recall that the transfer function of the zero-order hold is

$$G_{ZOH}(s) = \frac{1 - e^{-sT}}{s}$$

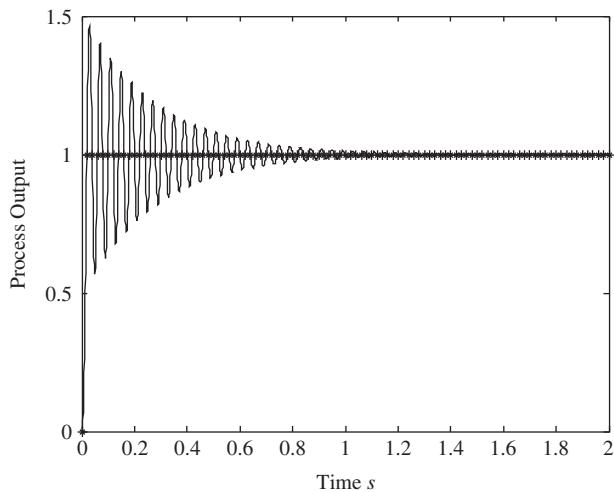


Figure 6.39

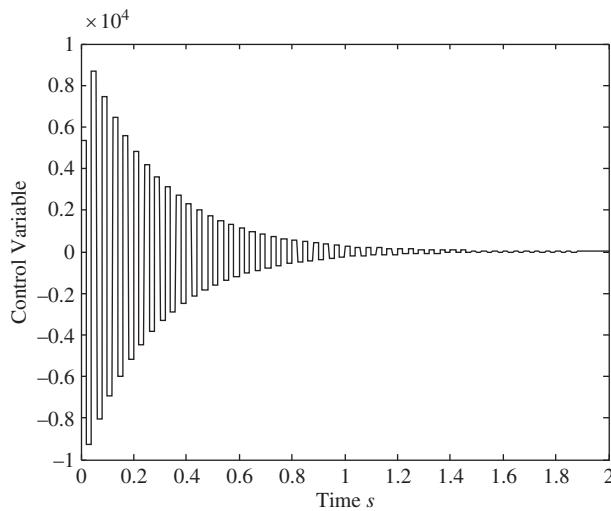
Sampled and analog step response for the deadbeat control of Example 6.22 ($T = 0.02$ s).

Using simple block diagram manipulation, we obtain

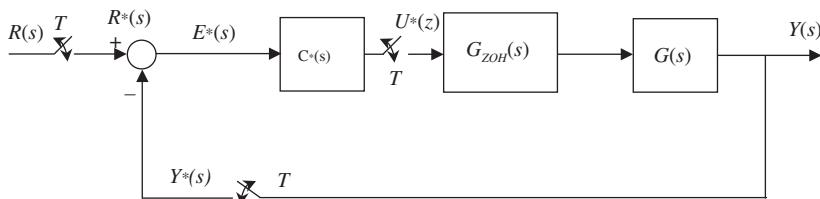
$$\begin{aligned} E^*(s) &= \frac{R^*(s)}{1 + (1 - e^{-sT})C^*(s)\left(\frac{G(s)}{s}\right)^*} \\ &= \frac{R^*(s)}{1 + C^*(s)G_{ZAS}^*(s)} \end{aligned} \quad (6.47)$$

Thus, the output is

$$Y(s) = \left(\frac{1 - e^{-sT}}{s} \right) G(s) C^*(s) \left[\frac{R^*(s)}{1 + C^*(s)G_{ZAS}^*(s)} \right] \quad (6.48)$$

Example 6.22—cont'd**Figure 6.40**

Control variable for the deadbeat control of Example 6.22 ($T = 0.02$ s).

**Figure 6.41**

Block diagram for finite settling time design with analog output.

Eqs. (6.47) and (6.48) are quite general and apply to any system described by the block diagram shown in Fig. 6.41. To obtain the specific output pertaining to our example, we substitute $z = e^{sT}$ in the relevant expressions to obtain

$$G_{ZAS}^*(s) = 1.8604 \times 10^{-4} e^{-sT} \frac{1 + 0.9293 e^{-sT}}{(1 - 0.8187 e^{-sT})(1 - 0.9802 e^{-sT})}$$

$$C^*(s) = 5375.0533 \frac{(1 - 0.9802 e^{-sT})(1 - 0.8187 e^{-sT})}{(1 - e^{-sT})(1 + 0.9293 e^{-sT})}$$

$$R^*(s) = \frac{1}{1 - e^{-sT}}$$

Then substituting in Eq. (6.48) and simplifying gives

$$Y(s) = 5375.0533 \frac{(1 - 0.9802 e^{-sT})(1 - 0.8187 e^{-sT})}{s(s + 1)(s + 10)(1 + 0.9293 e^{-sT})}$$

Example 6.22—cont'd

Expanding the denominator using the identity

$$\begin{aligned}(1 + 0.9293e^{-sT})^{-1} &= 1 + (-0.9293e^{-sT}) + (-0.9293e^{-sT})^2 + (-0.9293e^{-sT})^3 + \dots \\ &= 1 - 0.9293e^{-sT} + 0.8636e^{-2sT} - 0.8025e^{-3sT} + \dots\end{aligned}$$

we obtain

$$\begin{aligned}Y(s) &= \frac{5375.0533}{s(s+1)(s+10)}(1 - 2.782e^{-sT} + 3.3378e^{-2sT} \\ &\quad - 2.2993e^{-3sT} + 0.6930e^{-4sT} + \dots)\end{aligned}$$

Finally, we inverse Laplace transform to obtain the analog output

$$\begin{aligned}y(t) &= 5375.0533\left(\frac{1}{10} - \frac{1}{9}e^{-t} + \frac{1}{90}e^{-10t}\right)\mathbf{1}(t) + \\ &\quad 5375.0533\left(-0.2728 + 0.3031e^{-(t-T)} - 0.0303e^{-10(t-T)}\right)\mathbf{1}(t-T) + \\ &\quad 5375.0533\left(0.3338 - 0.3709e^{-(t-2T)} + 0.0371e^{-10(t-2T)}\right)\mathbf{1}(t-2T) + \\ &\quad 5375.0533\left(-0.2299 + 0.2555e^{-(t-3T)} - 0.0255e^{-10(t-3T)}\right)\mathbf{1}(t-3T) \\ &\quad 5375.0533\left(0.0693 - 0.0770e^{-(t-4T)} + 0.0077e^{-10(t-4T)}\right)\mathbf{1}(t-4T) + \dots\end{aligned}$$

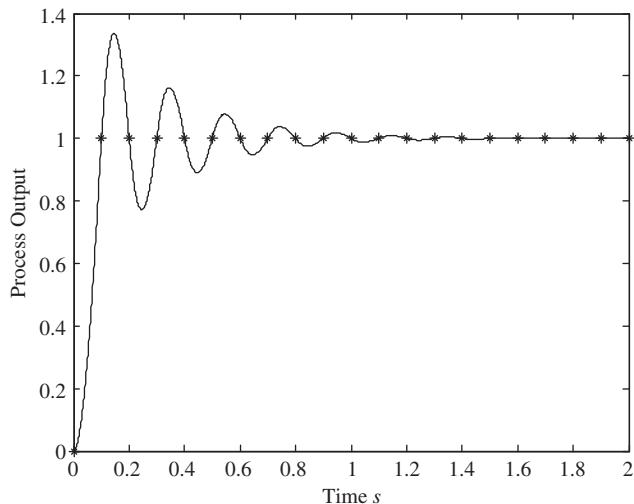
where $\mathbf{1}(t)$ is the unit step function. It is easy to evaluate that, at the sampling points, we have $y(0.02) = y(0.04) = y(0.06) = y(0.08) = \dots = 1$, but between samples the output oscillates wildly, as shown in Fig. 6.39. To reduce intersample oscillations, we set $T = 0.1$ s and obtain the transfer function

$$G_{ZAS}(z) = (1 - z^{-1})\mathcal{Z}\left\{\frac{G(s)}{s}\right\} = 35.501 \times 10^{-4} \frac{z + 0.6945}{(z - 0.9048)(z - 0.3679)}$$

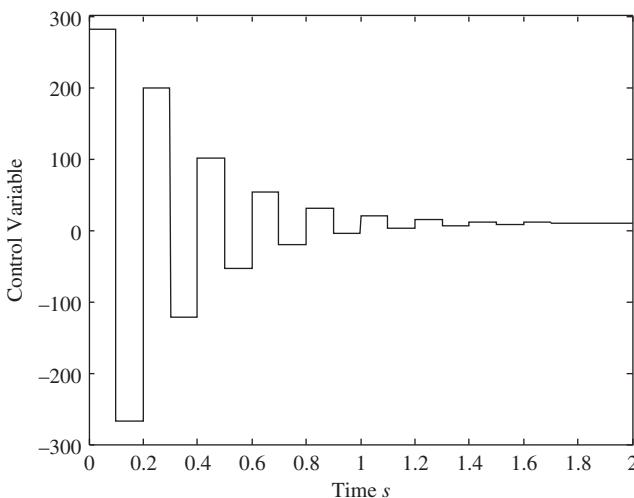
For $G_c(z) = z^{-1}$, we obtain

$$C(z) = \frac{281.6855(z - 0.9048)(z - 0.3679)}{(z - 1)(z + 0.6945)}$$

The resulting sampled and analog closed-loop step responses are shown in Fig. 6.42, whereas the corresponding control variable is shown in Fig. 6.43. The importance of the sampling period selection in the deadbeat controller design is obvious. Note that the oscillations and the amplitude of the control signal are both reduced, but it is also evident that the intersample oscillations cannot be avoided completely with this approach.

Example 6.22—cont'd**Figure 6.42**

Sampled and analog step response for the deadbeat control of Example 6.22 ($T = 0.1$ s).

**Figure 6.43**

Control variable for the deadbeat control of Example 6.22 ($T = 0.1$ s).

6.7.1 Eliminating intersample oscillation

To avoid intersample oscillations, we maintain the control variable constant after n samples, where n is the degree of the denominator of the discretized process. This is done at the expense of achieving the minimum settling time by considering the sampled process output as in the previous examples. The overall design requires a priori specification of the

reference signal. Considering a step signal for the sake of simplicity, a ripple-free response (i.e., a response without intersample oscillations) is achieved if, for any $l \geq n$ (see Eq. (6.45)), we have:

- $e(kT) = 0$
- $u(kT) = \text{constant}$.

Denoting the loop transfer function as $L(z) = C(z)G_{ZAS}(z) = N_L(z)/D_L(z)$, it can be demonstrated that these two conditions are satisfied if and only if

i.

$$N_L(z) + D_L(z) = z^l, l \geq n \quad (6.49)$$

- ii. $D_L(z)$ has a root at $z = 1$ so that the system with unity feedback will be type 1.
- iii. There is no unstable pole-zero cancellations between $C(z)$ and $G_{ZAS}(z)$ for closed-loop stability as shown in Chapter 4.

Using condition (ii), Eq. (6.49) makes the error transfer function

$$\frac{E(z)}{R(z)} = \frac{1}{1 + C(z)G_{ZAS}(z)} = \frac{D_L(z)}{N_L(z) + D_L(z)} = \frac{(z - 1)D_1(z)}{z^l}$$

Letting $D_1(z) = z^{n-1} + a_{n-2}z^{n-2} + \dots + a_1z + a_0$ gives

$$E(z) = (z - 1)(z^{n-1} + a_{n-2}z^{n-2} + \dots + a_1z + a_0)z^{-l} R(z)$$

For a step input, we have

$$E(z) = (z - 1)(z^{n-1} + a_{n-2}z^{n-2} + \dots + a_1z + a_0)z^{-l} \frac{z}{(z - 1)}$$

Simplifying and inverse transforming gives

$$e(k) = \delta(k - l + n) + a_{n-2}\delta(k - l + n - 1) + \dots + a_1\delta(k - l + 2) + a_0\delta(k - l + 1), \quad l \geq n$$

This shows that the error is zero after l time steps.

If we consider a process without poles and zeros outside the unit circle, if the system has no pole at $z = 1$, the controller can be determined by considering

$$G_{ZAS}(z) = \frac{a_{n-1}z^{n-1} + \dots + a_0}{D(z)} \quad (6.50)$$

and by expressing the controller as

$$C(z) = K \frac{D(z)}{(z - 1)(z^{n-1} + b_{n-2}z^{n-2} + \dots + b_0)} \quad (6.51)$$

where the coefficients K, b_{n-1}, \dots, b_0 (and a_{n-1}, \dots, a_0) are determined by solving Eq. (6.49), rewritten as

$$K(a_{n-1}z^{n-1} + \dots + a_0) + (z - 1)(z^{n-1} + \dots + b_0) = z^l, \quad l \geq n \quad (6.52)$$

If the process has a pole at $z = 1$, the same reasoning is applied with

$$G_{ZAS}(z) = \frac{a_{n-1}z^{n-1} + \dots + a_0}{(z - 1)D(z)} \quad (6.53)$$

to obtain again Eq. (6.52). Then, the controller is expressed as

$$C(z) = K \frac{D(z)}{z^{n-1} + b_{n-2}z^{n-2} \dots + b_0} \quad (6.54)$$

Example 6.23

Design a ripple-free deadbeat controller for the type 1 vehicle positioning system of Example 3.3 with transfer function

$$G(s) = \frac{1}{s(s+1)}$$

The sampling period is chosen as $T = 0.1$.

Solution

The discretized process transfer function is

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} \\ &= 4.8374 \times 10^{-3} \frac{z + 0.9672}{(z - 1)(z - 0.9048)} \end{aligned}$$

For this system, $a_0 = 0.004679$, $a_1 = 0.004837$, and $D(z) = z - 0.9048$. Thus, Eq. (6.52) can be written as

$$K(a_1z + a_0) + (z - 1)(z + b_0) = z^2$$

$$K(0.004679z + 0.004837) + (z - 1)(z + b_0) = z^2$$

Equating coefficient, we have

$$z^1: 0.004837K - 1 + b_0 = 0$$

$$z^0: 0.004679K - b_0 = 0$$

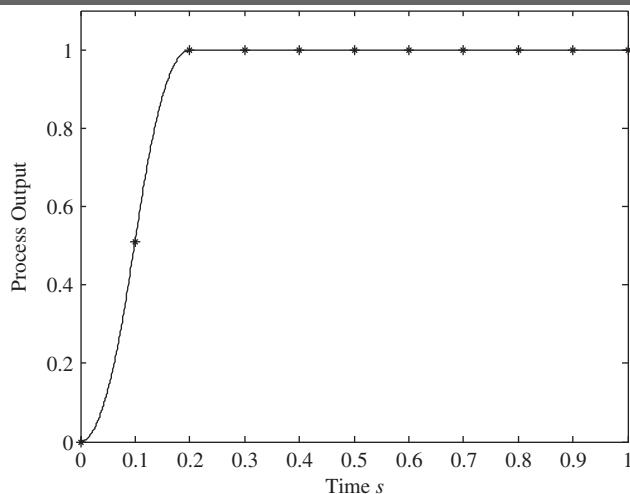
We solve for K and b_0

$$\begin{cases} b_0 = \frac{a_0}{a_0 + a_1} = 0.4917 \\ K = \frac{1}{a_0 + a_1} = 105.0833 \end{cases}$$

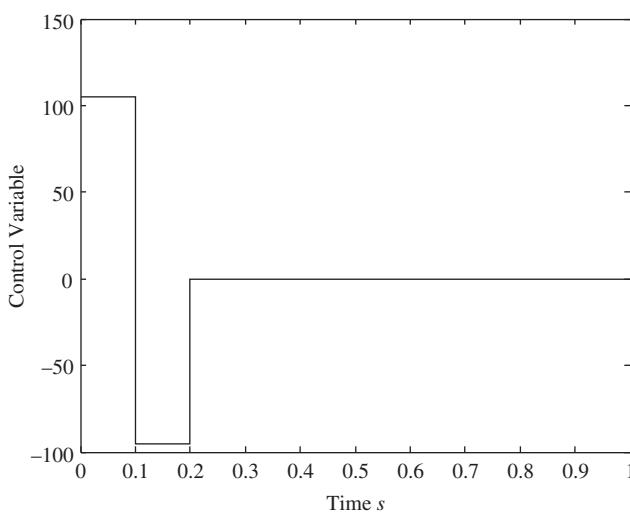
Because the process is type 1, we obtain the controller

$$C(z) = \frac{105.1z - 95.08}{z + 0.4917}$$

The resulting sampled and analog closed-loop step responses are shown in Fig. 6.44, and the corresponding control variable is shown in Fig. 6.45. Note that the control variable is constant after the second sample and that there is no intersample ripple.

Example 6.23—cont'd**Figure 6.44**

Sampled and analog step response for the deadbeat control of Example 6.23.

**Figure 6.45**

Control variable for the deadbeat control of Example 6.23.

Example 6.24

Design a ripple-free deadbeat controller for a DC motor speed control system where the (type 0) analog plant has the transfer function

$$G(s) = \frac{1}{(s+1)(s+10)}$$

Example 6.24—cont'd

The sampling period is chosen as $T = 0.1$ s.

Solution

For a sampling period of 0.01, the discretized process transfer function is

$$\begin{aligned} G_{ZAS}(z) &= (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} \\ &= \frac{0.00355z + 0.002465}{z^2 - 1.273z + 0.3329} \end{aligned}$$

As the process has no pole at unity, the controller transfer function is in the form of Eq. (6.51) and can be written as

$$C(z) = K \frac{z^2 - 1.273z + 0.3329}{(z - 1)(z + b_0)}$$

We write Eq. (6.52) as

$$K(a_1 z + a_0) + (z - 1)(z + b_0) = K(0.00355z + 0.002465) + (z - 1)(z + b_0) = z^2$$

and equate coefficients to obtain

$$0.00355K - 1 + b_0 = 0$$

$$0.002465K - b_0 = 0$$

which yields $b_0 = 0.4099$ and $K = 166.2352$ —that is,

$$C(z) = \frac{166.2352(z - 0.9048)(z - 0.3679)}{(z - 1)(z + 0.4099)}$$

The resulting sampled and analog closed-loop step responses are shown in Fig. 6.46 and the corresponding control variable is shown in Fig. 6.47. As in Example 6.23, there is no inter-sample ripple.

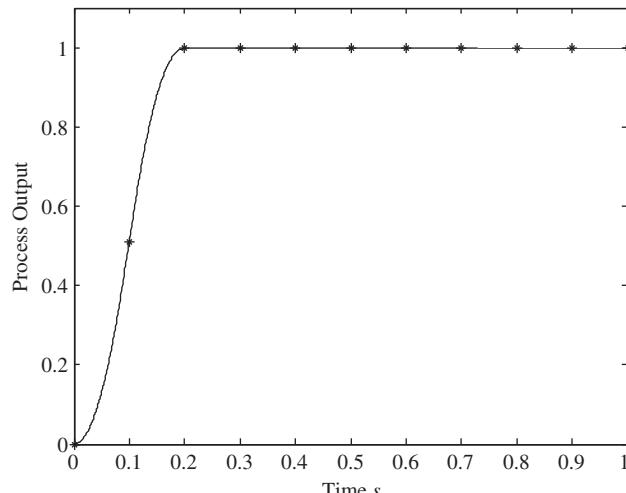
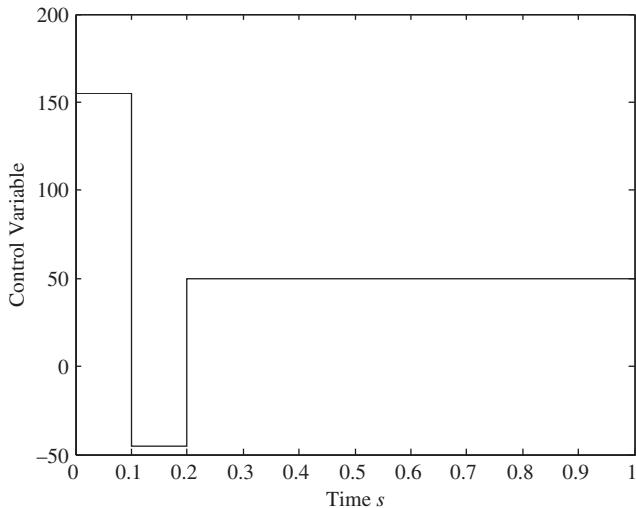


Figure 6.46

Sampled and analog step response for the deadbeat control of Example 6.24.

Example 6.24—cont'd**Figure 6.47**

Control variable for the deadbeat control of Example 6.24.

If the process has poles or zeros on or outside the unit circle, the finite settling time design procedure must be modified to include the constraints on the choice of $G_{cl}(z)$ outlined in [Section 6.6](#). Let the plant transfer function have the form

$$G_{ZAS}(z) = G_1(z) \frac{\prod_{i=1}^{h_z} (z - z_i)}{\prod_{j=1}^{h_p} (z - p_j)} \quad (6.55)$$

where $z_i, i = 1, \dots, h_z$, and $p_j, j = 1, \dots, h_p$, are the zeros and the poles that are on or outside the unit circle, and $G_1(z)$ is a transfer function with all its poles and zeros inside the unit circle.

As shown in [Section 6.6](#), for the closed-loop transfer function $G_{cl}(z)$ to be stable, it must include all the zeros of $G_{ZAS}(z)$ that are outside the unit circle. For causality, it must have (at least) the same pole-zero deficit as the transfer function $G_{ZAS}(z)$. Hence, the closed-loop transfer function must be in the form

$$G_{cl}(z) = \frac{(a_d z^d + \dots + a_1 z + a_0) \prod_{i=1}^{h_z} (z - z_i)}{z^l}, \quad l = n_z + h_z + d \quad (6.56)$$

where n_z is the zero-pole deficit of the process transfer function $G_{ZAS}(z)$, and $a_i, i = 1, \dots, d$, are coefficients to be determined later.

To avoid unstable pole-zero cancellation, $1 - G_{cl}(z)$ must include all the poles of $G_{ZAS}(z)$ that are on or outside the unit circle as zeros (see [Section 6.6](#)). In addition, $1 - G_{cl}(z)$ must

include as many zeros at unity as required for zero steady state error (the proof is left as an exercise). Hence, we require the form

$$1 - G_{cl}(z) = \frac{z^l - (a_d z^d + \dots + a_1 z + a_0) \prod_{i=1}^{h_z} (z - z_i)}{z^l}$$

$$= \frac{(z - 1)^{n_0} (b_f z^f + \dots + b_1 z + b_0) \prod_{i=1}^{h_p} (z - p_j)}{z^l} \quad (6.57)$$

where n_0 is the required number of poles at unity (integrators), and f satisfies

$$f = l - n_0 - h_p$$

to make the numerator of $1 - G_{cl}(z)$ have the same order as its denominator. For zero steady-state error due to a step input, the control loop must include at least one integrator. Therefore, we need to set $n_o = 1$ if $G_{ZAS}(z)$ has no poles at $z = 1$ and $n_o = 0$ otherwise.

The coefficients a_i and b_i are determined by equating the coefficients of the terms of the same order in the equation

$$z^l - (a_d z^d + \dots + a_1 z + a_0) \prod_{i=1}^{h_z} (z - z_i) = (z - 1)^{n_0} (b_f z^f + \dots + b_1 z + b_0) \prod_{i=1}^{h_p} (z - p_j) \quad (6.58)$$

Clearly, the expression shows that $b_f = 1$ and we only have $f + d + 1$ unknown coefficients.

Parameters d and f can be found, together with the order l of $G_{cl}(z)$, in order to meet the following conditions:

- The total number of unknown coefficients a_i and b_i must be equal to the order of $G_{cl}(z)$ —that is, to the number of equations available by equating coefficients of z^{-1} .

$$l = f + d + 1$$

- The order of the denominator polynomial $G_{cl}(z)$ and $1 - G_{cl}(z)$ must satisfy

$$l = n_z + h_z + d$$

The constraints can be written in matrix form as

$$\begin{bmatrix} -1 & 0 & 1 \\ 0 & -1 & 1 \\ 1 & 1 & -1 \end{bmatrix} \begin{bmatrix} d \\ f \\ g \end{bmatrix} = \begin{bmatrix} n_z + h \\ n_o + l \\ -1 \end{bmatrix}$$

$$\begin{bmatrix} -1 & 0 & 1 \\ 1 & 1 & -1 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} d \\ f \\ l \end{bmatrix} = \begin{bmatrix} n_z + h_z \\ -1 \\ n_0 + h_p \end{bmatrix}$$

We solve the equations to obtain

$$\begin{bmatrix} d \\ f \\ l \end{bmatrix} = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} n_z + h_z \\ -1 \\ n_0 + h_p \end{bmatrix} = \begin{bmatrix} n_0 + h_p - 1 \\ n_z + h_z - 1 \\ n_z + h_z + n_0 + h_p - 1 \end{bmatrix} \quad (6.59)$$

Example 6.25

Design a deadbeat controller for an isothermal chemical reactor whose transfer function is¹

$$G(s) = \frac{-1.1354(s - 2.818)}{(s + 2.672)(s + 2.047)}$$

with a sampling period $T = 0.01$.

Solution

For a sampling period of 0.01, the discretized process transfer function is

$$G_{ZAS}(z) = \frac{-0.010931(z - 1.029)}{(z - 0.9797)(z - 0.9736)}$$

The process has one zero and no poles outside the unit circle, and therefore we have $h_z = 1$ and $h_p = 0$. Further, we have the pole-zero deficit $n_z = 1$, and we select $n_o = 1$ because the process transfer function has no pole at $z = 1$. Substituting in Eq. (6.59), we obtain

$$\begin{bmatrix} d \\ f \\ l \end{bmatrix} = \begin{bmatrix} n_0 + h_p - 1 \\ n_z + h_z - 1 \\ n_z + h_z + n_0 + h_p - 1 \end{bmatrix} = \begin{bmatrix} 1 + 0 - 1 \\ 1 + 1 - 1 \\ 1 + 1 + 1 + 0 - 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$$

Eq. (6.58) is then written as

$$z^2 - a_0(z - 1.029) = (z - 1)(z + b_0)$$

$$z^2 - a_0z + 1.029 = z^2 + (b_0 - 1)z - b_0$$

and solved for $a_0 = -34.4828$ and $b_0 = 35.4828$. The resulting closed-loop transfer function

$$G_c(z) = -34.4828 \frac{z - 1.029}{z^2}$$

is substituted in Eq. (6.44) to obtain the controller transfer function

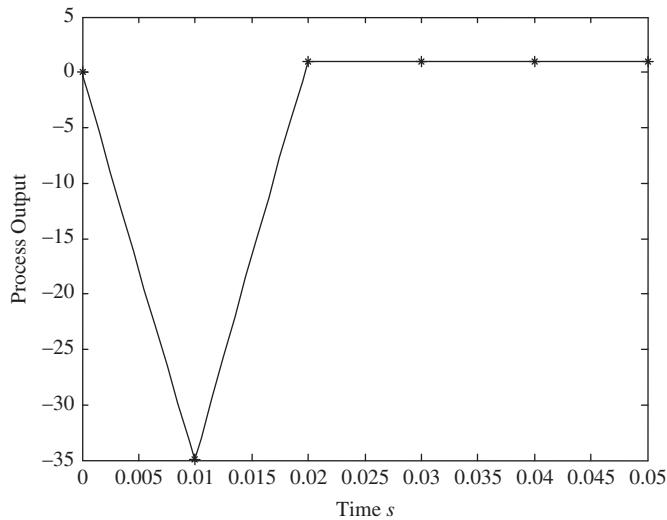
$$C(z) = 3154.4 \frac{(z - 0.9797)(z - 0.9736)}{(z + 35.4828)(z - 1)}$$

We observe that the controller has a pole outside the unit circle. In fact, there is no guarantee of the stability of the controller by using the finite settling time design if the process has poles or zeros outside the unit circle. Further, the transient response can exhibit unacceptable

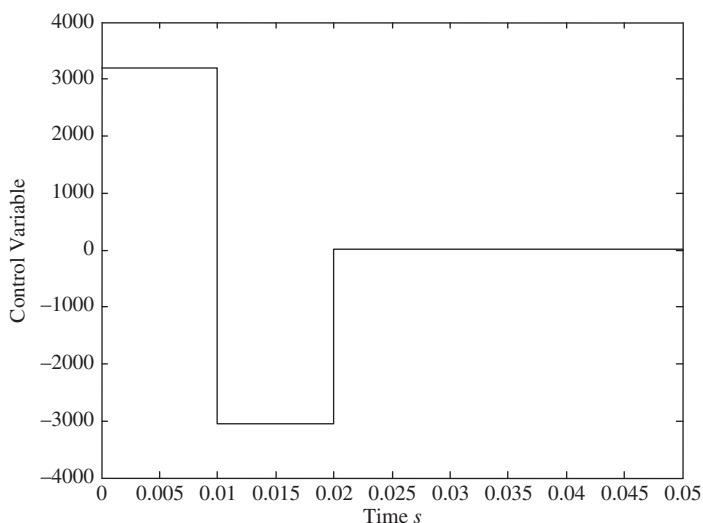
¹ Bequette, (B)W., 2003. Process Control—Modeling, Design, and Simulation. Prentice Hall, Upper Saddle River, NJ.

Example 6.25—cont'd

behavior, as in the step response shown in Figs. 6.48 and 6.49, where we can observe a large undershoot and large control magnitudes.

**Figure 6.48**

Sampled and analog step response for the deadbeat control of Example 6.25.

**Figure 6.49**

Control variable for the deadbeat control of Example 6.25.

Further reading

- Jacquot, R.G., 1981. Modern Digital Control Systems. Marcel Dekker, New York.
- Kuo, B.C., 1991. Automatic Control Systems. Prentice Hall, Englewood Cliffs, NJ.
- Kuo, B.C., 1992. Digital Control Systems. Saunders, Fort Worth, TX.
- Ogata, K., 1987. Digital Control Engineering. Prentice Hall, Englewood Cliffs, NJ.
- Oppenheim, A.V., Schafer, R.W., 1975. Digital Signal Processing. Prentice Hall, Englewood Cliffs, NJ.
- Ragazzini, J.R., Franklin, G.F., 1958. Sampled-Data Control Systems. McGraw-Hill, New York.

Problems

- 6.1 Sketch the z -domain root locus, and find the critical gain for the following systems:
- $G(z) = \frac{K}{z-0.4}$
 - $G(z) = \frac{K}{(z+0.9)(z-0.9)}$
 - $G(z) = \frac{Kz}{(z-0.2)(z-1)}$
 - $G(z) = \frac{K(z+0.9)}{(z-0.2)(z-0.8)}$
- 6.2 Prove that expression (6.6) describes a constant ω_n contour in the z -plane.
- 6.3 Verify that for the system with analog transfer function

$$G(s) = \frac{1}{(s+a)(s+b)}$$

the transfer function of the system with a zero-order hold and sampler has poles at e^{-aT} and e^{-bT} , where T is the sampling period.

- 6.4 Design proportional controllers for the systems of Problem 6.1 to meet the following specifications where possible. If the design specification cannot be met, explain why and suggest a more appropriate controller.
- A damping ratio of 0.7
 - A steady-state error of 10% due to a unit step
 - A steady-state error of 10% due to a unit ramp
- 6.5 Hold equivalence is a digital filter design approach that approximates an analog filter using

$$C(z) = \left(\frac{z-1}{z}\right) \mathcal{Z} \left\{ \mathcal{L}^{-1} \left[\frac{C_a(s)}{s} \right]^* \right\}$$

- a. Obtain the hold-equivalent digital filter for the PD, PI, and PID controllers. Modify the results as necessary to obtain a realizable filter with finite frequency response at the folding frequency.
- b. Why are the filters obtained using hold equivalence always stable?

6.6 Consider the system

$$G_a(s) = \frac{Y(s)}{U(s)} = \frac{1}{(s+1)^2}$$

Use pole-zero matching to determine a digital filter approximation with a sampling period of $T = 0.1$. Verify that not including a digital filter zero at unity results in a time delay by comparing the unit step response for the matched filter with and without a zero at -1 .

6.7 Consider the PI controller

$$C_a(s) = 10 \frac{(s+1)}{s}$$

Use pole-zero matching to determine a digital approximation with a sampling period of $T = 0.1$.

- 6.8 Show that the bilinear transformation of the PID controller expression (5.20) yields expression (6.33).
- 6.9 Show that the bilinear transformation of a PD controller with a high-frequency pole that makes the controller transfer function proper does not yield a pole at $z = -1$.
- 6.10 Design digital controllers to meet the desired specifications for the systems described in Problems 5.5, 5.8, and 5.9 by bilinearly transforming the analog designs.
- 6.11 Design a digital filter by applying the bilinear transformation to the analog (Butterworth) filter

$$C_a(s) = \frac{1}{s^2 + \sqrt{2}s + 1}$$

with $T = 0.1$ s. Then apply prewarping at the 3-dB frequency.

6.12 Design a digital PID controller with $T = 0.1$ for the plant

$$G(s) = \frac{1}{10s + 1} e^{-5s}$$

by applying the Ziegler–Nichols tuning rules presented in Table 5.1.

- 6.13 Design digital controllers to meet the desired specifications for the systems described in Problems 5.5, 5.8, and 5.9 in the z -domain directly.

- 6.14 In Example 4.9, we examined the closed-loop stability of the furnace temperature digital control system with proportional control and a sampling period of 0.01 s. We obtained the z -transfer function

$$G_{ZAS}(z) = 10^{-5} \frac{4.95z + 4.901}{z^2 - 1.97z + 0.9704}$$

Design a controller for the system to obtain zero steady-state error due to a step input without significant deterioration in the transient response.

- 6.15 Consider the DC motor position control system described in Example 3.6, where the (type 1) analog plant has the transfer function

$$G(s) = \frac{1}{s(s+1)(s+10)}$$

and design a digital controller by using frequency response methods to obtain a settling time of about 1 and an overshoot less than 5%.

- 6.16 Use direct control design for the system described in Problem 5.8 (with $T = 0.1$) to design a controller for the transfer function

$$G(s) = \frac{1}{(s+1)(s+5)}$$

to obtain zero steady-state error due to step, a settling time of less than 2 s, and an undamped natural frequency of 5 rad/s. Obtain the discretized and the analog output. Then apply the designed controller to the system

$$G(s) = \frac{1}{(s+1)(s+5)(0.1s+1)}$$

and obtain the discretized and the analog output to verify the robustness of the control system.

- 6.17 Design a deadbeat controller for the system of Problem 5.8 to obtain perfect tracking of a unit step in minimum finite time. Obtain the analog output for the system, and compare your design to that obtained in Problem 5.8. Then apply the controller to the process

$$G(s) = \frac{1}{(s+1)(s+5)(0.1s+1)}$$

to verify the robustness of the control system.

- 6.18 Find a solution for Problem 6.17 that avoids intersample ripple.

- 6.19 The loop gain of a system with unity feedback is for the form

$$L(z) = \frac{N(z)}{(z - 1)^{n_0} D(z)}$$

where n_0 is the number of integrators required to meet steady-state error requirements. Show that the transfer function $1 - G_{cl}(z)$ includes n_0 zeros at unity

- 6.20 Design a deadbeat controller for a steam-drum-level system whose transfer function is²

$$G(s) = \frac{0.25(-s + 1)}{s(2s + 1)}$$

by selecting a sampling period $T = 0.01$.

- 6.21 Determine a closed-form expression for the controller that avoids intersample ripple when the process has a pole at $z = 1$ and there are no poles and zeros outside the unit circle. Then, verify that the expression can be used for Example 6.23.
- 6.22 Find a solution for Example 6.25 by directly applying Eq. (6.51) and by avoiding the unstable pole-zero cancellation.
- 6.23 Solve Example 6.23 starting with the expression of the control variable $U(z)$

$$U(z) = \frac{Y(z)}{G_{ZAS}(z)} = \frac{Y(z)}{R(z)} \frac{R(z)}{G_{ZAS}(z)} = G_{cl}(z) \frac{R(z)}{G_{ZAS}(z)}$$

Hint: The control must be zero after two samples since the process is type 1.

Computer exercises

- 6.24 Write a MATLAB function to plot a constant damped natural frequency contour in the z -plane.
- 6.25 Write a MATLAB function to plot a time-constant contour in the z -plane.
- 6.26 Write a computer program that estimates a first-order-plus-dead-time transfer function with the tangent method and determines the digital PID parameters according to the Ziegler–Nichols formula. Apply the program to the system

$$G(s) = \frac{1}{(s + 1)^8}$$

and simulate the response of the digital control system (with $T = 0.1$) when a set point step change and a load disturbance step are applied. Compare the results with those of Exercise 5.14.

² Bequette, (B)W., 2003. Process Control—Modeling, Design, and Simulation. Prentice Hall, Upper Saddle River, NJ.

- 6.27 To examine the effect of the sampling period on the relative stability and transient response of a digital control system, consider the system

$$G(s) = \frac{1}{(s+1)(s+5)}$$

- a. Obtain the transfer function of the system, the root locus, and the critical gain for $T = 0.01$ s, 0.05 s, 0.1 s.
- b. Obtain the step response for each system at a gain of 2.
- c. Discuss the effect of the sampling period on the transient response and relative stability of the system based on your results from (a) and (b).

State-space representation

Objectives

After completing this chapter, the reader will be able to do the following:

1. Obtain a state-space model from the system transfer function or differential equation.
2. Determine a linearized model of a nonlinear system.
3. Determine the solution of linear (continuous-time and discrete-time) state-space equations.
4. Determine an input–output representation starting from a state–space representation.
5. Determine an equivalent state–space representation of a system by changing the basis vectors.

In this chapter, we discuss an alternative system representation in terms of the system **state variables**, known as the **state–space** representation or **realization**. We examine the properties, advantages, and disadvantages of this representation. We also show how to obtain an input–output representation from a state–space representation. Obtaining a state–space representation from an input–output representation is further discussed in Chapter 8.

The term **realization** arises from the fact that this representation provides the basis for implementing digital or analog filters. In addition, state–space realizations can be used to develop powerful controller design methodologies. Thus, state–space analysis is an important tool in the arsenal of today’s control system designer.

Chapter Outline

- 7.1 State variables** 254
- 7.2 State–space representation** 257
 - 7.2.1 State–space representation in MATLAB 259
 - 7.2.2 Linear versus nonlinear state–space equations 259
- 7.3 Linearization of nonlinear state equations** 262
- 7.4 The solution of linear state–space equations** 265
 - 7.4.1 The Leverrier algorithm 272
 - 7.4.1.1 *Leverrier algorithm* 273
 - 7.4.2 Sylvester’s expansion 277
 - 7.4.3 The state-transition matrix for a diagonal state matrix 278
 - 7.4.3.1 *Properties of constituent matrices* 282

7.4.4 Real form for complex conjugate eigenvalues 284

7.5 The transfer function matrix 285

7.5.1 MATLAB commands 287

7.6 Discrete-time state-space equations 289

7.6.1 MATLAB commands for discrete-time state-space equations 292

7.6.2 Complex conjugate eigenvalues 292

7.7 Solution of discrete-time state-space equations 293

7.7.1 z-transform solution of discrete-time state equations 295

7.8 z-transfer function from state-space equations 300

7.8.1 z-transfer function in MATLAB 303

7.9 Similarity transformation 303

7.9.1 Invariance of transfer functions and characteristic equations 306

Reference 307

Further reading 307

Problems 308

Computer exercises 315

7.1 State variables

Linear continuous-time single-input–single-output (SISO) systems are typically described by the input–output differential equation

$$\begin{aligned} \frac{d^n y}{dt^n} + a_{n-1} \frac{d^{n-1} y}{dt^{n-1}} + \dots + a_1 \frac{dy}{dt} + a_0 y \\ = c_n \frac{d^n u}{dt^n} + c_{n-1} \frac{d^{n-1} u}{dt^{n-1}} + \dots + c_1 \frac{du}{dt} + c_0 u \end{aligned} \quad (7.1)$$

where y is the system output, u is the system input, and a_i , $i = 0, 1, \dots, n-1$, c_j , $j = 0, 1, \dots, n$ are constants. The description is valid for time-varying systems if the coefficients a_i and c_j are explicit functions of time. For a **multi-input–multi-output** (MIMO) system, the representation is in terms of l input–output differential equations of the form Eq. (7.1), where l is the number of outputs. The representation can also be used for nonlinear systems if (7.1) is allowed to include nonlinear terms.

The solution of the differential Eq. (7.1) requires knowledge of the system input $u(t)$ for the period of interest as well as a set of constant initial conditions

$$y(t_0), dy(t_0)/dt, \dots, d^{n-1}y(t_0)/dt^{n-1}$$

where the notation signifies that the derivatives are evaluated at the initial time t_0 . The set of initial conditions is minimal in the sense that incomplete knowledge of this set would prevent the complete solution of Eq. (7.1). On the other hand, additional initial conditions are not needed to obtain the solution. The initial conditions provide a summary of the history of the system up to the initial time. This leads to the following definition.

Definition 7.1: System state

The state of a system is the minimal set of numbers $\{x_i(t_0), i = 1, 2, \dots, n\}$ needed together with the input $u(t)$, with t in the interval $[t_0, t_f]$, to uniquely determine the behavior of the system in the interval $[t_0, t_f]$. The number n is known as the order of the system.

As t increases, the state of the system evolves and each of the numbers $x_i(t)$ becomes a time variable. These variables are known as the **state variables**. In vector notation, the set of state variables form the **state vector**

$$\mathbf{x}(t) = [x_1 \quad x_2 \quad \dots \quad x_n]^T = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad (7.2)$$

The preceding equation follows standard notation in system theory where a column vector is bolded and a row vector is indicated by transposing a column.

State-space is an n -dimensional vector space where $\{x_i(t), i = 1, 2, \dots, n\}$ represent the coordinate axes. So for a second-order system, the state-space is two-dimensional and is known as the **state plane**. For the special case where the state variables are proportional to the derivatives of the output, the state plane is called the **phase plane** and the state variables are called the **phase variables**. Curves in state-space are known as the **state trajectories**, and a plot of state trajectories in the plane is the **state portrait** (or **phase portrait** for the phase plane).

Example 7.1

Consider the equation of motion of a point mass m driven by a force f

$$m\ddot{y} = f$$

where y is the displacement of the point mass. The solution of the differential equation is given by

$$y(t) = \frac{1}{m} \left\{ y(t_0) + \dot{y}(t_0)t + \int_{t_0}^t \int_{t_0}^t f(\tau) d\tau \right\}$$

Clearly, a complete solution can only be obtained, given the force, if the two initial conditions $\{y(t_0), \dot{y}(t_0)\}$ are known. Hence, these constants define the state of the system at time t_0 , and the system is second order. The state variables are

$$\begin{aligned} x_1(t) &= y(t) \\ x_2(t) &= \dot{y}(t) \end{aligned}$$

and the state vector is

$$\mathbf{x}(t) = [x_1 \quad x_2]^T = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

The state variables are phase variables in this case because the second is the derivative of the first.

Example 7.1—cont'd

The state variables are governed by the two first-order differential equations

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= u/m\end{aligned}$$

where $u = f$. The first of the two equations follows from the definitions of the two state variables. The second is obtained from the equation of motion for the point mass. The two differential equations together with the algebraic expression

$$y = x_1$$

are equivalent to the second-order differential equation because solving the first-order differential equations and then substituting in the algebraic expression yields the output y . For a force satisfying the state feedback law

$$\frac{u}{m} = -9x_1 - 3x_2$$

We have a second-order underdamped system with the solution depending only on the initial conditions. The solutions for different initial conditions can be obtained by repeatedly using the MATLAB commands **lsim** or **initial**. Each of these solutions yields position and velocity data for a phase trajectory, which is a plot of velocity versus position. A set of these trajectories corresponding to different initial states gives the phase portrait of Fig. 7.1. The time variable does not appear explicitly in the phase portrait and is an implicit parameter. Arrows indicate the direction of increasing time.

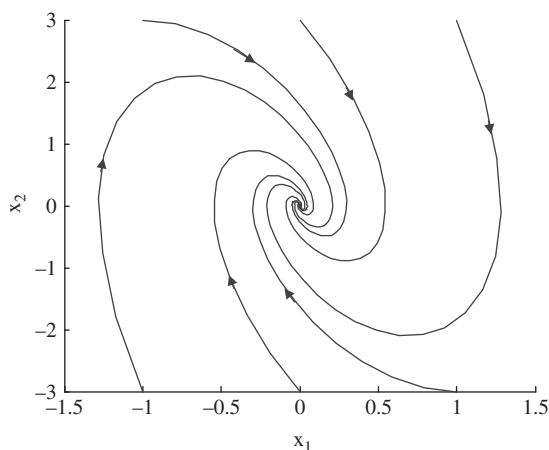


Figure 7.1
Phase portrait for a point mass.

Note that the choice of state variables is not unique. For example, one could use the displacement y and the sum of displacement and velocity as state variables (see Example 7.20). This choice has no physical meaning but nevertheless satisfies the definition of state variables. The freedom of choice is a general characteristic of state equations and is not restricted to this example. It allows us to represent a system so as to reveal its characteristics more clearly and is explored in later sections.

7.2 State-space representation

In Example 7.1, two first-order equations governing the state variables were obtained from the second-order input–output differential equation and the definitions of the state variables. These equations are known as **state equations**. In general, there are n state equations for an n^{th} -order system. State equations can be obtained for state variables of systems described by input–output differential equations, with the form of the equations depending on the nature of the system. For example, the equations are time varying for time-varying systems and nonlinear for nonlinear systems. State equations for linear time-invariant systems can also be obtained from their transfer functions.

The algebraic equation expressing the output in terms of the state variables is called the **output equation**. For multioutput systems, a separate output equation is needed to define each output. The state and output equations together provide a complete representation for the system described by the differential equation, which is known as the **state-space representation**. For linear systems, it is often more convenient to write the state equations as a single matrix equation referred to as the **state equation**. Similarly, the output equations can be combined in a single-output equation in matrix form. The matrix form of the state–space representation is demonstrated in Example 7.2.

Example 7.2

The state–space equations for the system of Example 7.1 in matrix form are

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1/m \end{bmatrix} u$$

$$y = [1 \quad 0] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

The general form of the state–space equations for linear systems is

$$\begin{aligned} \dot{\mathbf{x}}(t) &= A\mathbf{x}(t) + B\mathbf{u}(t) \\ \mathbf{y}(t) &= C\mathbf{x}(t) + D\mathbf{u}(t) \end{aligned} \tag{7.3}$$

where $\mathbf{x}(t)$ is an $n \times 1$ real vector, $\mathbf{u}(t)$ is an $m \times 1$ real vector, and $\mathbf{y}(t)$ is an $l \times 1$ real vector. The matrices in the equations are

Example 7.2—cont'd $A = n \times n$ state matrix $B = n \times m$ input or control matrix $C = l \times n$ output matrix $D = l \times m$ direct transmission matrix

The orders of the matrices are dictated by the dimensions of the vectors and the rules of vector-matrix multiplication. For example, in the single-input (SI) case, B is a column vector, and in the single-output (SO) case, both C and D are row vectors. For the SISO case, D is a scalar. The entries of the matrices are constant for time-invariant systems and functions of time for time-varying systems.

Example 7.3

The following are examples of state-space equations for linear systems.

1. A third-order 2-input–2-output (MIMO) linear time-invariant system:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} 1.1 & 0.3 & -1.5 \\ 0.1 & 3.5 & 2.2 \\ 0.4 & 2.4 & -1.1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 0.1 & 0 \\ 0 & 1.1 \\ 1.0 & 1.0 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

2. A second-order 2-output–single-input (SIMO) linear time-varying system:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} \sin(t) & \cos(t) \\ 1 & e^{-2t} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Here, the direct transmission matrix D is zero, and the input and output matrices are constant. But the system is time varying because the state matrix has some entries that are functions of time.

7.2.1 State-space representation in MATLAB

MATLAB has a special state-space representation obtained with the command `ss`. However, some state commands only operate on one or more of the matrices (A , B , C , D). To enter a matrix

$$A = \begin{bmatrix} 1 & 1 \\ -5 & -4 \end{bmatrix}$$

use the command

`>> A = [0, 1; -5, -4]`

If B , C , and D are similarly entered, we obtain the state-space quadruple \mathbf{p} with the command

`>> p = ss(A, B, C, D)`

We can also specify names for the input (torque, say) and output (position) using the `set` command

`>> set(p, 'input', 'output', 'torque', 'position')`

7.2.2 Linear versus nonlinear state-space equations

It is important to remember that the form Eq. (7.3) is only valid for linear state equations. Nonlinear state equations involve nonlinear functions and cannot be written in terms of the matrix quadruple (A , B , C , D).

Example 7.4

Obtain a state-space representation for the s -degree-of-freedom (s -D.O.F.) robotic manipulator from the equation of motion

$$M(\mathbf{q})\ddot{\mathbf{q}} + V(\mathbf{q}, \dot{\mathbf{q}})\dot{\mathbf{q}} + g(\mathbf{q}) = \tau$$

where.

\mathbf{q} = vector of generalized coordinates

$M(\mathbf{q})$ = $s \times s$ positive definite inertia matrix

$V(\mathbf{q}, \dot{\mathbf{q}})$ = $s \times s$ matrix of velocity-related terms

$g(\mathbf{q})$ = $s \times 1$ vector of gravitational terms

τ = vector of generalized forces

The output of the manipulator is the position vector \mathbf{q} .

Example 7.4—cont'd**Solution**

The system is of order $2s$, as $2s$ initial conditions are required to completely determine the solution. The most natural choice of state variables is the vector

$$\mathbf{x} = \text{col}\{\mathbf{x}_1, \mathbf{x}_2\} = \text{col}\{\mathbf{q}, \dot{\mathbf{q}}\}$$

where $\text{col}\{\cdot\}$ denotes a column vector. The associated state equations are

$$\begin{bmatrix} \dot{\mathbf{x}}_1 \\ \dot{\mathbf{x}}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{x}_2 \\ -M^{-1}(\mathbf{x}_1)\{V(\mathbf{x}_1, \mathbf{x}_2)\mathbf{x}_2 + \mathbf{g}(\mathbf{x}_1)\} \end{bmatrix} + \begin{bmatrix} 0 \\ M^{-1}(\mathbf{x}_1) \end{bmatrix} \mathbf{u}$$

with the generalized force now denoted by the symbol \mathbf{u} .

The output equation is

$$\mathbf{y} = \mathbf{x}_1$$

This equation is linear and can be written in the standard form

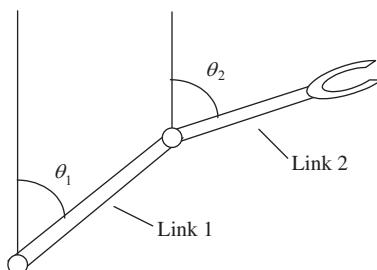
$$\mathbf{y} = [I_s \mid 0_{s \times s}] \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}$$

Example 7.5

Write the state-space equations for the 2-D.O.F. **anthropomorphic manipulator** in Fig. 7.2. The equations of motion of the manipulator are as in Example 7.4 with the definitions

$$M(\theta) = \begin{bmatrix} (m_1 + m_2)l_1^2 + m_2l_2^2 + 2m_2l_1l_2 \cos(\theta_2) & m_2l_2^2 + m_2l_1l_2 \cos(\theta_2) \\ m_2l_2^2 + m_2l_1l_2 \cos(\theta_2) & m_2l_2^2 \end{bmatrix}$$

$$V(\theta, \dot{\theta})\dot{\theta} = \begin{bmatrix} -m_2l_1l_2 \sin(\theta_2)\dot{\theta}_2(2\dot{\theta}_1 + \dot{\theta}_2) \\ m_2l_1l_2 \sin(\theta_2)\dot{\theta}_1^2 \end{bmatrix}$$

**Figure 7.2**

A 2-degree-of-freedom (2-D.O.F.) anthropomorphic manipulator.

Example 7.5—cont'd

$$\mathbf{g}(\boldsymbol{\theta}) = \begin{bmatrix} (m_1 + m_2)gl_1 \sin(\theta_1) + m_2gl_2 \sin(\theta_1 + \theta_2) \\ m_2gl_2 \sin(\theta_1 + \theta_2) \end{bmatrix}$$

where m_i , $i = 1, 2$ are the masses of the two links; l_i , $i = 1, 2$ are their lengths; and \mathbf{g} is the acceleration due to gravity. $(\boldsymbol{\theta}, \dot{\boldsymbol{\theta}})$ are the vectors of angular positions and angular velocities, respectively.

Solution

The state equations can be written using the results of Example 7.4 as

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ -L \left\{ \begin{bmatrix} -m_2l_1l_2\sin(\theta_2)\dot{\theta}_2(2\dot{\theta}_1 + \dot{\theta}_2) \\ m_2l_1l_2\sin(\theta_2)\dot{\theta}_1^2 \end{bmatrix} + \begin{bmatrix} (m_1 + m_2)gl_1\sin(\theta_1) + m_2gl_2\sin(\theta_1 + \theta_2) \\ m_2gl_2\sin(\theta_1 + \theta_2) \end{bmatrix} \right\} \end{bmatrix} + \begin{bmatrix} 0 \\ L \end{bmatrix} \mathbf{u}$$

where

$$L = \frac{1}{\det(D)} \begin{bmatrix} m_2l_2^2 & -(m_2l_2^2 + m_2l_1l_2 \cos(\theta_2)) \\ -(m_2l_2^2 + m_2l_1l_2 \cos(\theta_2)) & (m_1 + m_2)l_1^2 + m_2l_2^2 + 2m_2l_1l_2 \cos(\theta_2) \end{bmatrix}$$

$$\mathbf{x} = \text{col}\{\mathbf{x}_1, \mathbf{x}_2\} = \text{col}\{\boldsymbol{\theta}, \dot{\boldsymbol{\theta}}\}$$

The general form of nonlinear state-space equations is

$$\begin{aligned} \dot{\mathbf{x}} &= \mathbf{f}(\mathbf{x}, \mathbf{u}) \\ \mathbf{y} &= \mathbf{g}(\mathbf{x}, \mathbf{u}) \end{aligned} \tag{7.4}$$

where $\mathbf{f}(.)$ ($n \times 1$) and $\mathbf{g}(.)$ ($l \times 1$) are vectors of functions satisfying mathematical conditions that guarantee the existence and uniqueness of solution. But a form that is often encountered in practice and includes the equations of robotic manipulators is

$$\begin{aligned} \dot{\mathbf{x}} &= \mathbf{f}(\mathbf{x}) + \mathbf{B}(\mathbf{x})\mathbf{u} \\ \mathbf{y} &= \mathbf{g}(\mathbf{x}) + \mathbf{D}(\mathbf{x})\mathbf{u} \end{aligned} \tag{7.5}$$

The state Eq. (7.5) is said to be affine in the control because the RHS is affine (includes a constant vector) for constant \mathbf{x} .

7.3 Linearization of nonlinear state equations

Nonlinear state equations of the form Eq. (7.4) or (7.5) can be approximated by linear state equations of the form Eq. (7.3) for small ranges of the control and state variables.

The linear equations are based on the **first-order approximation**

$$f(x) = f(x_0) + \frac{df}{dx} \Big|_{x_0} \Delta x + O(\Delta^2 x) \quad (7.6)$$

where x_0 is a constant and $\Delta x = x - x_0$ is a perturbation from the constant. The error associated with the approximation is of order $\Delta^2 x$ and is therefore acceptable for small perturbations. For a function of n variables, Eq. (7.6) can be modified to

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \frac{\partial f}{\partial \mathbf{x}_1} \Big|_{\mathbf{x}_0} \Delta \mathbf{x}_1 + \dots + \frac{\partial f}{\partial \mathbf{x}_n} \Big|_{\mathbf{x}_0} \Delta \mathbf{x}_n + O(\|\Delta \mathbf{x}\|^2) \quad (7.7)$$

where \mathbf{x}_0 is the constant vector

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

and $\Delta \mathbf{x}$ denotes the perturbation vector

$$\Delta \mathbf{x} = [x_1 - x_{10} \quad x_2 - x_{20} \quad \dots \quad x_n - x_{n0}]^T$$

The term $\|\Delta \mathbf{x}\|^2$ denotes the sum of squares of the entries of the vector (i.e., its 2-norm), which is a measure of the length or “size” of the perturbation vector.¹

For nonlinear state-space equations of the form Eq. (7.4), let the i^{th} entry of the vector \mathbf{f} be f_i . Then applying Eq. (7.7) to f_i yields the approximation

$$\begin{aligned} f_i(\mathbf{x}, \mathbf{u}) &= f_i(\mathbf{x}_0, \mathbf{u}_0) + \frac{\partial f_i}{\partial \mathbf{x}_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \Delta \mathbf{x}_1 + \dots + \frac{\partial f_i}{\partial \mathbf{x}_n} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \Delta \mathbf{x}_n \\ &\quad + \frac{\partial f_i}{\partial \mathbf{u}_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \Delta \mathbf{u}_1 + \dots + \frac{\partial f_i}{\partial \mathbf{u}_m} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \Delta \mathbf{u}_m \end{aligned} \quad (7.8)$$

which can be rewritten as

¹ The error term dependent on this perturbation is assumed to be small and is neglected in the sequel.

$$f_i(\mathbf{x}, \mathbf{u}) - f_i(\mathbf{x}_0, \mathbf{u}_0) = \left[\frac{\partial f_i}{\partial x_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \cdots \frac{\partial f_i}{\partial x_n} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \right] \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_{n-1} \\ \Delta x_n \end{bmatrix} + \left[\frac{\partial f_i}{\partial u_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \cdots \frac{\partial f_i}{\partial u_m} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \right] \begin{bmatrix} \Delta u_1 \\ \Delta u_2 \\ \vdots \\ \Delta u_{m-1} \\ \Delta u_m \end{bmatrix} \quad (7.9)$$

In most situations where we seek a linearized model, the nominal state is an **equilibrium point**. This term refers to an initial state where the system remains unless perturbed. In other words, it is a system where the state's rate of change, as expressed by the RHS of the state equation, must be zero (see Definition 8.1). Thus, if the perturbation is about an equilibrium point, then the derivative of the state vector is zero at the nominal state; that is, $f_i(\mathbf{x}_0, \mathbf{u}_0)$ is zero and $\mathbf{f}(\mathbf{x}_0, \mathbf{u}_0)$ is a zero vector.

The i th entry g_i of the vector \mathbf{g} can be similarly expanded to yield the perturbation in the i th output

$$\begin{aligned} \Delta y_i &= g_i(\mathbf{x}, \mathbf{u}) - g_i(\mathbf{x}_0, \mathbf{u}_0) \\ &= \left[\frac{\partial g_i}{\partial x_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \cdots \frac{\partial g_i}{\partial x_n} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \right] \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_{n-1} \\ \Delta x_n \end{bmatrix} \\ &\quad + \left[\frac{\partial g_i}{\partial u_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \cdots \frac{\partial g_i}{\partial u_m} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \right] \begin{bmatrix} \Delta u_1 \\ \Delta u_2 \\ \vdots \\ \Delta u_{m-1} \\ \Delta u_m \end{bmatrix} \end{aligned} \quad (7.10)$$

We also note that the derivative of the perturbation vector is

$$\Delta \dot{\mathbf{x}} = \frac{d \Delta \mathbf{x}}{dt} = \frac{d(\mathbf{x} - \mathbf{x}_0)}{dt} = \dot{\mathbf{x}} \quad (7.11)$$

because the nominal state \mathbf{x}_0 is constant.

We now substitute the approximations Eqs. (7.9) and (7.10) in the state and output equations, respectively, to obtain the linearized equations

$$\begin{aligned} \Delta \dot{\mathbf{x}} &= \left[\begin{array}{ccc} \frac{\partial f_1}{\partial x_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial f_1}{\partial x_n} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial f_n}{\partial x_n} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \end{array} \right] \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_{n-1} \\ \Delta x_n \end{bmatrix} + \left[\begin{array}{ccc} \frac{\partial f_1}{\partial u_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial f_1}{\partial u_m} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial u_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial f_n}{\partial u_m} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \end{array} \right] \begin{bmatrix} \Delta u_1 \\ \Delta u_2 \\ \vdots \\ \Delta u_{m-1} \\ \Delta u_m \end{bmatrix} \\ \Delta \mathbf{y} &= \left[\begin{array}{ccc} \frac{\partial g_1}{\partial x_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial g_1}{\partial x_n} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n}{\partial x_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial g_n}{\partial x_n} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \end{array} \right] \cdot \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_{n-1} \\ \Delta x_n \end{bmatrix} + \left[\begin{array}{ccc} \frac{\partial g_1}{\partial u_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial g_1}{\partial u_m} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n}{\partial u_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial g_n}{\partial u_m} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \end{array} \right] \cdot \begin{bmatrix} \Delta u_1 \\ \Delta u_2 \\ \vdots \\ \Delta u_{m-1} \\ \Delta u_m \end{bmatrix} \end{aligned} \quad (7.12)$$

Dropping the Δ s reduces Eq. (7.12) to (7.13), with the matrices of the linear state equations defined as the **Jacobians**:

$$\begin{aligned} A &= \left[\begin{array}{ccc} \frac{\partial f_1}{\partial x_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial f_1}{\partial x_n} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial f_n}{\partial x_n} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \end{array} \right] \quad B = \left[\begin{array}{ccc} \frac{\partial f_1}{\partial u_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial f_1}{\partial u_m} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial u_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial f_n}{\partial u_m} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \end{array} \right] \\ C &= \left[\begin{array}{ccc} \frac{\partial g_1}{\partial x_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial g_1}{\partial x_n} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n}{\partial x_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial g_n}{\partial x_n} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \end{array} \right] \quad D = \left[\begin{array}{ccc} \frac{\partial g_1}{\partial u_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial g_1}{\partial u_m} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n}{\partial u_1} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} & \dots & \frac{\partial g_n}{\partial u_m} \Big|_{(\mathbf{x}_0, \mathbf{u}_0)} \end{array} \right] \end{aligned} \quad (7.13)$$

Example 7.6

Consider the equation of motion of the nonlinear spring-mass-damper system given by

$$m\ddot{y} + b(y)\dot{y} + k(y) = f$$

where y is the displacement, f is the applied force, m is a mass of 1 kg, $b(y)$ is a nonlinear damper constant, and $k(y)$ is a nonlinear spring force. Find the equilibrium position corresponding to a force f_0 in terms of the spring force, then linearize the equation of motion about this equilibrium.

Solution

The equilibrium of the system with a force f_0 is obtained by setting all the time derivatives equal to zero and solving for y to obtain

$$y_0 = k^{-1}(f_0)$$

where $k^{-1}(\bullet)$ denotes the inverse function. The equilibrium is therefore at zero velocity and the position y_0 .

The nonlinear state equation for the system with state vector $\mathbf{x} = [x_1, x_2]^T = [y, \dot{y}]^T$ is

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ -k(x_1) - b(x_1)x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix}u$$

where $u = f$. Then linearizing about the equilibrium, we obtain

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -\frac{dk(x_1)}{dx_1}|_{y_0} & -b(x_1)|_{y_0} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix}u$$

Clearly, the entries of the state matrix are constants whose values depend on the equilibrium position. In addition, terms that are originally linear do not change because of linearization.

7.4 The solution of linear state-space equations

The state-space Eq. (7.3) is linear and can therefore be Laplace transformed to obtain their solution. Clearly, once the state equation is solved for the state vector \mathbf{x} , substitution in the output equation easily yields the output vector \mathbf{y} . So we begin by examining the Laplace transform of the state equation. We recall that to solve a scalar differential equation of the form

$$\dot{x} = ax + bu,$$

we Laplace transform

$$X(s) - x(0) = aX(s) + bU(s)$$

then solve for $X(s)$

$$X(s) = \frac{x(0)}{s-a} + \frac{bU(s)}{s-a}$$

We inverse Laplace transform and use the convolution theorem to obtain

$$x(t) = e^{at}x(0) + \int_0^t e^{a(t-\tau)}bu(\tau)d\tau$$

We follow the same procedure to solve the state equation.

The state equation involves the derivative $\dot{\mathbf{x}}$ of the state vector \mathbf{x} . Because Laplace transformation is simply multiplication by a scalar followed by integration, the Laplace transform of this derivative is the vector of Laplace transforms of its entries. More specifically,

$$\begin{aligned}\mathcal{L}\{\dot{\mathbf{x}}(t)\} &= [s\mathbf{X}_i(s) - \mathbf{x}_i(0)] = s[\mathbf{X}_i(s)] - [\mathbf{x}_i(0)]\{\dot{\mathbf{x}}(t)\} \\ &= s\mathbf{X}(s) - \mathbf{x}(0)\end{aligned}\tag{7.14}$$

Using a similar argument, the Laplace transform of the product $A\mathbf{x}$ is

$$\begin{aligned}\mathcal{L}\{A\mathbf{x}(t)\} &= \mathcal{L}\left\{\left[\sum_{j=1}^n a_{ij}\mathbf{x}_j\right]\right\} = \left[\sum_{j=1}^n a_{ij}\mathcal{L}\{\mathbf{x}_j\}\right] = \left[\sum_{j=1}^n a_{ij}\mathbf{x}_j(s)\right] \\ &= A\mathbf{X}(s)\end{aligned}\tag{7.15}$$

Hence, the state equation

$$\dot{\mathbf{x}}(t) = A\mathbf{x}(t) + B\mathbf{u}(t), \quad \mathbf{x}(0)$$

has the Laplace transform

$$s\mathbf{x}(s) - \mathbf{x}(0) = A\mathbf{x}(s) + B\mathbf{U}(s)\tag{7.16}$$

Rearranging terms, we obtain

$$[sI_n - A]\mathbf{X}(s) = \mathbf{x}(0) + B\mathbf{U}(s)\tag{7.17}$$

Then premultiplying by the inverse of $[sI_n - A]$ gives

$$\mathbf{X}(s) = [sI_n - A]^{-1}[\mathbf{x}(0) + B\mathbf{U}(s)] \quad (7.18)$$

We now need to inverse Laplace transform to obtain the solution of the state equation. So we first examine the inverse known as the **resolvent matrix**

$$[sI_n - A]^{-1} = \frac{1}{s} \left[I_n - \frac{1}{s} A \right]^{-1} \quad (7.19)$$

This can be expanded as

$$\frac{1}{s} \left[I_n - \frac{1}{s} A \right]^{-1} = \frac{1}{s} \left\{ I_n + \frac{1}{s} A + \frac{1}{s^2} A^2 + \dots + \frac{1}{s^i} A^i + \dots \right\} \quad (7.20)$$

Then inverse Laplace transforming yields the series

$$\mathcal{L}^{-1} \left\{ [sI_n - A]^{-1} \right\} = I_n + At + \frac{(At)^2}{2!} + \dots + \frac{(At)^i}{i!} + \dots \quad (7.21)$$

This summation is a matrix version of the exponential function

$$e^{at} = 1 + at + \frac{(at)^2}{2!} + \dots + \frac{(at)^i}{i!} + \dots \quad (7.22)$$

It is therefore known as the **matrix exponential** and is written as

$$e^{At} = \sum_{i=0}^{\infty} \frac{(At)^i}{i!} = \mathcal{L}^{-1} \left\{ [sI_n - A]^{-1} \right\} \quad (7.23)$$

Returning to Eq. (7.18), we see that the first term can now be easily inverse-transformed using Eq. (7.23) to obtain the response due to the initial conditions, with zero input. This response is known as the **zero-input response** and is given by

$$\mathbf{x}_{ZI}(t) = e^{At} \mathbf{x}(0) \quad (7.24)$$

The second term requires the use of the convolution property of Laplace transforms, which states that multiplication of Laplace transforms is equivalent to convolution of their inverses. Hence, the solution of the state equation is given by

$$\mathbf{x}(t) = e^{At} \mathbf{x}(0) + \int_0^t e^{A(t-\tau)} B \mathbf{u}(\tau) d\tau \quad (7.25)$$

By superposition, the total response of the system due to both the initial state and the input is the sum of the zero-input response $\mathbf{x}_{ZI}(t)$ and the **zero-state response** $\mathbf{x}_{ZS}(t)$

$$\begin{aligned}\mathbf{x}(t) &= \mathbf{x}_{ZI}(t) + \mathbf{x}_{ZS}(t) \\ &= \underbrace{e^{At}\mathbf{x}(0)}_{\text{zero-input response}} + \underbrace{\int_0^t e^{A(t-\tau)}B\mathbf{u}(\tau)d\tau}_{\text{zero-state response}}\end{aligned}\quad (7.26)$$

The solution for nonzero initial time is obtained by simply shifting the time variable to get

$$\mathbf{x}(t) = e^{A(t-t_0)}\mathbf{x}(0) + \int_0^t e^{A(t-\tau)}B\mathbf{u}(\tau)d\tau \quad (7.27)$$

The zero-input response involves the change of the system state from the initial vector $\mathbf{x}(0)$ to the vector $\mathbf{x}(t)$ through multiplication by the matrix exponential. Hence, the matrix exponential is also called the **state-transition matrix**. This name is also given to a matrix that serves a similar function in the case of time-varying systems and depends on the initial as well as the final time and not just the difference between them. However, the matrix exponential form of the state-transition matrix is only valid for **linear time-invariant systems**.

To obtain the output of the system, we substitute Eq. (7.27) into the output equation

$$\mathbf{y}(t) = C\mathbf{x}(t) + D\mathbf{u}(t)$$

This gives the time response

$$\mathbf{y}(t) = C\left\{e^{A(t-t_0)}\mathbf{x}(t_0) + \int_{t_0}^t e^{A(t-\tau)}B\mathbf{u}(\tau)d\tau\right\} + D\mathbf{u}(t) \quad (7.28)$$

Example 7.7

The state equations of an armature-controlled DC motor are given by

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -10 & -11 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 10 \end{bmatrix} u$$

where x_1 is the angular position, x_2 is the angular velocity, and x_3 is the armature current. Find the following:

1. The state-transition matrix
2. The response due to an initial current of 10 mA with zero angular position and zero angular velocity
3. The response due to a unit step input
4. The response due to the initial condition in part 2 together with the input in part 3

Example 7.7—cont'd**Solution**

1. The state-transition matrix is the matrix exponential given by the inverse Laplace transform of the matrix

$$\begin{aligned}
 [sl_3 - A]^{-1} &= \begin{bmatrix} s & -1 & 0 \\ 0 & s & -1 \\ 0 & 10 & s+11 \end{bmatrix}^{-1} \\
 &= \frac{\begin{bmatrix} (s+1)(s+10) & s+11 & 1 \\ 0 & s(s+11) & s \\ 0 & -10s & s^2 \end{bmatrix}}{s(s+1)(s+10)} \\
 &= \begin{bmatrix} \frac{1}{s} & \frac{s+11}{s(s+1)(s+10)} & \frac{1}{s(s+1)(s+10)} \\ 0 & \frac{s+11}{(s+1)(s+10)} & \frac{1}{(s+1)(s+10)} \\ 0 & \frac{-10}{(s+1)(s+10)} & \frac{s}{(s+1)(s+10)} \end{bmatrix} \\
 &= \begin{bmatrix} \frac{1}{s} & \frac{1}{90} \left(\frac{99}{s} - \frac{100}{s+1} + \frac{1}{s+10} \right) & \frac{1}{90} \left(\frac{9}{s} - \frac{10}{s+1} + \frac{1}{s+10} \right) \\ 0 & \frac{1}{9} \left(\frac{10}{s+1} - \frac{1}{s+10} \right) & \frac{1}{9} \left(\frac{1}{s+1} - \frac{1}{s+10} \right) \\ 0 & -\frac{10}{9} \left(\frac{1}{s+1} - \frac{1}{s+10} \right) & \frac{1}{9} \left(\frac{10}{s+10} - \frac{1}{s+1} \right) \end{bmatrix}
 \end{aligned}$$

The preceding operations involve writing s in the diagonal entries of a matrix, subtracting entries of the matrix A , then inverting the resulting matrix. The inversion is feasible in this example but becomes progressively more difficult as the order of the system increases. The inverse matrix is obtained by dividing the adjoint matrix by the determinant because numerical matrix inversion algorithms cannot be used in the presence of the complex variable s . Next we inverse Laplace transform to obtain the state-transition matrix

Example 7.7—cont'd

$$e^{At} = \begin{bmatrix} 1 & \frac{1}{90}(99 - 100e^{-t} + e^{-10t}) & \frac{1}{90}(9 - 10e^{-t} + e^{-10t}) \\ 0 & \frac{1}{9}(10e^{-t} - e^{-10t}) & \frac{1}{9}(e^{-t} - e^{-10t}) \\ 0 & -\frac{10}{9}(e^{-t} - e^{-10t}) & \frac{1}{9}(10e^{-10t} - e^{-t}) \end{bmatrix}$$

The state-transition matrix can be decomposed as

$$e^{At} = \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \frac{e^0}{10} + \begin{bmatrix} 0 & -10 & -1 \\ 0 & 10 & 1 \\ 0 & -10 & -1 \end{bmatrix} \frac{e^{-t}}{9} + \begin{bmatrix} 0 & 1 & 1 \\ 0 & -10 & -10 \\ 0 & 100 & 100 \end{bmatrix} \frac{e^{-10t}}{90}$$

This last form reveals that a matrix exponential is nothing more than a matrix-weighted sum of scalar exponentials. The scalar exponentials involve the eigenvalues of the state matrix $\{0, -1, -10\}$ and are known as the **modes** of the system. The matrix weights have rank 1. This general property can be used to check the validity of the matrix exponential.

2. For an initial current of 10 mA and zero initial angular position and velocity, the initial state is

$$\mathbf{x}(0) = [0 \ 0 \ 0.01]^T$$

The zero-input response is

$$\begin{aligned} \mathbf{x}_{ZI}(t) &= e^{At}\mathbf{x}(0) \\ &= \left\{ \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \frac{e^0}{10} + \begin{bmatrix} 0 & -10 & -1 \\ 0 & 10 & 1 \\ 0 & -10 & -1 \end{bmatrix} \frac{e^{-t}}{9} + \begin{bmatrix} 0 & 1 & 1 \\ 0 & -10 & -10 \\ 0 & 100 & 100 \end{bmatrix} \frac{e^{-10t}}{90} \right\} \begin{bmatrix} 0 \\ 0 \\ 1/100 \end{bmatrix} \\ &= \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \frac{e^0}{1000} + \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} \frac{e^{-t}}{900} + \begin{bmatrix} 1 \\ -10 \\ 100 \end{bmatrix} \frac{e^{-10t}}{9000} \end{aligned}$$

By virtue of the decomposition of the state-transition matrix in part 1, the result is obtained using multiplication by constant matrices rather than ones with exponential entries.

Example 7.7—cont'd

3. The response due to a step input is easily evaluated starting in the s -domain to avoid the convolution integral. To simplify the matrix operations, the resolvent matrix is decomposed as

$$\begin{aligned}[sl_3 - A]^{-1} &= \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \frac{1}{10s} + \begin{bmatrix} 0 & -10 & -1 \\ 0 & 10 & 1 \\ 0 & -10 & 1 \end{bmatrix} \frac{1}{9(s+1)} \\ &\quad + \begin{bmatrix} 0 & 1 & 1 \\ 0 & -10 & -10 \\ 0 & 100 & 100 \end{bmatrix} \frac{1}{90(s+10)}\end{aligned}$$

The Laplace transform of the zero-state response is

$$\begin{aligned}\mathbf{X}_{zs}(s) &= [sl_3 - A]^{-1} B \mathbf{U}(s) \\ &= \left\{ \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \frac{1}{10s} + \begin{bmatrix} 0 & -10 & -1 \\ 0 & 10 & 1 \\ 0 & -10 & 1 \end{bmatrix} \frac{1}{9(s+1)} \right. \\ &\quad \left. + \begin{bmatrix} 0 & 1 & 1 \\ 0 & -10 & -10 \\ 0 & 100 & 100 \end{bmatrix} \frac{1}{90(s+10)} \right\} \begin{bmatrix} 0 \\ 0 \\ 10 \end{bmatrix} \frac{1}{s} \\ &= \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \frac{1}{s^2} + \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} \frac{(10/9)}{s(s+1)} + \begin{bmatrix} 1 \\ -10 \\ 100 \end{bmatrix} \frac{(1/9)}{s(s+10)} \\ &= \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \frac{1}{s^2} + \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} (10/9) \left[\frac{1}{s} - \frac{1}{s+1} \right] + \begin{bmatrix} 1 \\ -10 \\ 100 \end{bmatrix} (1/90) \left[\frac{1}{s} - \frac{1}{s+10} \right]\end{aligned}$$

Inverse Laplace transforming, we obtain the solution

Example 7.7—cont'd

$$\begin{aligned}\mathbf{x}_{ZS}(t) &= \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} t + \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} (10/9)[1 - e^{-t}] + \begin{bmatrix} 1 \\ -10 \\ 100 \end{bmatrix} (1/90)[1 - e^{-10t}] \\ &= \begin{bmatrix} -11 \\ 10 \\ 0 \end{bmatrix} (1/10) + \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} t + \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} (10/9)e^{-t} + \begin{bmatrix} -1 \\ 10 \\ -100 \end{bmatrix} (1/90)e^{-10t}\end{aligned}$$

4. The complete solution due to the initial conditions of part 2 and the unit step input of part 3 is simply the sum of the two responses obtained earlier. Hence,

$$\begin{aligned}\mathbf{x}(t) &= \mathbf{x}_{ZI}(t) + \mathbf{x}_{ZS}(t) \\ &= \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \frac{1}{1000} + \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} \frac{e^{-t}}{900} + \begin{bmatrix} 1 \\ -10 \\ 100 \end{bmatrix} \frac{e^{-10t}}{9000} + \begin{bmatrix} -11 \\ 10 \\ 0 \end{bmatrix} (1/10) \\ &\quad + \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} t + \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} (10/9)e^{-t} + \begin{bmatrix} -1 \\ 10 \\ -100 \end{bmatrix} (1/90)e^{-10t} \\ &= \begin{bmatrix} -1.099 \\ 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} 1.11e^{-t} + \begin{bmatrix} -1 \\ 10 \\ -100 \end{bmatrix} 1.1 \times 10^{-2}e^{-10t} + \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} t\end{aligned}$$

7.4.1 The Leverrier algorithm

The calculation of the resolvent matrix $[sI - A]^{-1}$ is clearly the bottleneck in the solution of state-space equations by Laplace transformation. The **Leverrier algorithm** is a convenient method to perform this calculation, which can be programmed on the current generation of handheld calculators that are capable of performing matrix arithmetic. The derivation of the algorithm is left as an exercise (see Problem 7.8).

We first write the resolvent matrix as

$$[sI_n - A]^{-1} = \frac{\text{adj}[sI_n - A]}{\det[sI_n - A]} \quad (7.29)$$

where $\text{adj}[\cdot]$ denotes the adjoint matrix and $\det[\cdot]$ denotes the determinant. Then we observe that, because its entries are determinants of matrices of order $n-1$, the highest

power of s in the adjoint matrix is $n-1$. Terms of the same powers of s can be separated with matrices of their coefficients, and the resolvent matrix can be expanded as

$$[sI_n - A]^{-1} = \frac{P_0 + P_1s + \cdots + P_{n-1}s^{n-1}}{a_0 + a_1s + \cdots + a_{n-1}s^{n-1} + s^n} \quad (7.30)$$

where P_i , $i = 1, 2, \dots, n-1$ are $n \times n$ constant matrices and a_i , $i = 1, 2, \dots, n-1$ are constant coefficients. The coefficients and matrices are calculated as follows.

7.4.1.1 Leverrier algorithm

1. Initialization: $k = n-1$

$$P_{n-1} = I_n \quad a_{n-1} = -\text{tr}\{A\} = -\sum_{i=1}^n a_{ii}$$

where $\text{tr}\{\cdot\}$ denotes the trace of the matrix.

2. Backward iteration: $k = n-2, \dots, 0$

$$P_k = P_{k+1}A + a_{k+1}I_n \quad a_k = -\frac{1}{n-k} \text{tr}\{P_k A\}$$

3. Check:

$$[0] = P_0A + a_0I_n$$

The algorithm requires matrix multiplication, matrix scalar multiplication, matrix addition, and trace evaluation. These are relatively simple operations available in many handheld calculators, with the exception of the trace operation. However, the trace operation can be easily programmed using a single repetition loop. The initialization of the algorithm is simple, and the backward iteration starts with the formulas

$$P_{n-2} = A + a_{n-1}I_n \quad a_{n-2} = -\frac{1}{2} \text{tr}\{P_{n-2}A\}$$

The Leverrier algorithm yields a form of the resolvent matrix that cannot be inverse Laplace transformed directly to obtain the matrix exponential. It is first necessary to expand the following s -domain functions into the partial fractions:

$$\begin{aligned} \frac{s^{n-1}}{\prod_{i=1}^n (s - \lambda_i)} &= \sum_{i=1}^n \frac{q_{i,n-1}}{s - \lambda_i} \cdot \frac{s^{n-2}}{\prod_{i=1}^n (s - \lambda_i)} \\ &= \sum_{i=1}^n \frac{q_{i,n-2}}{s - \lambda_i}, \dots, \frac{1}{\prod_{i=1}^n (s - \lambda_i)} = \sum_{i=1}^n \frac{q_{i,0}}{s - \lambda_i} \end{aligned} \quad (7.31)$$

where λ_i , $i = 1, \dots, n$ are the eigenvalues of the state matrix A defined by

$$\det[sI_n - A] = s^n + a_{n-1}s^{n-1} + \dots + a_1s + a_0 = (s - \lambda_1)(s - \lambda_2)\dots(s - \lambda_n) \quad (7.32)$$

To simplify the analysis, we assume that the eigenvalues in Eq. (7.32) are distinct (i.e., $\lambda_i \neq \lambda_j$, $i \neq j$). The repeated eigenvalue case can be handled similarly but requires higher-order terms in the partial fraction expansion.

Substituting in Eq. (7.30) and combining similar terms gives

$$\begin{aligned} [sI_n - A]^{-1} &= P_{n-1} \sum_{i=1}^n \frac{q_{i,n-1}}{s - \lambda_i} + \dots + P_1 \sum_{i=1}^n \frac{q_{i,1}}{s - \lambda_i} + P_0 \sum_{i=1}^n \frac{q_{i,0}}{s - \lambda_i} \\ &= \frac{1}{s - \lambda_1} \sum_{i=0}^{n-1} q_{1,j} P_j + \frac{1}{s - \lambda_2} \sum_{i=0}^{n-1} q_{2,j} P_j + \dots + \frac{1}{s - \lambda_n} \sum_{i=0}^{n-1} q_{n,j} P_j \\ &= \sum_{i=1}^n \frac{1}{s - \lambda_i} \left[\sum_{j=0}^{n-1} q_{i,j} P_j \right] \end{aligned}$$

Thus, we can write the resolvent matrix in the simple form

$$[sI_n - A]^{-1} = \sum_{i=1}^n \frac{1}{s - \lambda_i} Z_i \quad (7.33)$$

where the matrices Z_i , $i = 1, 2, \dots, n$ are given by

$$Z_i = \left[\sum_{j=0}^{n-1} q_{i,j} P_j \right] \quad (7.34)$$

Finally, we inverse the Laplace transform to obtain the matrix exponential

$$e^{At} = \sum_{i=1}^n Z_i e^{\lambda_i t} \quad (7.35)$$

This is the general form of the expansion used to simplify the computation in Example 7.7 and is the form we use throughout this chapter.

Example 7.8

Use the Leverrier algorithm to compute the matrix exponential for the state matrix of Example 7.7:

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -10 & -11 \end{bmatrix}$$

Example 7.8—cont'd**Solution**

1. Initialization:

$$P_2 = I_3 \quad a_2 = -\text{tr}\{A\} = 11$$

2. Backward iteration: $k = 1, 0$ a. $k = 1$

$$P_1 = A + a_2 I_3 = A + 11I_3 = \begin{bmatrix} 11 & 1 & 0 \\ 0 & 11 & 1 \\ 0 & -10 & 0 \end{bmatrix}$$

$$\begin{aligned} a_1 &= -\frac{1}{2}\text{tr}\{P_1 A\} = -\frac{1}{2}\text{tr}\left\{\begin{bmatrix} 11 & 1 & 0 \\ 0 & 11 & 1 \\ 0 & -10 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -10 & -11 \end{bmatrix}\right\} \\ &= -\frac{1}{2}\text{tr}\left\{\begin{bmatrix} 0 & 11 & 1 \\ 0 & -10 & 0 \\ 0 & 0 & -10 \end{bmatrix}\right\} = 10 \end{aligned}$$

b. $k = 0$

$$P_0 = P_1 A + a_1 I_3 = \begin{bmatrix} 11 & 1 & 0 \\ 0 & 11 & 1 \\ 0 & -10 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -10 & -11 \end{bmatrix} + \begin{bmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 10 \end{bmatrix} = \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$a_0 = -\frac{1}{3}\text{tr}\{P_0 A\} = -\frac{1}{3}\text{tr}\left\{\begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -10 & -11 \end{bmatrix}\right\} = -\frac{1}{3}\text{tr}\left\{\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}\right\} = 0$$

3. Check:

$$[0] = P_0 A + a_0 I_n = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -10 & -11 \end{bmatrix} \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} = [0]$$

Thus, we have

$$[sI_n - A]^{-1} = \frac{\begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 11 & 1 & 0 \\ 0 & 11 & 1 \\ 0 & -10 & 0 \end{bmatrix} s + \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} s^2}{10s + 11s^2 + s^3}$$

Example 7.8—cont'd

The characteristic polynomial of the system is

$$s^3 + 11s^2 + 10s = s(s+1)(s+10)$$

and the system eigenvalues are $\{0, -1, -10\}$. Next, we obtain the partial fraction expansions

$$\frac{s^2}{\prod_{i=1}^3 (s - \lambda_i)} = \sum_{i=1}^3 \frac{q_{i,2}}{s - \lambda_i} = \frac{0}{s} + \frac{(-1/9)}{s+1} + \frac{(10/9)}{s+10}$$

$$\frac{s}{\prod_{i=1}^3 (s - \lambda_i)} = \sum_{i=1}^3 \frac{q_{i,1}}{s - \lambda_i} = \frac{0}{s} + \frac{(1/9)}{s+1} + \frac{(-1/9)}{s+10}$$

$$\frac{1}{\prod_{i=1}^3 (s - \lambda_i)} = \sum_{i=1}^3 \frac{q_{i,0}}{s - \lambda_i} = \frac{(1/10)}{s} + \frac{(-1/9)}{s+1} + \frac{(1/90)}{s+10}$$

where some of the coefficients are zero due to cancellation. This allows us to evaluate the matrices

$$Z_1 = (1/10) \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$Z_2 = (-1/9) \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} + (1/9) \begin{bmatrix} 11 & 1 & 0 \\ 0 & 11 & 1 \\ 0 & -10 & 0 \end{bmatrix} + (-1/9) \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = (1/9) \begin{bmatrix} 0 & -10 & -1 \\ 0 & 10 & 1 \\ 0 & -10 & -1 \end{bmatrix}$$

$$\begin{aligned} Z_3 &= (1/90) \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} + (-1/9) \begin{bmatrix} 11 & 1 & 0 \\ 0 & 11 & 1 \\ 0 & -10 & 0 \end{bmatrix} + (10/9) \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \\ &= (1/90) \begin{bmatrix} 0 & 1 & 1 \\ 0 & -10 & -10 \\ 0 & 100 & 100 \end{bmatrix} \end{aligned}$$

Therefore, the state-transition matrix is

$$e^{At} = \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \frac{e^0}{10} + \begin{bmatrix} 0 & -10 & -1 \\ 0 & 10 & 1 \\ 0 & -10 & -1 \end{bmatrix} \frac{e^{-t}}{9} + \begin{bmatrix} 0 & 1 & 1 \\ 0 & -10 & -10 \\ 0 & 100 & 100 \end{bmatrix} \frac{e^{-10t}}{90}$$

Example 7.8—cont'd

Thus, we obtain the answer of Example 7.7 with far fewer partial fraction expansions but with some additional operations with constant matrices. Because these operations are easily performed using a calculator, this is a small price to pay for the resulting simplification.

7.4.2 Sylvester's expansion

The matrices Z_i , $i = 1, 2, \dots, n$, obtained in [Section 7.4.1](#) using the Leverrier algorithm, are known as the **constituent matrices** of A . The constituent matrices can also be calculated using Sylvester's formula as follows:

$$Z_i = \frac{\prod_{\substack{j=1 \\ j \neq i}}^n [A - \lambda_j I_n]}{\prod_{\substack{j=1 \\ j \neq i}}^n [\lambda_i - \lambda_j]} \quad (7.36)$$

where λ_i , $i = 1, 2, \dots, n$ are the eigenvalues of the matrix A .

Numerical computation of the matrix exponential using [Eq. \(7.36\)](#) can be problematic. For example, if two eigenvalues are almost equal, the scalar denominator in the equation is small, resulting in large computational errors. In fact, the numerical computation of the matrix exponential is not as simple as our presentation may suggest, with all known computational procedures failing in special cases. These issues are discussed in more detail in a well-known paper by [Moler and Van Loan \(1978\)](#).

Example 7.9

Calculate constituent matrices of the matrix A given in Examples 7.7 and 7.8 using Sylvester's formula.

Solution

$$Z_1 = \frac{\prod_{\substack{j=1 \\ j \neq 1}}^3 [A - \lambda_j I_3]}{\prod_{\substack{j=1 \\ j \neq 1}}^3 [0 - \lambda_j]} = \frac{(A + I_3)(A + 10I_3)}{(1)(10)} = (1/10) \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Example 7.9—cont'd

$$Z_2 = \frac{\prod_{\substack{j=1 \\ j \neq i}}^3 [A - \lambda_j I_n]}{\prod_{\substack{j=1 \\ j \neq i}}^3 [-1 - \lambda_j]} = \frac{A(A + 10I_n)}{(-1)(-1 + 10)} = (1/9) \begin{bmatrix} 0 & -10 & -1 \\ 0 & 10 & 1 \\ 0 & -10 & -1 \end{bmatrix}$$

$$Z_3 = \frac{\prod_{\substack{j=1 \\ j \neq i}}^3 [A - \lambda_j I_n]}{\prod_{\substack{j=1 \\ j \neq i}}^3 [-10 - \lambda_j]} = \frac{A(A + I_n)}{(-10)(-10 + 1)} = (1/90) \begin{bmatrix} 0 & 1 & 1 \\ 0 & -10 & -10 \\ 0 & 100 & 100 \end{bmatrix}$$

The constituent matrices can be used to define any analytic function of a matrix (i.e., a function possessing a Taylor series) and not just the matrix exponential. The following identity is true for any analytic function $f(\lambda)$:

$$f(A) = \sum_{i=1}^n Z_i f(\lambda_i) \quad (7.37)$$

This identity allows us to calculate the functions e^{At} and A^i , among others. Using the matrices Z_i described in Example 7.9 and Eq. (7.37), we obtain the state-transition matrix described in Example 7.7.

7.4.3 The state-transition matrix for a diagonal state matrix

For a state matrix Λ in the form

$$\Lambda = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\} \quad (7.38)$$

the resolvent matrix is

$$[sI_n - \Lambda]^{-1} = \text{diag}\left\{\frac{1}{s - \lambda_1}, \frac{1}{s - \lambda_2}, \dots, \frac{1}{s - \lambda_n}\right\} \quad (7.39)$$

The corresponding state-transition matrix is

$$\begin{aligned} e^{\Lambda t} &= \mathcal{L}\left\{[sI_n - \Lambda]^{-1}\right\} \\ &= \text{diag}\{e^{\lambda_1 t}, e^{\lambda_2 t}, \dots, e^{\lambda_n t}\} \\ &= \text{diag}\{1, 0, \dots, 0\}e^{\lambda_1 t} + \text{diag}\{0, 1, \dots, 0\}e^{\lambda_2 t} + \dots + \text{diag}\{0, 0, \dots, n\}e^{\lambda_n t} \end{aligned} \quad (7.40)$$

Thus, the i^{th} constituent matrix for the diagonal form is a diagonal matrix with unity entry (i, i) and all other entries equal to zero.

Example 7.10

Calculate the state-transition matrix if the state matrix A is

$$A = \text{diag}\{-1, -5, -6, -20\}$$

Solution

Using Eq. (7.40), we obtain the state-transition matrix

$$e^{At} = \text{diag}\{e^{-t}, e^{-5t}, e^{-6t}, e^{-20t}\}$$

Assuming distinct eigenvalues, the state matrix can in general be written in the form

$$A = V\Lambda V^{-1} = V\Lambda W \quad (7.41)$$

where

$$V = \begin{bmatrix} \mathbf{v}_1 & | & \mathbf{v}_2 & | & \cdots & | & \mathbf{v}_n \end{bmatrix} \quad W = \begin{bmatrix} \mathbf{w}_1^T \\ \vdots \\ \mathbf{w}_2^T \\ \vdots \\ \vdots \\ \mathbf{w}_n^T \end{bmatrix}$$

$\mathbf{v}_i, \mathbf{w}_i, i = 1, 2, \dots, n$ are the right and left eigenvectors of the matrix A , respectively. The fact that W is the inverse of V implies that their product is the identity matrix—that is,

$$WV = [w_i^T v_j] = I_n$$

Equating matrix entries gives

$$w_i^T v_j = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases} \quad (7.42)$$

The right eigenvectors are the usual eigenvectors of the matrix, whereas the left eigenvectors can be shown to be the eigenvectors of the matrix transpose.

The matrix A raised to any integer power i can be written in the form

$$A^i = (V\Lambda W)(V\Lambda W)\dots(V\Lambda W) = V\Lambda^i W \quad (7.43)$$

Substituting from Eq. (7.43) in the matrix exponential series expansion gives

$$\begin{aligned}
 e^{At} &= \sum_{i=0}^{\infty} \frac{(At)^i}{i!} \\
 &= \sum_{i=0}^{\infty} \frac{V(\Lambda t)^i W}{i!} \\
 &= V \left[\sum_{i=0}^{\infty} \frac{\text{diag}\left\{(\lambda_1 t)^i, (\lambda_2 t)^i, \dots, (\lambda_n t)^i\right\}}{i!} \right] W \\
 &= V \left[\text{diag} \left\{ \sum_{i=0}^{\infty} \frac{(\lambda_1 t)^i}{i!}, \sum_{i=0}^{\infty} \frac{(\lambda_2 t)^i}{i!}, \dots, \sum_{i=0}^{\infty} \frac{(\lambda_n t)^i}{i!}, \right\} \right] W
 \end{aligned}$$

That is,

$$e^{At} = Ve^{At}W \quad (7.44)$$

Thus, we have an expression for the matrix exponential of A in terms of the matrix exponential for the diagonal matrix Λ . The expression provides another method to calculate the state-transition matrix using the **eigenstructure** (eigenvalues and eigenvectors) of the state matrix. The drawback of the method is that the eigenstructure calculation is computationally costly.

Example 7.11

Obtain the constituent matrices of the state matrix using Eq. (7.44), then obtain the state-transition matrix for the system

$$\dot{x}(t) = \begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix} x(t)$$

Solution

The state matrix is in companion form, and its characteristic equation can be directly written with the coefficient obtained by negating the last row of the matrix

$$\lambda^2 + 3\lambda + 2 = (\lambda + 1)(\lambda + 2) = 0$$

The system eigenvalues are therefore $\{-1, -2\}$. The matrix of eigenvectors is the Van der Monde matrix

$$V = \begin{bmatrix} 1 & 1 \\ -1 & -2 \end{bmatrix} \quad W = V^{-1} = \frac{\begin{bmatrix} -2 & -1 \\ 1 & 1 \end{bmatrix}}{-2 + 1} = \begin{bmatrix} 2 & 1 \\ -1 & -1 \end{bmatrix}$$

Example 7.11—cont'd

The matrix exponential is

$$\begin{aligned}
 e^{At} = Ve^{At}W &= \begin{bmatrix} 1 & | & 1 \\ -1 & | & -2 \end{bmatrix} \begin{bmatrix} e^{-t} & 0 \\ \cdots & \cdots \\ 0 & e^{-2t} \end{bmatrix} \begin{bmatrix} 2 & 1 \\ -1 & -1 \end{bmatrix} \\
 &= \left\{ \begin{bmatrix} 1 \\ -1 \end{bmatrix} \begin{bmatrix} e^{-t} & | & 0 \\ -2 & | & \end{bmatrix} + \begin{bmatrix} 1 \\ -2 \end{bmatrix} \begin{bmatrix} 0 & | & e^{-2t} \\ -1 & | & \end{bmatrix} \right\} \begin{bmatrix} 2 & 1 \\ \cdots & \cdots \\ -1 & -1 \end{bmatrix} \\
 &= \begin{bmatrix} 1 \\ -1 \end{bmatrix} [2 \ 1] e^{-t} + \begin{bmatrix} 1 \\ -2 \end{bmatrix} [-1 \ -1] e^{-2t} \\
 &= \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} e^{-t} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} e^{-2t}
 \end{aligned}$$

As one may infer from the partitioned matrices used in Example 7.11, Eq. (7.44) can be used to write the state-transition matrix in terms of the constituent matrices. We first obtain the product

$$\begin{aligned}
 e^{At}W &= \begin{bmatrix} e^{\lambda_1 t} & | & 0 & | & \cdots & | & 0 \\ 0 & | & e^{\lambda_2 t} & | & \cdots & | & 0 \\ \vdots & | & \vdots & | & \ddots & | & \vdots \\ 0 & | & 0 & | & \cdots & | & e^{\lambda_n t} \end{bmatrix} \begin{bmatrix} \mathbf{w}_1^T \\ \vdots \\ \mathbf{w}_2^T \\ \vdots \\ \vdots \\ \vdots \\ \mathbf{w}_n^T \end{bmatrix} \\
 &= \begin{bmatrix} \mathbf{w}_1^T \\ \vdots \\ \mathbf{0}_{1 \times n} \\ \vdots \\ \mathbf{0}_{1 \times n} \end{bmatrix} e^{\lambda_1 t} + \begin{bmatrix} \mathbf{0}_{1 \times n} \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{0}_{1 \times n} \end{bmatrix} e^{\lambda_2 t} + \cdots + \begin{bmatrix} \mathbf{0}_{1 \times n} \\ \vdots \\ \mathbf{0}_{1 \times n} \\ \mathbf{w}_n^T \end{bmatrix} e^{\lambda_n t}
 \end{aligned}$$

Then we premultiply by the partition matrix V to obtain

$$Ve^{At}W = \left[\mathbf{v}_1 \mid \mathbf{v}_2 \mid \cdots \mid \mathbf{v}_n \right] \left\{ \begin{bmatrix} \mathbf{w}_1^T \\ \vdots \\ \mathbf{0}_{1 \times n} \\ \vdots \\ \mathbf{0}_{1 \times n} \end{bmatrix} e^{\lambda_1 t} + \begin{bmatrix} \mathbf{0}_{1 \times n} \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{0}_{1 \times n} \end{bmatrix} e^{\lambda_2 t} + \cdots + \begin{bmatrix} \mathbf{0}_{1 \times n} \\ \vdots \\ \mathbf{0}_{1 \times n} \\ \mathbf{w}_n^T \end{bmatrix} e^{\lambda_n t} \right\} \quad (7.45)$$

Substituting from Eq. (7.45) into Eq. (7.44) yields

$$e^{At} = \sum_{i=1}^n Z_i e^{\lambda_i t} = \sum_{i=1}^n \mathbf{v}_i \mathbf{w}_i^T e^{\lambda_i t} \quad (7.46)$$

Hence, the i^{th} constituent matrix of A is given by the product of the i^{th} right and left eigenvectors. The following properties of the constituent matrix can be proved using Eq. (7.46). The proofs are left as an exercise.

7.4.3.1 Properties of constituent matrices

1. Constituent matrices have rank 1.
2. The product of two constituent matrices is

$$Z_i Z_j = \begin{cases} Z_i, & i = j \\ 0, & i \neq j \end{cases}$$

Raising Z_i to any power gives the matrix Z_i . Z_i is said to be *idempotent*.

3. The sum of the n constituent matrices of an $n \times n$ matrix is equal to the identity matrix

$$\sum_{i=1}^n Z_i = I_n$$

Example 7.12

Obtain the constituent matrices of the state matrix of Example 7.11 using Eq. (7.46), and verify that they satisfy Properties 1 through 3. Then obtain the state-transition matrix for the system.

Solution

The constituent matrices are

$$Z_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix} [2 \quad 1] = \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix}$$

$$Z_2 = \begin{bmatrix} 1 \\ -2 \end{bmatrix} [-1 \quad -1] = \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix}$$

Both matrices have a second column equal to the first and clearly have rank 1. The product of the first and second matrices is

$$\begin{aligned} Z_1 Z_2 &= \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} \\ &= \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} = Z_2 Z_1 \end{aligned}$$

The squares of the matrices are

$$Z_1 Z_1 = \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix}$$

$$Z_2 Z_2 = \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix}$$

Example 7.12—cont'd

The state-transition matrix is

$$e^{At} = \sum_{i=1}^2 Z_i e^{\lambda_i t} = \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} e^{-t} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} e^{-2t}$$

Example 7.13

Obtain the state-transition matrix for the system with state matrix

$$A = \begin{bmatrix} 0 & 1 & 3 \\ 3 & 5 & 6 \\ 5 & 6 & 7 \end{bmatrix}$$

Solution

Using the MATLAB command **eig**, we obtain the matrices

$$V = \begin{bmatrix} 0.8283 & -0.2159 & 0.5013 \\ 0.1173 & -0.6195 & -0.7931 \\ -0.5479 & -0.7547 & 0.3459 \end{bmatrix} \quad W = \begin{bmatrix} 0.8355 & 0.3121 & -0.4952 \\ -0.4050 & -0.5768 & -0.7357 \\ 0.4399 & -0.7641 & 0.5014 \end{bmatrix}$$

and the eigenvalues $\{-1.8429, 13.3554, 0.4875\}$. Then multiplying each column of V by the corresponding row of W , we obtain the constituent matrices

$$Z_1 = \begin{bmatrix} 0.6920 & 0.2585 & -0.4102 \\ 0.0980 & 0.0366 & -0.0581 \\ -0.4578 & -0.1710 & 0.2713 \end{bmatrix} \quad Z_2 = \begin{bmatrix} 0.0874 & 0.1245 & 0.1588 \\ 0.2509 & 0.3573 & 0.4558 \\ 0.3057 & 0.4353 & 0.5552 \end{bmatrix}$$

$$Z_3 = \begin{bmatrix} 0.2205 & -0.3831 & 0.2513 \\ -0.3489 & 0.6061 & -0.3977 \\ -0.1521 & -0.2643 & 0.1734 \end{bmatrix}$$

The state-transition matrix is

$$e^{At} = \begin{bmatrix} 0.6920 & 0.2585 & -0.4102 \\ 0.0980 & 0.0366 & -0.0581 \\ -0.4578 & -0.1710 & 0.2713 \end{bmatrix} e^{-1.8429t} + \begin{bmatrix} 0.0874 & 0.1245 & 0.1588 \\ 0.2509 & 0.3573 & 0.4558 \\ 0.3057 & 0.4353 & 0.5552 \end{bmatrix} e^{13.3554t}$$

$$+ \begin{bmatrix} 0.2205 & -0.3831 & 0.2513 \\ -0.3489 & 0.6061 & -0.3977 \\ -0.1521 & -0.2643 & 0.1734 \end{bmatrix} e^{0.4875t}$$

7.4.4 Real form for complex conjugate eigenvalues

For the case of complex conjugate eigenvalues of a real matrix, the eigenvectors will also be complex conjugate. This is easily shown by taking the complex conjugate of the equation that defines the right eigenvectors

$$\bar{A}\bar{\mathbf{v}} = \bar{\lambda}\bar{\mathbf{v}} = A\bar{\mathbf{v}} = \bar{\lambda}\bar{\mathbf{v}} \quad (7.47)$$

where we used the fact that the real matrix A is identical to its conjugate. Similarly, for the left eigenvectors

$$\overline{\mathbf{w}^T A} = \overline{\lambda \mathbf{w}^T} = \overline{\mathbf{w}}^T A = \bar{\lambda} \overline{\mathbf{w}}^T \quad (7.48)$$

Thus, the constituent matrix for a complex eigenvalue λ as given by

$$Z = \mathbf{v}\mathbf{w}^T \quad (7.49)$$

is the conjugate of the constituent matrix of its complex conjugate eigenvalue $\bar{\lambda}$

$$\bar{Z} = \overline{\mathbf{v}\mathbf{w}^T} \quad (7.50)$$

Using this fact, we can simplify the expansion of the state-transition matrix with complex eigenvalues as follows:

$$e^{At} = Ze^{\lambda t} + \bar{Z}e^{\bar{\lambda}t} = 2Re\{Ze^{\lambda t}\} \quad (7.51)$$

For the eigenvalue $\lambda = \sigma + j\omega_d$ we have

$$e^{At} = 2Re\{Z\}e^{\sigma T} \cos(\omega_d t) - 2Im\{Z\}e^{\sigma T} \sin(\omega_d t) \quad (7.52)$$

Example 7.14

The RLC resonance circuit of Fig. 7.3 with capacitor voltage as output has the state-space model

$$\dot{\mathbf{x}} = \begin{bmatrix} 0 & 1 \\ -\frac{1}{LC} & -\frac{R}{L} \end{bmatrix} \mathbf{x} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u$$

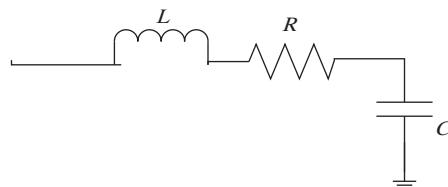


Figure 7.3
Series RLC circuit.

Example 7.14—cont'd

$$y = \begin{bmatrix} \frac{1}{LC} & 0 \end{bmatrix} \mathbf{x}$$

$$\mathbf{x} = [\mathbf{x}_1 \quad \mathbf{x}_2]^T$$

Assuming normalized component values such that $\frac{R}{L} = 2, \frac{1}{LC} = 10$, obtain a real form of the state-transition matrix of the system by combining complex conjugate terms.

Solution

The state matrix is in companion form, and the characteristic equation can be written by inspection:

$$\lambda^2 + \frac{R}{L}\lambda + \frac{1}{LC} = \lambda^2 + 2\lambda + 10 = 0$$

The eigenvalues of the system are

$$\lambda_{1,2} = -1 \pm j3$$

The right eigenvector matrix is the Van der Monde matrix

$$V = \begin{bmatrix} 1 & 1 \\ -1+j3 & -1-3j \end{bmatrix}$$

The left eigenvector matrix is

$$W = V^{-1} = \frac{1}{6} \begin{bmatrix} 3-j & -j \\ 3+j & j \end{bmatrix}$$

The constituent matrix for the first eigenvalue is

$$Z = \mathbf{v}\mathbf{w}^T = \begin{bmatrix} 1 \\ -1+j3 \end{bmatrix} \frac{1}{6} \begin{bmatrix} 3-j & -j \\ 3+j & j \end{bmatrix} = \frac{1}{6} \begin{bmatrix} 3-j & -j \\ 10j & 3+j \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \frac{j}{6} \begin{bmatrix} -1 & -1 \\ 10 & 1 \end{bmatrix}$$

The state transition matrix can be written as

$$e^{At} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} e^{-t} \cos(3t) - \frac{1}{3} \begin{bmatrix} -1 & -1 \\ 10 & 1 \end{bmatrix} e^{-t} \sin(3t)$$

7.5 The transfer function matrix

The transfer function matrix of a system can be derived from its state and output equations. We begin by Laplace transforming the state Eq. (7.3) with zero initial conditions to obtain

$$\mathbf{X}(s) = [sI_n - A]^{-1} B \mathbf{U}(s) \quad (7.53)$$

Then we Laplace transform the output equation and substitute from Eq. (7.53) to get

$$\mathbf{Y}(s) = C[sI_n - A]^{-1}B\mathbf{U}(s) + D\mathbf{U}(s)$$

The last equation can be rewritten in the form

$$\begin{aligned}\mathbf{Y}(s) &= H(s)\mathbf{U}(s) \xrightarrow{\mathcal{L}} \mathbf{y}(t) = H(t) * \mathbf{u}(t) \\ H(s) &= C[sI_n - A]^{-1}B + D \xrightarrow{\mathcal{L}} H(t) = Ce^{At}B + D\delta(t)\end{aligned}\quad (7.54)$$

where $H(s)$ is the transfer function matrix and $H(t)$ is the impulse response matrix. The equations emphasize the fact that the transfer function and the impulse response are Laplace transform pairs.

The preceding equation cannot be simplified further in the MIMO case because division by a vector $\mathbf{U}(s)$ is not defined. In the SI case, the transfer function can be expressed as a vector of ratios of outputs to inputs. This reduces to a single scalar ratio in the SISO case. In general, the i^{th} entry of the transfer function matrix denotes

$$h_{ij}(s) = \left. \frac{Y_i(s)}{U_j(s)} \right|_{\substack{U_l = 0, l \neq j \\ \text{zero initial conditions}}} \quad (7.55)$$

The transfer function can be rewritten in terms of the constituent matrices of A as

$$H(s) = C \left[\sum_{i=1}^n Z_i \frac{1}{s - \lambda_i} \right] B + D$$

Thus we have,

$$H(s) = \sum_{i=1}^n CZ_i B \frac{1}{s - \lambda_i} + D \xrightarrow{\mathcal{L}} H(t) = \sum_{i=1}^n CZ_i Be^{\lambda_i t} + D\delta(t) \quad (7.56)$$

This shows that the poles of the transfer function are the eigenvalues of the state matrix A . In some cases, however, one or both of the matrix products CZ_i , Z_iB are zero, and the eigenvalues do not appear in the reduced transfer function.

The evaluation of Eq. (7.54) requires the computation of the resolvent matrix and can be performed using the Leverrier algorithm or Sylvester's formula. Nevertheless, this entails considerable effort. Therefore, we should only use Eq. (7.54) to evaluate the transfer function if the state-space equations are given and the input-output differential equations are not. It is usually simpler to obtain the transfer function by Laplace transforming the input-output differential equation.

Example 7.15

Calculate the transfer function of the system of Example 7.7 with the angular position as output.

Solution

$$\begin{aligned}
 H(s) &= [1 \ 0 \ 0] \left\{ \begin{array}{l} \left[\begin{matrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \right] \frac{1}{10s} + \left[\begin{matrix} 0 & -10 & -1 \\ 0 & 10 & 1 \\ 0 & -10 & -1 \end{matrix} \right] \frac{1}{9(s+1)} + \\ \left[\begin{matrix} 0 & 1 & 1 \\ 0 & -10 & -10 \\ 0 & 100 & 100 \end{matrix} \right] \frac{1}{90(s+10)} \end{array} \right\} \begin{bmatrix} 0 \\ 0 \\ 10 \end{bmatrix} \\
 &= \frac{1}{s} - \frac{10}{9(s+1)} + \frac{1}{9(s+10)} \\
 &= \frac{10}{s(s+1)(s+10)}
 \end{aligned}$$

7.5.1 MATLAB commands

MATLAB obtains the transfer function for the matrices (A, B, C, D) with the commands

`>> g = tf(ss(A, B, C, D))`

For example, the matrices

$$\begin{aligned}
 A &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 1 \\ -3 & -4 & -2 \end{bmatrix} & B &= \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \\
 C &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix} & D &= \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}
 \end{aligned}$$

are entered as

```

>> A = [zeros(2,1), [1, 0; 1, 1]; -3, -4, -2];
>> B = [zeros(1,2); eye(2)];
>> C=zeros(2,1), [1, 0; 1, 1];
>> D = [zeros(2,1), ones(2,1)]

```

Then the transfer function for the first input is obtained with the transformation command

$\gg \mathbf{g} = \text{tf}(\text{ss}(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}))$

Transfer function from input 1 to output.

$s^2 + 2s$

#1: - - - - -

$s^3 + s^2 + 2s + 3$

$s^2 - 2s - 3$

#2: - - - - -

$s^3 + s^2 + 2s + 3$

Transfer function from input 2 to output

$s^3 + s^2 + 3s + 3$

#1: - - - - -

$s^3 + s^2 + 2s + 3$

$s^3 + 2s^2 + 2s + 3$

#2: - - - - -

$s^3 + s^2 + 2s + 3$

The first two terms are the first column of the transfer function matrix

$$H_1(s) = \frac{\begin{bmatrix} s^2 + 2s \\ s^2 - 2s - 3 \end{bmatrix}}{s^3 + s^2 + 2s + 3}$$

The next two terms are transfer function column corresponding to the second input

$$H_2(s) = \frac{\begin{bmatrix} s^3 + s^2 + 3s + 3 \\ s^3 + 2s^2 + 2s + 3 \end{bmatrix}}{s^3 + s^2 + 2s + 3}$$

The **tf** command can also be used to obtain the resolvent matrix by setting $B = C = I_n$ with D zero. The command takes the form

$\gg \mathbf{Resolvent} = \text{tf}(\text{ss}(\mathbf{A}, \text{eye}(\mathbf{n}), \text{eye}(\mathbf{n}), \mathbf{0}))$

To create a transfer function matrix, we define its numerator and denominator then use the command **tf** as in the case of a scalar transfer function. The numerator is created as a cell array with each entry including a vector of coefficient from the corresponding numerator polynomial. The denominator is created similarly using denominator polynomials. For example, to create the transfer function we obtained in our example.

```
>> num = {[1 2 0],[1 1 3 3];[1 -1 3],[1 2 2 3]};  
>> den = {[1 1 2 3],[1 1 2 3];[1 1 2 3],[1 1 2 3]};  
>> H = tf(num,den).
```

7.6 Discrete-time state-space equations

Given an analog system with piecewise constant inputs over a given sampling period, the system state variables at the end of each period can be related by a difference equation. The difference equation is obtained by examining the solution of the analog state equation derived in [Section 7.4](#), over a sampling period T . The solution is given by [Eq. \(7.27\)](#), which is repeated here for convenience:

$$\mathbf{x}(t_f) = e^{A(t_f-t_0)} \mathbf{x}(t_0) + \int_{t_0}^{t_f} e^{A(t_f-\tau)} B \mathbf{u}(\tau) d\tau \quad (7.57)$$

Let the initial time $t_0 = kT$ and the final time $t_f = (k+1)T$. Then [Eq. \(7.57\)](#) reduces to

$$\mathbf{x}(k+1) = e^{AT} \mathbf{x}(k) + \int_{kT}^{(k+1)T} e^{A((k+1)T-\tau)} B \mathbf{u}(\tau) d\tau \quad (7.58)$$

where $\mathbf{x}(k)$ denotes the vector at time kT . For a piecewise constant input

$$\mathbf{u}(t) = \mathbf{u}(k), \quad kT \leq t < (k+1)T$$

the input can be moved outside the integral. The integrand can then be simplified by changing the variable of integration to

$$\begin{aligned}\lambda &= (k+1)T - \tau \\ d\lambda &= -d\tau\end{aligned}$$

The integral now becomes

$$\left\{ \int_T^0 e^{A\lambda} B(-d\lambda) \right\} \mathbf{u}(k) = \left\{ \int_0^T e^{A\lambda} B d\lambda \right\} \mathbf{u}(k)$$

Substituting in [Eq. \(7.58\)](#), we obtain the discrete-time state equation

$$\mathbf{x}(k+1) = A_d \mathbf{x}(k) + B_d \mathbf{u}(k) \quad (7.59)$$

where

$$A_d = e^{AT} \quad (7.60)$$

$$B_d = \int_0^T e^{A\lambda} Bd\lambda \quad (7.61)$$

A_d is the discrete-time state matrix and B_d is the discrete input matrix, and they are clearly of the same orders as their continuous counterparts. The discrete-time state matrix is the state-transition matrix for the analog system evaluated at the sampling period T .

Eqs. (7.60) and (7.61) can be simplified further using properties of the matrix exponential. For invertible state matrix A , the integral of the matrix exponential is

$$\int e^{At} dt = A^{-1} [e^{AT} - I_n] = [e^{AT} - I_n] A^{-1} \quad (7.62)$$

This allows us to write B_d in the form

$$B_d = A^{-1} [e^{AT} - I_n] B = [e^{AT} - I_n] A^{-1} B \quad (7.63)$$

Using the expansion of the matrix exponential Eq. (7.35), we rewrite Eqs. (7.60) and (7.61) as

$$A_d = \sum_{i=1}^n Z_i e^{\lambda_i \tau} \quad (7.64)$$

$$B_d = \int_0^T \left(\sum_{i=1}^n Z_i e^{\lambda_i \tau} \right) Bd\tau \quad (7.65)$$

$$= \sum_{i=1}^n Z_i B \int_0^T e^{\lambda_i \tau} d\tau$$

The integrands in Eq. (7.65) are scalar functions, and the integral can be easily evaluated. Because we assume distinct eigenvalues, only one eigenvalue can be zero. Hence, we obtain the following expression for B_d :

$$B_d = \begin{cases} \sum_{i=1}^n Z_i B \left[\frac{1 - e^{\lambda_i T}}{-\lambda_i} \right] & \lambda_i \neq 0 \\ Z_1 B T + \sum_{i=2}^n Z_i B \left[\frac{1 - e^{\lambda_i T}}{-\lambda_i} \right] & \lambda_1 = 0 \end{cases} \quad (7.66)$$

The output equation evaluated at time kT is

$$\mathbf{y}(k) = C\mathbf{x}(k) + D\mathbf{u}(k) \quad (7.67)$$

The discrete-time state-space representation is given by Eqs. (7.59) and (7.67).

Eq. (7.59) is approximately valid for a general input vector $\mathbf{u}(t)$ provided that the sampling period T is sufficiently short. The equation can therefore be used to obtain the solution of the state equation in the general case.

Example 7.16

Obtain the discrete-time state equations for the system of Example 7.7

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -10 & -11 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 10 \end{bmatrix} u$$

for a sampling period $T = 0.01$ s.

Solution

From Example 7.7, the state-transition matrix of the system is

$$e^{At} = \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \frac{e^0}{10} + \begin{bmatrix} 0 & -10 & -1 \\ 0 & 10 & 1 \\ 0 & -10 & -1 \end{bmatrix} \frac{e^{-t}}{9} + \begin{bmatrix} 0 & 1 & 1 \\ 0 & -10 & -10 \\ 0 & 100 & 100 \end{bmatrix} \frac{e^{-10t}}{90}$$

Thus, the discrete-time state matrix is

$$A_d = e^{0.01 \times A} = \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \frac{e^0}{10} + \begin{bmatrix} 0 & -10 & -1 \\ 0 & 10 & 1 \\ 0 & -10 & 1 \end{bmatrix} \frac{e^{-0.01}}{9} + \begin{bmatrix} 0 & 1 & 1 \\ 0 & -10 & -10 \\ 0 & 100 & 100 \end{bmatrix} \frac{e^{-10 \times 0.01}}{90}$$

This simplifies to

$$A_d = \begin{bmatrix} 1.0 & 0.1 & 0.0 \\ 0.0 & 0.9995 & 0.0095 \\ 0.0 & -0.0947 & 0.8954 \end{bmatrix}$$

The discrete-time input matrix is

$$\begin{aligned} B_d &= Z_1 B(0.01) + Z_2 B(1 - e^{-0.01}) + Z_3 B(1 - e^{-10 \times 0.01}) / 10 \\ &= \begin{bmatrix} 0.01 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} (10/9)(1 - e^{-0.01}) + \begin{bmatrix} 1 \\ -10 \\ 100 \end{bmatrix} (1/90)(1 - e^{-10 \times 0.01}) \end{aligned}$$

This simplifies to

$$B_d = \begin{bmatrix} 1.622 \times 10^{-6} \\ 4.821 \times 10^{-4} \\ 9.468 \times 10^{-2} \end{bmatrix}$$

7.6.1 MATLAB commands for discrete-time state-space equations

The MATLAB command to obtain the discrete state-space quadruple **pd** from the continuous quadruple **p** with sampling period $T = 0.1$ is

$$\gg \mathbf{pd} = \mathbf{c2d}(\mathbf{p}, .1)$$

Alternatively, the matrices are obtained using Eqs. (7.60) and (7.63) and the MATLAB commands

$$\begin{aligned}\gg \mathbf{Ad} &= \mathbf{expm}(\mathbf{A} * \mathbf{0.1}) \\ \gg \mathbf{Bd} &= \mathbf{A} (\mathbf{Ad} - \mathbf{eye}(\mathbf{size}(\mathbf{A}))) * \mathbf{B}\end{aligned}$$

7.6.2 Complex conjugate eigenvalues

If an analog system with complex conjugate eigenvalues $\lambda_{1,2} = \sigma \pm j\omega_d$ is discretized with the analog input constant over each sampling period, then the resulting system has the complex conjugate eigenvalues $e^{\lambda_{1,2}T} = e^{\sigma \pm j\omega_d T}$. From Eq. (7.52), the constituent matrices for the discrete state matrix are the same as those for the analog system, and the discrete state matrix is in the form

$$e^{AT} = 2Re\{Z\}e^{\sigma T} \cos(\omega_d T) - 2Im\{Z\}e^{\sigma T} \sin(\omega_d T) \quad (7.68)$$

The discrete input matrix is in the form

$$\begin{aligned}B_d &= \int_0^T e^{A\tau} B d\tau = ZB \left(\frac{e^{\lambda T} - 1}{\lambda} \right) + \bar{Z}B \left(\frac{e^{\bar{\lambda}T} - 1}{\bar{\lambda}} \right) \\ &= 2Re\{Z\}BRe\left\{\frac{e^{\lambda T} - 1}{\lambda}\right\} - 2Im\{Z\}BIm\left\{\frac{e^{\lambda T} - 1}{\lambda}\right\} \\ \frac{e^{\lambda T} - 1}{\lambda} &= \frac{\sigma[e^{\sigma T} \cos(\omega_d T) - 1] + \omega_d e^{\sigma T} \sin(\omega_d T)}{|\lambda|^2} \\ &\quad - j \frac{\omega_d [e^{\sigma T} \cos(\omega_d T) - 1] - \sigma e^{\sigma T} \sin(\omega_d T)}{|\lambda|^2} \\ B_d &= 2Re\{Z\}B \frac{\sigma[e^{\sigma T} \cos(\omega_d T) - 1] + \omega_d e^{\sigma T} \sin(\omega_d T)}{|\lambda|^2} \\ &\quad + 2Im\{Z\}B \frac{\omega_d [e^{\sigma T} \cos(\omega_d T) - 1] - \sigma e^{\sigma T} \sin(\omega_d T)}{|\lambda|^2} \quad (7.69)\end{aligned}$$

For eigenvalues on the imaginary axis $\lambda = j\omega_d$, and we have the simpler forms

$$e^{AT} = 2Re\{Z\}\cos(\omega_d T) - 2Im\{Z\}\sin(\omega_d T) \quad (7.70)$$

$$B_d = 2Re\{Z\}B \frac{\sin(\omega_d T)}{\omega_d} + 2Im\{Z\}B \frac{\cos(\omega_d T) - 1}{\omega_d} \quad (7.71)$$

Example 7.17

Obtain the discrete state and input matrix for the series resonances circuit of Example 7.14 with a DAC and ADC and a sampling period T . Evaluate the matrices for $T = 0.05$ s and check your answer with MATLAB.

Solution

The state matrix is

$$A_d = e^{AT} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} e^{-T} \cos(3T) - \frac{1}{3} \begin{bmatrix} -1 & -1 \\ 10 & 1 \end{bmatrix} e^{-T} \sin(3T)$$

$$2ZB = 2 \left\{ \frac{1}{3} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \frac{j}{6} \begin{bmatrix} -1 & -1 \\ 10 & 1 \end{bmatrix} \right\} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \frac{j}{3} \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$B_d = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \frac{-[e^{-T} \cos(3T) - 1] + 3e^{-T} \sin(3T)}{10} + \frac{1}{3} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \frac{3[e^{-T} \cos(3T) - 1] + e^{-T} \sin(3T)}{10}$$

Evaluating the matrices at $T = 0.05$ s gives the same answer as the MATLAB commands

```
>> expm(A * .05)
ans =
    0.9879  0.0474
   -0.4738  0.8932

>> Bd = A (Ad - eye(size(A))) * B
Bd =
    0.0012
    0.0474
```

7.7 Solution of discrete-time state-space equations

We now seek an expression for the state at time k in terms of the initial condition vector $\mathbf{x}(0)$ and the input sequence $\mathbf{u}(k)$, $k = 0, 1, 2, \dots, k-1$. We begin by examining the discrete-time state Eq. (7.59):

$$\mathbf{x}(k+1) = A_d \mathbf{x}(k) + B_d \mathbf{u}(k)$$

At $k = 0, 1$, we have

$$\begin{aligned}\mathbf{x}(1) &= A_d \mathbf{x}(0) + B_d \mathbf{u}(0) \\ \mathbf{x}(2) &= A_d \mathbf{x}(1) + B_d \mathbf{u}(1)\end{aligned}\tag{7.72}$$

Substituting from the first into the second equation in Eq. (7.72) gives

$$\begin{aligned}\mathbf{x}(2) &= A_d[A_d \mathbf{x}(0) + B_d \mathbf{u}(0)] + B_d \mathbf{u}(1) \\ &= A_d^2 \mathbf{x}(0) + A_d B_d \mathbf{u}(0) + B_d \mathbf{u}(1)\end{aligned}\tag{7.73}$$

We then rewrite Eq. (7.73) as

$$\mathbf{x}(2) = A_d^2 \mathbf{x}(0) + \sum_{i=0}^{2-1} A_d^{2-i-1} B_d \mathbf{u}(i)\tag{7.74}$$

and observe that the expression generalizes to

$$\mathbf{x}(k) = A_d^k \mathbf{x}(0) + \sum_{i=0}^{k-1} A_d^{k-i-1} B_d \mathbf{u}(i)\tag{7.75}$$

This expression is in fact the general solution. Left as an exercise are details of the proof by induction where Eq. (7.75) is assumed to hold and it is shown that a similar form holds for $\mathbf{x}(k+1)$.

A more general expression for the solution with initial time k_0 and initial state $\mathbf{x}(k_0)$ is

$$\mathbf{x}(k) = A_d^{k-k_0} \mathbf{x}(k_0) + \sum_{i=k_0}^{k-1} A_d^{k-i-1} B_d \mathbf{u}(i)\tag{7.76}$$

Eq. (7.75) is the solution of the discrete-time state equation. The matrix A_d^k is known as the **state-transition matrix** for the discrete-time system, and it plays a role analogous to its continuous counterpart. A state-transition matrix can be defined for time-varying discrete-time systems, but it is not a matrix power, and it is dependent on both time k and initial time k_0 .

Eq. (7.75) includes two terms as in the continuous-time case. The first is the **zero-input response** due to nonzero initial conditions and zero input. The second is the **zero-state response** due to nonzero input and zero initial conditions. Because the system is linear, each term can be computed separately and then added to obtain the total response for a forced system with nonzero initial conditions.

Substituting from Eq. (7.75) in the discrete-time output Eq. (7.67) gives the output

$$\mathbf{y}(k) = C \left\{ A_d^k \mathbf{x}(0) + \sum_{i=0}^{k-1} A_d^{k-i-1} B_d \mathbf{u}(i) \right\} + D \mathbf{u}(k)\tag{7.77}$$

7.7.1 z-transform solution of discrete-time state equations

Eq. (7.75) can be obtained by z -transforming the discrete-time state Eq. (7.59). The z -transform is given by

$$z\mathbf{X}(z) - z\mathbf{x}(0) = A_d\mathbf{X}(z) + B_d\mathbf{U}(z) \quad (7.78)$$

Hence, $\mathbf{X}(z)$ is given by

$$\mathbf{X}(z) = [zI_n - A_d]^{-1}[z\mathbf{x}(0) + B_d\mathbf{U}(z)] \quad (7.79)$$

We therefore need to evaluate the inverse z -transform of the matrix $[zI_n - A_d]^{-1}z$. This can be accomplished by expanding the matrix in the form of the series

$$\begin{aligned} [zI_n - A_d]^{-1}z &= \left[I_n - \frac{1}{z}A_d \right]^{-1} \\ &= I_n + A_dz^{-1} + A_d^2z^{-2} + \dots + A_d^iz^{-i} + \dots \end{aligned} \quad (7.80)$$

The inverse z -transform of the series is

$$\mathcal{Z}\left\{[zI_n - A_d]^{-1}z\right\} = \left\{I_n, A_d, A_d^2, \dots, A_d^i, \dots\right\} \quad (7.81)$$

Hence, we have the z -transform pair

$$[zI_n - A_d]\mathcal{Z}^{-1} \leftrightarrow \left\{ A_d^k \right\}_{k=0}^{\infty} \quad (7.82)$$

This result is analogous to the scalar transform pair

$$\frac{z}{z - a_d} \xleftrightarrow{\mathcal{Z}} \left\{ a_d^k \right\}_{k=0}^{\infty}$$

The inverse matrix in Eq. (7.82) can be evaluated using the Leverrier algorithm of Section 7.4.1 to obtain the expression

$$[zI_n - A_d]^{-1}z = \frac{P_0z + P_1z^2 + \dots + P_{n-1}z^n}{a_0 + a_1z + \dots + a_{n-1}z^{n-1} + z^n} \quad (7.83)$$

Then, after denominator factorization and partial fraction expansion, we obtain

$$[zI_n - A_d]^{-1}z = \sum_{i=1}^n \frac{z}{z - \lambda_i} Z_i \quad (7.84)$$

where λ_i , $i = 1, 2, \dots, n$ are the eigenvalues of the discrete state matrix A_d . Finally, we inverse z -transform to obtain the discrete-time state-transition matrix

$$A_d^k = \sum_{i=1}^n Z_i \lambda_i^k \quad (7.85)$$

Writing Eq. (7.85) for $k = 1$ and using Eq. (7.60), we have

$$A_d = \sum_{i=1}^n Z_i \lambda_i = \sum_{i=1}^n Z_i(A) e^{\lambda_i(A)T} \quad (7.86)$$

where the parentheses indicate terms pertaining to the continuous-time state matrix A . Because the equality must hold for any sampling period T and any matrix A , we have the two equalities

$$\begin{aligned} Z_i &= Z_i(A) \\ \lambda_i &= e^{\lambda_i(A)T} \end{aligned} \quad (7.87)$$

In other words, the constituent matrices of the discrete-time state matrix are simply those of the continuous-time state matrix A , and its eigenvalues are exponential functions of the continuous-time characteristic values times the sampling period. This allows us to write the discrete-time state-transition matrix as

$$A_d^k = \sum_{i=1}^n Z_i e^{\lambda_i(A)kT} \quad (7.88)$$

For a state matrix with complex-conjugate eigenvalues $\lambda_{1,2} = \sigma \pm j\omega_d$, the state transition matrix for the discrete-time system can be written as

$$A_d^k = e^{AkT} = 2Re\{Z\}e^{\sigma kT} \cos(\omega_d kT) - 2Im\{Z\}e^{\sigma kT} \sin(\omega_d kT) \quad (7.89)$$

To complete the solution of the discrete-time state equation, we examine the zero-state response rewritten as

$$\mathbf{x}_{ZS}(z) = \left\{ [zI_n - A_d]^{-1} z \right\} z^{-1} B_d \mathbf{U}(z) \quad (7.90)$$

The term in braces in Eq. (7.90) has a known inverse transform, and multiplication by z^{-1} is equivalent to delaying its inverse transform by one sampling period. The remaining terms also have a known inverse transform. Using the convolution theorem, the inverse of the product is the convolution summation

$$\mathbf{x}_{ZS}(k) = \sum_{i=0}^{k-1} A_d^{k-i-1} B_d \mathbf{u}(i) \quad (7.91)$$

This completes the solution using z -transformation.

Using Eq. (7.88), the zero-state response can be written as

$$\mathbf{x}_{ZS}(k) = \sum_{i=0}^{k-1} \left[\sum_{j=1}^n Z_j e^{\lambda_j(A)(k-i-1)T} \right] B_d \mathbf{u}(i) \quad (7.92)$$

Then interchanging the order of summation gives

$$\mathbf{x}_{zs}(k) = \sum_{j=1}^n Z_j B_d e^{\lambda_j(A)(k-1)T} \left[\sum_{i=0}^{k-1} e^{-\lambda_j(A)iT} \mathbf{u}(i) \right] \quad (7.93)$$

This expression is useful in some special cases where the summation over i can be obtained in closed form. For example, if the system has no unity eigenvalues, and if the input is a vector of step functions $\mathbf{u}(k) = \mathbf{1} = [1 \dots 1]^T$, we can use the identity

$$\sum_{i=1}^{k-1} a^i = \frac{1-a^k}{1-a}, \quad a \neq 1$$

to obtain

$$\mathbf{x}_{zs}(k) = \sum_{j=1}^n Z_j B_d \mathbf{1} e^{\lambda_j(A)[k-1]T} \left[\frac{1 - e^{-\lambda_j(A)kT}}{1 - e^{-\lambda_j(A)T}} \right]$$

Example 7.18

Consider the state equation

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u$$

1. Solve the state equation for a unit step input and the initial condition vector $\mathbf{x}(0) = [1 \ 0]^T$.
2. Use the solution to obtain the discrete-time state equations for a sampling period of 0.1 s.
3. Solve the discrete-time state equations with the same initial conditions and input as in part 1, and verify that the solution is the same as that shown in part 1 evaluated at multiples of the sampling period T .

Solution

1. We begin by finding the state-transition matrix for the given system. The resolvent matrix is

$$\Phi(s) = \begin{bmatrix} s & -1 \\ 2 & s+3 \end{bmatrix}^{-1} = \frac{\begin{bmatrix} s+3 & 1 \\ -2 & s \end{bmatrix}}{s^2 + 3s + 2} = \frac{\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}s + \begin{bmatrix} 3 & 1 \\ -2 & 0 \end{bmatrix}}{(s+1)(s+2)}$$

To inverse Laplace transform, we need the partial fraction expansions

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

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

Example 7.18—cont'd

Thus, we reduce the resolvent matrix to

$$\begin{aligned}\Phi(s) &= \frac{-\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 3 & 1 \\ -2 & 0 \end{bmatrix}}{(s+1)} + \frac{2\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} 3 & 1 \\ -2 & 0 \end{bmatrix}}{(s+2)} \\ &= \frac{\begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix}}{(s+1)} + \frac{\begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix}}{(s+2)}\end{aligned}$$

The matrices in the preceding expansion are both rank 1 because they are the constituent matrices of the state matrix A . The expansion can be easily inverse Laplace transformed to obtain the state-transition matrix

$$\phi(t) = \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} e^{-t} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} e^{-2t}$$

The zero-input response of the system is

$$\begin{aligned}\mathbf{x}_{ZI}(t) &= \left\{ \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} e^{-t} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} e^{-2t} \right\} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 2 \\ -2 \end{bmatrix} e^{-t} + \begin{bmatrix} -1 \\ 2 \end{bmatrix} e^{-2t}\end{aligned}$$

For a step input, the zero-state response is

$$\begin{aligned}\mathbf{x}_{ZS}(t) &= \left\{ \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} e^{-t} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} e^{-2t} \right\} \begin{bmatrix} 0 \\ 1 \end{bmatrix} * 1(t) \\ &= \begin{bmatrix} 1 \\ -1 \end{bmatrix} e^{-t} * 1(t) + \begin{bmatrix} -1 \\ 2 \end{bmatrix} e^{-2t} * 1(t)\end{aligned}$$

where $*$ denotes convolution. Because convolution of an exponential and a step is a relatively simple operation, we do not need Laplace transformation. We use the identity

$$e^{-\alpha t} * 1(t) = \int_0^t e^{-\alpha \tau} d\tau = \frac{1 - e^{-\alpha t}}{\alpha}$$

to obtain

$$\mathbf{x}_{ZS}(t) = \begin{bmatrix} 1 \\ -1 \end{bmatrix} (1 - e^{-t}) + \begin{bmatrix} -1 \\ 2 \end{bmatrix} \frac{1 - e^{-2t}}{2} = \begin{bmatrix} 1/2 \\ 0 \end{bmatrix} - \begin{bmatrix} 1 \\ -1 \end{bmatrix} e^{-t} - \begin{bmatrix} -1 \\ 2 \end{bmatrix} \frac{e^{-2t}}{2}$$

Example 7.18—cont'd

The total system response is the sum of the zero-input and zero-state responses

$$\begin{aligned}\mathbf{x}(t) &= \mathbf{x}_{ZI}(t) + \mathbf{x}_{ZS}(t) \\ &= \begin{bmatrix} 2 \\ -2 \end{bmatrix} e^{-t} + \begin{bmatrix} -1 \\ 2 \end{bmatrix} e^{-2t} + \begin{bmatrix} 1/2 \\ 0 \end{bmatrix} - \begin{bmatrix} 1 \\ -1 \end{bmatrix} e^{-t} - \begin{bmatrix} -1 \\ 2 \end{bmatrix} \frac{e^{-2t}}{2} \\ &= \begin{bmatrix} 1/2 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ -1 \end{bmatrix} e^{-t} + \begin{bmatrix} -1 \\ 2 \end{bmatrix} \frac{e^{-2t}}{2}\end{aligned}$$

2. To obtain the discrete-time state equations, we use the state-transition matrix obtained in step 1 with t replaced by the sampling period. For a sampling period of 0.1, we have

$$A_d = \Phi(0.1) = \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} e^{-0.1} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} e^{-2 \times 0.1} = \begin{bmatrix} 0.9909 & 0.0861 \\ -0.1722 & 0.7326 \end{bmatrix}$$

The discrete-time input matrix B_d can be evaluated as shown earlier using Eq. (7.61). One may also observe that the continuous-time system response to a step input of the duration of one sampling period is the same as the response of a system due to a piecewise constant input discussed in Section 7.6. If the input is of unit amplitude, B_d can be obtained from the zero-state response of part 1 with t replaced by the sampling period $T = 0.1$. B_d is given by

$$\begin{aligned}B_d &= \left\{ \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} (1 - e^{-0.1}) + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} \frac{1 - e^{-2 \times 0.1}}{2} \right\} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} 1/2 \\ 0 \end{bmatrix} - \begin{bmatrix} 1 \\ -1 \end{bmatrix} e^{-0.1} - \begin{bmatrix} -1 \\ 2 \end{bmatrix} \frac{e^{-0.2}}{2} = \begin{bmatrix} 0.0045 \\ 0.0861 \end{bmatrix}\end{aligned}$$

3. The solution of the discrete-time state equations involves the discrete-time state-transition matrix

$$A_d^k = \Phi(0.1k) = \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} e^{-0.1k} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} e^{-0.2k}$$

The zero-input response is the product

$$\begin{aligned}\mathbf{x}_{ZI}(k) &= \Phi(0.1k)\mathbf{x}(0) \\ &= \left\{ \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} e^{-0.1k} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} e^{-0.2k} \right\} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 \\ -2 \end{bmatrix} e^{-0.1k} + \begin{bmatrix} -1 \\ 2 \end{bmatrix} e^{-0.2k}\end{aligned}$$

Comparing this result to the zero-input response of the continuous-time system reveals that the two are identical at all sampling points $k = 0, 1, 2, \dots$. Next, we z-transform the discrete-time state-transition matrix to obtain

$$\Phi(z) = \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} \frac{z}{z - e^{-0.1}} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} \frac{z}{z - e^{-0.2}}$$

Example 7.18—cont'd

Hence, the z-transform of the zero-state response for a unit step input is

$$\begin{aligned}\mathbf{X}_{zs}(z) &= \Phi(z)z^{-1}B_d\mathbf{U}(z) \\ &= \left\{ \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} \frac{z}{z - e^{-0.1}} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} \frac{z}{z - e^{-0.2}} \right\} \begin{bmatrix} 0.0045 \\ 0.0861 \end{bmatrix} \frac{z^{-1}z}{z - 1} \\ &= 9.5163 \times 10^{-2} \begin{bmatrix} 1 \\ -1 \end{bmatrix} \frac{z}{(z - e^{-0.1})(z - 1)} + 9.0635 \times 10^{-2} \begin{bmatrix} -1 \\ 2 \end{bmatrix} \frac{z}{(z - e^{-0.2})(z - 1)}\end{aligned}$$

We now need the partial fraction expansions

$$\begin{aligned}\frac{z}{(z - e^{-0.1})(z - 1)} &= \frac{1}{1 - e^{-0.1}} \left[\frac{z}{z - 1} + \frac{(-1)z}{z - e^{-0.1}} \right] = 10.5083 \left[\frac{z}{z - 1} + \frac{(-1)z}{z - 0.9048} \right] \\ \frac{z}{(z - e^{-0.2})(z - 1)} &= \frac{1}{1 - e^{-0.2}} \left[\frac{z}{z - 1} + \frac{(-1)z}{z - e^{-0.2}} \right] = 5.5167 \left[\frac{z}{z - 1} + \frac{(-1)z}{z - 0.8187} \right]\end{aligned}$$

Substituting in the zero-state response expression yields

$$\mathbf{X}_{zs}(z) = \begin{bmatrix} 0.5 \\ 0 \end{bmatrix} \frac{z}{z - 1} - \begin{bmatrix} 1 \\ -1 \end{bmatrix} \frac{z}{z - e^{-0.1}} - 0.5 \begin{bmatrix} -1 \\ 2 \end{bmatrix} \frac{z}{z - e^{-0.2}}$$

Then inverse-transforming gives the response

$$\mathbf{x}_{zs}(k) = \begin{bmatrix} 0.5 \\ 0 \end{bmatrix} - \begin{bmatrix} 1 \\ -1 \end{bmatrix} e^{-0.1k} - 0.5 \begin{bmatrix} -1 \\ 2 \end{bmatrix} e^{-0.2k}$$

This result is identical to the zero-state response for the continuous system at time $t = 0.1 k$, $k = 0, 1, 2, \dots$

Since the matrix A_d has distinct eigenvalues and no eigenvalues equal to unity, the zero-state response for a unit step input can also be obtained using Eq. (7.93)

$$\mathbf{x}_{zs}(k) = \left\{ \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} e^{-(k-1)} \begin{bmatrix} 1 - e^{kT} \\ 1 - e^T \end{bmatrix} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} e^{-2(k-1)} \begin{bmatrix} 1 - e^{2kT} \\ 1 - e^{2T} \end{bmatrix} \right\} \begin{bmatrix} 0.0045 \\ 0.0861 \end{bmatrix}$$

Substituting numerical values in the above expression yields the zero-state response obtained earlier. For example, substituting $k = 5$ in both expressions gives the vector

$$\mathbf{x}_{zs}(5) = [0.0774 \quad 0.2387]^T$$

7.8 z-transfer function from state-space equations

The z-transfer function can be obtained from the discrete-time state-space representation by z-transforming the output Eq. (7.67) to obtain

$$\mathbf{Y}(z) = \mathbf{C}\mathbf{X}(z) + \mathbf{D}\mathbf{U}(z) \tag{7.94}$$

The transfer function is derived under zero initial conditions. We therefore substitute the z -transform of the zero-state response Eq. (7.90) for $\mathbf{X}(z)$ to obtain

$$\mathbf{Y}(z) = C \left\{ [zI_n - A_d]^{-1} \right\} B_d \mathbf{U}(z) + D \mathbf{U}(z) \quad (7.95)$$

Thus, the z -transfer function matrix is defined by

$$\mathbf{Y}(z) = G(z) \mathbf{U}(z) \quad (7.96)$$

where $G(z)$ is the matrix

$$G(z) = C \left\{ [zI_n - A_d]^{-1} \right\} B_d + D \xleftrightarrow{\mathcal{Z}} G(k) = \begin{cases} CA_d^{k-1} B_d, & k \geq 1 \\ D, & k = 0 \end{cases} \quad (7.97)$$

and $G(k)$ is the impulse response matrix. The transfer function matrix and the impulse response matrix are z -transform pairs.

Substituting from Eq. (7.84) into Eq. (7.97) gives the alternative expression

$$G(z) = \sum_{i=1}^n CZ_i B_d \frac{1}{z - \lambda_i} + D \xleftrightarrow{\mathcal{Z}} G(k) = \begin{cases} \sum_{i=1}^n CZ_i B_d \lambda_i^{k-1}, & k \geq 1 \\ D, & k = 0 \end{cases} \quad (7.98)$$

Thus, the poles of the system are the eigenvalues of the discrete-time state matrix A_d . From Eq. (7.87), these are exponential functions of the eigenvalues $\lambda_i(A)$ of the continuous-time state matrix A . For a stable matrix A , the eigenvalues $\lambda_i(A)$ have negative real parts and the eigenvalues λ_i have magnitude less than unity. This implies that the discretization of Section 7.6 yields a stable discrete-time system for a stable continuous-time system.

Another important consequence of Eq. (7.98) is that the product CZ_iB can vanish and eliminate certain eigenvalues from the transfer function. This occurs if the product CZ_i is zero, the product Z_iB is zero, or both. If such cancellation occurs, the system is said to have an **output-decoupling zero** at λ_i , an **input-decoupling zero** at λ_i , or an **input-output-decoupling zero** at λ_i , respectively. The poles of the reduced transfer function are then a subset of the eigenvalues of the state matrix A_d . A state-space realization that leads to pole-zero cancellation is said to be **reducible** or **nonminimal**. If no cancellation occurs, the realization is said to be **irreducible** or **minimal**.

Clearly, in case of an output-decoupling zero at λ_i , the forced system response does not include the mode λ_i^k . In case of an input-decoupling zero, the mode is decoupled from or unaffected by the input. In case of an input-output-decoupling zero, the mode is decoupled from both the input and the output. These properties are related to the concepts of controllability and observability discussed in Chapter 8.

Example 7.19

Obtain the z-transfer function for the position control system described in [Example 7.18](#).

1. With x_1 as output
2. With x_1+x_2 as output

Solution

From [Example 7.18](#), we have

$$\Phi(z) = \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} \frac{z}{z - e^{-0.1}} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} \frac{z}{z - e^{-0.2}}$$

$$B_d = \begin{bmatrix} 1/2 \\ 0 \end{bmatrix} - \begin{bmatrix} 1 \\ -1 \end{bmatrix} e^{-0.1} - \begin{bmatrix} -1 \\ 2 \end{bmatrix} \frac{e^{-0.2}}{2} = \begin{bmatrix} 0.0045 \\ 0.0861 \end{bmatrix}$$

1. The output matrix C and the direct transmission matrix D for output x_1 are

$$C = [1 \ 0] \quad D = 0$$

Hence, the transfer function of the system is

$$G(z) = [1 \ 0] \left\{ \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} \frac{1}{z - e^{-0.1}} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} \frac{1}{z - e^{-0.2}} \right\} \begin{bmatrix} 0.0045 \\ 0.0861 \end{bmatrix}$$

$$= \frac{9.5163 \times 10^{-2}}{z - e^{-0.1}} - \frac{9.0635 \times 10^{-2}}{z - e^{-0.2}}$$

2. With x_1+x_2 as output, $C = [1, 1]$ and $D = 0$. The transfer function is

$$G(z) = [1 \ 1] \left\{ \begin{bmatrix} 2 & 1 \\ -2 & -1 \end{bmatrix} \frac{1}{z - e^{-0.1}} + \begin{bmatrix} -1 & -1 \\ 2 & 2 \end{bmatrix} \frac{1}{z - e^{-0.2}} \right\} \begin{bmatrix} 0.0045 \\ 0.0861 \end{bmatrix}$$

$$= \frac{0}{z - e^{-0.1}} + \frac{9.0635 \times 10^{-2}}{z - e^{-0.2}}$$

The system has an output-decoupling zero at $e^{-0.1}$ because the product CZ_1 is zero. The response of the system to any input does not include the decoupling term. For example, the step response is the inverse of the transform

$$Y(z) = \frac{9.0635 \times 10^{-2}z}{(z - e^{-0.2})(z - 1)}$$

$$= \frac{9.0635 \times 10^{-2}}{1 - e^{-0.2}} \left[\frac{z}{z - 1} + \frac{(-1)z}{z - e^{-0.2}} \right] = 0.5 \left[\frac{z}{z - 1} + \frac{(-1)z}{z - 0.8187} \right]$$

That is, the step response is

$$y(k) = 0.5 [1 - e^{-0.2k}]$$

7.8.1 z-transfer function in MATLAB

The expressions for obtaining z -domain and s -domain transfer functions differ only in that z in the former is replaced by s in the latter. The same MATLAB command is used to obtain s -domain and z -domain transfer functions. The transfer function for the matrices (A_d , B_d , C , D) is obtained with the commands

```
>> p = ss(Ad, Bd, C, D, T)
>> gd = tf(p)
```

where T is the sampling period. The poles and zeros of a transfer function are obtained with the command

```
>> zpk(gd)
```

For the system described in Example 7.19 (2), we obtain.

```
Zero/pole/gain:
0.090,635 (z-0.9048).
-----
(z-0.9048) (z-0.8187).
Sampling time: 0.1
```

The command reveals that the system has a zero at 0.9048 and poles at (0.9048, 0.8187) with a gain of 0.09035. With pole-zero cancellation, the transfer function is the same as that shown in Example 7.19 (2).

7.9 Similarity transformation

Any given linear system has an infinite number of valid state-space representations. Each representation corresponds to a different set of basis vectors in state-space. Some representations are preferred because they reveal certain system properties, whereas others may be convenient for specific design tasks. This section considers transformation from one representation to another.

Given a state vector $\mathbf{x}(k)$ with state-space representation of the form of Eq. (7.59), we define a new state vector $\mathbf{z}(k)$

$$\mathbf{x}(k) = T_r \mathbf{z}(k) \Leftrightarrow \mathbf{z}(k) = T_r^{-1} \mathbf{x}(k) \quad (7.99)$$

where the transformation matrix T_r is assumed invertible. Substituting for $\mathbf{x}(k)$ from Eq. (7.99) in the state Eq. (7.59) and the output Eq. (7.67) gives

$$\begin{aligned} T_r \mathbf{z}(k+1) &= A_d T_r \mathbf{z}(k) + B_d \mathbf{u}(k) \\ \mathbf{y}(k) &= C T_r \mathbf{z}(k) + D \mathbf{u}(k) \end{aligned} \quad (7.100)$$

Premultiplying the state equation by T_r^{-1} gives

$$\mathbf{z}(k+1) = T_r^{-1} A_d T_r \mathbf{z}(k) + T_r^{-1} B_d \mathbf{u}(k) \quad (7.101)$$

Hence, we have the state-space quadruple for the state vector $\mathbf{z}(k)$

$$(A, B, C, D) = (T_r^{-1} A_d T_r, T_r^{-1} B_d, C T_r, D) \quad (7.102)$$

Clearly, the quadruple for the state vector $\mathbf{x}(k)$ can be obtained from the quadruple of $\mathbf{z}(k)$ using the inverse of the transformation T_r (i.e., the matrix T_r^{-1}).

For the continuous-time system of Eq. (7.3), a state vector $\mathbf{z}(t)$ can be defined as

$$\mathbf{x}(t) = T_r \mathbf{z}(t) \Leftrightarrow \mathbf{z}(t) = T_r^{-1} \mathbf{x}(t) \quad (7.103)$$

and substitution in Eq. (7.3) yields Eq. (7.102). Thus, a discussion of similarity transformation is identical for continuous-time and discrete-time systems. We therefore drop the subscript d in the sequel.

Example 7.20

Consider the point mass m driven by a force f of Example 7.1, and determine the two state-space equations when $m = 1$ and

1. the displacement and the velocity are the state variables and
2. the displacement and the sum of displacement and velocity are the state variables.

Solution

1. From the first principle, as shown in Example 7.1, we have that the equation of motion is

$$m\ddot{y} = f$$

and in case 1, by considering x_1 as the displacement and x_2 as the velocity, we can write

$$\begin{aligned} \dot{x}_1(t) &= x_2(t) \\ \dot{x}_2(t) &= u(t) \\ y(t) &= x_1(t) \end{aligned}$$

Thus, we obtain the realization

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad C = [1 \quad 0] \quad D = 0$$

Example 7.20—cont'd

2. We express the new state variables, z_1 and z_2 , in terms of the state variable of part 1 as

$$\begin{aligned} z_1(t) &= x_1(t) \\ z_2(t) &= x_1(t) + x_2(t) \end{aligned}$$

This yields the inverse of the transformation matrix

$$T_r^{-1} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$$

Using Eq. (7.102) gives the realization

$$\begin{aligned} A &= T_r^{-1} A_d T_r = \begin{bmatrix} -1 & 1 \\ -1 & 1 \end{bmatrix} & B &= T_r^{-1} B_d = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\ C &= [1 \quad 0] & D &= 0 \end{aligned}$$

To obtain the transformation to diagonal form, we recall the expression

$$A = V \Lambda V^{-1} \Leftrightarrow \Lambda = V^{-1} A V \quad (7.104)$$

where V is the modal matrix of eigenvectors of A and $\Lambda = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ is the matrix of eigenvalues of A . Thus, for $A = \Lambda$ in Eq. (7.102), we use the modal matrix of A as the transformation matrix. The form thus obtained is not necessarily the same as the diagonal form obtained in Section 8.5.3 from the transfer function using partial fraction expansion. Even though all diagonal forms share the same state matrix, their input and output matrices may be different.

Example 7.21

Obtain the diagonal form for the state-space equations

$$\begin{aligned} \begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -0.04 & -0.5 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} u(k) \\ y(k) &= \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} \end{aligned}$$

Solution

The **eig** command of MATLAB yields the eigenvalues and the modal matrix whose columns have unity norm

$$\Lambda = \text{diag}\{0, -0.1, -0.4\} \quad V = \begin{bmatrix} 1 & -0.995 & 0.9184 \\ 0 & 0.0995 & -0.36741 \\ 0 & -0.00995 & 0.1469 \end{bmatrix}$$

Example 7.21—cont'd

The state matrix is in companion form and the modal matrix is also known to be the **Van der Monde** matrix (column norms need not be unity):

$$V = \begin{bmatrix} 1 & 1 & 1 \\ \lambda_1 & \lambda_2 & \lambda_3 \\ \lambda_1^2 & \lambda_2^2 & \lambda_3^2 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 0 & -0.1 & -0.4 \\ 0 & 0.01 & 0.16 \end{bmatrix}$$

The MATLAB command for similarity transformation is **ss2ss**. It requires the inverse T_r^{-1} of the similarity transformation matrix T_r . Two MATLAB commands transform to diagonal form. The first requires the modal matrix of eigenvectors V and the system to be transformed

```
>> s_diag = ss2ss(system, inv(v))
```

The second uses similarity transformation but does not require the transformation matrix

```
>> s_diag = canon(system, 'modal')
```

The two commands yield

$$\begin{aligned} A_t &= diag\{0, -0.1, -0.4\}, \quad B_t = [25 \quad 33.5012 \quad 9.0738]^T \\ C_t &= [1.0000 \quad -0.9950 \quad 0.9184] \quad D_t = 0 \end{aligned}$$

For complex conjugate eigenvalues, the command **canon** yields a real realization but its state matrix is not in diagonal form, whereas **ss2ss** will yield a diagonal but complex matrix.

7.9.1 Invariance of transfer functions and characteristic equations

Similar systems can be viewed as different representations of the same systems. This is justified by the following theorem.

Theorem 7.1

Similar systems have identical transfer functions and characteristic polynomials.

Proof

Consider the characteristic polynomials of similar realizations (A, B, C, D) and (A_1, B_1, C_1, D) :

$$\begin{aligned} \det(zI_n - A_1) &= \det(zI_n - T_r^{-1}AT_r) \\ &= \det[T_r^{-1}(zI_n - A)T_r] \\ &= \det[T_r^{-1}]\det(zI_n - A)\det[T_r] = \det(zI_n - A) \end{aligned}$$

where we used the identity $\det[T_r^{-1}] \times \det[T_r] = 1$.

Theorem 7.1—cont'd

The transfer function matrix is

$$\begin{aligned}G_1(s) &= C_1[zI_n - A_1]B_1 + D_1 \\&= CT_r[zI_n - T_r^{-1}AT_r]T_r^{-1}B + D \\&= C[T_r(zI_n - T_r^{-1}AT_r)T_r^{-1}]B + D \\&= C[zI_n - A]B + D = G(s)\end{aligned}$$

where we used the identity $(A \ B \ C)^{-1} = C^{-1}B^{-1}A^{-1}$.

Clearly, not all systems with the same transfer function are similar, due to the possibility of pole-zero cancellation. Systems that give rise to the same transfer function are said to be **equivalent**.

Example 7.22

Show that the following system is equivalent to the system shown in Example 7.19 (2).

$$\begin{aligned}x(k+1) &= 0.8187x(k) + 9.0635 \times 10^{-2}u(k) \\y(k) &= x(k)\end{aligned}$$

Solution

The transfer function of the system is

$$G(z) = \frac{9.0635 \times 10^{-2}}{z - 0.8187}$$

which is identical to the reduced transfer function of Example 7.19 (2).

Reference

Moler, C.B., Van Loan, C.F., 1978. Nineteen dubious ways to calculate the exponential of a matrix. SIAM Rev. 20, 801–836.

Further reading

Belanger, P.R., 1995. Control Engineering: A Modern Approach. Saunders, Fort Worth, TX.

- Brogan, W.L., 1985. Modern Control Theory. Prentice Hall, Englewood Cliff
- Chen, C.T., 1984. Linear System Theory and Design. HRW, New York.
- Gupta, S.C., Hasdorff, L., 1970. Fundamentals of Automatic Control. Wiley, New York.
- Friedland, B., 1986. Control System Design: An Introduction to State—Space Methods. McGraw-Hill, New York.
- Gupta, S.C., Hasdorff, L., 1970. Fundamentals of Automatic Control. Wiley, New York.
- Hou, S.-H., 1998. A simple proof of the Leverrier-Faddeev characteristic polynomial algorithm. SIAM Rev. 40 (3), 706–709.
- Kailath, T., 1980. Linear Systems. Prentice Hall, Englewood Cliffs, NJ.
- Sinha, N.K., 1988. Control Systems. HRW, New York.

Problems

7.1 Classify the state—space equations regarding linearity and time variance:

a. $\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} \sin(t) & 1 \\ 0 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 1 \\ 2 \end{bmatrix} u$
 $y = [1 \quad 1] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

b. $\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 2 & 0 \\ 1 & 5 & 7 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} u$
 $y = [1 \quad 1 \quad 2] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$

c. $\dot{x} = -2x^2 + 7x + xu$

$y = 3x$

d. $\dot{x} = -7x + u$

$y = 3x^2$

7.2 The equations of motion of a 2-D.O.F. manipulator are

$$M\ddot{\theta} + D(\dot{\theta}) + g(\theta) = \begin{bmatrix} T \\ f \end{bmatrix}$$

$$M = \begin{bmatrix} m_{11} & m_{12} \\ m_{12} & m_{22} \end{bmatrix} \quad d(\dot{\theta}) = \begin{bmatrix} 0 \\ D_2\dot{\theta}_2 \end{bmatrix} \quad g(\theta) = \begin{bmatrix} g_1(\theta) \\ g_2(\theta) \end{bmatrix}$$

$$M_a = M^{-1} = \begin{bmatrix} m_{a11} & m_{a12} \\ m_{a12} & m_{a22} \end{bmatrix}$$

where $\theta = [\theta_1, \theta_2]^T$ is a vector of joint angles. The entries of the positive definite inertia matrix M depend on the robot coordinates θ . D_2 is a damping constant. The terms g_i , $i = 1, 2$, are gravity-related terms that also depend on the coordinates. The right-hand side is a vector of generalized forces.

- a. Obtain a state-space representation for the manipulator, and then linearize it in the vicinity of a general operating point $(\mathbf{x}_0, \mathbf{u}_0)$.
 - b. Obtain the linearized model in the vicinity of zero coordinates, velocities, and inputs.
 - c. Show that, if the entries of the state matrix are polynomials, the answer to (b) can be obtained from (a) by letting all nonlinear terms go to zero.
- 7.3 Obtain the matrix exponentials for the state matrices using four different approaches:
- a. $A = \text{diag}\{-3, -5, -7\}$
 - b. $A = \begin{bmatrix} 0 & 0 & 1 \\ 0 & -1 & 0 \\ -6 & 0 & 0 \end{bmatrix}$
 - c. $A = \begin{bmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & -6 & -5 \end{bmatrix}$
 - d. A is a block diagonal matrix with the matrices of (b) and (c) on its diagonal.
- 7.4 Obtain the zero-input responses of the systems of Problem 7.3 due to the initial condition vectors:
- (a), (b), (c), $[1, 1, 0]^T$ and $[1, 0, 0]^T$
 - (d) $[1, 1, 0, 1, 0, 0]^T$
- 7.5 Determine the discrete-time state equations for the systems of Problem 7.3 (a), (b), and (c), with $\mathbf{b} = [0, 0, 1]^T$ in terms of the sampling period T .
- 7.6 Prove that the (right) eigenvectors of the matrix A^T are the left eigenvectors of A and that A^T 's eigenvalues are the eigenvalues of A using

$$A = V \Lambda V^{-1} = V \Lambda W$$

- 7.7 Prove that for any square matrix A with distinct eigenvalues λ_i and any analytic function $f(\cdot)$

$$f(A) = \sum_{i=1}^n Z_i f(\lambda_i)$$

where Z_i , $i=1,\dots,n$ are the constituent matrices of A .

Hint: Use the identity $A^k = \sum_{i=1}^n Z_i \lambda_i^k$, $k = 0, 1, \dots$

- 7.8 Prove the properties of the constituent matrices given in Section 7.4.3 using Eq. (7.46).
- 7.9 a. Derive the expressions for the terms of the adjoint matrix used in the Leverrier algorithm. *Hint:* Multiply both sides of Eq. (7.30) by the matrix $[sI - A]$ and equate coefficients.
- b. Derive the expressions for the coefficients of the characteristic equations used in the Leverrier algorithm. *Hint:* Laplace transform the derivative expression for the matrix exponential to obtain

$$s\mathcal{L}\{e^{At}\} - I = \mathcal{L}\{e^{At}\}A$$

Take the trace, and then use the identity

$$\text{tr}([sI_n - A]^{-1}) = \frac{a_1 + 2a_2s + \cdots + (n-1)a_{n-1}s^{n-2} + ns^{n-1}}{a_0 + a_1s + \cdots + a_{n-1}s^{n-1} + s^n}$$

- 7.10 The biological component of a fishery system is assumed to be governed by the population dynamics equation

$$\frac{dx(t)}{dt} = rx(t)(1 - x(t)/K) - h(t)$$

where r is the intrinsic growth rate per unit time, K is the environment carrying capacity, $x(t)$ is the stock biomass, and $h(t)$ is the harvest rate in weight.²

- a. Determine the harvest rate for a sustainable fish population $x_0 < K$.
- b. Linearize the system in the vicinity of the fish population x_0 .
- c. Obtain a discrete-time model for the linearized model with a fixed average yearly harvest rate $h(k)$ in the k th year.
- d. Obtain a condition for the stability of the fish population from your discrete-time model, and comment on the significance of the condition.
- 7.11 The following differential equations represent a simplified model of an overhead crane:³

$$(m_L + m_C)\ddot{x}_1(t) + m_L l(\ddot{x}_3(t)\cos x_3(t) - \dot{x}_3^2(t)\sin x_3(t)) = u$$

$$m_L \ddot{x}_1(t) + \cos x_3(t) + m_L l \ddot{x}_3(t) = -m_L g \sin x_3(t)$$

where m_C is the mass of the trolley, m_L is the mass of the hook/load, l is the rope length, g is the gravity acceleration, u is the force applied to the trolley, x_1 is the

² Clark, C.W., 1990. Mathematical Bioeconomics: The Optimal Management of Renewable Resources. Wiley, New York.

³ Piazz A., Visioli, A., 2002. Optimal dynamic-inversion-based control of an overhead crane. IEE Proceedings: Control Theory and Applications 149(5), 405–411.

position of the trolley, and x_3 is the rope angle. Consider the position of the load $y = x_1 + l \sin x_3$ as the output.

- Determine a linearized state-space model of the system about the equilibrium point $\mathbf{x} = 0$ with state variables x_1, x_3 , the first derivative of x_1 , and the first derivative of x_3 .
 - Determine a second state-space model when the sum of the trolley position and of the rope angle is substituted for the rope angle as a third state variable.
- 7.12 An unmanned autonomous vehicle, shown in Fig. P7.12, is governed by the nonlinear equations⁴

$$\begin{aligned} m\ddot{x} &= -u(\cos(\psi)\sin(\theta)\cos(\phi) + \sin(\psi)\sin(\phi)) \\ m\ddot{y} &= -u(\sin(\psi)\sin(\theta)\cos(\phi) - \cos(\psi)\sin(\phi)) \\ m\ddot{z} &= -u(\cos(\theta)\cos(\phi)) + mg \\ \ddot{\psi} &= \tau_\psi \\ \ddot{\theta} &= \tau_\theta \\ \ddot{\phi} &= \tau_\phi \end{aligned}$$

where $\xi = [x, y, z]^T$ is the translation vector of the vehicle and $\eta = [\psi, \theta, \phi]^T$ is its rotation vector, with the rotation angles referred to as yaw, pitch and roll, respectively. The forces applied to the system are the torques $\tau_\psi, \tau_\theta, \tau_\phi$, and the control input u .

- (a) Show that the $\eta = 0$ is an equilibrium point of the translational dynamics of the vehicle and find the corresponding control input u .

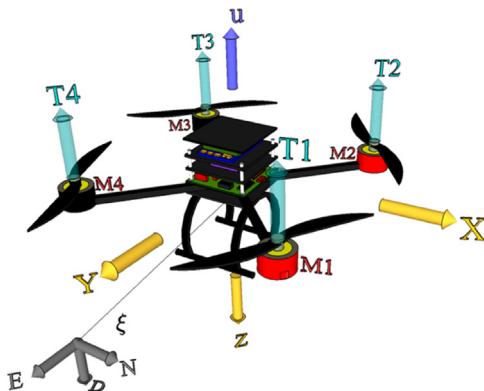


Figure P7.12

Unmanned autonomous vehicle. From Garcia-Carrillo et al., used with permission.

⁴ L.R. García Carrillo, A. Dzul, R. Lozano and C. Pégard, Quad Rotorcraft Control: Vision-Based Hovering and Navigation, Springer, September 2012.

- (b) Linearize the translational dynamic about the equilibrium $\eta = 0$.
- 7.13 Obtain the diagonal form for the discretized armature-controlled DC motor system of Example 7.16 with the motor angular position as output.
- 7.14 A system whose state and output responses are always nonnegative for any nonnegative initial conditions and any nonnegative input is called a **positive system**.⁵ Positive systems arise in many applications where the system variable can never be negative, including chemical processes, biological systems, and economics, among others. Show that the SISO discrete-time system (A, b, c^T) is positive if and only if all the entries of the state, input, and output matrix are positive.
- 7.15 To monitor river pollution, we need to model the concentration of biodegradable matter contained in the water in terms of biochemical oxygen demand for its degradation. We also need to model the dissolved oxygen deficit defined as the difference between the highest concentration of dissolved oxygen and the actual concentration in mg/l. If the two variables of interest are the state variables x_1 and x_2 , respectively, then an appropriate model is given by

$$\dot{\mathbf{x}} = \begin{bmatrix} -k_1 & 0 \\ k_1 & -k_2 \end{bmatrix} \mathbf{x}$$

where k_1 is a biodegradation constant and k_2 is a reaeration constant, and both are positive. Assume that the two positive constants are unequal. Obtain a discrete-time model for the system with sampling period T , and show that the system is positive.

- 7.16 Autonomous underwater vehicles (AUVs) are robotic submarines that can be used for a variety of studies of the underwater environment. The vertical and horizontal dynamics of the vehicle must be controlled to remotely operate the AUV. The INFANTE (Fig. P7.16) is a research AUV operated by the Instituto Superior Técnico of Lisbon, Portugal.⁶ The variables of interest in horizontal motion are the sway speed and the yaw angle. A linearized model of the horizontal plane motion of the vehicle is given by

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} -0.14 & -0.69 & 0.0 \\ -0.19 & -0.048 & 0.0 \\ 0.0 & 1.0 & 0.0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 0.056 \\ -0.23 \\ 0.0 \end{bmatrix} u$$

⁵ Farina, L., Rinaldi, S., 2000. Positive Linear Systems: Theory & Applications. Wiley-Interscience, New York.

⁶ Silvestre, C., Pascoa, A., 2004. Control of the INFANTEAUV using gainscheduledstaticoutput feedback. Control Engineering Practice 12, 1501–1509.



Figure P7.16

The INFANTE autonomous underwater vehicle (AUV). *From Silvestre and Pascoa, 2004; used with permission.*

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

where x_1 is the sway speed, x_2 is the yaw angle, x_3 is the yaw rate, and u is the rudder deflection. Obtain the discrete state-space model for the system with a sampling period of 50 ms.

- 7.17 A simplified linear state-space model of the ingestion of a drug in the bloodstream is given by⁷

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} -K_1 & 0 \\ K_1 & -K_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} u$$

where x_1 = the mass of the absorbed drug in mg; x_2 = the mass of the drug in the bloodstream in mg; u = the rate at which the drug is ingested in mg/min

- Find the state-transition matrix of the system.
 - Obtain a discrete-time state-space model for the system with a sampling period T in terms of the model parameters.
- 7.18 A typical assumption in most mathematical models is that the system differential equations or transfer functions have real coefficients. In a few applications, this assumption is not valid. In models of rotating machines, the evolution of the vectors governing the system with time depends on their space orientation relative to fixed inertial axes. For an induction motor, assuming symmetry, two stator fixed axes are used as the reference frame: the direct axis (d) in the horizontal direction and the quadrature axis (q) in the vertical direction. The terms in the quadrature direction are identified with a (j) coefficient that is absent from the direct axis terms. The two axes are shown in Fig. P7.18. We write the equations for the electrical subsystem of the

⁷ McClamroch, N.H., 1980. State Models of Dynamic Systems: A Case Study Approach. Springer-Verlag.

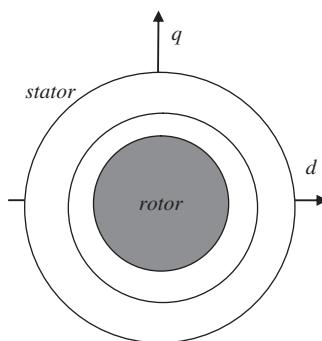


Figure P7.18
Stator frame for the induction motor.

motor in terms of the stator and rotor currents and voltages. Each current is decomposed into a direct axis component and quadrature component, with the latter identified with the term (*j*). The s-domain equations of the motor are obtained from its equivalent circuit using Kirchhoff's laws. The equations relative to the stator axes and including complex terms are

$$\begin{bmatrix} v_s \\ v_r \end{bmatrix} = \begin{bmatrix} R_s + sL_s & sL_m \\ (s - j\omega_r)L_m & R_r + (s - j\omega_r)L_r \end{bmatrix} \begin{bmatrix} i_s \\ i_r \end{bmatrix}$$

where

- a. Write the state equations for the induction motor.
- b. Without obtaining the eigenvalues of the state matrix, show that the two eigenvalues are not complex conjugate.

R_s (R_r) = stator (rotor) resistance

L_s , L_r , L_m = stator, rotor, mutual inductance, respectively

ω_r = rotor angular velocity

- 7.19 Throughout the text, we assumed uniform sampling with sampling period T . In some situations, it is advantageous to vary the sampling period with time. Show that, when the sampling period varies, the discrete-time state-space model of linear time invariant system with state-space matrices (A, B, C, D) and piece-wise constant input $\mathbf{u}(t) = \mathbf{u}(t_k)$, $t \in [t_k, t_{k+1})$, is time varying and is given by

$$\begin{aligned} \mathbf{x}(k+1) &= A(k)\mathbf{x}(k) + B(k)\mathbf{u}(k) \\ \mathbf{y}(k) &= C\mathbf{x}(k) + D\mathbf{u}(k) \end{aligned}$$

with $A(k) = e^{A\Delta t_k}$, $B(k) = \int_0^{\Delta t_k} e^{A\tau} B d\tau$, $\Delta t_k = t_{k+1} - t_k$, $\mathbf{x}(k) = \mathbf{x}(t_k)$, and t_k is the k^{th} sampling point.

- 7.20 In many practical applications, the output sampling in a digital control system is not exactly synchronized with the input transition. Show that the output equation corresponding to output sampling at $t_k = kT + \Delta_k, k = 0, 1, 2$, is

$$\mathbf{y}(k) = C(k)\mathbf{x}(k) + D(k)\mathbf{u}(k)$$

where

$$C(k) = Ce^{A\Delta_k}, \quad D(k) = CB_d(\Delta_k) + D, \quad B_d(\Delta_k) = \int_0^{\Delta_k} e^{At} B d\tau$$

- If the direct transmission matrix D is zero, does the input directly influence the sampled output?
- When is the resulting input–output model time-invariant?

- 7.21 Show that the solution of the state equation for the time-varying discrete time system

$$\mathbf{x}(k+1) = A(k)\mathbf{x}(k) + B(k)\mathbf{u}(k)$$

is given by

$$\begin{aligned}\mathbf{x}(k) &= \phi(k, k_0)\mathbf{x}(k_0) + \sum_{i=k_0}^{k-1} \phi(k, i+1)B(i)\mathbf{u}(i) \\ \phi(k, k_0) &= A(k-1)A(k-2)\dots A(k_0), \\ \phi(k, k) &= I_n\end{aligned}$$

Computer exercises

- 7.22 Write a computer program to simulate the systems of Problem 7.1 for various initial conditions with zero input, and discuss your results referring to the solutions of Problem 7.1. Obtain plots of the phase trajectories for any second-order system.
- 7.23 The Fitzhugh–Nagumo model is a simplified model of the electric potential across cell membranes. The model is given by the circuit of Fig. P7.23. In the figure, the membrane is represented by a nonlinear parallel R–C circuit with three branches. The first branch is a voltage-dependent membrane conductance $g(v)$ where v is the transmembrane voltage. The second is an inductance L with a series resistance R and bias voltage source b . The third branch is a capacitance C . A current I is applied to the membrane.

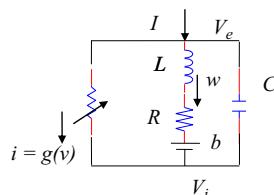


Figure P7.23
Circuit model of the cell membrane.

To write state equations for the circuit, we use the physical state variable v (capacitor voltage) and w (inductor current). The circuit equations for the model are given by

$$\begin{aligned} C \frac{dv}{dt} &= -g(v) - w + I \\ \frac{dw}{dt} &= \frac{1}{L}(v - R w - b) \end{aligned}$$

The nonlinear conductance is governed by the formula

$$g(v) = v(v - v_0)(v - 1)$$

The capacitance value for the model is $C = 1$, while other values will change resulting in drastically different responses.

- (a) Find the equilibrium points of the system with no applied current and zero bias b .
- (b) For the parameter values $L = 40$, $v_0 = 0.2$, $R = 0.5$, simulate the system using SIMULINK and plot the phase plane trajectory of the system.

7.24 Problem 7.23 investigates the behavior of the Fitzhugh–Nagumo model. We consider the same model with different parameter values:

$$\begin{aligned} \frac{dv}{dt} &= -v(v - 0.2)(v - 1) - w \\ \frac{dw}{dt} &= v - 0.25w \end{aligned}$$

- (a) Use the forward differencing approach of Section 6.3.1 to obtain a discretized Fitzhugh–Nagumo model
- (b) Simulate the discretized Fitzhugh–Nagumo model with a sampling period $T = 0.5$ s using SIMULINK, and obtain state plane trajectories for the system. The behavior of the discretized system is repetitive but not periodic with a small variation in the behavior with every cycle. This type of behavior is called **chaotic**.
- (c) Simulate the original system with the same parameter values to show that it is not chaotic and that, in this case, chaos is a result of the approximation associated with the discretization and not due to the natural behavior of the system.

- 7.25 Write a program to obtain the state-transition matrix using the Leverrier algorithm.
- 7.26 Simulate the systems of Problem 7.3 (a–c) with the initial conditions of Problem 7.4, and obtain state–trajectory plots with one state variable fixed for each system.
- 7.27 Repeat Problem 7.5 using a computer-aided design (CAD) package for two acceptable choices of the sampling period and compare the resulting systems.
- 7.28 Simulate the river pollution system of Problem 7.15 for the normalized parameter values of $k_1 = 1$, $k_2 = 2$, with a sampling period $T = 0.01$ s for the initial conditions $\mathbf{x}^T(0) = [1, 0], [0, 1], [1, 1]$, and plot all the results together.

- 7.29 Repeat Problem 7.16 using a CAD package.
- 7.30 By evaluating the derivatives of the matrix exponential at $t = 0$, show that for any matrix A with distinct eigenvalues the constituent matrices can be obtained by solving the equation

$$\begin{bmatrix} I_n & I_n & \dots & I_n \\ \lambda_1 I_n & \lambda_2 I_n & \dots & \lambda_n I_n \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^{n-1} I_n & \lambda_2^{n-1} I_n & \dots & \lambda_n^{n-1} I_n \end{bmatrix} \begin{bmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_n \end{bmatrix} = \begin{bmatrix} I_n \\ A \\ \vdots \\ A^{n-1} \end{bmatrix}$$

Write a MATLAB function that evaluates the constituent matrices and save them in a cell array.

- 7.31 The Cayley–Hamilton theorem states that, if $f(\lambda) = \lambda^n + a_{n-1}\lambda^{n-1} + \dots + a_1\lambda + a_0$ is the characteristic polynomial of a matrix A , then $f(A) = A^n + a_{n-1}A^{n-1} + \dots + a_1A + a_0I_n = \mathbf{0}$. This allows us to express the n^{th} power as

$$A^n = -[a_{n-1}A^{n-1} + \dots + a_1A + a_0I_n]$$

- a. Verify the validity of the theorem for the matrix

$$A = \begin{bmatrix} 3 & -2 & 0 \\ 8 & -3 & -4 \\ 0 & 4 & -9 \end{bmatrix}$$

using the MATLAB commands **poly** and **polyvalm**. Note that the answer you get will not be exact because of numerical errors.

- b. Use the Cayley–Hamilton theorem to show that the matrix exponential of a matrix A can be written as

$$e^{At} = \sum_{i=0}^{n-1} \alpha(t)A^i$$

- c. Show that the initial value of the vector of time functions

$$\alpha(t) = [\alpha_0(t) \quad \dots \quad \alpha_{n-1}(t)]^T$$

is given by

$$\alpha(\mathbf{0}) = [1 \quad \mathbf{0}_{1 \times (n-1)}]^T$$

- d. By differentiating the expression of (b), use the Cayley–Hamilton theorem to show that the vector of time functions satisfies

$$\dot{\alpha}(t) = \begin{bmatrix} \mathbf{0}^T \\ I_{n-1} \end{bmatrix} \left| -\mathbf{a}^T \right] \alpha(t)$$

- e. The transpose of the matrix of (d) is in a companion form whose eigenvector decomposition can be written by inspection. Use this fact to obtain the equation

$$\begin{bmatrix} 1 & \lambda_1 & \dots & \lambda_1^{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \lambda_n & \dots & \lambda_n^{n-1} \end{bmatrix} \alpha(t) = \begin{bmatrix} e^{\lambda_1 t} \\ \vdots \\ e^{\lambda_n t} \end{bmatrix}$$

- f. Write a MATLAB function to calculate the matrix exponential using the Cayley–Hamilton theorem and the results of (b–e).

- 7.32 Write a MATLAB script to calculate the zero state response of the linear time invariant system of Example 7.18 with unit step input at multiples of the sampling period, and the zero-state response of the discrete-time system, and show that at $k = 5$ the responses are both $\mathbf{x}_{zs}(5) = [0.0774 \quad 0.2387]^T$.

Properties of state-space models

Objectives

After completing this chapter, the reader will be able to do the following:

1. Determine the equilibrium point of a discrete-time system.
2. Determine the asymptotic stability of a discrete-time system.
3. Determine the input–output stability of a discrete-time system.
4. Determine the poles and zeros of multivariable systems.
5. Determine controllability and stabilizability.
6. Determine observability and detectability.
7. Obtain canonical state–space representations from an input–output representation of a system.

In Chapter 7, we described state–space models of linear discrete time systems and how these models can be obtained from transfer functions or input–output differential equations. We also obtained the solutions to continuous time and discrete-time state equations. Now, we examine some properties of these models that play an important role in system analysis and design.

We examine controllability, which determines the effectiveness of state feedback control; observability, which determines the possibility of state estimation from the output measurements; and stability. These three properties are independent, so a system can be unstable but controllable, uncontrollable but stable, and so on. However, systems whose uncontrollable dynamics are stable are stabilizable, and systems whose unobservable dynamics are stable are called detectable. Stabilizability and detectability are more likely to be encountered in practice than controllability and observability. Finally, we show how state–space representations of a system in several canonical forms can be obtained from its input–output representation.

To simplify our notation, the subscript d used in Chapter 7 with discrete-time state and input matrices is dropped if the discussion is restricted to discrete-time systems. In sections involving both continuous-time and discrete-time systems, the subscript d is retained. We begin with a discussion of stability.

Chapter Outline**8.1 Stability of state-space realizations 320**

- 8.1.1 Asymptotic stability 320
- 8.1.2 Bounded-Input–Bounded-Output stability 325

8.2 Controllability and stabilizability 329

- 8.2.1 MATLAB commands for controllability testing 337
- 8.2.2 Controllability of systems in normal form 337
- 8.2.3 Stabilizability 338

8.3 Observability and detectability 343

- 8.3.1 MATLAB commands 347
- 8.3.2 Observability of systems in normal form 347
- 8.3.3 Detectability 348

8.4 Poles and zeros of multivariable systems 350

- 8.4.1 Poles and zeros from the transfer function matrix 351
- 8.4.2 Zeros from state-space models 355

8.5 State-space realizations 357

- 8.5.1 Controllable canonical realization 358
 - 8.5.1.1 *Systems with no input differencing* 358
 - 8.5.1.2 *Systems with input differencing* 360
- 8.5.2 Controllable form in MATLAB 363
- 8.5.3 Parallel realization 363
 - 8.5.3.1 *Parallel realization for multiinput-multioutput systems* 366
- 8.5.4 Observable form 369

8.6 Duality 370**8.7 Hankel realization 372****8.8 Realizations for continuous-time systems 377**

Further reading 378

Problems 379

Computer exercises 385

8.1 Stability of state-space realizations

The concepts of **asymptotic stability** and **bounded-input–bounded-output (BIBO) stability** of transfer functions are discussed in Chapter 4. Here, we give more complete coverage using state-space models. We first discuss the asymptotic stability of state-space realizations. Then we discuss the BIBO stability of their input–output responses.

8.1.1 Asymptotic stability

The natural response of a linear system because of its initial conditions may

1. Converge to the origin asymptotically

2. Remain in a bounded region in the vicinity of the origin
3. Grow unbounded

In the first case, the system is said to be **asymptotically stable**; in the second, the system is **marginally stable**; and in the third, it is **unstable**. Clearly, physical variables are never actually unbounded even though linear models suggest this. Once a physical variable leaves a bounded range of values, the linear model ceases to be valid and the system must be described by a nonlinear model. Critical to the understanding of stability of both linear and nonlinear systems is the concept of an **equilibrium state**.

Definition 8.1: Equilibrium

An equilibrium point or state is an initial state from which the system never departs unless perturbed.

For the state equation,

$$\mathbf{x}(k+1) = \mathbf{f}[\mathbf{x}(k)] \quad (8.1)$$

all equilibrium states \mathbf{x}_e satisfy the condition

$$\begin{aligned} \mathbf{x}(k+1) &= \mathbf{f}[\mathbf{x}(k)] \\ &= \mathbf{f}[\mathbf{x}_e] = \mathbf{x}_e \end{aligned} \quad (8.2)$$

For linear state equations, equilibrium points satisfy the condition

$$\begin{aligned} \mathbf{x}(k+1) &= A\mathbf{x}(k) \\ &= A\mathbf{x}_e = \mathbf{x}_e \Leftrightarrow [A - I_n]\mathbf{x}_e = 0 \end{aligned} \quad (8.3)$$

For an invertible matrix $A - I_n$, Eq. (8.3) has the unique solution $\mathbf{0}$, and the linear system has one equilibrium state at the origin. Later we show that invertibility of $A - I_n$ is a necessary condition for asymptotic stability.

Unlike linear systems, nonlinear models may have several equilibrium states. A trajectory leaving one equilibrium state, around which the linear model is valid, may drive the system to another equilibrium state where another linear model is required. Nonlinear systems may also exhibit other more complex phenomena that are not discussed in this chapter but are discussed in Chapter 11.

The equilibrium of a system with constant input \mathbf{u} can be derived from Definition 8.1 by first substituting the constant input value in the state equation to obtain the form of Eq. (8.2).

Example 8.1

Find the equilibrium points of the following systems:

1. $x(k+1) = x(k)[x(k) - 0.5]$

2. $x(k+1) = 2x(k)$

3. $\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 0.1 & 0 \\ 1 & 0.9 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$

Solution

1. At equilibrium, we have

$$x_e = x_e[x_e - 0.5]$$

We rewrite the equilibrium condition as

$$x_e[x_e - 1.5] = 0$$

Hence, the system has the two equilibrium states

$$x_e = 0 \quad \text{and} \quad x_e = 1.5$$

2. The equilibrium condition is

$$x_e = 2x_e$$

The system has one equilibrium point at $x_e = 0$.

3. The equilibrium condition is

$$\begin{bmatrix} x_{1e}(k) \\ x_{2e}(k) \end{bmatrix} = \begin{bmatrix} 0.1 & 0 \\ 1 & 0.9 \end{bmatrix} \begin{bmatrix} x_{1e}(k) \\ x_{2e}(k) \end{bmatrix} \Leftrightarrow \begin{bmatrix} 0.1 - 1 & 0 \\ 1 & 0.9 - 1 \end{bmatrix} \begin{bmatrix} x_{1e}(k) \\ x_{2e}(k) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Because the coefficient matrix is invertible, the system has a unique equilibrium state at $x_e = [x_{1e}, x_{2e}]^T = [0, 0]^T$.

Although the systems in Example 8.1 all have equilibrium points, convergence to an equilibrium point is not guaranteed. This additional property defines asymptotic stability for linear systems assuming that the necessary condition of a unique equilibrium at the origin is satisfied.

Definition 8.2: Asymptotic stability

A linear system is said to be asymptotically stable if all its trajectories converge to the origin—that is, for any initial state $x(k_0)$, $x(k) \rightarrow 0$ as $k \rightarrow \infty$.

Theorem 8.1

A discrete-time linear system is asymptotically (**Schur**) **stable** if and only if all the eigenvalues of its state matrix are inside the unit circle.

Proof

To prove the theorem, we examine the zero-input response in Eq. (7.75) with the state-transition matrix given by Eq. (7.88). Substitution yields

$$\begin{aligned}\mathbf{x}_{ZI}(k) &= A^{k-k_0} \mathbf{x}(k_0) \\ &= \sum_{i=1}^n Z_i \mathbf{x}(k_0) \lambda_i^{k-k_0}\end{aligned}\quad (8.4)$$

Sufficiency

The response decays to zero if the eigenvalues are all inside the unit circle. Hence, the condition is sufficient.

Necessity

To prove necessity, we assume that the system is stable but that one of its eigenvalues λ_j is outside the unit circle. Let the initial condition vector $\mathbf{x}(k_0) = \mathbf{v}_j$, the j^{th} eigenvector of A . Then (see [Section 7.4.3](#)) the system response is

$$\begin{aligned}\mathbf{x}_{ZI}(k) &= \sum_{i=1}^n Z_i \mathbf{x}(k_0) \lambda_i^{k-k_0} \\ &= \sum_{i=1}^n \mathbf{v}_i \mathbf{w}_i^T \mathbf{v}_j \lambda_i^{k-k_0} = \mathbf{v}_j \lambda_j^{k-k_0}\end{aligned}$$

which is clearly unbounded as $k \rightarrow \infty$. Hence, the condition is also necessary by contradiction.

Remark

For some nonzero initial states, the product $Z_j \mathbf{x}(k_0)$ is zero because the matrix Z_j is rank 1 ($\mathbf{x}(k_0) = \mathbf{v}_i, i \neq j$). The response to those initial states may converge to zero even if the corresponding λ_j has magnitude greater than unity. However, the solution grows unbounded for an initial state in the one direction for which the product is nonzero. Therefore, not *all* trajectories converge to the origin for $|\lambda_j| > 1$. Hence, the condition of Theorem 8.1 is necessary.

We now discuss the necessity of an invertible matrix $A - I_n$ for asymptotic stability. The matrix can be decomposed as

$$A - I_n = V[A - I_n]V^{-1} = V \text{ diag}\{\lambda_i - 1\}V^{-1}$$

An asymptotically stable matrix A has no unity eigenvalues, and $A - I_n$ is invertible.

Example 8.2

Determine the stability of the systems of Example 8.1 (2) and 8.1 (3) using basic principles, and verify your results using Theorem 8.1.

Solution

8.1 (2): Consider any initial state $x(0)$. Then the response of the system is

$$\begin{aligned}x(1) &= 2x(0) \\x(2) &= 2x(1) = 2 \times 2x(0) = 4x(0) \\x(3) &= 2x(2) = 2 \times 4x(0) = 8x(0) \\&\vdots \\x(k) &= 2^k x(0)\end{aligned}$$

Clearly, the response is unbounded as $k \rightarrow \infty$. Thus, the system is unstable. The system has one eigenvalue at $2 > 1$, which violates the stability condition of Theorem 8.1.

8.1 (3): The state equation is

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 0.1 & 0 \\ 1 & 0.9 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$$

Using Sylvester's expansion Eq. (7.34), we obtain the constituent matrices

$$\begin{bmatrix} 0.1 & 0 \\ 1 & 0.9 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ -1/0.8 & 0 \end{bmatrix} 0.1 + \begin{bmatrix} 0 & 0 \\ 1/0.8 & 1 \end{bmatrix} 0.9$$

From Eq. (7.73), the response of the system due to an arbitrary initial state $x(0)$ is

$$\begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} = \left\{ \begin{bmatrix} 1 & 0 \\ -1/0.8 & 0 \end{bmatrix} (0.1)^k + \begin{bmatrix} 0 & 0 \\ 1/0.8 & 1 \end{bmatrix} (0.9)^k \right\} \begin{bmatrix} x_1(0) \\ x_2(0) \end{bmatrix}$$

This decays to zero as $k \rightarrow \infty$. Hence, the system is asymptotically stable. Both system eigenvalues (0.1 and 0.9) are inside the unit circle, and the system satisfies the conditions of Theorem 8.1.

Note that the response due to any initial state with $x_1(0) = 0$ does not include the first eigenvalue. Had this eigenvalue been unstable, the response due to this special class of initial conditions would remain bounded. However, the response due to other initial conditions would be unbounded, and the system would be unstable.

In Chapter 4, asymptotic stability was studied in terms of the poles of the system with the additional condition of no pole-zero cancellation. The condition for asymptotic stability was identical to the condition imposed on the eigenvalues of a stable state matrix in Theorem 8.1. This is because the eigenvalues of the state matrix A and the poles of the transfer function are identical in the absence of pole-zero cancellation. Next, we examine a stability definition based on the input-output response.

8.1.2 Bounded-Input–Bounded-Output stability

For an input–output system description, the system output must remain bounded for any bounded-input function. To test the boundedness of an n -dimensional output vector, a measure of the vector length or size known as the **norm** of the vector must be used (see Appendix III). A vector \mathbf{x} is bounded if it satisfies the condition

$$\|\mathbf{x}\| < b_x < \infty \quad (8.5)$$

for some finite constant b_x where $\|\mathbf{x}\|$ denotes any vector norm.

For real or complex $n \times 1$ vectors, all norms are equivalent in the sense that if a vector has a bounded norm $\|\mathbf{x}\|_a$, then any other norm $\|\mathbf{x}\|_b$, is also bounded. This is true because for any two norms, $\|\mathbf{x}\|_a$ and $\|\mathbf{x}\|_b$ there exist finite positive constants k_1 and k_2 such that

$$k_1\|\mathbf{x}\|_a \leq \|\mathbf{x}\|_b \leq \|\mathbf{x}\|_a k_2 \quad (8.6)$$

Because a change of norm merely results in a scaling of the constant b_x , boundedness is independent of which norm is actually used in Eq. (8.5). Stability of the input–output behavior of a system is defined as follows.

Definition 8.3: Bounded-input–bounded-output stability

A system is BIBO stable if its output is bounded for any bounded input. That is,

$$\|\mathbf{u}(k)\| < b_u < \infty \Rightarrow \|\mathbf{y}(k)\| < b_y < \infty \quad (8.7)$$

To obtain a necessary and sufficient condition for BIBO stability, we need the concept of the norm of a matrix (see Appendix III). Matrix norms can be defined based on a set of axioms or properties. Alternatively, we define the norm of a matrix in terms of the vector norm as

$$\begin{aligned} \|A\| &= \max_{\mathbf{x}} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} \\ &= \max_{\|\mathbf{x}\|=1} \|A\mathbf{x}\| \end{aligned} \quad (8.8)$$

Multiplying by the norm of \mathbf{x} , we obtain the inequality

$$\|A\mathbf{x}\| \leq \|A\| \cdot \|\mathbf{x}\| \quad (8.9)$$

which applies to any matrix norm induced from a vector norm, as well some other norms, and is known as the submultiplicative norm property. Next, we give a condition for BIBO stability in terms of the norm of the impulse response matrix.

Theorem 8.2

A system is BIBO stable if and only if the norm of its impulse response matrix is **absolutely summable**. That is,

$$\sum_{k=0}^{\infty} \|G(k)\| < \infty \quad (8.10)$$

Proof**Sufficiency**

The output of a system with impulse response matrix $G(k)$ and input $\mathbf{u}(k)$ is

$$\mathbf{y}(i) = \sum_{k=0}^i G(k) \mathbf{u}(i-k)$$

The norm of $\mathbf{y}(i)$ satisfies

$$\begin{aligned} \|\mathbf{y}(i)\| &= \left\| \sum_{k=0}^i G(k) \mathbf{u}(i-k) \right\| \\ &\leq \sum_{k=0}^i \|G(k) \mathbf{u}(i-k)\| \leq \sum_{k=0}^i \|G(k)\| \|\mathbf{u}(i-k)\| \end{aligned}$$

For a bounded input, $\mathbf{u}(k)$ satisfies

$$\|\mathbf{u}(k)\| < b_u < \infty$$

Hence, the output is bounded by

$$\|\mathbf{y}(i)\| \leq b_u \sum_{k=0}^i \|G(k)\| \leq b_u \sum_{k=0}^{\infty} \|G(k)\|$$

which is finite if condition (8.10) is satisfied.

Necessity

The proof is by contradiction. We assume that the system is BIBO stable but that Eq. (8.10) is violated. We write the output $\mathbf{y}(k)$ in the form

$$\begin{aligned} \mathbf{y}(k) &= G(k) \mathbf{u} \\ &= [G(0) \mid G(1) \mid \cdots \mid G(k)] \begin{bmatrix} \mathbf{u}(k) \\ \vdots \\ \mathbf{u}(0) \end{bmatrix} \end{aligned}$$

Proof—cont'd

The norm of the output vector can be written as

$$\begin{aligned}\|\mathbf{y}(k)\| &= \max_s |y_s(k)| \\ &= \max_s |\mathbf{g}_s^T(k) \mathbf{u}|\end{aligned}$$

where $\mathbf{g}_s^T(k)$ is the s row of the matrix $G(k)$. Select the vectors $\mathbf{u}(i)$ to have the r^{th} entry with unity magnitude and with sign opposite to that of $g_{sr}(i)$, where $g_{sr}(i)$ denotes the sr^{th} entry of the impulse response matrix $G(i)$. Using the definition of a matrix norm as the maximum row sum, this gives the output norm

$$\|\mathbf{y}(k)\| = \max_s \sum_{i=0}^k \sum_{r=1}^m |g_{sr}(i)|$$

For BIBO stability, this sum remains finite. But for the violation of Eq. (8.10),

$$\sum_{i=0}^k \|G(k)\| = \sum_{i=0}^k \sum_{s=1}^l \sum_{r=1}^m |g_{sr}(i)| \rightarrow \infty \quad \text{as } k \rightarrow \infty$$

This contradicts the assumption of BIBO stability.

Example 8.3

Determine the BIBO stability of the system with difference equations

$$\begin{aligned}y_1(k+2) + 0.1y_2(k+1) &= u(k) \\ y_2(k+2) + 0.9y_1(k+1) &= u(k)\end{aligned}$$

Solution

To find the impulse response matrix, we first obtain the transfer function matrix

$$\begin{aligned}\begin{bmatrix} Y_1(z) \\ Y_2(z) \end{bmatrix} &= \begin{bmatrix} z^2 & 0.1z \\ 0.9z & z^2 \end{bmatrix}^{-1} U(z) \\ &= \frac{\begin{bmatrix} z^2 & -0.1z \\ -0.9z & z^2 \end{bmatrix}}{z^2(z^2 - 0.09)} U(z) = \frac{\begin{bmatrix} z & -0.1 \\ -0.9 & z \end{bmatrix}}{z(z - 0.3)(z + 0.3)} U(z)\end{aligned}$$

Inverse z-transforming the matrix gives the impulse response

$$G(k) = \begin{bmatrix} 1.6667\{(0.3)^{k-1} - (-0.3)^{k-1}\} & 1.111 \delta(k-1) - 0.5556\{(-0.3)^{k-1} + (0.3)^{k-1}\} \\ 10 \delta(k-1) - 5\{(-0.3)^{k-1} + (0.3)^{k-1}\} & 1.6667\{(0.3)^{k-1} - (-0.3)^{k-1}\} \end{bmatrix},$$

$k \geq 1$, and zero elsewhere.

Example 8.3—cont'd

The entries of the impulse response matrix are all absolutely summable because

$$\sum_{k=0}^{\infty} (0.3)^k = \frac{1}{1 - 0.3} = (0.7)^{-1} < \infty$$

Thus, any induced matrix norm is also summable. For example, if we use the maximum row sum (infinity norm, see Appendix III), we have the norm

$$\|G(k)\| = |10 \delta(k-1) - 5\{(-0.3)^{k-1} + (0.3)^{k-1}\}| + |1.6667\{(0.3)^{k-1} - (-0.3)^{k-1}\}|, k \geq 1$$

The sum is

$$\begin{aligned} \sum_{k=0}^{\infty} \|G(k)\| &= \sum_{k=1}^{\infty} |10 \delta(k-1) - 5\{(-0.3)^{k-1} + (0.3)^{k-1}\}| + |1.6667\{(0.3)^{k-1} - (-0.3)^{k-1}\}| \\ &< 10 + \frac{40}{3} \sum_{k=0}^{\infty} (0.3)^k = 10 + \frac{40}{3 \times 0.7} < \infty \end{aligned}$$

Similar results are obtained using any other norm by virtue of Eq. (8.6). Hence, the impulse response satisfies Eq. (8.10), and the system is BIBO stable.

Although it is possible to test the impulse response for BIBO stability, it is much easier to develop tests based on the transfer function. Theorem 8.3 relates BIBO stability to asymptotic stability and sets the stage for such a test.

Theorem 8.3

If a discrete-time linear system is asymptotically stable, then it is BIBO stable. Furthermore, in the absence of unstable pole-zero cancellation, the system is asymptotically stable if it is BIBO stable.

Proof

Substituting from Eq. (7.98) into Eq. (8.10) gives

$$\begin{aligned} \sum_{k=0}^{\infty} \|G(k)\| &= \|D\| + \sum_{k=0}^{\infty} \left\| \sum_{j=1}^n CZ_j B \lambda_j^k \right\| \\ &\leq \|D\| + \sum_{k=0}^{\infty} \sum_{j=1}^n \|CZ_j B\| |\lambda_j|^k \end{aligned}$$

Proof—cont'd

For an asymptotically stable system, all poles are inside the unit circle. Therefore,

$$\sum_{k=0}^{\infty} \|G(k)\| \leq \|D\| + \sum_{j=1}^n \|CZ_jB\| \frac{1}{1 - |\lambda_j|}$$

which is finite provided that the entries of the matrices are finite. Thus, whenever the product CZ_jB is nonzero for all j , BIBO and asymptotic stability are equivalent. However, the system is BIBO stable but not asymptotically stable if the product is zero for some λ_j with magnitude greater than unity. In other words, a system with unstable input-decoupling, output-decoupling, or input-output-decoupling zeros and all other modes stable is BIBO stable but not asymptotically stable. In the absence of unstable pole-zero cancellation, BIBO stability and asymptotic stability are equivalent.

Using Theorem 8.3, we can test BIBO stability by examining the poles of the transfer function matrix. Recall that, in the absence of pole-zero cancellation, the eigenvalues of the state matrix are the system poles. If the poles are all inside the unit circle, we conclude that the system is BIBO stable. However, we cannot conclude asymptotic stability except in the absence of unstable pole-zero cancellation. We reexamine this issue later in this chapter.

Example 8.4

Test the BIBO stability of the transfer function of the system of Example 8.3.

Solution

The transfer function is

$$G(z) = \frac{\begin{bmatrix} z & -0.1 \\ -0.9 & z \end{bmatrix}}{z(z - 0.3)(z + 0.3)}$$

The system has poles at the origin, 0.3, and -0.3 , all of which are inside the unit circle. Hence, the system is BIBO stable.

8.2 Controllability and stabilizability

When obtaining z -transfer functions from discrete-time state-space equations, we discover that it is possible to have modes that do not affect the input–output relationship.

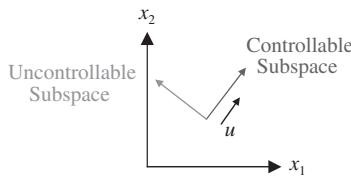


Figure 8.1
Controllable and uncontrollable subspaces.

This occurs as a result of pole-zero cancellation in the transfer function. We now examine this phenomenon more closely to assess its effect on control systems.

The time response of a system is given by the sum of the system modes weighted by eigenvectors. It is therefore important that each system mode be influenced by the input so as to be able to select a time response that meets design specifications. Systems possessing this property are said to be **completely controllable** or simply **controllable**. Modes that cannot be influenced by the control are called **uncontrollable modes**. A more precise definition of controllability is given next.

Definition 8.4: Controllability

A system is said to be controllable if for any initial state $\mathbf{x}(k_0)$ there exists a control sequence $\mathbf{u}(k)$, $k = k_0, k_0+1, \dots, k_f-1$, such that an arbitrary final state $\mathbf{x}(k_f)$ can be reached in finite k_f .

Controllability can be given a geometric interpretation if we consider a second-order system. Fig. 8.1 shows the decomposition of the state plane into a controllable subspace and an uncontrollable subspace. The controllable subspace represents states that are reachable using the control input, whereas the uncontrollable subspace represents states that are not reachable. Some authors prefer to use the term **reachability** instead of “controllability.” Reachable systems are ones where every point in state space is reachable from the origin.¹ This definition is equivalent to our definition of controllability.

Theorem 8.4 establishes that the previous controllability definition is equivalent to the ability of the system input to influence all its modes.

¹ Traditionally, controllable systems are defined as systems where the origin can be reached from every point in the state space. This is equivalent to our definition for continuous-time systems but not for discrete-time systems. For simplicity, we adopt the more practically relevant Definition 8.4.

Theorem 8.4: Controllability condition

A linear time-invariant system is completely controllable if and only if the products $\mathbf{w}_i^T B_d, i = 1, 2, \dots, n$ are all nonzero where \mathbf{w}_i^T is the i^{th} left eigenvector of the state matrix. Furthermore, modes for which the product $\mathbf{w}_i^T B_d$ is zero are uncontrollable.

Proof

Necessity and uncontrollable modes

Definition 8.4 with finite final time k_f guarantees that all system modes are influenced by the control. To see this, we examine the zero-input response

$$\mathbf{x}_{Z_i}(k) = \sum_{i=1}^n Z_i \mathbf{x}(k_0) \lambda_i^k$$

For any eigenvalue λ_i inside the unit circle, the corresponding mode decays to zero exponentially regardless of the influence of the input. However, to go to the zero state in finite time, it is necessary that the input influence all modes. For a zero product $Z_i B_d$, the i^{th} mode is not influenced by the control and can only go to zero asymptotically. Therefore, the i^{th} mode is uncontrollable for zero $Z_i B_d$. In Section 7.4.3, we showed that Z_i is a matrix of rank 1 given by the product of the i^{th} right eigenvector \mathbf{v}_i of the state matrix, and its i^{th} left eigenvector is \mathbf{w}_i . Hence, the product $Z_i B_d$ is given by

$$Z_i B_d = \mathbf{v}_i \mathbf{w}_i^T B_d = \begin{bmatrix} v_{i1} \\ v_{i2} \\ \vdots \\ v_{in} \end{bmatrix} \mathbf{w}_i^T B_d = \begin{bmatrix} v_{i1} \mathbf{w}_i^T \\ v_{i2} \mathbf{w}_i^T \\ \vdots \\ v_{in} \mathbf{w}_i^T \end{bmatrix} B_d$$

where $v_{ij}, j = 1, \dots, n$ are the entries of the i^{th} eigenvector of the state matrix. Therefore, the product $Z_i B_d$ is zero if and only if the product $\mathbf{w}_i^T B_d$ is zero.

Sufficiency

We examine the total response for any initial and terminal state. From Eq. (7.75), we have

$$\mathbf{x} = \mathbf{x}(k) - A_d^k \mathbf{x}(k_0) = \sum_{i=k_0}^{k-1} A_d^{k-i-1} B_d \mathbf{u}(i) \quad (8.11)$$

Using the Cayley–Hamilton theorem, it can be shown that for $k > n$, no “new” terms are added to the preceding summation. By that we mean that the additional terms will be linearly dependent and are therefore superfluous. Thus, if the vector \mathbf{x} cannot be obtained in n sampling periods by proper control sequence selection, it cannot be obtained over a longer duration. We therefore assume that $k = n$ and use the expansion of Eq. (7.92) to obtain

Proof—cont'd

$$\mathbf{x} = \sum_{i=0}^{n-1} \left[\sum_{j=1}^n Z_j \lambda_j^{n-i-1} \right] B_d \mathbf{u}(i) \quad (8.12)$$

The outer summation in Eq. (8.12) can be written in matrix form to obtain

$$\begin{aligned} \mathbf{x} &= \left[\begin{array}{c|c|c|c|c} \sum_{j=1}^n Z_j B_d \lambda_j^{n-1} & \sum_{j=1}^n Z_j B_d \lambda_j^{n-2} & \dots & \sum_{j=1}^n Z_j B_d \lambda_j & \sum_{j=1}^n Z_j B_d \\ \hline \end{array} \right] \begin{bmatrix} \mathbf{u}(0) \\ \vdots \\ \mathbf{u}(1) \\ \vdots \\ \mathbf{u}(n-2) \\ \vdots \\ \mathbf{u}(n-1) \end{bmatrix} \\ &= L\mathbf{u} \end{aligned} \quad (8.13)$$

The matrix B_d is $n \times m$ and is assumed of full rank m . We need to show that the $n \times m.n$ matrix L has full row rank. The constituent matrices Z_i , $i = 1, 2, \dots, n$ are of rank 1, and each has a column linearity independent of the other constituent matrices. Therefore, provided that the individual products $Z_i B_d$ are all nonzero, each submatrix of L has rank 1, and L has n linearly independent columns (i.e., L has full rank n). Because the rank of a matrix cannot exceed the number of rows for fewer rows than columns, we have for any vector \mathbf{x}

$$\text{rank}\{L|\mathbf{x}\} = \text{rank}\{L\} = n \quad (8.14)$$

This condition guarantees the existence of a solution \mathbf{u} to Eq. (8.12). However, the solution is, in general, nonunique because the number of elements of \mathbf{u} (unknowns) is $m.n > n$. A unique solution exists only in the single-input case ($m = 1$). It follows that if the products $Z_i B_d$, $i = 1, 2, \dots, n$ are all nonzero, then one can solve for the vector of input sequences needed to reach any $\mathbf{x}(k_f)$ from any $\mathbf{x}(k_0)$. As discussed in the proof of necessity, the product $Z_i B_d$ is zero if and only if the product $\mathbf{w}_i^T B_d$ is zero.

Theorem 8.5 gives a simpler controllability test, but, as first stated, it does not provide a means of determining which of the system modes is uncontrollable.

Theorem 8.5: Controllability rank condition

A linear time-invariant system is completely controllable if and only if the $n \times m.n$ controllability matrix

$$\mathcal{C} = [B_d \mid A_d B_d \mid \dots \mid A_d^{n-1} B_d] \quad (8.15)$$

has rank n .

Proof

We first write Eq. (8.13) in the form (recall the Cayley–Hamilton theorem)

$$\mathbf{x} = \left[B_d \mid A_d B_d \mid \dots \mid A_d^{n-1} B_d \right] \begin{bmatrix} \mathbf{u}(n-1) \\ \vdots \\ \mathbf{u}(n-2) \\ \vdots \\ \cdot \\ \vdots \\ \mathbf{u}(1) \\ \vdots \\ \mathbf{u}(0) \end{bmatrix} = \mathcal{C}\mathbf{u} \quad (8.16)$$

We now have a system of linear equations for which Eq. (8.15) is a necessary and sufficient condition for a solution \mathbf{u} to exist for any \mathbf{x} . Hence, Eq. (8.15) is a necessary and sufficient condition for controllability.

If the controllability matrix \mathcal{C} is rank deficient, then the rank deficiency is equal to the number of linearly independent row vectors that are mapped to zero on multiplication by \mathcal{C} . These row vectors are the uncontrollable states of the system as well as the left eigenvectors of the state matrix A . The vectors are also the transpose of the right eigenvectors of A^T . Thus, the rank test can be used to determine the uncontrollable modes of the system if the eigenstructure of A^T is known.

Example 8.5

Determine the controllability of the following state equation:

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} -2 & -2 & 0 \\ 0 & 0 & 1 \\ 0 & -0.4 & -0.5 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \mathbf{u}(k)$$

Solution

The controllability matrix for the system is

$$\mathcal{C} = \left[B_d \mid A_d B_d \mid A_d^2 B_d \right]$$

$$= \left[\begin{array}{c|cc|cc|cc} 1 & 0 & -2 & -2 & 2 & 2 \\ 0 & 1 & 1 & 1 & -0.5 & -0.9 \\ 1 & 1 & -0.5 & -0.9 & -0.15 & 0.05 \end{array} \right]$$

Example 8.5—cont'd

The matrix has rank 3, which implies that the third-order system is controllable. In fact, the first three columns of the matrix are linearly independent, and the same conclusion can be reached without calculating the entire controllability matrix. In general, one can gradually compute more columns until n linearly independent columns are obtained to conclude controllability for an n^{th} -order system.

Although the controllability tests are given here for discrete-time systems, they are applicable to continuous-time systems. This fact is used in the following two examples.

Example 8.6

Show that the following state equation is uncontrollable, and determine the uncontrollable mode. Then obtain the discrete-time state equation, and verify that it is also uncontrollable.

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

Solution

The system has a zero eigenvalue with multiplicity 2. The controllability matrix for the system is

$$\begin{aligned} \mathcal{C} &= \left[\begin{array}{c|cc} B & AB \end{array} \right] \\ &= \left[\begin{array}{c|cc} -1 & 0 \\ 0 & 0 \end{array} \right] \end{aligned}$$

The matrix has rank 1 (less than 2), which implies that the second-order system has one (2–1) uncontrollable mode. In fact, integrating the state equation reveals that the second state variable is governed by the equation

$$x_2(t) = x_2(t_0), \quad t \geq t_0$$

Therefore, the second state variable cannot be altered using the control and corresponds to an uncontrollable mode. The first state variable is influenced by the input and can be arbitrarily controlled.

The state-transition matrix for this system can be obtained directly using the series expansion because the matrix A^i is zero for $i > 1$. Thus, we have the discrete-time state matrix

$$A_d = e^{AT} = I_2 + AT = \begin{bmatrix} 1 & T \\ 0 & 1 \end{bmatrix}$$

Example 8.6—cont'd

The discrete-time input matrix is

$$B_d = \int_0^T e^{A\tau} Bd\tau = \int_0^T [I_2 + A\tau] \begin{bmatrix} -1 \\ 0 \end{bmatrix} d\tau = \begin{bmatrix} -T \\ 0 \end{bmatrix}$$

The controllability matrix for the discrete-time system is

$$\mathcal{C} = \begin{bmatrix} B_d & A_d B_d \end{bmatrix} = \begin{bmatrix} -T & -T \\ 0 & 1 \end{bmatrix}$$

As with the continuous time system, the matrix has rank 1 (less than 2), which implies that the second-order system is uncontrollable.

The solution of the difference equation for the second state variable gives

$$x_2(k) = x_2(k_0), \quad k \geq k_0$$

with initial time k_0 . Thus, the second state variable corresponds to an uncontrollable unity mode (e^0), whereas the first state variable is controllable as with the continuous system.

Example 8.7

Show that the state equation for the motor system

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -10 & -11 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} u$$

with DAC and ADC is uncontrollable, and determine the uncontrollable modes. Obtain the transfer functions of the continuous-time and discrete-time systems, and relate the uncontrollable modes to the poles of the transfer function.

Solution

From Example 7.7 the state-transition matrix of the system is

$$e^{At} = \begin{bmatrix} 10 & 11 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \frac{e^0}{10} + \begin{bmatrix} 0 & -10 & -1 \\ 0 & 10 & 1 \\ 0 & -10 & -1 \end{bmatrix} \frac{e^{-t}}{9} + \begin{bmatrix} 0 & 1 & 1 \\ 0 & -10 & -10 \\ 0 & 100 & 100 \end{bmatrix} \frac{e^{-10t}}{90}$$

Example 8.7—cont'd

Multiplying by the input matrix, we obtain

$$e^{At}B = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} e^0 + \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \frac{e^{-t}}{9} + \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \frac{e^{-10t}}{90} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \frac{e^0}{10}$$

Thus, the modes e^{-t} and e^{-10t} are both uncontrollable for the analog subsystem. The two modes are also uncontrollable for the digital system, including DAC and ADC.

Using Eq. 7.66, the input matrix for the discrete-time system is

$$B_d = Z_1 BT + Z_2 B(1 - e^{-T}) + Z_3 B(1 - e^{-10T}) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} T$$

and the matrix $A_d = e^{AT}$. The controllability matrix for the discrete-time system is

$$\begin{aligned} \mathcal{C} &= \left[B_d \mid A_d B_d \mid A_d^2 B_d \right] \\ &= \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} T \end{aligned}$$

which clearly has rank 1, indicating that there are 2 (i.e., 3–1) uncontrollable modes. The left eigenvectors corresponding to the eigenvalues e^{-T} and e^{-10T} have zero first entry, and their product with the controllability matrix \mathcal{C} is zero. Hence, the two corresponding modes are not controllable. The zero eigenvalue has the left eigenvector $[1, 0, 0]$, and its product with \mathcal{C} is nonzero. Thus, the corresponding mode is controllable.

The transfer function of the continuous-time system with output x_1 is

$$G(s) = C[sI_3 - A]^{-1}B = \frac{1}{s}$$

This does not include the uncontrollable modes (input-decoupling modes) because they cancel when the resolvent matrix and the input matrix are multiplied.

The z-transfer function corresponding to the system with DAC and ADC is

$$G_{ZAS}(z) = \frac{z-1}{z} \mathcal{Z} \left\{ \mathcal{Z}^{-1} \left[\frac{G(s)}{s} \right] \right\} = \frac{T}{z-1}$$

The reader can easily verify that the transfer function is identical to that obtained using the discrete-time state-space representation. It includes the pole at e^0 but has no poles at e^{-T} or e^{-10T} .

8.2.1 MATLAB commands for controllability testing

The MATLAB commands to calculate the controllability matrix and determine its rank are

```
>> c = ctrb(A, C)
>> rank(c)
```

8.2.2 Controllability of systems in normal form

Checking controllability is particularly simple if the system is given in normal form—that is, if the state matrix is in the form

$$A = \text{diag}\{\lambda_1(A), \lambda_2(A), \dots, \lambda_n(A)\}$$

The corresponding state-transition matrix is

$$e^{At} = \mathcal{L}\left\{[sI_n - A]^{-1}\right\} = \text{diag}\left\{e^{\lambda_1(A)t}, e^{\lambda_2(A)t}, \dots, e^{\lambda_n(A)t}\right\} \quad (8.17)$$

The discrete-time state matrix is

$$\begin{aligned} A_d &= e^{AT} = \text{diag}\left\{e^{\lambda_1(A)T}, e^{\lambda_2(A)T}, \dots, e^{\lambda_n(A)T}\right\} \\ &= \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\} \end{aligned} \quad (8.18)$$

Theorem 8.6 establishes controllability conditions for a system in normal form.

Theorem 8.6: Controllability of systems in normal form

A system in normal form is controllable if and only if its input matrix has no zero rows. Furthermore, if the input matrix has a zero row, then the corresponding mode is uncontrollable.

Proof

Necessity and uncontrollable modes

For a system in normal form, the state equations are in the form

$$x_i(k+1) = \lambda_i x_i(k) + \mathbf{b}_i^T \mathbf{u}(k), \quad i = 1, 2, \dots, n$$

where \mathbf{b}_i^T is the i^{th} row of the input matrix B_d . For a zero row, the system is unforced and can only converge to zero asymptotically for $|\lambda_i|$ inside the unit circle. Because controllability requires convergence to any final state (including the origin) in finite time, the i^{th} mode is uncontrollable for zero \mathbf{b}_i^T .

Proof—cont'd**Sufficiency**

From the sufficiency proof of Theorem 8.4, we obtain Eq. (8.13) to be solved for the vector \mathbf{u} of controls over the period $k = 0, 1, \dots, n-1$. The solution exists if the matrix L in Eq. (8.13) has rank n . For a system in normal form, the state-transition matrix is in the form

$$A_d^k = \text{diag}\{\lambda_i^k\}$$

and the $n \times n.m$ matrix L is in the form

$$L = [\text{diag}\{\lambda_j^{n-1}\}B_d \mid \text{diag}\{\lambda_j^{n-2}\}B_d \mid \dots \mid \text{diag}\{\lambda_j\}B_d \mid B_d]$$

Substituting for B_d in terms of its rows and using the rules for multiplying partitioned matrices, we obtain

$$L = \begin{bmatrix} \lambda_1^{n-1}\mathbf{b}_1^T & \lambda_1^{n-2}\mathbf{b}_1^T & \dots & \mathbf{b}_1^T \\ \lambda_2^{n-1}\mathbf{b}_2^T & \lambda_2^{n-2}\mathbf{b}_2^T & \dots & \mathbf{b}_2^T \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_n^{n-1}\mathbf{b}_n^T & \lambda_1^{n-2}\mathbf{b}_n^T & \dots & \mathbf{b}_n^T \end{bmatrix}$$

For a matrix B_d with no zero rows, the rows of L are linearly independent and the matrix has full rank n . This guarantees the existence of a solution \mathbf{u} to Eq. (8.13).

Example 8.8

Determine the controllability of the system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} -2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -0.5 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \mathbf{u}(k)$$

Solution

The system is in normal form, and its input matrix has no zero rows. Hence, the system is completely controllable.

8.2.3 Stabilizability

The system in Example 8.8 is controllable but has one unstable eigenvalue (outside the unit circle) at (-2) . This clearly demonstrates that controllability implies the ability to control the modes of the system regardless of their stability. If the system has

uncontrollable modes, then these modes cannot be influenced by the control and may or may not decay to zero asymptotically. Combining the independent concepts of stability and controllability gives the following definition.

Definition 8.5: Stabilizability

A system is said to be stabilizable if all its uncontrollable modes are asymptotically stable.

Physical systems are often stabilizable rather than controllable. This poses no problem provided that the uncontrollable dynamics decay to zero sufficiently fast so as not to excessively slow down the system.

Example 8.9

Determine the controllability and stabilizability of the system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} -2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -0.5 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 1 & 1 \end{bmatrix} \mathbf{u}(k)$$

Solution

The system is in normal form, and its input matrix has one zero row corresponding to its zero eigenvalue. Hence, the system is uncontrollable. However, the uncontrollable mode at the origin is asymptotically stable, and the system is therefore stabilizable.

An alternative procedure to determine the stabilizability of the system is to transform it to a form that partitions the state into a controllable part and an uncontrollable part, and then determine the asymptotic stability of the uncontrollable part. This can be done by exploiting the following theorem.

Theorem 8.7: Standard form for uncontrollable systems

Consider the pair (A_d, B_d) with n_c controllable modes and $n-n_c$ uncontrollable modes and the transformation matrix

$$T_c = \begin{bmatrix} \mathbf{q}_1^T \\ \vdots \\ \mathbf{q}_{n_c}^T \\ \hline \mathbf{q}_{n_c+1}^T \\ \vdots \\ \mathbf{q}_n^T \end{bmatrix}^{-1} = \begin{bmatrix} Q_1 \\ \hline Q_2 \end{bmatrix}^{-1} \quad (8.19)$$

Theorem 8.7: Standard form for uncontrollable systems—cont'd

where $\{\mathbf{q}_i^T, i = n_c + 1, \dots, n\}$ are linearly independent vectors in the null space of the controllability matrix \mathcal{C} and $\{\mathbf{q}_i^T, i = 1, \dots, n_c\}$ are arbitrary linearly independent vectors selected to make T_c nonsingular.

The transformed state-transition matrix \mathbf{A} and the transformed input matrix \mathbf{B} have the following form:

$$\mathbf{A} = \begin{bmatrix} A_c & | & A_{uc} \\ \hline 0 & | & A_{uc} \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} B_c \\ \hline 0 \end{bmatrix} \quad (8.20)$$

where A_c is an $n_c \times n_c$ matrix, B_c is an $n_c \times m$ matrix, and the pair (A_c, B_c) is controllable.

Proof

If the rank of the controllability matrix \mathcal{C} is n_c , then the dimension of its null space is $n - n_c$ and there exist linearly independent vectors $\{\mathbf{q}_i^T, i = n_c + 1, \dots, n\}$ satisfying $\mathbf{q}_i^T \mathcal{C} = \mathbf{0}^T$. The transformation T_c can be obtained by adding n_c linearly independent row vectors. The null space of the controllability matrix is the unreachable space of the system and is spanned by the set of left eigenvectors corresponding to uncontrollable modes. Hence, the vectors $\{\mathbf{q}_i^T, i = n_c + 1, \dots, n\}$ are linear combinations of the left eigenvectors of the matrix A_d corresponding to uncontrollable modes. We can therefore write the transformation matrix in the form

$$T_c = \begin{bmatrix} Q_1 \\ \hline \overline{Q}_2 W_{uc} \end{bmatrix}^{-1}$$

where the matrix of coordinates \overline{Q}_2 is nonsingular and W_{uc} is a matrix whose rows are the left eigenvectors of A , $\{\mathbf{w}_i^T, i = 1, \dots, n - n_c\}$, corresponding to uncontrollable modes. By the eigenvector controllability test, if the i^{th} mode is uncontrollable, we have

$$\mathbf{w}_i^T B_d = \mathbf{0}^T, i = 1, \dots, n - n_c$$

Therefore, transforming the input matrix gives

$$\mathbf{B} = T_c^{-1} B_d = \begin{bmatrix} Q_1 \\ \hline \overline{Q}_2 W_{uc} \end{bmatrix}^{-1} B_d = \begin{bmatrix} B_c \\ \hline \mathbf{0} \end{bmatrix}$$

Next, we show that the transformation yields a state matrix in the form of Eq. (8.20). We first recall that the left and right eigenvectors can be scaled to satisfy the condition

$$\mathbf{w}_i^T \mathbf{v}_i = 1, \quad i = 1, \dots, n - n_c$$

Proof—cont'd

If we combine the condition for all the eigenvectors corresponding to uncontrollable modes, we have

$$W_{uc}V_{uc} = [\mathbf{w}_i^T \mathbf{v}_j] = I_{n-n_c}$$

We write the transformation matrix in the form

$$T_c = [R_1 \mid R_2] = [R_1 \mid V_{uc}\bar{R}_2]$$

which satisfies the condition

$$T_c^{-1}T_c = \begin{bmatrix} \frac{Q_1}{\bar{Q}_2 W_{uc}} \\ \bar{Q}_2 W_{uc} \end{bmatrix} [R_1 \mid V_{uc}\bar{R}_2] = \begin{bmatrix} I_{n_c} & 0 \\ 0 & I_{n-n_c} \end{bmatrix}$$

Because \bar{Q}_2 is nonsingular, we have the equality

$$W_{uc}R_1 = 0$$

The similarity transformation gives

$$\begin{aligned} T_c^{-1}A_dT_c &= \begin{bmatrix} \frac{Q_1}{\bar{Q}_2 W_{uc}} \\ \bar{Q}_2 W_{uc} \end{bmatrix} A_d [R_1 \mid V_{uc}\bar{R}_2] \\ &= \begin{bmatrix} Q_1 A_d R_1 & Q_1 A_d V_{uc} \bar{R}_2 \\ \bar{Q}_2 W_{uc} A_d R_1 & \bar{Q}_2 W_{uc} A_d V_{uc} \bar{R}_2 \end{bmatrix} \end{aligned} \quad (8.21)$$

Because premultiplying the matrix A_d by a left eigenvector gives the same eigenvector scaled by an eigenvalue, we have the condition

$$W_{uc}A_d = A_u W_{uc}, A_u = \text{diag}\{\lambda_{u1}, \dots, \lambda_{u,n-n_c}\}$$

with λ_u denoting an uncontrollable eigenvalue. This simplifies our matrix equality Eq. (8.21) to

$$\begin{aligned} T_c^{-1}A_dT_c &= \begin{bmatrix} Q_1 A_d R_1 & Q_1 A_d V_{uc} \bar{R}_2 \\ \bar{Q}_2 A_u W_{uc} R_1 & \bar{Q}_2 A_u W_{uc} V_{uc} \bar{R}_2 \end{bmatrix} = \begin{bmatrix} Q_1 A_d R_1 & Q_1 A_d V_{uc} \bar{R}_2 \\ \bar{Q}_2 A_u W_{uc} R_1 & \bar{Q}_2 A_u \bar{R}_2 \end{bmatrix} \\ &= \begin{bmatrix} Q_1 A_d R_1 & Q_1 A_d V_{uc} \bar{R}_2 \\ \mathbf{0} & \bar{Q}_2 A_u \bar{R}_2 \end{bmatrix} = \begin{bmatrix} A_c & A_{uc} \\ \mathbf{0} & A_{uc} \end{bmatrix} \end{aligned}$$

Because all the uncontrollable modes correspond to the eigenvalues of the matrix A_{uc} , the remaining modes are all controllable, and the pair (A_c, B_c) is completely controllable.

Because the pair (A_d, B_d) has n_c controllable modes, the rank of the controllability matrix is n_c and the transformation matrix T_c is guaranteed to exist. The columns $\{\mathbf{q}_i^T, i = n_c + 1, \dots, n\}$ of the inverse transformation matrix have a geometric interpretation. They form a basis set for the unobservable subspace of the system and can be selected as the eigenvectors of the uncontrollable modes of the system. The remaining columns of T_c form a basis set of the controllable subspace and can be conveniently selected as linearly independent columns $\{\mathbf{q}_i^T, i = 1, \dots, n\}$ of the controllability matrix. Clearly, these vectors are not in the null space of the controllability matrix as they satisfy the equation $\mathbf{q}_i^T \mathcal{C} \neq 0, i = 1, \dots, n$. Starting with these vectors, the transformation matrix can be formed by selecting $n - n_c$ additional linearly independent vectors.

The similarity transformation is a change of basis that allows us to separate the controllable subspace from the uncontrollable subspace. After transformation, we can check the stabilizability of the system by verifying that all the eigenvalues of the matrix A_{uc} are inside the unit circle.

Example 8.10

Determine the controllability, stability, and stabilizability of the system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} -1.85 & 4.2 & -0.15 \\ -0.3 & 0.1 & 0.3 \\ -1.35 & 4.2 & -0.65 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 4 & 3 \\ 1 & 1 \\ 2 & 1 \end{bmatrix} \mathbf{u}(k)$$

Solution

We first find the controllability matrix of the system

$$\begin{aligned} \mathcal{C} &= \left[\begin{array}{c|cc} B_d & A_d B_d & A_d^2 B_d \end{array} \right] \\ &= \left[\begin{array}{cc|cc|cc} 4 & 3 & -3.5 & -1.5 & 4.75 & 0.75 \\ 1 & 1 & -0.5 & -0.5 & 0.25 & 0.25 \\ 2 & 1 & -2.5 & -0.5 & 4.25 & 0.25 \end{array} \right] \end{aligned}$$

The matrix has rank 2, and the first two columns are linearly independent. The dimension of the null space is $3 - 2 = 1$. We form a transformation matrix of the two first columns of the controllability matrix and a third linearly independent column, giving

$$T_c = \left[\begin{array}{cc|c} 4 & 3 & 0 \\ 1 & 1 & 0 \\ 2 & 1 & 1 \end{array} \right]$$

Example 8.10—cont'd

The system is transformed to

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \left[\begin{array}{cc|c} -2 & 0 & -1.05 \\ 1.5 & -0.5 & 1.35 \\ 0 & 0 & 0.1 \end{array} \right] \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{u}(k)$$

The system has one uncontrollable mode corresponding to the eigenvalue at 0.1, which is inside the unit circle. The system is therefore stabilizable but not controllable. The two other eigenvalues are at -2 , outside the unit circle, and at -0.5 , inside the unit circle, and both corresponding modes are controllable. The system is unstable because one eigenvalue is outside the unit circle.

8.3 Observability and detectability

To effectively control a system, it may be advantageous to use the system state to select the appropriate control action. Typically, the output or measurement vector for a system includes a combination of some but not all the state variables. The system state must then be estimated from the time history of the measurements and controls. However, state estimation is only possible with proper choice of the measurement vector. Systems whose measurements are so chosen are said to be **completely observable** or simply **observable**. Modes that cannot be detected using the measurements and the control are called **unobservable modes**. A more precise definition of observability is given next.

Definition 8.6: Observability

A system is said to be observable if any initial state $\mathbf{x}(k_0)$ can be estimated from the control sequence $\mathbf{u}(k)$, $k = k_0, k_0+1, \dots, k_f-1$ and the measurements $\mathbf{y}(k)$, $k = k_0, k_0+1, \dots, k_f$.

As in the case of controllability, observability can be given a geometric interpretation if we consider a second-order system. Fig. 8.2 shows the decomposition of the state plane into an observable subspace and an unobservable subspace. The observable subspace includes all initial states that can be identified using the measurement history, whereas the unobservable subspace includes states that are indistinguishable.

The following theorem establishes that this observability definition is equivalent to the ability to estimate all system modes using its controls and measurements.

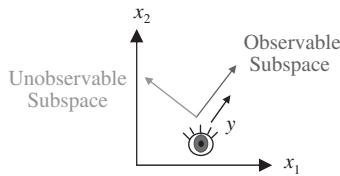


Figure 8.2
Observable and unobservable subspaces.

Theorem 8.8: Observability

A system is observable if and only if $C\mathbf{v}_i$ is nonzero for $i = 1, 2, \dots, n$, where \mathbf{v}_i is the i^{th} eigenvector of the state matrix. Furthermore, if the product $C\mathbf{v}_i$ is zero, then the i^{th} mode is unobservable.

Proof

We first define the vector

$$\begin{aligned}\bar{\mathbf{y}}(k) &= \mathbf{y}(k) - C \sum_{i=k_0}^{k-1} A_d^{k-i-1} B_d \mathbf{u}(i) - D \mathbf{u}(k) \\ &= CA_d^{k-k_0} \mathbf{x}(k_0) = \sum_{i=1}^n CZ_i \lambda_i^{k-k_0} \mathbf{x}(k_0)\end{aligned}$$

Sufficiency

Stack the output vectors $\bar{\mathbf{y}}(i), i = 1, \dots, k$ to obtain

$$\begin{bmatrix} \bar{\mathbf{y}}(0) \\ \vdots \\ \bar{\mathbf{y}}(1) \\ \vdots \\ \bar{\mathbf{y}}(n-2) \\ \vdots \\ \bar{\mathbf{y}}(n-1) \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^n CZ_j \\ \vdots \\ \sum_{j=1}^n CZ_j \lambda_j \\ \vdots \\ \sum_{j=1}^n CZ_j \lambda_j^{n-2} \\ \vdots \\ \sum_{j=1}^n CZ_j \lambda_j^{n-1} \end{bmatrix} \mathbf{x}(k_0) = L \mathbf{x}(k_0) \quad (8.22)$$

The products CZ_i are

$$\begin{aligned}CZ_i &= C\mathbf{v}_i w_i^T = C\mathbf{v}_i [w_{i1} \quad w_{i2} \quad \dots \quad w_{in}] \\ &= C [w_{i1} \mathbf{v}_i \quad w_{i2} \mathbf{v}_i \quad \dots \quad w_{in} \mathbf{v}_i]\end{aligned}$$

where $w_{ij}, j = 1, \dots, n$ are the entries of the i^{th} left eigenvector, and \mathbf{v}_i is the i^{th} right eigenvector of the state matrix A_d (see Section 7.4.3). Therefore, if the products are nonzero, the matrix L has n linearly independent columns (i.e., has rank n). An $n \times n$ submatrix of L

Proof—cont'd

can be formed by discarding the dependent rows. This leaves us with n equations in the n unknown entries of the initial condition vector. A unique solution can thus be obtained.

Necessity

Let $\mathbf{x}(k_0)$ be in the direction of \mathbf{v}_i , the i^{th} eigenvector of A_d , and let CZ_i be zero. Then the vector $\mathbf{y}(k)$ is zero for any k regardless of the amplitude of the initial condition vector (recall that $Z\mathbf{v}_j = \mathbf{0}$ whenever $i \neq j$). Thus, all initial vectors in the direction \mathbf{v}_i are indistinguishable using the output measurements, and the system is not observable.

The next theorem establishes an observability rank test that can be directly checked using the state and output matrices.

Theorem 8.9: Observability rank condition

A linear time-invariant system is completely observable if and only if the $l.n \times n$ observability matrix

$$\mathcal{O} = \begin{bmatrix} C \\ - \\ - \\ CA_d \\ - \\ - \\ \vdots \\ - \\ - \\ CA_d^{n-1} \end{bmatrix} \quad (8.23)$$

has rank n .

Proof

We first write Eq. (8.22) in the form

$$\begin{bmatrix} \bar{\mathbf{y}}(0) \\ - \\ - \\ \bar{\mathbf{y}}(1) \\ - \\ - \\ \vdots \\ - \\ - \\ \bar{\mathbf{y}}(n-2) \\ - \\ - \\ \bar{\mathbf{y}}(n-1) \end{bmatrix} = \begin{bmatrix} C \\ - \\ - \\ CA_d \\ - \\ - \\ \vdots \\ - \\ - \\ CA_d^{n-1} \end{bmatrix} \mathbf{x}(k_0) = \mathcal{O} \mathbf{x}(k_0) \quad (8.24)$$

We now have a system of linear equations that can include, at most, n linearly independent equations and, by the Cayley–Hamilton theorem, adding more entries to the stacked output vector can only result in linearly dependent rows. If n independent equations exist, their

Proof—cont'd

coefficients can be used to form an $n \times n$ invertible matrix. The rank condition Eq. (8.23) is a necessary and sufficient condition for a unique solution $\mathbf{x}(k_0)$ to exist. Thus, Eq. (8.23) is a necessary and sufficient condition for observability.

If the observability matrix \mathcal{O} is rank deficient, then the rank deficiency is equal to the number of linearly independent column vectors that are mapped to zero on multiplication by \mathcal{O} . This number equals the number of unobservable states of the system, and the column vectors mapped to zero are the right eigenvectors of the state matrix. Thus, the rank test can be used to determine the unobservable modes of the system if the eigenstructure of the state matrix is known.

Example 8.11

Determine the observability of the system using two different tests:

$$A = \begin{bmatrix} & \mathbf{0}_{2 \times 1} & I_2 \\ 0 & & -3 & 4 \end{bmatrix} \quad C = [0 \ 0 \ 1]$$

If the system is not completely observable, determine the unobservable modes.

Solution

Because the state matrix is in companion form, its characteristic equation is easily obtained from its last row. The characteristic equation is

$$\lambda^3 - 4\lambda^2 + 3\lambda = \lambda(\lambda - 1)(\lambda - 3) = 0$$

Hence, the system eigenvalues are $\{0, 1, 3\}$.

The companion form of the state matrix allows us to write the modal matrix of eigenvectors as the Van der Monde matrix:

$$V = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 3 \\ 0 & 1 & 9 \end{bmatrix}$$

Observability is tested using the product of the output matrix and the modal matrix:

$$CV = [0 \ 0 \ 1] \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 3 \\ 0 & 1 & 9 \end{bmatrix} = [0 \ 1 \ 9]$$

The product of the output matrix and the eigenvector for the zero eigenvalue is zero. We conclude that the system has an output-decoupling zero at zero (i.e., one unobservable mode).

Example 8.11—cont'd

The observability matrix of the system is

$$\mathcal{O} = \begin{bmatrix} C \\ CA_d \\ CA_d^2 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & -3 & 4 \\ 0 & -12 & 13 \end{bmatrix}$$

which has rank $2 = 3 - 1$. Hence, the system has one unobservable mode. The product of the observability matrix and the eigenvector for the zero eigenvalue is zero. Hence, it corresponds to the unobservable mode of the system.

8.3.1 MATLAB commands

The MATLAB commands to test observability are

```
>> o = obsv(A, C)%Obtain the observability matrix
>> rank(o)
```

8.3.2 Observability of systems in normal form

The observability of a system can be easily checked if the system is in normal form by exploiting Theorem 8.10.

Theorem 8.10: Observability of normal form systems

A system in normal form is observable if and only if its output matrix has no zero columns. Furthermore, if the input matrix has a zero column, then the corresponding mode is unobservable.

Proof

The proof is similar to that of the Controllability Theorem 8.6 and is left as an exercise.

8.3.3 Detectability

If a system is not observable, it is preferable that the unobservable modes be stable. This property is called **detectability**.

Definition 8.7: Detectability

A system is detectable if all its unobservable modes decay to zero asymptotically.

Example 8.12

Determine the observability and detectability of the system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} -0.4 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & -2 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 1 & 1 \end{bmatrix} \mathbf{u}(k)$$

$$\mathbf{y}(k) = \begin{bmatrix} 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}$$

Solution

The system is in normal form, and its output matrix has one zero column corresponding to its eigenvalue -2 . Hence, the system is unobservable. The unobservable mode at $-2, |-2| > 1$, is unstable, and the system is therefore not detectable.

Similar to the concepts described for stabilizability, we can determine the detectability of the system by transforming it to a form that partitions the state into an observable part and an unobservable part, and then determining the asymptotic stability of the unobservable part. We do this with Theorem 8.11.

Theorem 8.11: Standard form for unobservable systems

Consider the pair (A_d, C_d) with n_0 observable modes and $n-n_0$ unobservable modes and the $n \times n$ transformation matrix

$$T_o = \left[\begin{array}{c|c} T_{o1} & T_{o2} \end{array} \right] = \left[\begin{array}{cccc|cccc} \mathbf{t}_1 & \dots & \mathbf{t}_{n-n_0} & \mathbf{t}_{n-n_0+1} & \dots & \mathbf{t}_n \end{array} \right]$$

where $\{\mathbf{t}_i, i = 1, \dots, n-n_0\}$ are linearly independent vectors in the null space of the observability matrix \mathcal{O} and $\{\mathbf{t}_i, i = n-n_0+1, \dots, n\}$ are arbitrary vectors selected to make

Theorem 8.11: Standard form for unobservable systems—cont'd

To invertible. The transformed state-transition matrix \mathbf{A} and the transformed input matrix \mathbf{C} have the following form:

$$\mathbf{A} = \begin{bmatrix} A_u & | & A_{uo} \\ \hline 0_{n_0 \times n-n_0} & | & A_o \end{bmatrix} \quad \mathbf{C} = \begin{bmatrix} \mathbf{0}_{l \times n-n_0} & | & C_o \end{bmatrix}$$

where A_o is an $n_0 \times n_0$ matrix, C_o is an $l \times n_0$ matrix, and the pair (A_o, C_o) is observable.

Proof

The proof is similar to that of Theorem 8.7, and it is left as an exercise.

Because the pair (A_d, C_d) has n_o observable modes, the rank of the observability matrix is n_o and the transformation matrix T_o is guaranteed to exist. The columns $\{\mathbf{t}_i, i = 1, \dots, n - n_o\}$ of the transformation matrix have a geometric interpretation. They form a basis set for the unobservable subspace of the system and can be selected as the eigenvectors of the unobservable modes of the system. The remaining columns of T_o form a basis set of the observable subspace and can be selected as the transpose of linearly independent rows \mathbf{r}^T of the observability matrix. Clearly, these vectors are not in the null space of the observability matrix as they satisfy the equation $\mathcal{O} \mathbf{r} \neq \mathbf{0}$. Starting with these vectors, the transformation matrix can be formed by selecting $n - n_o$ additional linearly independent vectors.

The similarity transformation is a change of basis that allows us to separate the observable subspace from the unobservable subspace. After transformation, we can check the detectability of the system by verifying that all the eigenvalues of the matrix A_u are inside the unit circle.

Example 8.13

Determine the observability and detectability of the system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} 2.0 & 4.0 & 2.0 \\ -1.1 & -2.5 & -1.15 \\ 2.6 & 6.8 & 2.8 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} \mathbf{u}(k)$$

$$y(k) = \begin{bmatrix} 2 & 10 & 3 \\ 1 & 8 & 3 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}$$

Example 8.13—cont'd**Solution**

We first find the observability matrix of the system

$$\mathcal{O} = \begin{bmatrix} C \\ CA_d \\ CA_d^2 \end{bmatrix} = \begin{bmatrix} 2 & 10 & 3 \\ 1 & 8 & 3 \\ 0.8 & 3.4 & 0.9 \\ 1.0 & 4.4 & 1.2 \\ 0.2 & 0.82 & 0.21 \\ 0.28 & 0.16 & 0.3 \end{bmatrix}$$

The matrix has rank 2, and the system is not completely observable. The dimension of the null space is $3 - 2 = 1$, and there is only one vector satisfying $\mathcal{O} \mathbf{v}_1 = \mathbf{0}$. It is the eigenvector \mathbf{v}_1 corresponding to the unobservable mode and is given by

$$\mathbf{v}_1 = \begin{bmatrix} 1 \\ -0.5 \\ 1 \end{bmatrix}$$

The transformation matrix T_o can then be completed by adding any two linearly independent columns—for example,

$$T_o = \begin{bmatrix} 1 & | & 1 & 0 \\ -0.5 & | & 0 & 1 \\ 1 & | & 0 & 0 \end{bmatrix}$$

The system is transformed to

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} 2 & | & 2.6 & 6.8 \\ 0 & | & -0.6 & -2.8 \\ 0 & | & 0.2 & 0.9 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ 1 & -1 \\ 1 & 1.5 \end{bmatrix} \mathbf{u}(k)$$

$$y(k) = \begin{bmatrix} 0 & | & 2 & 10 \\ 0 & | & 1 & 8 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}$$

The system has one unobservable mode corresponding to the unstable eigenvalue at 2, which is outside the unit circle. The system is therefore not detectable.

8.4 Poles and zeros of multivariable systems

As for single-input–single-output (SISO) systems, poles and zeros determine the stability, controllability, and observability of a multivariable system. For SISO systems, zeros are the zeros of the numerator polynomial of the scalar transfer function, whereas poles are

the zeros of the denominator polynomial. For multi-input–multi-output (MIMO) systems, the transfer function is not scalar, and it is no longer sufficient to determine the zeros and poles of individual entries of the transfer function matrix. In fact, element zeros do not play a major role in characterizing multivariable systems and their properties beyond their effect on the shape of the time response of the system. In our discussion of transfer function matrices in [Section 7.8](#), we mention three important types of multivariable zeros:

1. Input-decoupling zeros, which correspond to uncontrollable modes.
2. Output-decoupling zeros, which correspond to unobservable modes.
3. Input-output-decoupling zeros, which correspond to modes that are neither controllable nor observable.

We can obtain those zeros and others, as well as the system poles, from the transfer function matrix.

8.4.1 Poles and zeros from the transfer function matrix

We first define system poles and zeros based on a transfer function matrix of a MIMO system.

Definition 8.8: Poles

Poles are the roots of the least common denominator of all nonzero minors of all orders of the transfer function matrix.

The least common denominator of the preceding definition is known as the **pole polynomial** and is essential for the determination of the poles. The pole polynomial is in general not equal to the least common denominator of all nonzero elements, known as the **minimal polynomial**. The MATLAB command **pole** gives the poles of the system based on its realization. Unfortunately, the command may not use the minimal realization of the transfer function of the system. Thus, the command may give values for the pole that need not be included in a minimal realization of the transfer function.

From the definition, the reader may guess that the poles of the system are the same as the poles of the elements. Although this is correct, it is not possible to guess the multiplicity of the poles by inspection, as the following example demonstrates.

Example 8.14

Determine the poles of the transfer function matrix

$$G(z) = \begin{bmatrix} \frac{1}{z-1} & \frac{1}{(z-1)(z-0.5)} \\ \frac{z-0.1}{(z-0.2)(z-0.5)} & \frac{1}{z-0.2} \end{bmatrix}$$

Solution

The least common denominator of the matrix entries is

$$(z-0.2)(z-0.5)(z-1)$$

The determinant of the matrix is

$$\det[G(z)] = \frac{1}{(z-1)(z-0.2)} - \frac{z-0.1}{(z-0.2)(z-0.5)^2(z-1)}$$

The denominator of the determinant of the matrix is

$$(z-0.2)(z-0.5)^2(z-1)$$

The least common denominator of all the minors is

$$(z-0.2)(z-0.5)^2(z-1)$$

Thus, the system has poles at $\{0.2, 0.5, 0.5, 1\}$. Note that the pole at 0.2 and the pole at 0.5 are poles of more than one element of the transfer function but are not repeated poles of the system.

For this relatively simple transfer function, we can obtain the minimal number of blocks needed for realization, as shown in Fig. 8.3. This clearly confirms our earlier determination of the system poles.

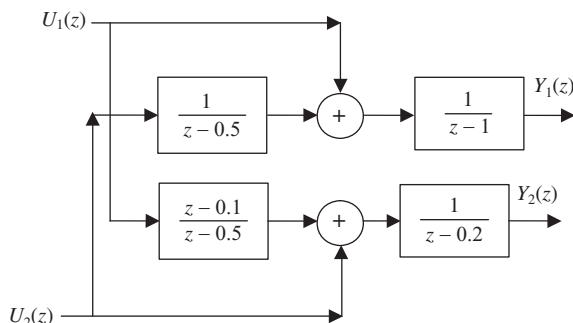


Figure 8.3
Block diagram of the system described in Example 8.14.

Example 8.14—cont'd

The MATLAB command **pole** gives

```
>> pole(g)
ans =
    1.0000
    0.5000
    0.2000
    1.0000
    0.5000
    0.2000
```

This includes two poles that are not needed for the minimal realization of the transfer function.

The definition of the multivariable zero is more complex than that of the pole, and several types of zeros can be defined. We define zeros as follows.

Definition 8.9: System zeros

Consider an $l \times m$ transfer function matrix written such that each entry has a denominator equal to the pole polynomial. Zeros are values that make the transfer function matrix rank deficient—that is, values z_0 that satisfy either of the two equations

$$G(z_0)\mathbf{w}_{in} = 0, \mathbf{w}_{in} \neq \mathbf{0} \quad (8.25)$$

$$\mathbf{w}_{out}^T G(z_0) = 0, \mathbf{w}_{out} \neq \mathbf{0} \quad (8.26)$$

for some nonzero real vector \mathbf{w}_{in} known as the **input zero direction** or a nonzero real vector \mathbf{w}_{out} known as the **output zero direction**.

For any matrix, the rank is equal to the order of the largest nonzero minor. Thus, provided that the terms have the appropriate denominator, the zeros are the divisors of all minors of order equal to the minimum of the pair (l, m) . For a square matrix, zeros of the determinant rewritten with the pole polynomial as its denominator are the zeros of the system. These definitions are examined in the following examples.

Example 8.15

Determine the z-transfer function of a digitally controlled single-axis milling machine with a sampling period of 40 ms if its analog transfer function is given by²

$$G(s) = \text{diag} \left\{ \frac{3150}{(s+35)(s+150)}, \frac{1092}{(s+35)(s+30)} \right\}$$

Find the poles and zeros of the transfer function.

Solution

Using the MATLAB command **c2d**, we obtain the transfer function matrix

$$G_{ZAS}(s) = \text{diag} \left\{ \frac{0.4075(z + 0.1067)}{(z - 0.2466)(z - 0.2479 \times 10^{-2})}, \frac{0.38607(z + 0.4182)}{(z - 0.2466)(z - 0.3012)} \right\}$$

Because the transfer function is diagonal, the determinant is the product of its diagonal terms. The least common denominator of all nonzero minors is the denominator of the determinant of the transfer function

$$(z - 0.2479 \times 10^{-2})(z - 0.2466)^2(z - 0.3012)$$

We therefore have poles at $\{0.2479 \times 10^{-2}, 0.2466, 0.2466, 0.3012\}$. The system is stable because all the poles are inside the unit circle.

To obtain the zeros of the system, we rewrite the transfer function in the form

$$G_{ZAS}(z) = \frac{(z - 0.2466)}{(z - 0.2466)^2(z - 0.2479 \times 10^{-2})(z - 0.3012)} \\ \times \begin{bmatrix} 0.4075(z + 0.1067)(z - 0.3012) & 0 \\ 0 & 0.38607(z + 0.4182)(z - 0.2479 \times 10^{-2}) \end{bmatrix}$$

For this square 2-input-2-output system, the determinant of the transfer function is

$$\frac{(z + 0.1067)(z - 0.3012)(z + 0.4182)(z - 0.2479 \times 10^{-2})(z - 0.2466)^2}{(z - 0.2466)^4(z - 0.2479 \times 10^{-2})^2(z - 0.3012)^2} \\ = \frac{(z + 0.1067)(z + 0.4182)}{(z - 0.2466)^2(z - 0.2479 \times 10^{-2})(z - 0.3012)}$$

The roots of the numerator are the zeros $\{-0.1067, -0.4182\}$.

The same answer is obtained using the MATLAB script.

```
Ga = [zpk([], [-35 -150], 3150), 0; 0, zpk([], [-35 -30], 1092)];  
gd = c2d(Ga, 0.04) % Hold equivalence  
pole(gd)  
tzero(gd)
```

² Rober, S.J., Shin, Y.C., 1995. Modeling and control of CNC machines using a PC-based open architecture controller, Mechatronics 5 (4), 401–420.

Example 8.16

Determine the zeros of the transfer function matrix

$$G(z) = \begin{bmatrix} \frac{1}{z-1} & \frac{1}{(z-1)(z-0.5)} \\ \frac{z-0.1}{(z-0.2)(z-0.5)} & \frac{1}{z-0.2} \end{bmatrix}$$

Solution

From Example 8.14, we know that the poles of the system are at $\{0.2, 0.5, 0.5, 1\}$. We rewrite the transfer function matrix in the form

$$G(z) = \begin{bmatrix} (z-0.2)(z-0.5)^2 & (z-0.2)(z-0.5) \\ (z-0.1)(z-0.5)(z-1) & (z-0.5)^2(z-1) \end{bmatrix} / (z-0.2)(z-0.5)^2(z-1)$$

Zeros are the roots of the greatest common divisor of all minors of order equal to 2, which in this case is simply the determinant of the transfer function matrix. These values make the transfer function matrix rank deficient so that we can find a nonzero vector whose product with the matrix is zero. The determinant of the matrix with cancellation to reduce the denominator to the characteristic polynomial is

$$\det[G(z)] = \frac{(z-0.5)^2 - (z-0.1)}{(z-0.2)(z-0.5)^2(z-1)} = \frac{z^2 - 2z + 0.35}{(z-0.2)(z-0.5)^2(z-1)}$$

The roots of the numerator yield zeros at $\{1.8062, 0.1938\}$.

The same values are obtained using MATLAB but MATLAB gives the additional values $\{1, 0.2\}$.

8.4.2 Zeros from state-space models

System zeros can be obtained from a state-space realization using the definition of the zero. However, this is complicated by the possibility of pole-zero cancellation if the system is not minimal. We rewrite the zero condition Eq. (8.25) in terms of the state-space matrices as

$$\begin{aligned} G(z_0)\mathbf{w} &= C[z_0I_n - A]^{-1}B\mathbf{w} + D\mathbf{w} \\ &= C\mathbf{x}_w + D\mathbf{w} = 0 \\ \mathbf{x}_w &= [z_0I_n - A]^{-1}B\mathbf{w} \end{aligned} \tag{8.27}$$

where \mathbf{x}_w is the state response for the system with input \mathbf{w} . We can now rewrite Eq. (8.27) in the following form, which is more useful for numerical solution:

$$\begin{bmatrix} -(z_0 I_n - A) & B \\ C & D \end{bmatrix} \begin{bmatrix} \mathbf{x}_w \\ \mathbf{w} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad (8.28)$$

The matrix in Eq. (8.28) is known as **Rosenbrock's system matrix**. Eq. (8.28) can be solved for the unknowns provided that a solution exists.

Note that the zeros are invariant under similarity transformation because

$$\begin{bmatrix} -(z_0 I_n - T_r^{-1} A T_r) & T_r^{-1} B \\ C T_r & D \end{bmatrix} \begin{bmatrix} \mathbf{x}_w \\ \mathbf{w} \end{bmatrix} = \begin{bmatrix} T_r^{-1} & 0 \\ 0 & I_l \end{bmatrix} \begin{bmatrix} -(z_0 I_n - A) & B \\ C & D \end{bmatrix} \begin{bmatrix} T_r \mathbf{x}_w \\ \mathbf{w} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}$$

In Chapter 9, we will discuss state feedback where the control input is obtained as a linear combination of the measured state variables. We show that the zeros obtained from Eq. (8.28) are invariant under state feedback. Hence, the zeros are known as **invariant zeros**.

The invariant zeros include decoupling zeros if the system is not minimal. If these zeros are removed from the set of invariant zeros, then the remaining zeros are known as **transmission zeros**. In addition, not all decoupling zeros can be obtained from Eq. (8.28). We therefore have the following relation:

$$\begin{aligned} \{\text{system zeros}\} &= \{\text{transmission zeros}\} + \{\text{input - decoupling zeros}\} \\ &\quad + \{\text{output - decoupling zeros}\} - \{\text{input - output - decoupling zeros}\} \end{aligned}$$

Note that some authors refer to invariant zeros as transmission zeros and do not use the term “invariant zeros.”

MATLAB calculates zeros from Rosenbrock's matrix of a state-space model \mathbf{p} using the command **tzero**. The MATLAB manual identifies the result as transmission zeros, but in our terminology this refers to invariant zeros. Although the command accepts a transfer function matrix, its results are based on a realization that need not be minimal and may include superfluous zeros that would not be obtained using the procedure described in Section 8.4.1.

Example 8.17

Consider the system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} -0.4 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & -2 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 1 & 1 \end{bmatrix} \mathbf{u}(k)$$

Example 8.17—cont'd

$$y(k) = [1 \ 1 \ 0] \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}$$

$$\left[\begin{array}{c|c} -(z_0 I_n - A) & B \\ \hline C & D \end{array} \right] \begin{bmatrix} \mathbf{x}_w \\ \mathbf{w} \end{bmatrix} = \left[\begin{array}{ccc|cc} -0.4 - z_0 & 0 & 0 & 1 & 0 \\ 0 & 3 - z_0 & 0 & 0 & 0 \\ 0 & 0 & -2 - z_0 & 1 & 1 \\ \hline 0 & 0 & -2 & 0 & 0 \end{array} \right] \begin{bmatrix} x_{w1} \\ x_{w2} \\ x_{w3} \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

From the second row, we have the invariant zero $z_0 = 3$, since this value makes the matrix rank deficient. This is also an input-decoupling zero because it corresponds to an uncontrollable mode because the second row of the input matrix is zero. The system has an output-decoupling zero at -2 because the output matrix has a zero third column, but this cannot be determined from Rosenbrock's system matrix.

8.5 State-space realizations

State-space realizations can be obtained from input-output time domain or z -domain models. Because every system has infinitely many state-space realizations, we only cover a few standard canonical forms that have special desirable properties. These realizations play an important role in digital filter design or controller implementation. Because difference equations are easily transformable to z -transfer functions and vice versa, we avoid duplication by obtaining realizations from either the z -transfer functions or the difference equations.

In this section, we mainly discuss SISO transfer functions of the form

$$G(z) = \frac{\mathbf{c}_n z^n + \mathbf{c}_{n-1} z^{n-1} + \dots + \mathbf{c}_1 z + \mathbf{c}_0}{z^n + a_{n-1} z^{n-1} + a_{n-2} z^{n-2} + \dots + a_1 z + a_0} \quad (8.29)$$

$$= c_n + \frac{c_{n-1} z^{n-1} + c_{n-2} z^{n-2} + \dots + c_1 z + c_0}{z^n + a_{n-1} z^{n-1} + a_{n-2} z^{n-2} + \dots + a_1 z + a_0}$$

where the leading numerator coefficients \mathbf{c}_n and c_{n-1} can be zero and $\mathbf{c}_n = c_n$, or the corresponding difference equation

$$y(k+n) + a_{n-1} y(k+n-1) + \dots + a_1 y(k+1) + a_0 y(k) \quad (8.30)$$

$$= \mathbf{c}_n u(k+n) + \mathbf{c}_{n-1} u(k+n-1) + \dots + \mathbf{c}_1 u(k+1) + \mathbf{c}_0 u(k)$$

8.5.1 Controllable canonical realization

The **controllable canonical realization** is so called because it possesses the property of controllability. The controllable form is also known as the **phase variable** form or as the first controllable form. A second controllable form, known as **controller form**, is identical in structure to the first but with the state variables numbered backward. We begin by examining the special case of a difference equation whose RHS is the forcing function at time k . This corresponds to a transfer function with unity numerator. We then use our results to obtain realizations for a general SISO linear discrete-time system described by Eq. (8.29) or Eq. (8.30).

8.5.1.1 Systems with no input differencing

Consider the special case of a system whose difference equation includes the input at time k only. The difference equation considered is of the form

$$y(k+n) + a_{n-1}y(k+n-1) + \dots + a_1y(k+1) + a_0y(k) = u(k) \quad (8.31)$$

Define the state vector

$$\begin{aligned} \mathbf{x}(k) &= [x_1(k) \ x_2(k) \ \dots \ x_{n-1}(k) \ x_n(k)]^T \\ &= [y(k) \ y(k+1) \ \dots \ y(k+n-2) \ y(k+n-1)]^T \end{aligned} \quad (8.32)$$

Hence, the difference Eq. (8.31) can be rewritten as

$$x_n(k+1) = -a_{n-1}x_n(k) - \dots - a_1x_2(k) - a_0x_1(k) + u(k) \quad (8.33)$$

Using the definitions of the state variables and Eq. (8.33), we obtain the matrix state equation

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ \vdots \\ x_{n-1}(k+1) \\ x_n(k+1) \end{bmatrix} = \begin{bmatrix} 0 & 1 & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 0 & 1 \\ -a_0 & -a_1 & \dots & -a_{n-2} & -a_{n-1} \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ \vdots \\ x_{n-1}(k) \\ x_n(k) \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} u(k) \quad (8.34)$$

and the output equation

$$y(k) = [1 \ 0 \ \dots \ 0 \ 0] \begin{bmatrix} x_1(k) \\ x_2(k) \\ \vdots \\ x_{n-1}(k) \\ x_n(k) \end{bmatrix} \quad (8.35)$$

The state-space equations can be written more concisely as

$$\begin{aligned}\mathbf{x}(k+1) &= \left[\begin{array}{c|ccc} \mathbf{0}_{n-1 \times 1} & | & I_{n-1} \\ \hline -a_0 & | & -a_1 & \cdots & \cdots & -a_{n-1} \end{array} \right] \mathbf{x}(k) + \begin{bmatrix} \mathbf{0}_{n-1 \times 1} \\ \vdots \\ 1 \end{bmatrix} u(k) \\ y(k) &= [1 \mid \mathbf{0}_{1 \times n-1}] \mathbf{x}(k)\end{aligned}\quad (8.36)$$

Clearly, the state-space equations can be written by inspection from the transfer function

$$G(z) = \frac{1}{z^n + a_{n-1}z^{n-1} + \dots + a_1z + a_0} \quad (8.37)$$

or from the corresponding difference equation because either includes all the needed coefficients, a_i , $i = 0, 1, 2, \dots, n-1$.

Example 8.18

Obtain the controllable canonical realization of the difference equation

$$y(k+3) + 0.5y(k+2) + 0.4y(k+1) - 0.8y(k) = u(k)$$

using basic principles; then show how the realization can be written by inspection from the transfer function or the difference equation.

Solution

Select the state vector

$$\begin{aligned}\mathbf{x}(k) &= [x_1(k) \ x_2(k) \ x_3(k)]^T \\ &= [y(k) \ y(k+1) \ y(k+2)]^T\end{aligned}$$

and rewrite the difference equation as

$$x_3(k+1) = -0.5x_3(k) - 0.4x_2(k) + 0.8x_1(k) + u(k)$$

Using the definitions of the state variables and the difference equation, we obtain the state-space equations

$$\begin{aligned}\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -0.8 & -0.4 & -0.5 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} u(k) \\ y(k) &= [1 \mid 0 \ 0] \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}\end{aligned}$$

because the system is of order $n = 3$, $n-1 = 2$. Hence, the upper half of the state matrix includes a 2×1 zero vector next to a 2×2 identity matrix. The input matrix is a column vector because the system is single input (SI), and it includes a 2×1 zero vector and a unity last entry. The output matrix is a row vector because the system is single output (SO), and it includes a 1×2 zero vector and a unity first entry. With the exception of the last row of the

Example 8.18—cont'd

state matrix, all matrices in the state-space description are completely determined by the order of the system. The last row of the state matrix has entries equal to the coefficients of the output terms in the difference equation with their signs reversed. The same coefficients appear in the denominator of the transfer function

$$G(z) = \frac{1}{z^3 + 0.5z^2 + 0.4z - 0.8}$$

Therefore, the state-space equations for the system can be written by inspection from the transfer function or input-output difference equation.

8.5.1.2 Systems with input differencing

We now use the results from the preceding section to obtain the controllable canonical realization of the transfer function of Eq. (8.29). We assume that the constant term has been extracted from the transfer function, if necessary, and that we are left with the form

$$G(z) = \frac{Y(z)}{U(z)} = c_n + G_d(z) \quad (8.38)$$

where

$$G_d(z) = \frac{c_{n-1}z^{n-1} + c_{n-2}z^{n-2} + \dots + c_1z + c_0}{z^n + a_{n-1}z^{n-1} + \dots + a_1z + a_0} \quad (8.39)$$

We next consider a transfer function with the same numerator as G_d but with unity numerator, and we define a new variable $p(k)$ whose z -transform satisfies

$$\frac{P(z)}{U(z)} = \frac{1}{z^n + a_{n-1}z^{n-1} + \dots + a_1z + a_0} \quad (8.40)$$

The state equation of a system with the preceding transfer function can be written by inspection, with the state variables chosen as

$$\begin{aligned} \mathbf{x}(k) &= [x_1(k) \quad x_2(k) \quad \dots \quad x_{n-1}(k) \quad x_n(k)]^T \\ &= [p(k) \quad p(k+1) \quad \dots \quad p(k+n-2) \quad p(k+n-1)]^T \end{aligned} \quad (8.41)$$

This choice of state variables is valid because none of the variables can be written as a linear combination of the others. However, we have used neither the numerator of the transfer function nor the constant c_n . Nor have we related the state variables to the output. So we multiply Eq. (8.39) by $U(z)$ and use Eq. (8.40) to obtain

$$\begin{aligned}
Y(z) &= c_n U(z) + G_d(z)U(z) \\
&= c_n U(z) + \frac{c_{n-1}z^{n-1} + c_{n-2}z^{n-2} + \dots + c_1z + c_0}{z^n + a_{n-1}z^{n-1} + \dots + a_1z + a_0} U(z) \\
&= c_n U(z) + [c_{n-1}z^{n-1} + c_{n-2}z^{n-2} + \dots + c_1z + c_0] P(z)
\end{aligned} \tag{8.42}$$

Then we inverse z -transform and use the definition of the state variables to obtain

$$\begin{aligned}
y(k) &= c_n u(k) + c_{n-1}p(k+n-1) + c_{n-2}p(k+n-2) + \dots + c_1p(k+1) + c_0p(k) \\
&= c_n u(k) + c_{n-1}x_n(k) + c_{n-2}x_{n-1}(k) + \dots + c_1x_2(k) + c_0x_1(k)
\end{aligned} \tag{8.43}$$

Finally, we write the output equation in the matrix form

$$y(k) = [c_0 \quad c_1 \quad \dots \quad c_{n-2} \quad c_{n-1}] \begin{bmatrix} x_1(k) \\ x_2(k) \\ \vdots \\ x_{n-1}(k) \\ x_n(k) \end{bmatrix} + du(k) \tag{8.44}$$

where $d = c_n$.

As in the preceding section, the state-space equations can be written by inspection from the difference equation or the transfer function.

A simulation diagram for the system is shown in Fig. 8.4. The simulation diagram shows how the system can be implemented in terms of summer, delay, and scaling operations. The number of delay elements needed for implementation is equal to the order of the system. In addition, two summers and at most $2n+1$ gains are needed. These operations can be easily implemented using a microprocessor or digital signal processing chip.

Example 8.19

Write the state-space equations in controllable canonical form for the following transfer functions:

$$1. \quad G(z) = \frac{0.5(z-0.1)}{z^3+0.5z^2+0.4z-0.8}$$

$$2. \quad G(z) = \frac{z^4+0.1z^3+0.7z^2+0.2z}{z^4+0.5z^2+0.4z-0.8}$$

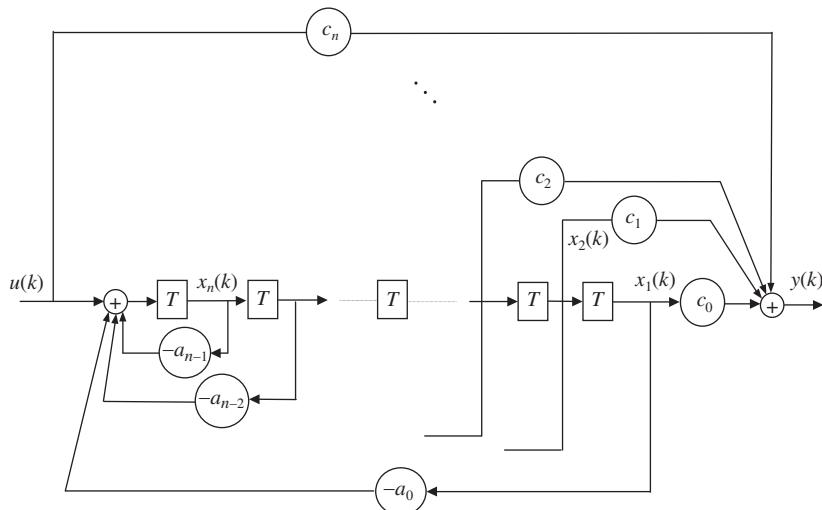


Figure 8.4
Simulation diagram for the controllable canonical realization.

Example 8.19—cont'd

Solution

1. The transfer function has the same denominator as that shown in Example 8.18. Hence, it has the same state equation. The numerator of the transfer function can be expanded as $(-0.05 + 0.5z + 0z^2)$, and the output equation is of the form

$$y(k) = [-0.05 \quad 0.5 \quad 0] \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}$$

2. The transfer function has the highest power of z equal to 4, in both the numerator and denominator. We therefore begin by extracting a constant equal to the numerator z^4 coefficient to obtain

$$\begin{aligned} G(z) &= 1 + \frac{0.1z^3 + (0.7 - 0.5)z^2 + (0.2 - 0.4)z - (-0.8)}{z^4 + 0.5z^2 + 0.4z - 0.8} \\ &= 1 + \frac{0.1z^3 + 0.2z^2 - 0.2z + 0.8}{z^4 + 0.5z^2 + 0.4z - 0.8} \end{aligned}$$

Now we can write the state-space equations by inspection as

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \\ x_4(k+1) \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0.8 & -0.4 & -0.5 & 0 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \\ x_4(k) \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} u(k)$$

Example 8.19—cont'd

$$y(k) = [0.8 \quad -0.2 \quad 0.2 \quad 0.1] \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \\ x_4(k) \end{bmatrix} + u(k)$$

Theorem 8.5 can be used to show that any system in controllable form is actually controllable. The proof is straightforward and is left as an exercise (see Problem 8.15).

8.5.2 Controllable form in MATLAB

MATLAB gives a canonical form that is almost identical to the controllable form of this section with the command

>> [A,B,C,D] = tf2ss(num,den)

Using the command with the system described in Example 8.19(2) gives the state-space equations

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \\ x_4(k+1) \end{bmatrix} = \begin{bmatrix} 0 & -0.5 & -0.4 & 0.8 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \\ x_4(k) \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} u(k)$$

$$y(k) = [0.1 \quad 0.2 \quad -0.2 \quad 0.8] \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \\ x_4(k) \end{bmatrix} + u(k)$$

The equations are said to be in **controller form**. By drawing a simulation diagram for this system and then for the system described in Example 8.19 (2), we can verify that the two systems are identical. The state variables for the form in MATLAB are simply numbered in the reverse of the order used in this text (see Fig. 8.4); that is, the variables x_i in the MATLAB model are none other than the variables x_{n-i+1} , $i = 1, 2, \dots, n$.

8.5.3 Parallel realization

The **parallel realization** of a transfer function is based on the partial fraction expansion

$$G_d(z) = d + \frac{c_{n-1}z^{n-1} + c_{n-2}z^{n-2} + \dots + c_1z + c_0}{z^n + a_{n-1}z^{n-1} + \dots + a_1z + a_0} = d + \sum_{i=1}^n \frac{K_i}{z + p_i} \quad (8.45)$$

The expansion is represented by the block diagram of Fig. 8.5. The summation in Eq. (8.45) gives the parallel configuration shown in the figure, which justifies the name **parallel realization**.

A simulation diagram can be obtained for the parallel realization by observing that the z -transfer functions in each of the parallel branches can be rewritten as

$$\frac{1}{z + p_i} = \frac{z^{-1}}{1 - (-p_i z^{-1})}$$

which can be represented by a positive feedback loop with forward transfer function z^{-1} and feedback gain $-p_i$, as shown in Fig. 8.6. Recall that z^{-1} is simply a time delay so that a physical realization of the transfer function in terms of constant gains and fixed time delays is now possible.

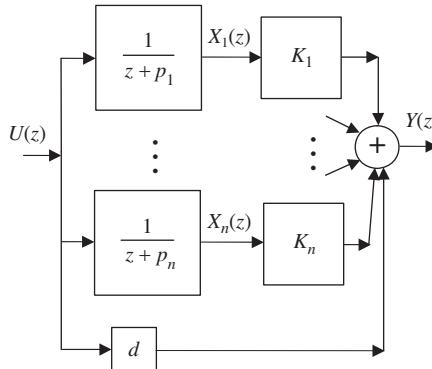


Figure 8.5
Block diagram for parallel realization.

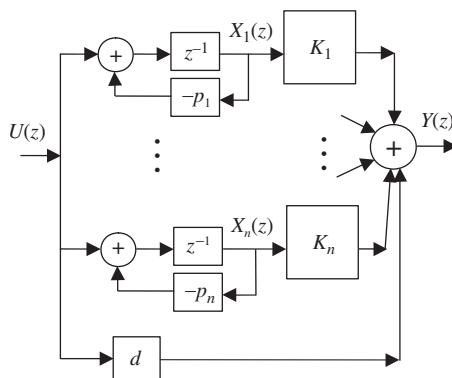


Figure 8.6
Simulation diagram for parallel realization.

We define the state variables as the outputs of the first-order blocks and inverse z -transform to obtain the state equations

$$x_i(k+1) = -p_i x_i(k) + u(k), \quad i = 1, 2, \dots, n \quad (8.46)$$

The output is given by the summation

$$y(k) = \sum_{i=1}^n k_i x_i(k) + du(k) \quad (8.47)$$

Eqs. (8.46) and (8.47) are equivalent to the state-space representation

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ \vdots \\ x_{n-1}(k+1) \\ x_n(k+1) \end{bmatrix} = \begin{bmatrix} -p_1 & 0 & \dots & 0 & 0 \\ 0 & -p_2 & \dots & 0 & 0 \\ \vdots & \vdots & \dots & \vdots & \vdots \\ 0 & 0 & \dots & -p_{n-1} & 0 \\ 0 & 0 & \dots & 0 & -p_n \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ \vdots \\ x_{n-1}(k) \\ x_n(k) \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{bmatrix} u(k) \quad (8.48)$$

$$y(k) = [K_1 \ K_2 \ \dots \ K_{n-1} \ K_n] \begin{bmatrix} x_1(k) \\ x_2(k) \\ \vdots \\ x_{n-1}(k) \\ x_n(k) \end{bmatrix} + du(k)$$

Example 8.20

Obtain a parallel realization for the transfer function

$$G(z) = \frac{2z^2 + 2z + 1}{z^2 + 5z + 6}$$

Solution

We first write the transfer function in the form

$$\begin{aligned} G(z) &= 2 + \frac{2z^2 + 2z + 1 - 2(z^2 + 5z + 6)}{z^2 + 5z + 6} \\ &= 2 - \frac{8z + 11}{(z + 2)(z + 3)} \end{aligned}$$

Then we obtain the partial fraction expansion

$$G(z) = 2 + \frac{5}{z + 2} - \frac{13}{z + 3}$$

Example 8.20—cont'd

Finally, we have the state-space equations

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} -2 & 0 \\ 0 & -3 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} u(k)$$

$$y(k) = [5 \quad -13] \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + 2u(k)$$

The block diagram for the parallel realization is shown in Fig. 8.7, and the simulation diagram is shown in Fig. 8.8. Clearly, the system is unstable with two eigenvalues outside the unit circle.

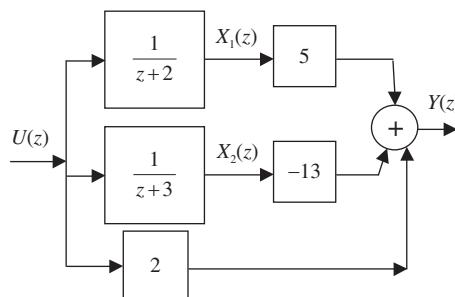


Figure 8.7
Block diagram for Example 8.20.

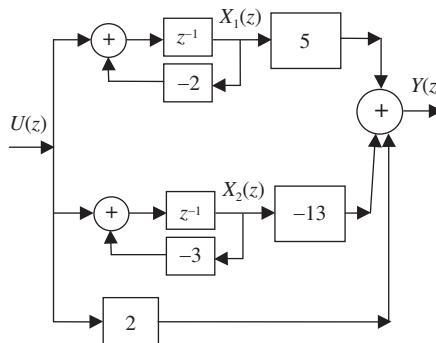


Figure 8.8
Simulation diagram for Example 8.20.

8.5.3.1 Parallel realization for multiinput-multioutput systems

For MIMO systems, a parallel realization can be obtained by partial fraction expansion as in the SISO case. The method presented here requires that the minimal polynomial of the system have no repeated roots. The partial fractions in this case are constant matrices.

The partial fraction expansion is in the form

$$G_d(z) = D + \frac{P_{n-1}z^{n-1} + P_{n-2}z^{n-2} + \dots + P_1z + P_0}{z^n + a_{n-1}z^{n-1} + \dots + a_1z + a_0} = D + \sum_{i=1}^n \frac{K_i}{z + p_i} \quad (8.49)$$

where D , P_i , and K_i , $i = 1, \dots, n-1$ are $l \times m$ matrices. Each of the matrices K_i must be decomposed into the product of two full-rank matrices: an input component matrix B_i and an output component matrix C_i . We write the partial fraction matrices in the form

$$K_i = C_i B_i, \quad i = 1, \dots, n \quad (8.50)$$

For full-rank component matrices, their order is dictated by the rank of the matrix K_i . This follows from the fact that the rank of the product is at most equal to the minimum dimension of the components. For $\text{rank}(K_i) = r_i$, we have C_i as an $l \times r_i$ matrix and B_i as an $r_i \times m$ matrix. The rank of the matrix K_i represents the minimum number of poles p_i needed for the parallel realization. In fact, it can be shown that this realization is indeed minimal. The parallel realization is given by the quadruple

$$A = \begin{bmatrix} -I_{r_1}p_1 & 0 & \dots & 0 & 0 \\ 0 & -I_{r_2}p_2 & \dots & 0 & 0 \\ \vdots & \vdots & \dots & \vdots & \vdots \\ 0 & 0 & \dots & -I_{r_{n-1}}p_{n-1} & 0 \\ 0 & 0 & \dots & 0 & -I_{r_n}p_n \end{bmatrix} \quad B = \begin{bmatrix} B_1 \\ B_2 \\ \vdots \\ B_{n-1} \\ B_n \end{bmatrix} \quad (8.51)$$

$$C = [C_1 \ C_2 \ \dots \ C_{n-1} \ C_n]D$$

Example 8.21

Obtain a parallel realization for the transfer function matrix of Example 8.14:

$$G(z) = \begin{bmatrix} \frac{1}{z-1} & \frac{1}{(z-1)(z-0.5)} \\ \frac{z-0.1}{(z-0.2)(z-0.5)} & \frac{1}{z-0.2} \end{bmatrix}$$

Solution

The minimal polynomial $(z-0.2)(z-0.5)(z-1)$ has no repeated roots. The partial fraction expansion of the matrix is

Example 8.21—cont'd

$$\begin{aligned}
 G(z) &= \begin{bmatrix} \frac{1}{z-1} & \frac{1}{(z-1)(z-0.5)} \\ \frac{z-0.1}{(z-0.2)(z-0.5)} & \frac{1}{z-0.2} \end{bmatrix} \\
 &= \frac{\begin{bmatrix} 0 & 0 \\ -\frac{1}{3} & 1 \end{bmatrix}}{z-0.2} + \frac{\begin{bmatrix} 0 & -2 \\ \frac{4}{3} & 0 \end{bmatrix}}{z-0.5} + \frac{\begin{bmatrix} 1 & 2 \\ 0 & 0 \end{bmatrix}}{z-1}
 \end{aligned}$$

The ranks of the partial fraction coefficient matrices given with the corresponding poles are (1, 0.2), (2, 0.5), and (1, 1). The matrices can be factorized as

$$\begin{bmatrix} 0 & 0 \\ -\frac{1}{3} & 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \begin{bmatrix} -\frac{1}{3} & 1 \end{bmatrix}$$

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

$$\begin{bmatrix} 1 & 2 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} [1 \quad 2]$$

The parallel realization is given by the quadruple

$$A = \begin{bmatrix} 0.2 & 0 & 0 & 0 \\ 0 & 0.5 & 0 & 0 \\ 0 & 0 & 0.5 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad B = \begin{bmatrix} -\frac{1}{3} & 1 \\ 0 & -2 \\ \frac{4}{3} & 0 \\ 1 & 2 \end{bmatrix}$$

$$C = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \end{bmatrix} \quad D = 0_{2 \times 2}$$

Note that the realization is minimal as it is fourth order, and in Example 8.14, the system was determined to have four poles.

8.5.4 Observable form

The **observable realization** of a transfer function can be obtained from the controllable realization using the following steps:

1. Transpose the state matrix A .
2. Transpose and interchange the input matrix B and the output matrix C .

Clearly, the coefficients needed to write the state-space equations all appear in the transfer function, and all equations can be written in observable form directly. Using Eqs. (8.36) and (8.44), we obtain the realization

$$\begin{aligned} \mathbf{x}(k+1) &= \begin{bmatrix} \mathbf{0}_{1 \times n-1} & -a_0 \\ & -a_1 \\ I_{n-1} & \vdots \\ & -a_{n-1} \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ c_{n-1} \end{bmatrix} u(k) \\ y(k) &= [\mathbf{0}_{1 \times n-1} \quad 1] \mathbf{x}(k) + du(k) \end{aligned} \quad (8.52)$$

The simulation diagram for the observable realization is shown in Fig. 8.9. Note that it is possible to renumber the state variables in the simulation diagram to obtain an equivalent realization known as **observer form**. However, the resulting realization will clearly have different matrices from those of (8.52). The second observable realization can be obtained from the controller form by transposing matrices as in the two preceding steps.

The two observable realizations can also be obtained from basic principles; however, the derivations are omitted in this text and are left as an exercise. Theorem 8.5 can be used to show that any system in observable form is actually observable. The proof is straightforward and is left as an exercise (see Problem 8.16).

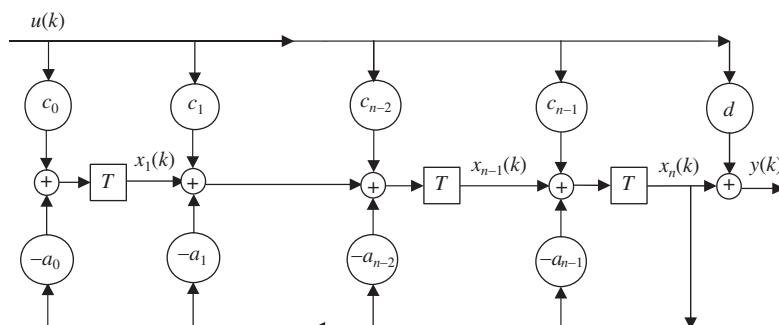


Figure 8.9
Simulation diagram for the observable canonical realization.

Example 8.22

Write the state-space equations in observable canonical form for the transfer function of Example 8.19 (2).

$$G(z) = \frac{z^4 + 0.1z^3 + 0.7z^2 + 0.2z}{z^4 + 0.5z^2 + 0.4z - 0.8}$$

Solution

The transfer function can be written as a constant plus a transfer function with numerator order less than the denominator order as in Example 8.19 (2). Then, using Eq. (8.52), we obtain the following state-space equations:

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \\ x_4(k+1) \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 & 0.8 \\ 1 & 0 & 0 & -0.4 \\ 0 & 1 & 0 & -0.5 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \\ x_4(k) \end{bmatrix} + \begin{bmatrix} 0.8 \\ -0.2 \\ 0.2 \\ 0.1 \end{bmatrix} u(k)$$

$$y(k) = [0 \quad 0 \quad 0 \quad 1] \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \\ x_4(k) \end{bmatrix} + u(k)$$

The same realization is obtained from the controllable realization of Example 8.19 (2) by transposing the state matrix and transposing, and then interchanging the matrices B and C .

8.6 Duality

The concepts of controllability (stabilizability) and observability (detectability) are often referred to as **duals**. This term is justified by Theorem 8.12.

Theorem 8.12

The pair (A, B) is controllable (stabilizable) if and only if (A^T, B^T) is observable (detectable). System (A, C) is observable (detectable) if and only if (A^T, C^T) is controllable (stabilizable).

Proof

The relevant controllability and observability matrices are related by the equations

$$\mathcal{C}(A, B) = \begin{bmatrix} B & | & AB & | & \dots & | & A^{n-1}B \end{bmatrix} = \begin{bmatrix} B^T \\ - \\ B^T A^T \\ - \\ \vdots \\ - \\ B^T (A^T)^{n-1} \end{bmatrix}^T = \mathcal{O}^T(A^T, B^T)$$

$$\mathcal{C}(A^T, B^T) = \begin{bmatrix} B^T & | & A^T B^T & | & \dots & | & (A^T)^{n-1} B^T \end{bmatrix} = \begin{bmatrix} B \\ - \\ BA \\ - \\ \vdots \\ - \\ BA^{n-1} \end{bmatrix}^T = \mathcal{O}^T(A, B)$$

The proof follows from the equality of the rank of any matrix and the rank of its transpose. The statements regarding detectability and stabilizability are true because (i) a matrix and its transpose have the same eigenvalues and (ii) the right eigenvectors of the transpose of the matrix when transposed give the left eigenvectors of the original matrix. Details of the proof are left as an exercise (Problem 8.32).

Example 8.23

Show that the reducible transfer function

$$G(z) = \frac{0.3(z - 0.5)}{(z - 1)(z - 0.5)}$$

has a controllable but unobservable realization and an observable but uncontrollable realization.

Solution

The transfer function can be written as

$$G(z) = \frac{0.3z - 0.15}{z^2 - 15z + 0.5}$$

The controllable realization for this system is

Example 8.23—cont'd

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -0.5 & 1.5 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \mathbf{u}(k)$$

$$y = [-0.15 \quad 0.3] \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$$

The observability matrix for this realization is

$$\mathcal{O} = \begin{bmatrix} C \\ CA_d \\ CA_d^2 \end{bmatrix} = \begin{bmatrix} -0.15 & 0.3 \\ -0.15 & 0.3 \\ -0.15 & 0.3 \end{bmatrix}$$

The observability matrix has rank 1. The rank deficit is $2 - 1 = 1$, corresponding to one unobservable mode.

Transposing the state, output, and input matrices and interchanging the input and output matrices gives the observable realization

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 0 & -0.5 \\ 1 & 1.5 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} -0.15 \\ 0.3 \end{bmatrix} \mathbf{u}(k)$$

$$y(k) = [0 \quad 1] \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$$

By duality, the realization is observable but has one uncontrollable mode.

8.7 Hankel realization

In this section we obtain a realization for a MIMO system using singular value decomposition (see Appendix III). Consider a transfer function matrix expanded in the form

$$G(z) = \sum_{k=0}^{\infty} G(k)z^{-k} \quad (8.53)$$

The terms of the infinite expansion are known as the **Markov parameter matrices** of the system. The terms can be obtained using

$$G(0) = \lim_{z \rightarrow \infty} G(z)$$

$$G(i) = \lim_{z \rightarrow \infty} z^i \left\{ G(z) - \sum_{k=0}^{i-1} G(k)z^{-k} \right\}, i = 1, 2, \dots, \infty \quad (8.54)$$

From Eq. (7.97) we have the direct transmission matrix

$$D = G(0) \quad (8.55)$$

and the remaining terms of the expansion must be used to obtain a realization for the system. To this end, we first define the **Hankel matrix** as

$$H(p,p) = \begin{bmatrix} G(1) & \dots & G(p) \\ \vdots & \ddots & \vdots \\ G(p) & \dots & G(2p-1) \end{bmatrix} \quad (8.56)$$

We also need the shifted Hankel matrix

$$H_s(p,p) = \begin{bmatrix} G(2) & \dots & G(p+1) \\ \vdots & \ddots & \vdots \\ G(p+1) & \dots & G(2p) \end{bmatrix} \quad (8.57)$$

Let p be the degree of the least common denominator of the transfer function matrix. Then we know that the order of the minimal realization is $n \leq p$. For the $n \times n$ block Hankel matrix, we see from the impulse response expansion Eq. (7.97) that

$$H(n,n) = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{n-1} \end{bmatrix} [B \ AB \ \dots \ A^{n-1}B] = \mathcal{OC} \quad (8.58)$$

$$H_s(n,n) = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{n-1} \end{bmatrix} A [B \ AB \ \dots \ A^{n-1}B] = \mathcal{O}A\mathcal{C} \quad (8.59)$$

So that the Hankel matrix of the appropriate size is the product of the observability matrix and the controllability matrix of the minimal realization. The problem is that we do not know the order of this realization in advance. However, with the condition $n \leq p$, we can obtain the singular value decomposition of the Hankel matrix

$$H(p,p) = U \begin{bmatrix} \Sigma & 0_{n \times (p-n)} \\ 0_{(p-n) \times n} & 0_{(p-n) \times (p-n)} \end{bmatrix} V^* \quad (8.60)$$

where U and V are orthogonal matrices and $\Sigma = \text{diag}\{\sigma_1, \dots, \sigma_n\}$ is a diagonal matrix with positive entries. The $n \times n$ block Hankel matrix is the submatrix given by

$$H(n,n) = U_{1:n} \Sigma^{1/2} \Sigma^{1/2} V_{1:n}^* \quad (8.61)$$

With $U_{1:n}$ denoting the first n columns of the matrix and $\Sigma^{1/2} = \text{diag}\left\{\sigma_1^{1/2} \cdots \sigma_n^{1/2}\right\}$, we can now select the matrices

$$\mathcal{C} = \Sigma^{1/2} V_{ln}^* \quad (8.62)$$

$$\mathcal{O} = U_{ln} \Sigma^{1/2} \quad (8.63)$$

From the above matrices we have

$$B = \text{first } m \text{ columns of } \mathcal{C} \quad (8.64)$$

$$C = \text{first } l \text{ rows of } \mathcal{O} \quad (8.65)$$

$$A = \mathcal{O}^+ H_s(n, n) \mathcal{C}^+ \quad (8.66)$$

where \mathcal{C}^+ denotes the pseudoinverse of the matrix.

We now summarize the procedure to obtain the Hankel realization.

Procedure 8.1

- Find the order p of the least common denominator of all the entries of the transfer function matrix.
- Obtain the first $2p$ Markov parameter matrices $G(k), k = 0, 1, \dots, 2p$ and form the $p \times p$ block Hankel matrix $H(p, p)$.
- Obtain the singular value decomposition of the matrix $H(p, p)$ and determine the order of the minimal realization n .
- Calculate the controllability and observability matrices of the minimal realization using the first n right and left singular vectors and the nonzero singular values.
- Obtain the $n \times n$ block Hankel matrix $H_s(n, n)$.
- Obtain a minimal realization using Eqs. (8.55) and (8.64)–(8.66).

Example 8.24

Obtain the Hankel realization of the system of Example 8.14.

Solution

The transfer function matrix

$$G(z) = \begin{bmatrix} \frac{1}{z-1} & \frac{1}{(z-1)(z-0.5)} \\ \frac{z-0.1}{(z-0.2)(z-0.5)} & \frac{1}{z-0.2} \end{bmatrix}$$

Example 8.24—cont'd

has the least common denominator

$$(z - 0.2)(z - 0.5)(z - 1)$$

and thus $p = 3$ and $2p = 6$. The first six Markov parameter matrices are

$$G(0) = \lim_{z \rightarrow \infty} G(z) = \mathbf{0}_{2 \times 2}$$

$$G(1) = \lim_{z \rightarrow \infty} z\{G(z) - G(0)\} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$$

$$G(2) = \lim_{z \rightarrow \infty} z^2\{G(z) - z^{-1}G(1) - G(0)\} = \begin{bmatrix} 1 & 1 \\ 0.6 & 0.2 \end{bmatrix}$$

$$G(3) = \lim_{z \rightarrow \infty} z^3\{G(z) - z^{-2}G(2) - z^{-1}G(1) - G(0)\} = \begin{bmatrix} 1 & 15 \\ 0.32 & 0.04 \end{bmatrix}$$

$$G(4) = \lim_{z \rightarrow \infty} z^4\{G(z) - z^{-3}G(3) - z^{-2}G(2) - z^{-1}G(1) - G(0)\} = \begin{bmatrix} 1 & 1.75 \\ 0.164 & 0.008 \end{bmatrix}$$

$$G(5) = \lim_{z \rightarrow \infty} z^5\{G(z) - z^{-4}G(4) - z^{-3}G(3) - z^{-2}G(2) - z^{-1}G(1) - G(0)\} = \begin{bmatrix} 1 & 1.875 \\ 0.0828 & 0.0016 \end{bmatrix}$$

We form the Hankel matrix using the Markov parameters

$$H(3,3) = \begin{bmatrix} G(1) & G(2) & G(3) \\ G(2) & G(3) & G(4) \\ G(3) & G(4) & G(5) \end{bmatrix} = \begin{bmatrix} 1 & 0 & \cdot & 1 & 1 & \cdot & 1 & 1.5 \\ 1 & 1 & \cdot & 0.6 & 0.2 & \cdot & 0.32 & 0.04 \\ \cdot & \cdot \\ 1 & 1 & \cdot & 1 & 1.5 & \cdot & 1 & 1.75 \\ 0.6 & 0.2 & \cdot & 0.32 & 0.04 & \cdot & 0.164 & 0.008 \\ \cdot & \cdot \\ 1 & 1.5 & \cdot & 1 & 1.75 & \cdot & 1 & 1.875 \\ 0.32 & 0.04 & \cdot & 0.164 & 0.008 & \cdot & 0.0828 & 0.0016 \end{bmatrix}$$

Then we determine its rank and singular value decomposition using the MATLAB command.

n = rank(Hank); % Find the rank of the Hankel matrix.

[L,Sig,R] = svd(Hank); % Singular value decomposition Hank = L * Sig * R';

ans = 4

Example 8.24—cont'd

Using the singular value decomposition of the Hankel matrix, we obtain the controllability matrix

$$\mathcal{C} = \Sigma^{1/2} V_{tn}^* \\ = \begin{bmatrix} -0.7719 & -0.8644 & -0.7463 & -1.1275 & -0.7294 & -1.2712 & -0.7202 & -1.3455 \\ -0.4768 & -0.9415 & -0.1550 & 0.0267 & 0.0588 & 0.3521 & 0.1763 & 0.4831 \\ 0.7110 & -0.4946 & 0.4158 & -0.3515 & 0.2322 & 0.1721 & 0.1332 & -0.0608 \\ -0.1441 & 0.0909 & 0.0867 & -0.1812 & 0.1150 & -0.0560 & 0.1118 & 0.0588 \end{bmatrix}$$

and the observability matrix

$$\mathcal{O} = U_{tn} \Sigma^{1/2} = \begin{bmatrix} -1.0926 & 0.6558 & 0.6740 & 0.0689 \\ -0.3560 & -0.9411 & 0.4206 & 0.1569 \\ -1.3429 & 0.1463 & 0.0480 & 0.0064 \\ -0.1453 & -0.3235 & 0.4292 & -0.1969 \\ -1.4680 & -0.1085 & -0.2651 & -0.0248 \\ -0.0661 & -0.1323 & 0.2584 & -0.1535 \\ -1.5306 & -0.2358 & -0.4216 & -0.0404 \\ -0.0318 & -0.0603 & 0.1380 & -0.0878 \end{bmatrix}$$

We also need the shifted Hankel matrix

$$H_s(3, 3) = \begin{bmatrix} G(2) & G(3) & G(4) \\ G(3) & G(4) & G(5) \\ G(4) & G(5) & G(6) \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1.5 & 1 & 1.75 \\ 0.6 & 0.2 & 0.32 & 0.04 & 0.164 & 0.008 \\ 1 & 1.5 & 1 & 1.75 & 1 & 1.875 \\ 0.32 & 0.04 & 0.164 & 0.008 & 0.0828 & 0.0016 \\ 1 & 1.75 & 1 & 1.875 & 1 & 1.9375 \\ 0.164 & 0.008 & 0.0828 & 0.0016 & 0.0416 & 0.0003 \end{bmatrix}$$

The Hankel realization is given by the matrices

$$A = \mathcal{O}^+ H_s(n, n) \mathcal{C}^+ = \begin{bmatrix} 1.064 & 0.1216 & 0.1943 & 0.03386 \\ -0.2599 & 0.3485 & -0.2312 & 0.1736 \\ -0.009736 & 0.06056 & 0.5896 & -0.1247 \\ -0.04398 & 0.1419 & 0.2094 & 0.1976 \end{bmatrix}$$

$$B = \text{first 2 columns of } \mathcal{C} = \begin{bmatrix} -0.7719 & -0.8644 \\ -0.4768 & -0.9415 \\ 0.7110 & -0.4946 \\ -0.1441 & 0.0909 \end{bmatrix}$$

Example 8.24—cont'd

$$C = \text{first 2 rows of } \mathcal{O} = \begin{bmatrix} -1.0926 & 0.6558 & 0.6740 & 0.0689 \\ -0.3560 & -0.9411 & 0.4206 & 0.1569 \end{bmatrix}$$

$$D = G(0) = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

8.8 Realizations for continuous-time systems

It is important to note that all the methods presented in this chapter to obtain realization of discrete-time systems are applicable to continuous time systems with time advance replaced by differentiation and time delay is replaced by integration. Thus, the simulation diagrams for the continuous-time systems can be obtained from the simulation diagrams for discrete-time systems by replacing all the time delay blocks with integrators. As for discrete-time systems, realizations can be obtained from the transfer functions by inspection as in the following example.

Example 8.25

For the following transfer function, obtain (a) the controllable form, (b) the observable form, (c) the normal form

$$G(s) = \frac{2s^2 + 2s + 1}{s^2 + 5s + 6}$$

Solution

We observe that this is identical in form to the system of Example 8.20 with z replaced by s . We first use long division to write the transfer function in the form

$$G(s) = 2 - \frac{8s + 11}{s^2 + 5s + 6} = 2 - \frac{8s + 11}{(s + 2)(s + 3)}$$

(a) By inspection, the controllable form is

$$\begin{bmatrix} \dot{x}_1(t) \\ x_2(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -6 & -5 \end{bmatrix} \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u(t)$$

Example 8.25—cont'd

$$y(t) = [-11 \quad -8]x(t) + 2u(t)$$

(b) We obtain the observable form by inspection or from the controllable form

$$\begin{bmatrix} \dot{x}_1(t) \\ x_2(t) \end{bmatrix} = \begin{bmatrix} 0 & -6 \\ 1 & -5 \end{bmatrix} \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} + \begin{bmatrix} -11 \\ -8 \end{bmatrix} u(t)$$

$$y(t) = [0 \quad 1]x(t) + 2u(t)$$

(c) For the parallel realization, we obtain the partial fraction expansion

$$G(s) = 2 + \frac{5}{s+2} - \frac{13}{s+3}$$

We have the state-space equations

$$\begin{bmatrix} \dot{x}_1(t) \\ x_2(t) \end{bmatrix} = \begin{bmatrix} -2 & 0 \\ 0 & -3 \end{bmatrix} \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} u(t)$$

$$y(t) = [5 \quad -13]x(t) + 2u(t)$$

Further reading

- Antsaklis, P.J., Michel, A.N., 1997. Linear Systems. McGraw-Hill, New York.
- Belanger, P.R., 1995. Control Engineering. A Modern Approach. Saunders, Fort Worth, TX.
- Chen, C.T., 1984. Linear System Theory and Design. HRW, New York.
- D'Azzo, J.J., Houpis, C.H., 1988. Linear Control System Analysis and Design. McGraw-Hill, New York.
- Delchamps, D.F., 1988. State Space and Input-Output Linear Systems. Springer-Verlag, New York.
- Gupta, S.C., Hasdorff, L., 1970. Fundamentals of Automatic Control. Wiley, New York.
- Kailath, T., 1980. Linear Systems. Prentice Hall, Englewood Cliffs, NJ.
- Moler, C.B., Van Loan, C.F., 1978. Nineteen dubious ways to calculate the exponential of a matrix. SIAM Rev. 20, 801e836.
- Patel, R.V., Munro, N., 1982. Multivariable System Theory and Design. Pergamon Press.
- Sinha, N.K., 1988. Control Systems. HRW, New York.
- Skogestad, S., Postlethwaite, I., 2005. Multivariable Feedback Control: Analysis and Design. Wiley, Chichester.

Problems

8.1 Find the equilibrium state and the corresponding output for the system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -0.5 & -0.1 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u(k)$$

$$y(k) = [1 \quad 1] \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$$

when

- a. $u(k)=0$
- b. $u(k)=1$

8.2 A mechanical system has the state-space equations

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -0.5 & -a_1 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \mathbf{u}(k)$$

$$y(k) = [1 \quad 0] \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$$

where a_1 is dependent on viscous friction.

- a. Using the results of Chapter 4, determine the range of the parameter a_1 for which the system is internally stable.
- b. Predict the dependence of the parameter a_1 on viscous friction, and use physical arguments to justify your prediction. (*Hint:* Friction dissipates energy and helps the system reach its equilibrium.)

8.3 Determine the internal stability and the input-output stability of the following linear systems:

$$\text{a. } \begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 0.1 & 0 \\ 1 & 0.2 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 0 \\ 0.2 \end{bmatrix} \mathbf{u}(k)$$

$$y(k) = [1 \quad 1] \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$$

b.

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} -0.2 & 0.2 & 0 \\ 0 & 1 & 0.1 \\ 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 1 & 1 \end{bmatrix} \mathbf{u}(k)$$

$$y = [1 \ 0 \ 0] \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}$$

c.

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 0.1 & 0.3 \\ 1 & 0.2 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 0 \\ 0.2 \end{bmatrix} \mathbf{u}(k)$$

$$y(k) = [1 \ 1] \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$$

d.

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} -0.1 & -0.3 & 0 \\ 0.1 & 1 & 0.1 \\ 0.3 & 0 & -1 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{u}(k)$$

$$y = [1 \ 0 \ 1] \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}$$

- 8.4 Determine the stable, marginally stable, and unstable modes for each of the unstable systems presented in Problem 8.3.
- 8.5 Determine the controllability and stabilizability of the systems presented in Problem 8.3.
- 8.6 Transform the following system to standard form for uncontrollable systems, and use the transformed system to determine if it is stabilizable:

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} 0.05 & 0.09 & 0.1 \\ 0.05 & 1.1 & -1 \\ 0.05 & -0.9 & 1 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \mathbf{u}(k)$$

$$y = [1 \ 0 \ 1] \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}$$

- 8.7 Transform the system to the standard form for unobservable systems, and use the transformed system to determine if it is detectable:

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} -0.2 & -0.08 \\ 0.125 & 0 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} \mathbf{u}(k)$$

$$y(k) = [1 \quad 0.8] \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$$

- 8.8 Determine the controllability and stabilizability of the systems presented in Problem 8.3 with the input matrices changed to the following:

- a. $B = [1 \quad 0]^T$
- b. $B = [1 \quad 1 \quad 0]^T$
- c. $B = [1 \quad 0]^T$
- d. $B = [1 \quad 0 \quad 1]^T$

- 8.9 An engineer is designing a control system for a chemical process with reagent concentration as the sole control variable. After determining that the system is not controllable, why is it impossible for him to control all the modes of the system by an innovative control scheme using the same control variables? Explain, and suggest an alternative solution to the engineer's problem.

- 8.10 The engineer introduced in Problem 8.9 examines the chemical process more carefully and discovers that all the uncontrollable modes with concentration as control variable are asymptotically stable with sufficiently fast dynamics. Why is it possible for the engineer to design an acceptable controller with reagent concentration as the only control variable? If such a design is possible, give reasons for the engineer to prefer it over a design requiring additional control variables.

- 8.11 Determine the observability and detectability of the systems of Problem 8.3.

- 8.12 Repeat Problem 8.11 with the following output matrices:

- a. $C = [0 \quad 1]$
- b. $C = [0 \quad 1 \quad 0]$
- c. $C = [1 \quad 0]$
- d. $C = [1 \quad 0 \quad 0]$

- 8.13 Consider the system

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -0.005 & -0.11 & -0.7 \end{bmatrix} \quad b = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad c = [0.5 \quad 1 \quad 0] \quad d = 0$$

- a. Can we design a controller for the system that can influence all its modes?
- b. Can we design a controller for the system that can observe all its modes?

Justify your answers using the properties of the system.

- 8.14 Consider system (A, B, C) and the family of systems $(\alpha A, \beta B, \gamma C)$ with each of (α, β, γ) nonzero.
- Show that, if λ is an eigenvalue of A with right eigenvector \mathbf{v} and left eigenvector \mathbf{w}^T , then $\alpha\lambda$ is an eigenvalue of αA with right eigenvector \mathbf{v}/α and left eigenvector \mathbf{w}^T/α .
 - Show that (A, B) is controllable if and only if $(\alpha A, \beta B)$ is controllable for any nonzero constants (α, β) .
 - Show that system (A, C) is observable if and only if $(\alpha A, \gamma C)$ is observable for any nonzero constants (α, γ) .
- 8.15 Show that any system in controllable form is controllable.
- 8.16 Show that any system in observable form is observable.
- 8.17 The realization (A, B, C, D) with $B=C=I_n$, $D = \mathbf{0}_{n \times n}$, allows us to use the MATLAB commands **zpk** or **tf** to obtain the resolvent matrix of A .
- Show that the realization is minimal (no pole-zero cancellation) for any state matrix A (i) using the rank test, and then (ii) using the eigenvector test.
 - Obtain the resolvent matrix of Example 7.7 using the MATLAB commands **zpk** and **tf**.
- 8.18 Obtain state-space representations for the following linear systems:
- In controllable form
 - In observable form
 - In diagonal form
 - $G(z) = 3 \frac{z+0.5}{(z-0.1)(z+0.1)}$
 - $G(z) = 5 \frac{z(z+0.5)}{(z-0.1)(z+0.1)(z+0.8)}$
 - $G(z) = \frac{z^2(z+0.5)}{(z-0.4)(z-0.2)(z+0.8)}$
 - $G(z) = \frac{z(z-0.1)}{z^2-0.9z+0.8}$
- 8.19 Obtain the controller form that corresponds to a renumbering of the state variables of the controllable realization (also known as phase variable form) from basic principles.
- 8.20 Obtain the transformation matrix to transform a system in phase variable form to controller form. Prove that the transformation matrix will also perform the reverse transformation.
- 8.21 Use the two steps of [Section 8.5.4](#) to obtain a second observable realization from controller form. What is the transformation that will take this form to the first observable realization of [Section 8.5.4](#)?
- 8.22 Show that the observable realization obtained from the phase variable form realizes the same transfer function.

- 8.23 Show that the transfer functions of the following systems are identical, and give a detailed explanation.

$$A = \begin{bmatrix} 0 & 1 \\ -0.02 & 0.3 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \mathbf{c}^T = [0 \ 1] \quad \mathbf{d} = 0$$

$$A = \begin{bmatrix} 0.1 & 0 \\ 0 & 0.2 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \mathbf{c}^T = [-1 \ 2] \quad \mathbf{d} = 0$$

- 8.24 Obtain a parallel realization for the transfer function matrix presented in Example 8.15.
- 8.25 Consider the transfer function of a single-input–multioutput (SIMO) system

$$\mathbf{g}(z) = [g_1(z) \ g_2(z) \ \cdots \ g_l(z)]^T$$

We can obtain a controllable realization for the system using the approach used for SISO systems. To show this we focus on the following example

$$\mathbf{g}(z) = \left[\frac{z}{(z-0.2)(z-0.4)} \quad \frac{9z(z-0.9)}{(z-1)(z-0.1)} \right]^T$$

- (a) Use long division to obtain the direct transmission matrix D and subtract it from the transfer function.
- (b) Write state-space equations for each entry of the transfer function vector in controllable form.
- (c) Combine the two state equations to obtain the overall state equation for the system.
- (d) Obtain the transfer function from the state space realization and verify that it is the original transfer function given in the problem.
- (e) Find the minimal characteristic polynomial and the poles of the system. Is the realization obtained minimal?

- 8.26 Consider the column transfer function

$$\mathbf{g}(z) = \left[\frac{z}{(z-0.1)(z-0.4)} \quad \frac{z-0.9}{(z-1)(z-0.1)} \right]^T$$

- (a) Obtain the minimal characteristic polynomial of the system.
- (b) Write the state equation for a transfer function with unity numerator and with the denominator as the minimal polynomial.
- (c) Rewrite the entries of the transfer function vector with a common denominator and use the numerators to write the output equation of the system.

- (d) Verify that the state-space model with the state equation of (b) the output equation of (c) has the transfer function given in the problem. Is this a minimal realization?
- 8.27 Use the controllable realization of Problem 8.26 to obtain an observable realization for the row transfer function

$$\mathbf{g}(z)^T = \left[\frac{z}{(z-0.1)(z-0.4)} \quad \frac{z-0.9}{(z-1)(z-0.1)} \right]$$

- (i) Transpose A , B , C , then interchange B and C , then write the state-space equations.
- (ii) Transpose the expression for the transfer function of the column in terms of its state-space matrices to show that this procedure is valid for any row transfer function.
- 8.28 Find the poles and zeros of the following transfer function matrices:

a. $G(z) = \begin{bmatrix} \frac{z-0.1}{(z+0.1)^2} & \frac{1}{z+0.1} \\ 0 & \frac{1}{z-0.1} \end{bmatrix}$

b. $G(z) = \begin{bmatrix} 0 & \frac{1}{z-0.1} \\ \frac{z-0.2}{(z-0.1)(z-0.3)} & \frac{1}{z-0.3} \\ 0 & \frac{2}{z-0.3} \end{bmatrix}$

- 8.29 Obtain a parallel realization for the system of Problem 8.28(b).
- 8.30 The linearized model of the horizontal plane motion of the INFANTE autonomous underwater vehicle (AUV) of Problem 7.16 is given by

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} = \begin{bmatrix} -0.14 & -0.69 & 0.0 \\ -0.19 & -0.048 & 0.0 \\ 0.0 & 1.0 & 0.0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 0.056 \\ -0.23 \\ 0.0 \end{bmatrix} u$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

where x_1 is the sway speed, x_2 is the yaw angle, x_3 is the yaw rate, and u is the rudder deflection. For the discrete-time system with a sampling period $T = 50$ ms:

- Obtain the zeros of the system using Rosenbrock's system matrix.
 - Determine the system decoupling zeros by testing the system's controllability and observability.
- 8.31 The terminal composition of a binary distillation column uses reflux and steam flow as the control variables. The 2-input-2-output system is governed by the transfer function

$$G(s) = \begin{bmatrix} \frac{12.8e^{-s}}{16.7s + 1} & \frac{-18.9e^{-3s}}{21.0s + 1} \\ \frac{6.6e^{-7s}}{10.9s + 1} & \frac{-19.4e^{-3s}}{14.4s + 1} \end{bmatrix}$$

Find the discrete transfer function of the system with DAC, ADC, and a sampling period of one time unit; then determine the poles and zeros of the discrete-time system.

- 8.32 Show that the λ_i is an uncontrollable mode of the pair (A, B) if and only if it is an unobservable mode of its dual (A^T, B^T) .

Computer exercises

- 8.33 Write computer programs to simulate the second-order systems described in Problem 8.3 for various initial conditions. Obtain state plane plots and discuss your results, referring to the solutions presented in Examples 8.1 and 8.2.
- 8.34 Repeat Problem 8.23 using a computer-aided design (CAD) package. Comment on any discrepancies between CAD results and solution by hand.
- 8.35 Write a MATLAB function that determines the equilibrium state of the system with the state matrix A , the input matrix B , and a constant input u as input parameters.
- 8.36 Select a second-order state equation in diagonal form for which some trajectories converge to the origin and others diverge. Simulate the system using Simulink, and obtain plots for one diverging trajectory and one converging trajectory with suitable initial states.
- 8.37 Write a MATLAB program to find the Hankel realization of a given transfer function and use it to obtain the results of Example 8.24.

State feedback control

Objectives

After completing this chapter, the reader will be able to do the following:

1. Design state feedback control using pole placement.
2. Design servomechanisms using state-space models.
3. Analyze the behavior of multi variable zeros under state feedback.
4. Design state estimators (observers) for state-space models.
5. Design controllers using observer state feedback.
6. Assign system poles using transfer functions with output feedback.

State variable feedback allows the flexible selection of linear system dynamics. Often, not all state variables are available for feedback, and the remainder of the state vector must be estimated. This chapter includes an analysis of state feedback and its limitations. It also includes the design of state estimators for use when some state variables are not available and the use of state estimates in feedback control.

Throughout this chapter, we assume that the state vector \mathbf{x} is $n \times 1$, the control vector \mathbf{u} is $m \times 1$, and the output vector \mathbf{y} is $l \times 1$. We drop the subscript d for discrete system matrices. With minor changes, the design methodologies are applicable to continuous-time systems.

Chapter Outline

9.1 State and output feedback 388

9.2 Pole placement 389

- 9.2.1 Pole placement by transformation to controllable form 393
- 9.2.2 Pole placement using a matrix polynomial 395
- 9.2.3 Choice of the closed-loop eigenvalues 397
- 9.2.4 MATLAB commands for pole placement 402
- 9.2.5 Pole placement for multi-input systems 403
- 9.2.6 Pole placement by output feedback 406

9.3 Servo problem 407

9.4 Invariance of system zeros 411

9.5 State estimation 413

- 9.5.1 Full-order observer 414

9.5.2 Reduced-order observer 417

9.6 Observer state feedback 421

9.6.1 Choice of observer eigenvalues 423

9.7 Pole assignment using transfer functions 429**Further reading 434****Problems 435****Computer exercises 439****9.1 State and output feedback**

State feedback involves the use of the state vector to compute the control action for specified system dynamics. Fig. 9.1 shows a linear system (A, B, C) with constant state feedback gain matrix K . Using the rules for matrix multiplication, we deduce that the matrix K is $m \times n$ so that for a single-input (SI) system, K is a row vector.

The equations for the linear system and the feedback control law are, respectively, given by:

$$\begin{aligned}\mathbf{x}(k+1) &= A\mathbf{x}(k) + B\mathbf{u}(k) \\ \mathbf{y}(k) &= C\mathbf{x}(k)\end{aligned}\tag{9.1}$$

$$\mathbf{u}(k) = -K\mathbf{x}(k) + \mathbf{v}(k)\tag{9.2}$$

where $\mathbf{v}(k)$ is the reference input vector.

The two equations can be combined to yield the closed-loop state equation

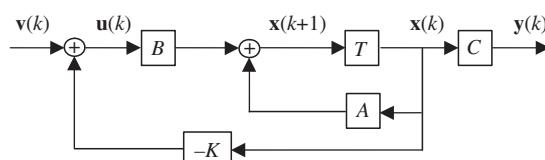
$$\begin{aligned}\mathbf{x}(k+1) &= A\mathbf{x}(k) + B[-K\mathbf{x}(k) + \mathbf{v}(k)] \\ &= [A - BK]\mathbf{x}(k) + B\mathbf{v}(k)\end{aligned}\tag{9.3}$$

We define the closed-loop state matrix as

$$A_{cl} = A - BK\tag{9.4}$$

and rewrite the closed-loop system state-space equations in the form

$$\begin{aligned}\mathbf{x}(k+1) &= A_{cl}\mathbf{x}(k) + B\mathbf{v}(k) \\ \mathbf{y}(k) &= C\mathbf{x}(k)\end{aligned}\tag{9.5}$$

**Figure 9.1**

Block diagram of constant state feedback control.

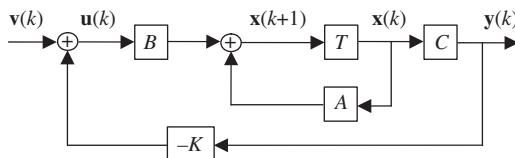


Figure 9.2
Block diagram of constant output feedback control.

The dynamics of the closed-loop system depend on the **eigenstructure** (eigenvalues and eigenvectors) of the matrix A_{cl} . Thus, the desired system dynamics can be chosen with appropriate choice of the gain matrix K . Limitations of this control scheme are addressed in the next section.

For many physical systems, it is either prohibitively costly or impossible to measure all the state variables. The output measurements \mathbf{y} must then be used to obtain the control \mathbf{u} as shown in Fig. 9.2.

The feedback control for output feedback is

$$\begin{aligned}\mathbf{u}(k) &= -K_y \mathbf{y}(k) + \mathbf{v}(k) \\ &= -K_y C \mathbf{x}(k) + \mathbf{v}(k)\end{aligned}\tag{9.6}$$

Substituting in the state equation gives the closed-loop system

$$\begin{aligned}\mathbf{x}(k+1) &= A\mathbf{x}(k) + B[-K_y C \mathbf{x}(k) + \mathbf{v}(k)] \\ &= [A - BK_y C]\mathbf{x}(k) + B\mathbf{v}(k)\end{aligned}\tag{9.7}$$

The corresponding state matrix is

$$A_y = A - BK_y C\tag{9.8}$$

Intuitively, less can be accomplished using output feedback than state feedback because less information is used in constituting the control law. In addition, the postmultiplication by the C matrix in Eq. (9.8) restricts the choice of closed-loop dynamics. However, output feedback is a more general design problem because state feedback is the special case where C is the identity matrix.

9.2 Pole placement

Using output or state feedback, the poles or eigenvalues of the system can be assigned subject to system-dependent limitations. This is known as **pole placement**, **pole assignment**, or **pole allocation**. We state the problem as follows.

Definition 9.1: Pole placement

Choose the gain matrix K or K_y to assign the system eigenvalues to an arbitrary set $\{\lambda_i, i = 1, \dots, n\}$.

The following theorem gives conditions that guarantee a solution to the pole-placement problem with state feedback.

Theorem 9.1: State feedback

If the pair (A, B) is controllable, then there exists a feedback gain matrix K that arbitrarily assigns the system poles to any set $\{\lambda_i, i = 1, \dots, n\}$. Furthermore, if the pair (A, B) is stabilizable, then the controllable modes can all be arbitrarily assigned.

Proof**Necessity**

We first show that controllability is invariant under state feedback. If the system is not controllable, then by Theorem 8.4 we have $\mathbf{w}_i^T B = 0^T$ for some left eigenvector \mathbf{w}_i^T .

Premultiplying the closed-loop state matrix by \mathbf{w}_i^T gives

$$\mathbf{w}_i^T A_{cl} = \mathbf{w}_i^T (A - BK) = \lambda_i \mathbf{w}_i^T - \mathbf{w}_i^T BK = \lambda_i \mathbf{w}_i^T$$

Thus, λ_i is an eigenvalue and \mathbf{w}_i^T is a left eigenvector of the closed-loop system for any state feedback gain matrix K . Hence, the i^{th} **eigenpair** (eigenvalue and eigenvector) of A is unchanged by state feedback and cannot be arbitrarily assigned. Therefore, controllability is a necessary condition for arbitrary pole assignment.

Sufficiency

We first give the proof for the SI case where \mathbf{b} is a column matrix. For a controllable pair (A, B) , we can assume without loss of generality that the pair is in controllable form. We rewrite Eq. (9.4) as

$$\mathbf{b} \mathbf{k}^T = A - A_{cl}$$

Proof—cont'd

Substituting the controllable form matrices gives

$$\begin{bmatrix} \mathbf{0}_{n-1 \times 1} \\ 1 \end{bmatrix} \begin{bmatrix} k_1 & k_2 & \cdots & k_n \end{bmatrix} = \begin{bmatrix} \mathbf{0}_{n-1 \times 1} & | & I_{n-1} \\ -a_0 & -a_1 & \cdots & \cdots & -a_{n-1} \end{bmatrix}$$

$$-\begin{bmatrix} \mathbf{0}_{n-1 \times 1} & | & I_{n-1} \\ -a_0^d & -a_1^d & \cdots & \cdots & -a_{n-1}^d \end{bmatrix}$$

That is,

$$\begin{bmatrix} \mathbf{0}_{n-1 \times n} \\ \cdots \\ k_1 & k_2 & \cdots & k_n \end{bmatrix} = \begin{bmatrix} \mathbf{0}_{n-1 \times n} \\ \cdots \\ a_0^d - a_0 & a_1^d - a_1 & \cdots & \cdots & a_{n-1}^d - a_{n-1} \end{bmatrix}$$

Equating the last rows of the matrices yields

$$[k_1 \ k_2 \ \dots \ k_n] = [a_0^d - a_0 \ a_1^d - a_1 \ \dots \ \dots \ a_{n-1}^d - a_{n-1}] \quad (9.9)$$

which is the control that yields the desired characteristic polynomial coefficients.

Sufficiency Proof 2

We now give a more general proof by contraposition—that is, we assume that the result is not true and prove that the assumption is not true. So we assume that the eigenstructure of the j^{th} mode is unaffected by any state feedback for any choice of the matrix K and prove that the system is uncontrollable. In other words, we assume that the j^{th} eigenpair is the same for the open-loop and closed-loop systems. Then we have

$$\mathbf{w}_j^T B K = \mathbf{w}_j^T \{A - A_{cl}\} = \left\{ \lambda_j \mathbf{w}_j^T - \lambda_j \mathbf{w}_j^T \right\} = 0^T$$

Assuming K full rank, we have $\mathbf{w}_j^T B = \mathbf{0}$. By Theorem 8.4, the system is not controllable.

For a controllable SI system, the matrix (row vector) K has n entries with n eigenvalues to be assigned, and the pole placement problem has a unique solution. Clearly, with more inputs, the K matrix has more rows and consequently more unknowns than the n equations dictated by the n specified eigenvalues. This freedom can be exploited to obtain a solution that has desirable properties in addition to the specified eigenvalues. For example, the eigenvectors of the closed-loop state matrix can be selected subject to constraints. There is a rich literature that covers eigenstructure assignment, but it is beyond the scope of this text.

The sufficiency proof of Theorem 9.1 provides a method to obtain the feedback matrix K to assign the poles of a controllable system for the SI case with the system in controllable form. This approach is explored later. We first give a simple procedure applicable to low-order systems.

Procedure 9.1: Pole placement by equating coefficients

- Evaluate the desired characteristic polynomial from the specified eigenvalues λ_i^d , $i = 1, \dots, n$ using the expression

$$\Delta_c^d(\lambda) = \prod_{i=1}^n (\lambda - \lambda_i^d) \quad (9.10)$$

- Evaluate the closed-loop characteristic polynomial using the expression

$$\det\{\lambda I_n - (A - BK)\} \quad (9.11)$$

- Equate the coefficients of the two polynomials to obtain n equations to be solved for the entries of the matrix K .

Example 9.1: Pole assignment

Assign the eigenvalues $\{0.3 \pm j0.2\}$ to the pair

$$A = \begin{bmatrix} 0 & 1 \\ 3 & 4 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

Solution

For the given eigenvalues, the desired characteristic polynomial is

$$\Delta_c^d(\lambda) = (\lambda - 0.3 - j0.2)(\lambda - 0.3 + j0.2) = \lambda^2 - 0.6\lambda + 0.13$$

The closed-loop state matrix is

$$\begin{aligned} A - \mathbf{b}\mathbf{k}^T &= \begin{bmatrix} 0 & 1 \\ 3 & 4 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \end{bmatrix} [k_1 \quad k_2] \\ &= \begin{bmatrix} 0 & 1 \\ 3 - k_1 & 4 - k_2 \end{bmatrix} \end{aligned}$$

The closed-loop characteristic polynomial is

$$\begin{aligned} \det\{\lambda I_2 - (A - \mathbf{b}\mathbf{k}^T)\} &= \det \begin{bmatrix} \lambda & -1 \\ -(3 - k_1) & \lambda - (4 - k_2) \end{bmatrix} \\ &= \lambda^2 - (4 - k_2)\lambda - (3 - k_1) \end{aligned}$$

Equating coefficients gives the two equations

Example 9.1: Pole assignment—cont'd

1. $4 - k_2 = 0.6 \Rightarrow k_2 = 3.4$
2. $-3 + k_1 = 0.13 \Rightarrow k_1 = 3.13$

that is,

$$\mathbf{k}^T = [3.13, 3.4]$$

Because the system is in controllable form, the same result can be obtained as the coefficients of the open-loop characteristic polynomial minus those of the desired characteristic polynomial using Eq. (9.9).

9.2.1 Pole placement by transformation to controllable form

Any controllable single-input-single-output (SISO) system can be transformed into controllable form using the transformation

$$T_c = \mathcal{C}\mathcal{C}_c^{-1} = [\mathbf{b} \mid A\mathbf{b} \mid \dots \mid A^{n-1}\mathbf{b}] \begin{bmatrix} a_1 & a_2 & \cdots & a_{n-1} & 1 \\ a_2 & a_3 & \cdots & 1 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{n-1} & 1 & \cdots & 0 & 0 \\ 1 & 0 & \cdots & 0 & 0 \end{bmatrix}$$

$$T_c^{-1} = \mathcal{C}_c \mathcal{C}^{-1} = \begin{bmatrix} 0 & 0 & 0 & \cdots & 1 \\ 0 & 0 & 0 & \cdots & t_{2,n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 1 & \cdots & t_{n-2,n} \\ 0 & 1 & t_{2,n} & \cdots & t_{n-1,n} \\ 1 & t_{2,n} & t_{3,n} & \cdots & t_{n,n} \end{bmatrix} [\mathbf{b} | A\mathbf{b} | \dots | A^{n-1}\mathbf{b}]^{-1} \quad (9.13)$$

where \mathcal{C} is the controllability matrix, the subscript c denotes the controllable form, and the entries of the matrix T_c^{-1} are given by

$$t_{2,n} = a_{n-1}$$

$$t_{j+1,n} = - \sum_{i=0}^{j-1} a_{n-i-1} t_{j-i,n}, \quad j = 2, \dots, n-1$$

The forms of the matrix T_c and its inverse given above can be derived by induction and the proof is left as an exercise. The state feedback for a system in controllable form is

$$u = -\mathbf{k}_c^T \mathbf{x}_c = -\mathbf{k}_c^T (T_c^{-1} \mathbf{x}) \quad (9.14)$$

$$\mathbf{k}^T = [a_0^d - a_0 \quad a_1^d - a_1 \quad \dots \quad \dots \quad a_{n-1}^d - a_{n-1}] T_c^{-1} \quad (9.15)$$

We now have the following pole placement procedure.

Procedure 9.2

1. Obtain the characteristic polynomial of the pair (A, B) using the Leverrier algorithm described in [Section 7.4.1](#).
2. Obtain the transformation matrix T_c^{-1} using the coefficients of the polynomial from step 1.
3. Obtain the desired characteristic polynomial coefficients from the given eigenvalues using [Eq. \(9.10\)](#).
4. Compute the state feedback matrix using [Eq. \(9.15\)](#).

Procedure 9.2 requires the numerically troublesome inversion of the controllability matrix to obtain T . However, it does reveal an important characteristic of state feedback. From [Eq. \(9.15\)](#), we observe that the feedback gains tend to increase as the change from the open-loop to the closed-loop polynomial coefficients increases. Procedure 9.2, like Procedure 9.1, works well for low-order systems but can be implemented more easily using a computer-aided design (CAD) program.

Example 9.2

Design a feedback controller for the pair

$$A = \begin{bmatrix} 0.1 & 0 & 0.1 \\ 0 & 0.5 & 0.2 \\ 0.2 & 0 & 0.4 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0.01 \\ 0 \\ 0.005 \end{bmatrix}$$

to obtain the eigenvalues $\{0.1, 0.4 \pm j 0.4\}$.

Solution

The characteristic polynomial of the state matrix is

$$\lambda^3 - \lambda^2 + 0.27\lambda - 0.01 \quad \text{i.e.,} \quad a_2 = -1, a_1 = 0.27, \quad a_0 = -0.01$$

The transformation matrix T_c^{-1} is

$$T_c^{-1} \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 0.73 \end{bmatrix} \times 10^3 \begin{bmatrix} 10 & 1.5 & 0.6 \\ 0 & 1 & 1.3 \\ 5 & 4 & 1.9 \end{bmatrix}^{-1} = 10^3 \begin{bmatrix} 0.1923 & 1.25 & -0.3846 \\ -0.0577 & 0.625 & 0.1154 \\ 0.0173 & 0.3125 & 0.1654 \end{bmatrix}$$

The desired characteristic polynomial is

$$\lambda^3 - 0.9\lambda^2 + 0.4\lambda - 0.032 \quad \text{i.e.,} \quad a_2^d = -0.9, a_1^d = 0.4, a_0^d = -0.032$$

Hence, we have the feedback gain vector

Example 9.2—cont'd

$$\begin{aligned}
 \mathbf{k}^T &= [a_0^d - a_0 \quad a_1^d - a_1 \quad a_2^d - a_2] T_c^{-1} \\
 &= [-0.032 + 0.01 \quad 0.4 - 0.27 \quad -0.9 + 1] \times 10^3 \begin{bmatrix} 0.1923 & 1.25 & -0.3846 \\ -0.0577 & 0.625 & 0.1154 \\ 0.0173 & 0.3125 & 0.1654 \end{bmatrix} \\
 &= [-10 \quad 85 \quad 40]
 \end{aligned}$$

9.2.2 Pole placement using a matrix polynomial

The gain vector for pole placement can be expressed in terms of the desired closed-loop characteristic polynomial. The expression, known as **Ackermann's formula**, is

$$\mathbf{k}^T = \mathbf{t}_1^T \Delta_c^d(A) \quad (9.16)$$

where \mathbf{t}_1^T is the first row of the matrix T_c^{-1} of Eq. (9.13) and $\Delta_c^d(\lambda)$ is the desired closed-loop characteristic polynomial. From Theorem 9.1, we know that state feedback can arbitrarily place the eigenvalues of the closed-loop system for any controllable pair (A, \mathbf{b}) . In addition, any controllable pair can be transformed into controllable form (A_c, \mathbf{b}_c) . By the Cayley–Hamilton theorem, the state matrix satisfies its own characteristic polynomial $\Delta(\lambda)$ but not that corresponding to the desired pole locations. That is,

$$\begin{aligned}
 \Delta_c(A) &= \sum_{i=0}^{n-1} a_i A^i = \sum_{i=0}^{n-1} a_i A_c^i = 0, \quad a_n = 1 \\
 \Delta_c^d(A) &= \sum_{i=0}^n a_i^d A^i \neq 0, \quad a_n^d = 1
 \end{aligned}$$

Subtracting and using the identity $A_c = T_c^{-1} A T_c$ gives

$$T_c^{-1} \Delta_c^d(A) T_c = \sum_{i=0}^{n-1} (a_i^d - a_i) A_c^i \quad (9.17)$$

The state matrix in controllable form possesses an interesting property, which we use in this proof. If the matrix raised to power i , with $i = 1, 2, \dots, n-1$, is premultiplied by the first elementary vector

$$\mathbf{e}_1^T = [1 \quad 0_{n-1 \times 1}^T]$$

the result is the $(i+1)$ th elementary vector—that is,

$$\mathbf{e}_1^T A_c^i = \mathbf{e}_{i+1}^T, \quad i = 0, \dots, n-1 \quad (9.18)$$

Premultiplying Eq. (9.17) by the elementary vector \mathbf{e}_1^T , and then using Eq. (9.18), we obtain

$$\begin{aligned} \mathbf{e}_1^T T_c^{-1} \Delta_c^d(A) T_c &= \sum_{i=0}^{n-1} (a_i^d - a_i) \mathbf{e}_1^T A_c^i \\ &= \sum_{i=0}^{n-1} (a_i^d - a_i) \mathbf{e}_{i+1}^T \\ &= [a_0^d - a_0 \quad a_1^d - a_1 \quad \cdots \quad a_{n-1}^d - a_{n-1}] \end{aligned}$$

Using Eq. (9.15), we obtain

$$\begin{aligned} \mathbf{e}_1^T T_c^{-1} \Delta_c^d(A) T_c &= [a_0^d - a_0 \quad a_1^d - a_1 \quad \cdots \quad a_{n-1}^d - a_{n-1}] \\ &= \mathbf{k}_c^T T_c \end{aligned}$$

Postmultiplying by T_c^{-1} and observing that the first row of the inverse is $\mathbf{t}_1^T = \mathbf{e}_1^T T_c^{-1}$, we obtain Ackermann's formula (9.16).

Minor modifications in Procedure 9.2 allow pole placement using Ackermann's formula. The formula requires the evaluation of the first row of the matrix T_c^{-1} rather than the entire matrix. However, for low-order systems, it is often simpler to evaluate the inverse and then use its first row. The following example demonstrates pole placement using Ackermann's formula.

Example 9.3

Obtain the solution described in Example 9.2 using Ackermann's formula.

Solution

The desired closed-loop characteristic polynomial is

$$\Delta_c^d(\lambda) = \lambda^3 - 0.9\lambda^2 + 0.4\lambda - 0.032 \quad \text{i.e.,} \quad d_2^d = -0.9, d_1^d = 0.4, d_0^d = -0.032$$

The first row of the inverse transformation matrix is

$$\begin{aligned} \mathbf{t}_1^T &= \mathbf{e}_1^T T_c^{-1} = 10^3 [1 \quad 0 \quad 0] \begin{bmatrix} 0.1923 & 1.25 & -0.3846 \\ -0.0577 & 0.625 & 0.1154 \\ 0.0173 & 0.3125 & 0.1654 \end{bmatrix} \\ &= 10^3 [0.1923 \quad 1.25 \quad -0.3846] \end{aligned}$$

Example 9.3—cont'd

We use Ackermann's formula to compute the gain vector

$$\begin{aligned}
 \mathbf{k}^T &= \mathbf{t}_1^T \Delta_c^d(A) \\
 &= \mathbf{t}_1^T \{A^3 - 0.9A^2 + 0.4A - 0.032I_3\} \\
 &= 10^3 [0.1923 \quad 1.25 \quad -0.3846] \times 10^{-3} \begin{bmatrix} 6 & 0 & 18 \\ 4 & 68 & 44 \\ 36 & 0 & 48 \end{bmatrix} \\
 &= [-10 \quad 85 \quad 40]
 \end{aligned}$$

9.2.3 Choice of the closed-loop eigenvalues

Procedures 9.1 and 9.2 yield the feedback gain matrix once the closed-loop eigenvalues have been arbitrarily selected. The desired locations of the eigenvalues are directly related to the desired transient response of the system. In this context, considerations similar to those made in [Section 6.6](#) can be applied to select the desired time response. However, the designer must take into account that poles associated with fast modes will lead to high gains for the state feedback matrix and consequently to a high **control effort**. High gains may also lead to performance degradation due to nonlinear behavior such as analog-to-digital converter (ADC) or actuator saturation. If all the desired closed-loop eigenvalues are selected at the origin of the complex plane, the deadbeat control strategy is implemented (see [Section 6.7](#)), and the closed-loop characteristic polynomial is chosen as

$$\Delta_c^d(\lambda) = \lambda^n \quad (9.19)$$

Substituting in Ackermann's [formula \(9.16\)](#) gives the feedback gain matrix

$$\mathbf{k}^T = \mathbf{t}_1^T A^n \quad (9.20)$$

The resulting control law will drive all the states to zero in at most n sampling intervals starting from any initial condition. However, the limitations of deadbeat control discussed in [Section 6.7](#) apply—namely, the control variable can assume unacceptably high values, and undesirable intersample oscillations can occur.

Example 9.4

Determine the gain vector \mathbf{k} using Ackermann's formula for the discretized state-space model of the armature-controlled DC motor (see Example 7.16) where

$$A_d = \begin{bmatrix} 1.0 & 0.01 & 0.0 \\ 0.0 & 0.9995 & 0.0095 \\ 0.0 & -0.09470 & 0.8954 \end{bmatrix}, \quad B_d = \begin{bmatrix} 1.622 \times 10^{-6} \\ 4.821 \times 10^{-4} \\ 9.468 \times 10^{-2} \end{bmatrix}$$

for the following choices of closed-loop eigenvalues:

1. $\{0.1, 0.4 \pm j0.4\}$
2. $\{0.4, 0.6 \pm j0.33\}$
3. $\{0, 0, 0\}$ (deadbeat control)

Simulate the system in each case to obtain the zero-input response starting from the initial condition $\mathbf{x}(0) = [1, 1, 1]$, and discuss the results.

Solution

The characteristic polynomial of the state matrix is

$$\Delta(\lambda) = \lambda^3 - 2.895\lambda^2 + 2.791\lambda - 0.896$$

That is,

$$a_2 = -2.895, \quad a_1 = 2.791, \quad a_0 = -0.896$$

The controllability matrix of the system is

$$\mathcal{C} = 10^{-3} \begin{bmatrix} 0.001622 & 0.049832 & 0.187964 \\ 0.482100 & 1.381319 & 2.185571 \\ 94.6800 & 84.37082 & 75.73716 \end{bmatrix}$$

Using Eq. (9.13) gives the transformation matrix

$$T_c^{-1} = 10^4 \begin{bmatrix} 1.0527 & -0.0536 & 0.000255 \\ -1.9948 & 0.20688 & -0.00102 \\ 0.94309 & -0.14225 & 0.00176 \end{bmatrix}$$

1. The desired closed-loop characteristic polynomial is

$$\Delta_c^d(\lambda) = \lambda^3 - 0.9\lambda^2 + 0.4\lambda - 0.032$$

That is,

$$a_2^d = -0.9, \quad a_1^d = 0.4, \quad a_0^d = -0.032$$

By Ackermann's formula, the gain vector is

$$\begin{aligned} \mathbf{k}^T &= \mathbf{t}_1^T \Delta_c^d(A) \\ &= \mathbf{t}_1^T \{ A^3 - 0.9A^2 + 0.4A - 0.032I_3 \} \\ &= 10^4 [1.0527 \quad -0.0536 \quad 0.000255] \times \{ A^3 - 0.9A^2 + 0.4A - 0.032I_3 \} \\ &= 10^3 [4.9268 \quad 1.4324 \quad 0.0137] \end{aligned}$$

Example 9.4—cont'd

The zero-input response for the three states and the corresponding control variable u are shown in Fig. 9.3.

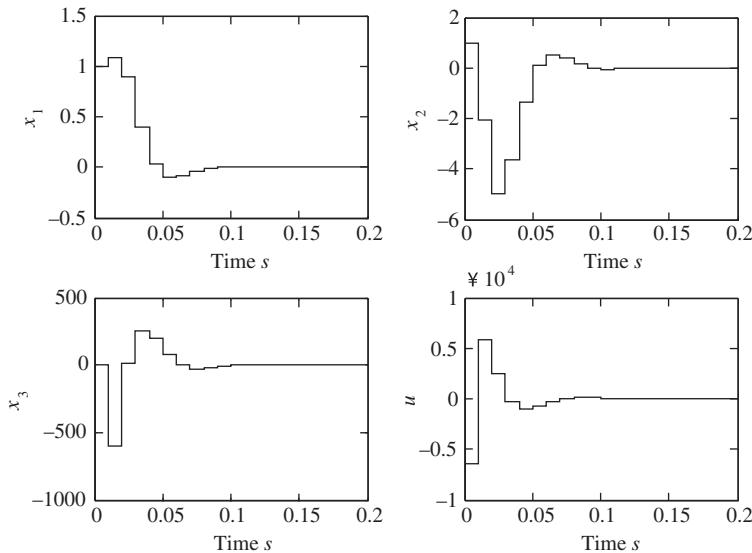


Figure 9.3

Zero-input state response and control variable for case 1 of Example 9.4.

2. The desired closed-loop characteristic polynomial is

$$\Delta_c^d(\lambda) = \lambda^3 - 1.6\lambda^2 + 0.9489\lambda - 0.18756$$

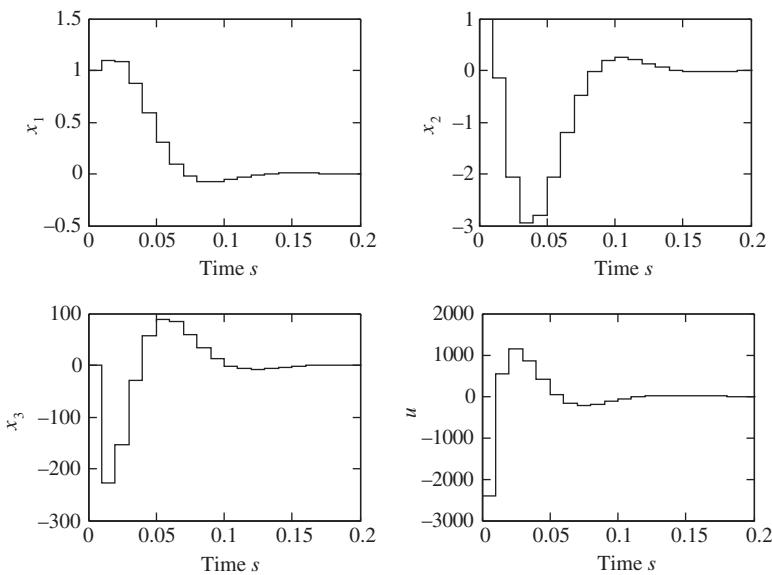
That is,

$$a_2^d = -1.6, \quad a_1^d = 0.9489, \quad a_0^d = -0.18756$$

And the gain vector is

$$\begin{aligned} \mathbf{k}^T &= \mathbf{t}_1^T \Delta_c^d(A) \\ &= \mathbf{t}_1^T \{A^3 - 1.6A^2 + 0.9489A - 0.18756I_3\} \\ &= 10^4 [1.0527 \quad -0.0536 \quad 0.000255] \times \{A^3 - 1.6A^2 + 0.9489A - 0.18756I_3\} \\ &= 10^3 [1.6985 \quad 0.70088 \quad 0.01008] \end{aligned}$$

The zero-input response for the three states and the corresponding control variable u are shown in Fig. 9.4.

Example 9.4—cont'd**Figure 9.4**

Zero-input state response and control variable for case 2 of Example 9.4.

3. The desired closed-loop characteristic polynomial is

$$\Delta_c^d(\lambda) = \lambda^3 \quad \text{i.e.,} \quad d_2^d = 0, \quad d_1^d = 0, \quad d_0^d = 0$$

and the gain vector is

$$\begin{aligned} \mathbf{k}^T &= \mathbf{t}_1^T \Delta_c^d(A) \\ &= \mathbf{t}_1^T \{A^3\} \\ &= 10^4 [1.0527 \quad -0.0536 \quad 0.000255] \times \{A^3\} \\ &= 10^4 [1.0527 \quad 0.2621 \quad 0.0017] \end{aligned}$$

The zero-input response for the three states and the corresponding control variable u are shown in Fig. 9.5.

We observe that when eigenvalues associated with faster modes are selected, higher gains are required for the state feedback and the state variables have transient responses with larger oscillations. Specifically, for the deadbeat control of case 3, the gain values are one order of magnitude larger than those of cases 1 and 2, and the magnitude of its transient oscillations is much larger. Further, for deadbeat control, the zero state is reached in $n = 3$ sampling intervals, as predicted by the theory. However, unpredictable transient intersample behavior occurs between 0.01 and 0.02 in the motor velocity x_2 . This is shown in Fig. 9.6, where the analog velocity and the sampled velocity of the motor are plotted.

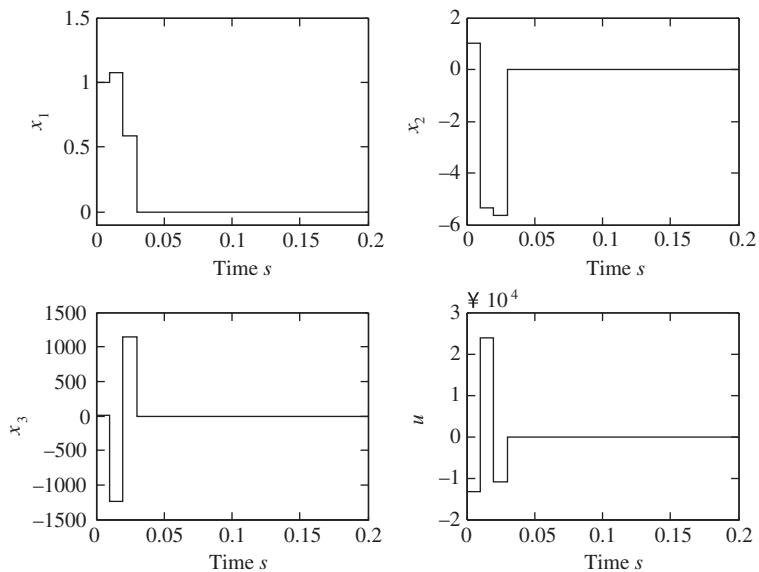
Example 9.4—cont'd

Figure 9.5
Zero-input state response and control variable for case 3 of Example 9.4.

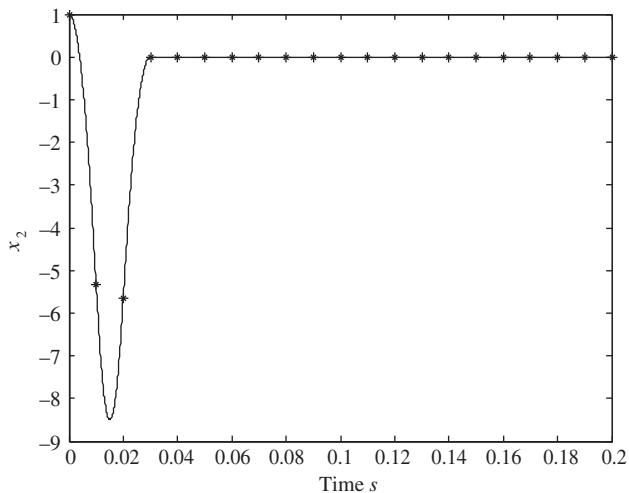


Figure 9.6
Sampled and analog velocity of the motor with deadbeat control.

Another consideration in selecting the desired closed-loop poles is the **robustness** of the system to modeling uncertainties. Although we have so far assumed that the state-space model of the system is known perfectly, this is never true in practice. Thus, it is desirable that control systems have poles with low sensitivity to perturbations in their state-space matrices. It is well known that low pole sensitivity is associated with designs that minimize the condition number of the modal matrix of eigenvectors of the state matrix; that is, the ratio of the maximum to the minimum singular value (see Appendix III for a discussion of singular values).

For a SISO system, the eigenvectors are fixed once the eigenvalues are chosen. Thus, the choice of eigenvalues determines the pole sensitivity of the system. For example, selecting the same eigenvalues for the closed-loop system leads to a high condition number of the eigenvector matrix and therefore to a closed-loop system that is very sensitive to coefficient perturbations. For multi-input–multi-output (MIMO) systems, the feedback gains for a given choice of eigenvalues are nonunique. This allows for more than one choice of eigenvectors and can be exploited to obtain robust designs. However, robust pole placement for MIMO systems is beyond the scope of this text and is not discussed further.

9.2.4 MATLAB commands for pole placement

The pole placement command is **place**. The following example illustrates the use of the command:

```
>> A = [0, 1; 3, 4];
>> B = [0; 1];
>> poles = [0.3+j*.2, 0.3-j*.2];
>> K = place(A, B, poles)
    place: ndigits = 16
    K = 3.1300   3.4000
```

`ndigits` is a measure of the accuracy of pole placement.

Alternatively, we can obtain the same answer using Ackermann's formula and the command

```
>> K = acker(A, B, poles)
```

It is also possible to compute the state feedback gain matrix of Eqs. (9.15) or (9.16) using basic MATLAB commands as follows:

1. Generate the characteristic polynomial of the state matrix to use Eq. (9.15).

```
>> poly(A)
```

2. Obtain the coefficients of the desired characteristic polynomial from a set of desired eigenvalues given as the entries of vector **poles**.

$\gg \text{desired} = \text{poly}(\text{poles})$

The vector **desired** contains the desired coefficients in descending order.

3. For Eq. (9.16), generate the polynomial matrix for a matrix **A** corresponding to the desired polynomial:

$\gg \text{polyvalm}(\text{desired}, \text{A})$

9.2.5 Pole placement for multi-input systems

For multi-input systems, the solution to the pole placement problem for a specified set of eigenvalues is not unique. However, there are constraints on the eigenvectors of the closed-loop matrix. Using the singular value decomposition of the input matrix (see Appendix III), we can obtain a complete characterization of the eigenstructure of the closed-loop matrix that is assignable by state feedback.

Theorem 9.2

For any controllable pair (A, B) , there exists a state feedback matrix K that assigns the eigen-pairs $\{(\lambda_i, \mathbf{v}_{di}), i = 1, \dots, n\}$ if and only if

$$U_1^T [A - \lambda_i I_n] \mathbf{v}_{di} = \mathbf{0}_{n \times 1} \quad (9.21)$$

where

$$B = U \begin{bmatrix} Z \\ \mathbf{0}_{n-m \times m} \end{bmatrix} \quad (9.22)$$

with

$$U = [U_0 | U_1], U^{-1} = U^T \quad (9.23)$$

The state feedback gain is given by

$$K = -Z^{-1} U_0^T [V_d A V_d^{-1} - A] \quad (9.24)$$

Proof

The singular value decomposition of the full-rank input matrix is

$$B = U \begin{bmatrix} \Sigma_B \\ \mathbf{0}_{n-m \times m} \end{bmatrix} V^T = U \begin{bmatrix} Z \\ \mathbf{0}_{n-m \times m} \end{bmatrix}$$

$$Z = \Sigma_B V^T$$

where the matrices of singular vectors satisfy $U^{-1} = U^T$, $V^{-1} = V^T$ and the inverse of the matrix Z is

$$Z^{-1} = V \Sigma_B^{-1}$$

The closed-loop state matrix is

$$A_{cl} = A - BK = V_d \Lambda V_d^{-1}$$

$$\text{with } \Lambda = \text{diag}\{\lambda_1, \dots, \lambda_n\}, V_d = [\mathbf{v}_{d1}, \dots, \mathbf{v}_{dn}]$$

We write the matrix of left singular vectors U as in Eq. (9.23), and then multiply by the closed-loop state matrix

$$U^T [A - V_d \Lambda V_d^{-1}] = \begin{bmatrix} U_0^T \\ U_1^T \end{bmatrix} [A - V_d \Lambda V_d^{-1}] = U^T BK$$

Using the singular value decomposition of B , we have

$$\begin{bmatrix} U_0^T (A - V_d \Lambda V_d^{-1}) \\ U_1^T (A - V_d \Lambda V_d^{-1}) \end{bmatrix} = \begin{bmatrix} ZK \\ \mathbf{0}_{n-m \times m} \end{bmatrix}$$

The lower part of the matrix equality is equivalent to Eq. (9.21), while the upper part gives Eq. (9.24).

Condition (9.21) implies that the eigenvectors that can be assigned by state feedback must be in the null space of the matrix U_1^T ; that is, they must be mapped to zero by the matrix. This shows the importance of the singular value decomposition of the input matrix and its influence on the choice of available eigenvectors. Recall that the input matrix must also make the pair (A, B) controllable. While the discussion of this section helps us characterize state feedback in the multivariable case, it does not really provide an easy recipe for the choice of state feedback control. In fact, no such recipe exists, although there are alternative methods to choose the eigenvectors to endow the closed-loop system with desirable properties, these methods are beyond the scope of this text.

Example 9.5

For the pair

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -0.005 & -0.11 & -0.7 \end{bmatrix} \quad B = \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$$

and the eigenvalues {0.1, 0.2, 0.3},

1. Determine a basis set for the space from which each eigenvector can be chosen.
2. Select a set of eigenvectors and determine the state feedback gain matrix that assigns the desired eigenstructure. Verify that the closed-loop state matrix has the correct eigenstructure.

Solution

The singular value decomposition of the input matrix is obtained using MATLAB

```
>> [Ub, sigb, Vb] = svd(B)
```

$$\begin{aligned} \text{Ub} = \\ -0.5000 & \quad -0.5000 & \quad -0.7071 \\ -0.5000 & \quad -0.5000 & \quad 0.7071 \\ -0.7071 & \quad 0.7071 & \quad 0 \end{aligned}$$

$$\begin{aligned} \text{sigb} = \\ 1.8478 & \quad 0 \\ 0 & \quad 0.7654 \\ 0 & \quad 0 \end{aligned}$$

$$\begin{aligned} \text{Vb} = \\ -0.3827 & \quad 0.9239 \\ -0.9239 & \quad -0.38 \end{aligned}$$

The system is third order with two inputs—that is, $n = 3$, $m = 2$, $n-m = 1$ —and the matrix of left singular vectors is partitioned as

$$U_b = [U_0 | U_1] = \left[\begin{array}{cc|c} -0.5 & -0.5 & -0.7071 \\ -0.5 & -0.5 & 0.7071 \\ -0.7071 & 0.7071 & 0 \end{array} \right]$$

For the eigenvalue $\lambda_1 = 0.1$, Eq. (9.21) gives

$$U_1^T [A - 0.1I_3] \mathbf{v}_{d1} = [0.0707 \quad -0.7778 \quad 0.7071] \mathbf{v}_{d1} = \mathbf{0}_{3 \times 1}$$

Example 9.5—cont'd

The corresponding eigenvector is

$$\mathbf{v}_{d1} = [0.7383 \quad 0.4892 \quad 0.4643]^T$$

Similarly, we determine the two remaining eigenvectors and obtain the modal matrix of eigenvectors

$$V_d = \begin{bmatrix} 0.7383 & 0.7620 & 0.7797 \\ 0.4892 & 0.4848 & 0.4848 \\ 0.4643 & 0.4293 & 0.3963 \end{bmatrix}$$

The closed-loop state matrix is

$$A_{cl} = V_d \Lambda V_d^{-1} = \begin{bmatrix} 0.2377 & 2.0136 & -2.3405 \\ 0.2377 & 1.0136 & -1.3405 \\ 0.6511 & -0.2695 & -0.6513 \end{bmatrix}$$

We also need the matrix

$$\Sigma_b V_b = \begin{bmatrix} -0.7071 & -1.7071 \\ 0.7071 & -0.2929 \\ 0 & 0 \end{bmatrix}$$

The state feedback gain matrix is

$$K = -Z^{-1} U_0^T [A_{cl} - A] = \begin{bmatrix} -0.4184 & 1.1730 & -2.3892 \\ -0.2377 & -1.0136 & 2.3405 \end{bmatrix}$$

The MATLAB command **place** gives the solution

$$K_m = \begin{bmatrix} -0.0346 & 0.3207 & -2.0996 \\ -0.0499 & 0.0852 & 0.7642 \end{bmatrix}$$

that also assigns the desired eigenvalues but with a different choice of eigenvectors.

9.2.6 Pole placement by output feedback

As one would expect, using output feedback limits our ability to assign the eigenvalues of the state system relative to what is achievable using state feedback. It is, in general, not possible to arbitrarily assign the system poles even if the system is completely controllable and completely observable. Only l poles can be arbitrarily assigned if the system has l linearly independent outputs. It is possible to arbitrarily assign the controllable dynamics of the system using dynamic output feedback, and a satisfactory solution can be obtained if the system is stabilizable and detectable. Several approaches are available for the design

of such a dynamic controller. One solution is to obtain an estimate of the state using the output and input of the system and use it in state feedback as explained in [Section 9.6](#).

9.3 Servo problem

The schemes shown in [Figs. 9.1 and 9.2](#) are regulators that drive the system state to zero starting from any initial condition capable of rejecting impulse disturbances. In practice, it is often necessary to track a constant reference input \mathbf{r} with zero steady-state error. For this purpose, a possible approach is to use the **two degree-of-freedom control scheme** in [Fig. 9.7](#), so called because we now have two matrices to select: the feedback gain matrix K and the reference gain matrix F .

The reference input of [Eq. \(9.2\)](#) becomes $\mathbf{v}(k) = F\mathbf{r}(k)$, and the control law is chosen as

$$\mathbf{u}(k) = -K\mathbf{x}(k) + F\mathbf{r}(k) \quad (9.25)$$

with $\mathbf{r}(k)$ the reference input to be tracked. The corresponding closed-loop system equations are

$$\begin{aligned} \mathbf{x}(k+1) &= A_{cl}\mathbf{x}(k) + BF\mathbf{r}(k) \\ y(k) &= C\mathbf{x}(k) \end{aligned} \quad (9.26)$$

where the closed-loop state matrix is

$$A_{cl} = A - BK$$

Using the formula for the transfer function, the z -transform of the corresponding output is given by (see [Section 7.8](#))

$$\mathbf{Y}(z) = C[zI_n - A_{cl}]^{-1}BF\mathbf{R}(z)$$

The steady-state tracking error for a unit step input is given by

$$\begin{aligned} \lim_{z \rightarrow 1} (z-1)\{\mathbf{Y}(z) - \mathbf{R}(z)\} &= \lim_{z \rightarrow 1} \left\{ C[zI_n - A_{cl}]^{-1}BF - I \right\} \\ &= C[I_n - A_{cl}]^{-1}BF - I \end{aligned}$$

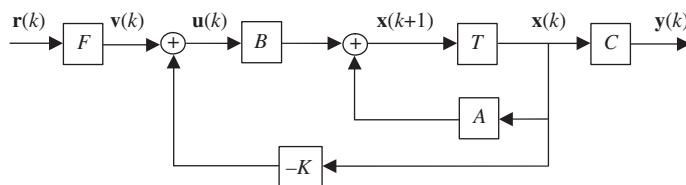


Figure 9.7
Block diagram of the two degree-of-freedom controller.

For zero steady-state error, we require the condition

$$C[I_n - A_{cl}]^{-1}BF = I_n \quad (9.27)$$

If the system is square ($m = l$) and A_{cl} is stable (no unity eigenvalues), we solve for the reference gain

$$F = [C(I_n - A_{cl})^{-1}B]^{-1} \quad (9.28)$$

Example 9.6

Design a state-space controller for the discretized state-space model of the DC motor speed control system described in Example 6.9 (with $T = 0.02$) to obtain zero steady-state error due to a unit step, a damping ratio of 0.7, and a settling time of about 1 s.

Solution

The discretized transfer function of the system with digital-to-analog converter (DAC) and ADC is

$$G_{ZAS}(z) = (1 - z^{-1}) \mathcal{Z}\left\{\frac{G(s)}{s}\right\} = 1.8604 \times 10^{-4} \frac{z + 0.9293}{(z - 0.8187)(z - 0.9802)}$$

The corresponding state-space model, computed with the MATLAB command `ss(G)`, is

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 1.799 & -0.8025 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 0.01563 \\ 0 \end{bmatrix} u(k)$$

$$y(k) = [0.01191 \quad 0.01107] \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$$

The desired eigenvalues of the closed-loop system are selected as $\{0.9 \pm j0.09\}$, as in Example 6.9. This yields the feedback gain vector

$$K = [-0.068517 \quad 0.997197]$$

and the closed-loop state matrix

$$A_{cl} = \begin{bmatrix} 1.8 & -0.8181 \\ 1 & 0 \end{bmatrix}$$

The feedforward gain is

$$F = [C(I_n - A_{cl})^{-1}B]^{-1}$$

$$= \left[[0.01191 \quad 0.01107] \left(\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} 1.8 & -0.8181 \\ 1 & 0 \end{bmatrix} \right)^{-1} \begin{bmatrix} 0.01563 \\ 0 \end{bmatrix} \right]^{-1} = 50.42666$$

The response of the system to a step reference input r is shown in Fig. 9.8. The system has a settling time of about 0.84 s and percentage overshoot of about 4%, with a peak time of about 1 s. All design specifications are met.

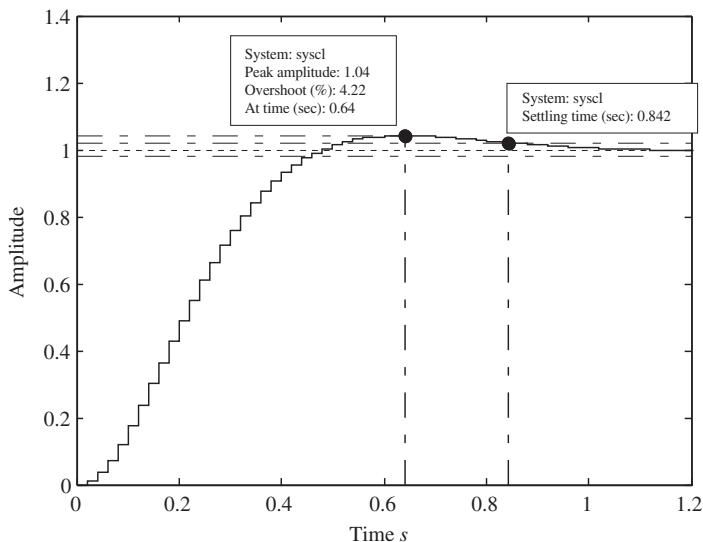
Example 9.6—cont'd

Figure 9.8
Step response of the closed-loop system of Example 9.6.

The control law (9.25) is equivalent to a **feedforward action** determined by F to yield zero steady-state error for a constant reference input \mathbf{r} . Because the forward action does not include any form of feedback, this approach is not robust to modeling uncertainties. Thus, modeling errors (which always occur in practice) will result in nonzero steady-state error. To eliminate such errors, we introduce the **integral control** shown in Fig. 9.9, with a new state added for each control error integrated.

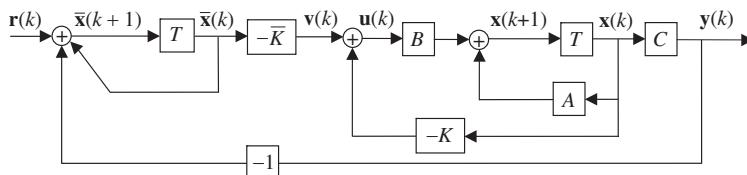


Figure 9.9
Control scheme with integral control.

The resulting state-space equations are

$$\begin{aligned}\mathbf{x}(k+1) &= A\mathbf{x}(k) + B\mathbf{u}(k) \\ \bar{\mathbf{x}}(k+1) &= \bar{\mathbf{x}}(k) + \mathbf{r}(k) - \mathbf{y}(k) \\ \mathbf{y}(k) &= C\mathbf{x}(k) \\ \mathbf{u}(k) &= K\mathbf{x}(k) - \bar{K}\bar{\mathbf{x}}(k)\end{aligned}\tag{9.29}$$

where $\bar{\mathbf{x}}$ is $l \times 1$. The state-space equations can be combined and rewritten in terms of an augmented state vector $\mathbf{x}_a(k) = [\mathbf{x}(k) \ \bar{\mathbf{x}}(k)]^T$ as

Example 9.6—cont'd

$$\begin{bmatrix} \mathbf{x}(k+1) \\ \bar{\mathbf{x}}(k+1) \end{bmatrix} = \begin{bmatrix} A & 0 \\ -C & I_I \end{bmatrix} \begin{bmatrix} \mathbf{x}(k) \\ \bar{\mathbf{x}}(k) \end{bmatrix} - \begin{bmatrix} B \\ 0 \end{bmatrix} \begin{bmatrix} K & \bar{K} \end{bmatrix} \begin{bmatrix} \mathbf{x}(k) \\ \bar{\mathbf{x}}(k) \end{bmatrix} + \begin{bmatrix} 0 \\ I_I \end{bmatrix} \mathbf{r}(k)$$

$$\mathbf{y}(k) = [C \quad 0] \begin{bmatrix} \mathbf{x}(k+1) \\ \bar{\mathbf{x}}(k+1) \end{bmatrix}$$

That is,

$$\begin{aligned} \mathbf{x}_d(k+1) &= (\tilde{A} - \tilde{B}\tilde{K})\mathbf{x}_d(k) + \begin{bmatrix} 0 \\ I_I \end{bmatrix} \mathbf{r}(k) \\ \mathbf{y}(k) &= [C \quad 0] \mathbf{x}_d(k) \end{aligned} \quad (9.30)$$

where

$$\tilde{A} = \begin{bmatrix} A & 0 \\ -C & I_I \end{bmatrix} \quad \tilde{B} = \begin{bmatrix} B \\ 0 \end{bmatrix} \quad \tilde{K} = [K \quad \bar{K}] \quad (9.31)$$

It can be shown that the system of Eq. (9.31) is controllable if and only if the original system is controllable. The eigenvalues of the closed-loop system state matrix $A_{cl} = (\tilde{A} - \tilde{B}\tilde{K})$ can be arbitrarily assigned by computing the gain matrix \tilde{K} using any of the procedures for the regulator problem as described in Section 9.2.

Example 9.7

Solve the design problem presented in Example 9.6 using integral control.

Solution

The state-space matrices of the system are

$$\begin{aligned} A &= \begin{bmatrix} 1.799 & -0.8025 \\ 1 & 0 \end{bmatrix} & B &= \begin{bmatrix} 0.01563 \\ 0 \end{bmatrix} \\ C &= [0.01191 \quad 0.01107] \end{aligned}$$

Adding integral control, we obtain

$$\tilde{A} = \begin{bmatrix} A & 0 \\ -C & 1 \end{bmatrix} = \begin{bmatrix} 1.799 & -0.8025 & 0 \\ 1 & 0 & 0 \\ -0.01191 & -0.01107 & 1 \end{bmatrix} \quad \tilde{B} = \begin{bmatrix} B \\ 0 \end{bmatrix} = \begin{bmatrix} 0.01563 \\ 0 \\ 0 \end{bmatrix}$$

In Example 9.6, the eigenvalues were selected as $\{0.9 \pm j0.09\}$. Using integral control increases the order of the system by one, and an additional eigenvalue must be selected. The desired eigenvalues are selected as $\{0.9 \pm j0.09, 0.2\}$, and the additional eigenvalue at 0.2 is chosen for its negligible effect on the overall dynamics. This yields the feedback gain vector

Example 9.7—cont'd

$$\tilde{K} = [51.1315 \quad -40.4431 \quad -40.3413]$$

The closed-loop system state matrix is

$$A_{cl} = \begin{bmatrix} 1 & -0.1706 & 0.6303 \\ 1 & 0 & 0 \\ -0.0119 & -0.0111 & 1 \end{bmatrix}$$

The response of the system to a unit step reference signal r is shown in Fig. 9.10. The figure shows that the control specifications are satisfied. The settling time of 0.87 is well below the specified value of 1 s, and the percentage overshoot is about 4.2%, which is less than the value corresponding to $\zeta = 0.7$ for the dominant pair.

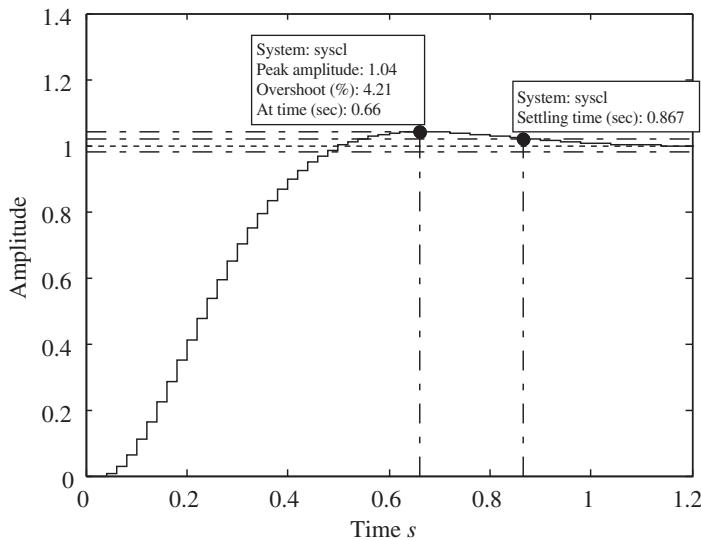


Figure 9.10
Step response of the closed-loop system of Example 9.7.

9.4 Invariance of system zeros

A severe limitation of the state-feedback control scheme is that it cannot change the location of the zeros of the system, which significantly affect the transient response. To show this, we consider the z -transform of the system including the direct transmission matrix D :

$$\mathbf{x}(k+1) = (A - BK)\mathbf{x}(k) + B\mathbf{v}(k)$$

$$\mathbf{y}(k) = C\mathbf{x}(k) - DK\mathbf{x}(k) - D\mathbf{v}(k)$$

The z-transform is given by

$$\begin{aligned} -(zI_n - A - BK)\mathbf{X}(z) + B\mathbf{V}(z) &= 0 \\ C\mathbf{X}(z) + D\mathbf{V}(z) &= \mathbf{Y}(z) \end{aligned} \quad (9.32)$$

If $z = z_0$ is a zero of the system, then $\mathbf{Y}(z_0)$ is zero with $\mathbf{V}(z_0)$ and $\mathbf{X}(z_0)$ nonzero. Thus, for $z = z_0$, the state-space Eq. (9.32) can be rewritten as

$$\begin{bmatrix} -(z_0 I_n - A + BK) & B \\ C & D \end{bmatrix} \begin{bmatrix} \mathbf{X}(z_0) \\ \mathbf{V}(z_0) \end{bmatrix} = \mathbf{0} \quad (9.33)$$

Using the vector $\mathbf{V}(z_0) + K\mathbf{X}(z_0)$ in place of $\mathbf{V}(z_0)$, we rewrite Eq. (9.33) in terms of the state-space matrices of the open-loop system as

$$\begin{bmatrix} -(z_0 I_n - A + BK) & B \\ C & D \end{bmatrix} \begin{bmatrix} \mathbf{X}(z_0) \\ \mathbf{V}(z_0) + K\mathbf{X}(z_0) \end{bmatrix} = \begin{bmatrix} -(z_0 I_n - A) & B \\ C & D \end{bmatrix} \begin{bmatrix} \mathbf{X}(z_0) \\ \mathbf{V}(z_0) \end{bmatrix} = 0$$

We observe that with the state feedback $\mathbf{u}(k) = -K\mathbf{x}(k) + \mathbf{r}(k)$, the open-loop state-space quadruple (A, B, C, D) becomes

$$(A - BK, B, C - DK, D)$$

Thus, the zeros of the closed-loop system are the same as those of the plant and are invariant under state feedback. Similar reasoning can establish the same result for the system of Eq. (9.26).

Example 9.8

Consider the following continuous-time system:

$$G(s) = \frac{s+1}{(2s+1)(3s+1)}$$

Obtain a discrete model for the system with digital control and a sampling period $T = 0.02$, and then design a state-space controller with integral control and with the same closed-loop eigenvalues as in Example 9.7.

Solution

The analog system with DAC and ADC has the transfer function

$$G_{ZAS}(z) = (1 - z^{-1}) \mathcal{Z} \left\{ \frac{G(s)}{s} \right\} = 33.338 \times 10^{-4} \frac{z - 0.9802}{(z - 0.9934)(z - 0.99)}$$

with an open-loop zero at 0.9802. The corresponding state-space model (computed with the MATLAB command `ss(G)`) is

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 1.983 & -0.983 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 0.0625 \\ 0 \end{bmatrix} u(k)$$

Example 9.8—cont'd

$$y(k) = [0.0534 \quad 0.0524] \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix}$$

The desired eigenvalues of the closed-loop system are selected as $\{0.9 \pm j0.09\}$ with an additional eigenvalue at 0.2, and this yields the feedback gain vector

$$\tilde{K} = [15.7345 \quad -24.5860 \quad -2.19016]$$

The closed-loop system state matrix is

$$A_{cl} = \begin{bmatrix} 1 & -0.55316 & 13.6885 \\ 1 & 0 & 0 \\ 0.05342 & -0.05236 & 1 \end{bmatrix}$$

The response of the system to a unit step reference signal r , shown in Fig. 9.11, has a huge peak overshoot due to the closed-loop zero at 0.9802. The closed-loop control cannot change the location of the zero.

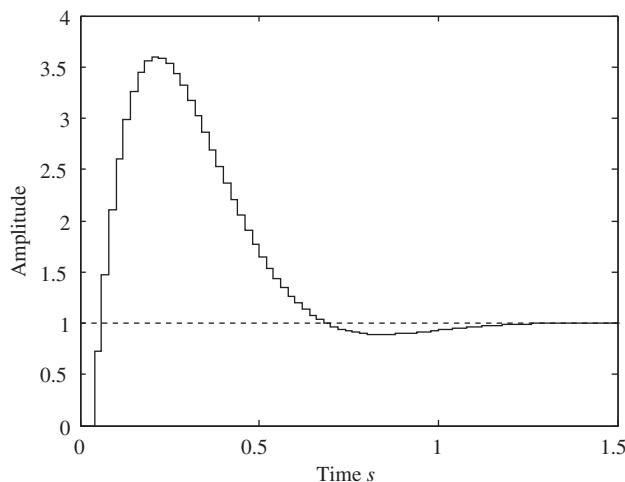


Figure 9.11
Step response of the closed-loop system of Example 9.8.

9.5 State estimation

In most applications, measuring the entire state vector is impossible or prohibitively expensive. To implement state feedback control, an estimate $\hat{x}(k)$ of the state vector can be used. The state vector can be estimated from the input and output measurements by using a **state estimator** or **observer**.

9.5.1 Full-order observer

To estimate all the states of the system, one could in theory use a system with the same state equation as the plant to be observed. In other words, one could use the open-loop system

$$\hat{\mathbf{x}}(k+1) = A\hat{\mathbf{x}}(k) + B\mathbf{u}(k)$$

However, this open-loop estimator assumes perfect knowledge of the system dynamics and lacks the feedback needed to correct the errors that are inevitable in any implementation. The limitations of this observer become obvious on examining its error dynamics. We define the estimation error as $\tilde{\mathbf{x}} = \mathbf{x} - \hat{\mathbf{x}}$. We obtain the error dynamics by subtracting the open-loop observer dynamics from the system dynamics Eq. (9.1).

$$\tilde{\mathbf{x}}(k+1) = A\tilde{\mathbf{x}}(k)$$

The error dynamics are determined by the state matrix of the system and cannot be chosen arbitrarily. For an unstable system, the observer will be unstable and cannot track the state of the system.

A practical alternative is to feed back the difference between the measured and the estimated output of the system, as shown in Fig. 9.12. This yields to the following observer:

$$\hat{\mathbf{x}}(k+1) = A\hat{\mathbf{x}}(k) + B\mathbf{u}(k) + L[y(k) - C\hat{\mathbf{x}}(k)] \quad (9.34)$$

Subtracting the observer state equation from the system dynamics yields the estimation error dynamics

$$\tilde{\mathbf{x}}(k+1) = (A - LC)\tilde{\mathbf{x}}(k) \quad (9.35)$$

The error dynamics are governed by the eigenvalues of the observer matrix $A_0 = A - LC$. We transpose the matrix to obtain

$$A_0^T = A^T - C^T L^T \quad (9.36)$$

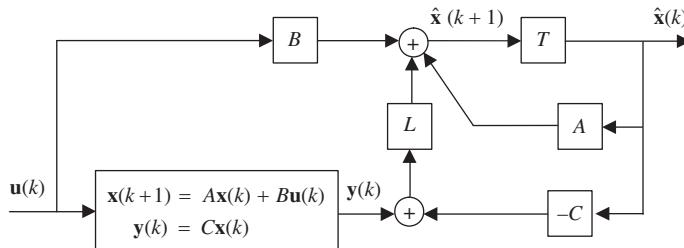


Figure 9.12
Block diagram of the full-order state estimator.

which has the same eigenvalues as the observer matrix. We observe that Eq. (9.36) is identical to the controller design Eq. (9.4) with the pair (A, B) replaced by the pair (A^T, C^T) . We therefore have the following theorem.

Theorem 9.3: State estimation

If the pair (A, C) is observable, then there exists a feedback gain matrix L that arbitrarily assigns the observer poles to any set $\{\lambda_i, i = 1, \dots, n\}$. Furthermore, if the pair (A, C) is detectable, then the observable modes can all be arbitrarily assigned.

Proof

Based on Theorem 8.12, system (A, C) is observable (detectable) if and only if (A^T, C^T) is controllable (stabilizable). Therefore, Theorem 9.3 follows from Theorem 9.1.

Based on Theorem 9.3, the matrix gain L can be determined from the desired observer poles, as discussed in Section 9.2. Hence, we can arbitrarily select the desired observer poles or the associated characteristic polynomial. From Eq. (9.36), it follows that the MATLAB command for the solution of the observer pole placement problem is

```
>> L = place(A', C', poles)'
```

Example 9.9

Determine the observer gain matrix L for the discretized state-space model of the armature-controlled DC motor of Example 9.4 with the observer eigenvalues selected as $\{0.1, 0.2 \pm j0.2\}$.

Solution

Recall that the system matrices are

$$A = \begin{bmatrix} 1.0 & 0.1 & 0.0 \\ 0.0 & 0.995 & 0.0095 \\ 0.0 & -0.0947 & 0.8954 \end{bmatrix} \quad B = \begin{bmatrix} 1.622 \times 10^{-6} \\ 4.821 \times 10^{-4} \\ 9.468 \times 10^{-2} \end{bmatrix} \quad C = [1 \quad 0 \quad 0]$$

The MATLAB command **place** gives the observer gain

$$L = \begin{bmatrix} 2.3949 \\ 18.6734 \\ 4.3621 \end{bmatrix}$$

Expression (9.34) represents a **prediction observer**, because the estimated state vector (and any associated control action) at a given sampling instant does not depend on the current measured value of the system output. Alternatively, a **filtering observer** estimates the state vector based on the current output (assuming negligible computation time). The error correction term for the filtering observer uses the difference between the current output $\mathbf{y}(k+1)$ and its estimate and has the form

$$\begin{aligned}\widehat{\mathbf{x}}(k+1) &= A\widehat{\mathbf{x}}(k) + B\mathbf{u}(k) + L[\mathbf{y}(k+1) - C\widehat{\mathbf{x}}(k+1)] \\ &= A\widehat{\mathbf{x}}(k) + B\mathbf{u}(k) + L[\mathbf{y}(k+1) - C(A\widehat{\mathbf{x}}(k) + B\mathbf{u}(k))]\end{aligned}\quad (9.37)$$

The error dynamics are now represented by

$$\widetilde{\mathbf{x}}(k+1) = (A - LCA)\widetilde{\mathbf{x}}(k) \quad (9.38)$$

Expression (9.38) is the same as Expression (9.35) with the matrix product CA substituted for C . The observability matrix $\overline{\mathcal{O}}$ of system (A, CA) is

$$\overline{\mathcal{O}} = \begin{bmatrix} CA \\ \vdash \vdash \vdash \\ CA^2 \\ \vdash \vdash \vdash \\ \vdots \\ \vdash \vdash \vdash \\ CA^n \end{bmatrix} = \mathcal{O}A$$

where \mathcal{O} is the observability matrix of the pair (A, C) (see Section 8.3). Thus, if the pair (A, C) is observable, the pair (A, CA) is observable unless A has one or more zero eigenvalues. If A has zero eigenvalues, the pair (A, CA) is detectable because the zero eigenvalues are associated with stable modes. Further, the zero eigenvalues are associated with the fastest modes, and the design of the observer can be completed by selecting a matrix L that assigns suitable values to the remaining eigenvalues of $A - LCA$.

Example 9.10

Determine the filtering observer gain matrix L for the system of Example 9.9.

Solution

Using the MATLAB command **place**,

$$\gg L = \text{place}(A', (C * A)', \text{poles})'$$

we obtain the observer gain

$$\begin{aligned}L = \\ 1.0e + 022* \\ 0.009910699530206 \\ 0.140383004697920 \\ 4.886483453094875\end{aligned}$$

9.5.2 Reduced-order observer

A full-order observer is designed so that the entire state vector is estimated from knowledge of the input and output of the system. The reader might well ask, why estimate n state variables when we already have l measurements that are linear functions of the same variables? Would it be possible to estimate $n-l$ variables only and use them with the measurements to estimate the entire state?

This is precisely what is done in the design of a **reduced-order observer**. A reduced-order observer is generally more efficient than a full-order observer. However, a full-order observer may be preferable in the presence of significant measurement noise. In addition, the design of the reduced-order observer is more complex.

We consider the linear time-invariant system Eq. (9.1) where the input matrix B and the output matrix C are assumed full rank. Then the entries of the output vector $\mathbf{y}(k)$ are linearly independent and form a partial state vector of length l , leaving $n-l$ variables to be determined. We thus have the state vector

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{z} \end{bmatrix} = Q_o^{-1} \mathbf{x} = \begin{bmatrix} C \\ M \end{bmatrix} \mathbf{x} \quad (9.39)$$

where M is a full-rank $n-l \times n$ matrix with rows that are linearly independent of those of C , and \mathbf{z} is the unknown partial state.

The state-space matrices for the transformed state variables are

$$C_t = CQ_o = \left[I_l \mid \mathbf{0}_{l \times n-l} \right] \quad (9.40)$$

$$B_t = Q_o^{-1} B = \left[\begin{array}{c|c} \overline{B}_1 & \\ \hline \cdots & \\ \overline{B}_2 & \end{array} \right] \} \begin{array}{c} l \\ n-l \end{array} \quad (9.41)$$

$$A_t = Q_o^{-1} A Q_o = \left[\begin{array}{c|c} \overline{A}_1 & \overline{A}_2 \\ \hline \overline{A}_3 & \overline{A}_4 \\ \hline l & n-l \end{array} \right] \} \begin{array}{c} l \\ n-l \end{array} \quad (9.42)$$

The identity in the output matrix indicates that for the transformed system the measurement vector $\mathbf{y}(k)$ is not part of the state vector. The state equation for the unknown partial state is

$$\mathbf{z}(k+1) = \overline{A}_3 \mathbf{y}(k) + \overline{A}_4 \mathbf{z}(k) + \overline{B}_2 \mathbf{u}(k) \quad (9.43)$$

Using the state equation for the known partial state, that is, the measurements, we define an output variable to form a state-space model with Eq. (9.43) as

$$\mathbf{y}_z(k) = \mathbf{y}(k+1) - \bar{\mathbf{A}}_1\mathbf{y}(k) + \bar{\mathbf{B}}_1\mathbf{u}(k) = \bar{\mathbf{A}}_2\mathbf{z}(k) \quad (9.44)$$

This output represents the portion of the known partial state $\mathbf{y}(k+1)$ that is computed using the unknown partial state. The observer dynamics, including the error in computing \mathbf{y}_z , are assumed linear time invariant of the form

$$\begin{aligned} \hat{\mathbf{z}}(k+1) &= \bar{\mathbf{A}}_3\mathbf{y}(k) + \bar{\mathbf{A}}_4\hat{\mathbf{z}}(k) + \bar{\mathbf{B}}_2\mathbf{u}(k) + L[\mathbf{y}_z(k) - \bar{\mathbf{A}}_2\hat{\mathbf{z}}(k)] \\ &= (\bar{\mathbf{A}}_4 - L\bar{\mathbf{A}}_2)\hat{\mathbf{z}}(k) + \bar{\mathbf{A}}_3\mathbf{y}(k) + L[\mathbf{y}(k+1) - \bar{\mathbf{A}}_1\mathbf{y}(k) - \bar{\mathbf{B}}_1\mathbf{u}(k)] + \bar{\mathbf{B}}_2\mathbf{u}(k) \end{aligned} \quad (9.45)$$

where $\hat{\mathbf{z}}$ denotes the estimate of the partial state vector \mathbf{z} . Unfortunately, the observer Eq. (9.45) includes the term $\mathbf{y}(k+1)$, which is not available at time k . Moving the term to the LHS reveals that its use can be avoided by estimating the variable

$$\bar{\mathbf{x}}(k) = \hat{\mathbf{z}}(k) - Ly(k) \quad (9.46)$$

Using the observer dynamics Eq. (9.45) and the definition (9.46), we obtain the observer

$$\begin{aligned} \bar{\mathbf{x}}(k+1) &= \hat{\mathbf{z}}(k+1) - Ly(k+1) \\ &= [\bar{\mathbf{A}}_4 - L\bar{\mathbf{A}}_2](\bar{\mathbf{x}}(k) + Ly(k)) + [\bar{\mathbf{A}}_3 - L\bar{\mathbf{A}}_1]\mathbf{y}(k) + [\bar{\mathbf{B}}_2 - L\bar{\mathbf{B}}_1]\mathbf{u}(k) \end{aligned}$$

The observer for the unknown partial state is given by

$$\bar{\mathbf{x}}(k+1) = A_o\bar{\mathbf{x}}(k) + A_y\mathbf{y}(k) + B_o\mathbf{u}(k) \quad (9.47)$$

where

$$\begin{aligned} A_o &= \bar{\mathbf{A}}_4 - L\bar{\mathbf{A}}_2 \\ A_y &= A_oL + \bar{\mathbf{A}}_3 - L\bar{\mathbf{A}}_1 \\ B_o &= \bar{\mathbf{B}}_2 - L\bar{\mathbf{B}}_1 \end{aligned} \quad (9.48)$$

The block diagram of the reduced-order observer is shown in Fig. 9.13.

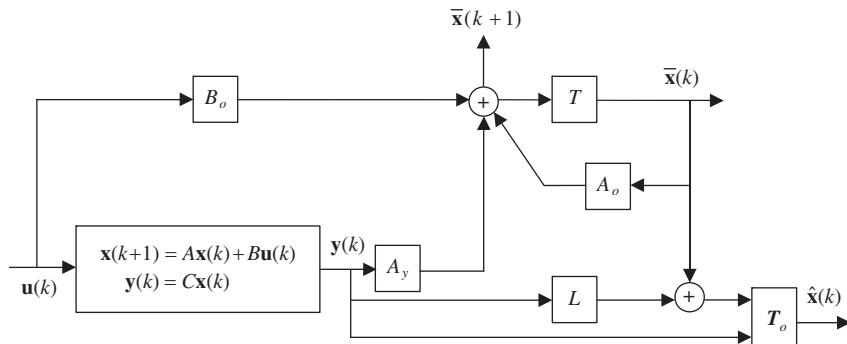


Figure 9.13
Block diagram of the reduced-order observer.

The dynamic of the reduced-order observer Eq. (9.47) is governed by the matrix A_o . The eigenvalues of A_o must be selected inside the unit circle and must be sufficiently fast to track the state of the observed system. This reduces observer design to the solution of Eq. (9.48) for the observer gain matrix L . Once L is obtained, the other matrices in Eq. (9.48) can be computed and the state vector $\hat{\mathbf{x}}$ can be obtained using the equation

$$\begin{aligned}\hat{\mathbf{x}}(k) &= Q_o \begin{bmatrix} \mathbf{y}(k) \\ \hat{\mathbf{z}}(k) \end{bmatrix} \\ &= Q_o \begin{bmatrix} I_l & 0_{l \times n-l} \\ L & I_{n-l} \end{bmatrix} \begin{bmatrix} \mathbf{y}(k) \\ \bar{\mathbf{x}}(k) \end{bmatrix} = T_o \begin{bmatrix} \mathbf{y}(k) \\ \bar{\mathbf{x}}(k) \end{bmatrix}\end{aligned}\quad (9.49)$$

where the transformation matrix Q_o is defined in Eq. (9.39).

Transposing the state matrix of Eq. (9.48) yields

$$A_o^T = \bar{A}_4^T - \bar{A}_2^T L^T \quad (9.50)$$

Because Eq. (9.50) is identical in form to the controller design Eq. (9.4), it can be solved as discussed in Section 9.2. We recall that the poles of the matrix A_o^T can be arbitrarily assigned provided that the pair $(\bar{A}_4^T, \bar{A}_2^T)$ is controllable. From the duality concept discussed in Section 8.6, this is equivalent to the observability of the pair (\bar{A}_4, \bar{A}_2) . Theorem 9.4 gives a necessary and sufficient condition for the observability of the pair.

Theorem 9.4

The pair (\bar{A}_4, \bar{A}_2) is observable if and only if system (A, C) is observable.

Proof

The proof is left as an exercise.

Example 9.11

Design a reduced-order observer for the discretized state-space model of the armature-controlled DC motor of Example 9.4 with the observer eigenvalues selected as $\{0.2 \pm j0.2\}$.

Solution

Recall that the system matrices are

Example 9.11—cont'd

$$A = \begin{bmatrix} 1.0 & 0.1 & 0.0 \\ 0.0 & 0.995 & 0.0095 \\ 0.0 & -0.0947 & 0.8954 \end{bmatrix} \quad B = \begin{bmatrix} 1.622 \times 10^{-6} \\ 4.821 \times 10^{-4} \\ 9.468 \times 10^{-2} \end{bmatrix} \quad C = [1 \ 0 \ 0]$$

The output matrix C is in the required form, and there is no need for similarity transformation. The second and third state variables must be estimated. The state matrix is partitioned as

$$A = \left[\begin{array}{c|cc} a_1 & \mathbf{a}_2^T \\ \hline \mathbf{a}_3 & A_4 \end{array} \right] = \left[\begin{array}{c|ccc} 1.0 & & 0.1 & & 0.0 \\ \hline 0.0 & & 0.9995 & & 0.0095 \\ 0.0 & & -0.0947 & & 0.8954 \end{array} \right]$$

The similarity transformation can be selected as an identity matrix; that is,

$$Q_o = Q_o^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

and therefore we have $A_t = A$, $B_t = B$, and $C_t = C$. Hence, we need to solve the linear equation

$$A_o = \begin{bmatrix} 0.9995 & 0.0095 \\ -0.0947 & 0.8954 \end{bmatrix} - \bar{\mathbf{I}}[0.1 \ 0]$$

to obtain the observer gain

$$\bar{\mathbf{I}} = [14.949 \ 550.191]^T$$

The corresponding observer matrices are

$$A_o = \begin{bmatrix} -0.4954 & 0.0095 \\ -55.1138 & 0.8954 \end{bmatrix}$$

$$\mathbf{b}_o = \bar{\mathbf{b}}_2 - \bar{\mathbf{I}}\bar{\mathbf{b}}_1 = \begin{bmatrix} 4.821 \times 10^{-4} \\ 9.468 \times 10^{-2} \end{bmatrix} - \begin{bmatrix} 14.949 \\ 550.191 \end{bmatrix} \times 1.622 \times 10^{-6} = \begin{bmatrix} 0.04579 \\ 9.37876 \end{bmatrix} \times 10^{-2}$$

$$\begin{aligned} \mathbf{a}_y &= A_o \mathbf{I} + \mathbf{a}_3 - \mathbf{I} \mathbf{a}_1 \\ &= \begin{bmatrix} -0.4954 & 0.0095 \\ -55.1138 & 0.8954 \end{bmatrix} \begin{bmatrix} 14.949 \\ 550.191 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 14.949 \\ 550.191 \end{bmatrix} \times 1 = \begin{bmatrix} -0.1713 \\ -8.8145 \end{bmatrix} \times 10^2 \end{aligned}$$

The state estimate can be computed using

$$\hat{\mathbf{x}}(k) = Q_o \begin{bmatrix} 1 & \mathbf{0}_{1 \times 2}^T \\ \mathbf{I} & I_2 \end{bmatrix} \begin{bmatrix} \mathbf{y}(k) \\ \bar{\mathbf{x}}(k) \end{bmatrix} = \begin{bmatrix} 1 & \mathbf{0}_{1 \times 2}^T \\ \frac{1}{14.949} & I_2 \\ \frac{550.191}{14.949} & \mathbf{I}_2 \end{bmatrix} \begin{bmatrix} \mathbf{y}(k) \\ \bar{\mathbf{x}}(k) \end{bmatrix}$$

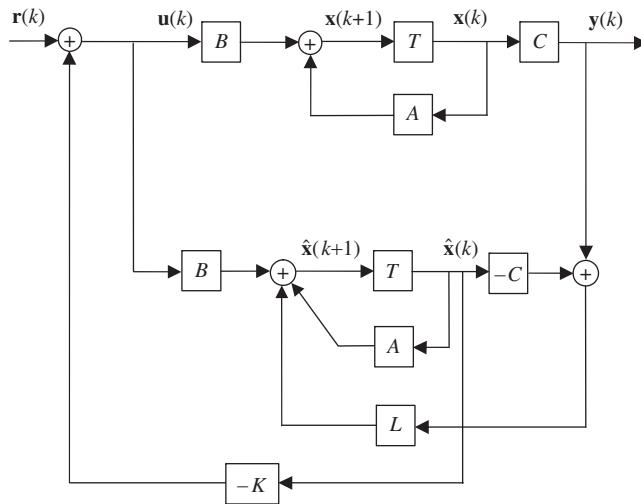


Figure 9.14

Block diagram of a system with observer state feedback.

9.6 Observer state feedback

If the state vector is not available for feedback control, a state estimator can be used to generate the control action as shown in Fig. 9.14. The corresponding control vector is

$$\mathbf{u}(k) = -K\hat{\mathbf{x}}(k) + \mathbf{v}(k) \quad (9.51)$$

Substituting in the state Eq. (9.1) gives

$$\mathbf{x}(k+1) = \mathbf{A}\mathbf{x}(k) - \mathbf{B}K\hat{\mathbf{x}}(k) + \mathbf{B}\mathbf{v}(k) \quad (9.52)$$

Adding and subtracting the term $BK\mathbf{x}(k)$, we rewrite Eq. (9.52) in terms of the estimation error $\tilde{\mathbf{x}} = \mathbf{x} - \hat{\mathbf{x}}$ as

$$\mathbf{x}(k+1) = (\mathbf{A} - BK)\mathbf{x}(k) + BK\tilde{\mathbf{x}}(k) + B\mathbf{v}(k) \quad (9.53)$$

If a full-order (predictor) observer is used, by combining Eq. (9.53) with Eq. (9.35), we obtain

$$\begin{bmatrix} \mathbf{x}(k+1) \\ \tilde{\mathbf{x}}(k+1) \end{bmatrix} = \begin{bmatrix} \mathbf{A} - BK & BK \\ 0 & \mathbf{A} - LC \end{bmatrix} \begin{bmatrix} \mathbf{x}(k) \\ \tilde{\mathbf{x}}(k) \end{bmatrix} + \begin{bmatrix} \mathbf{B} \\ 0 \end{bmatrix} \mathbf{v}(k) \quad (9.54)$$

The state matrix of Eq. (9.54) is block triangular and its characteristic polynomial is

$$\begin{aligned} \Delta_{cl}(\lambda) &= \det[\lambda I - (\mathbf{A} - BK)]\det[\lambda I - (\mathbf{A} - LC)] \\ &= \Delta_c(\lambda)\Delta_o(\lambda) \end{aligned} \quad (9.55)$$

Thus, the eigenvalues of the closed-loop system can be selected separately from those of the observer. This important result is known as the **separation theorem** or the **uncertainty equivalence principle**.

Analogously, if a reduced-order observer is employed, the estimation error $\tilde{\mathbf{x}}$ can be expressed in terms of the errors in estimating \mathbf{y} and \mathbf{z} as

$$\begin{aligned}\tilde{\mathbf{x}}(k) &= Q_o \begin{bmatrix} \tilde{\mathbf{y}}(k) \\ \tilde{\mathbf{z}}(k) \end{bmatrix} \\ \tilde{\mathbf{y}}(k) &= \mathbf{y}(k) - \hat{\mathbf{y}}(k) \quad \tilde{\mathbf{z}}(k) = \mathbf{z}(k) - \hat{\mathbf{z}}(k)\end{aligned}\tag{9.56}$$

We partition the matrix Q_o into an $n \times l$ matrix Q_y and an $n \times n-l$ matrix Q_z to allow the separation of the two error terms and rewrite the estimation error as

$$\tilde{\mathbf{x}}(k) = [Q_y \mid Q_z] \begin{bmatrix} \tilde{\mathbf{y}}(k) \\ \tilde{\mathbf{z}}(k) \end{bmatrix} = Q_y \tilde{\mathbf{y}}(k) + Q_z \tilde{\mathbf{z}}(k)$$

Assuming negligible measurement error $\tilde{\mathbf{y}}$, the estimation error reduces to

$$\tilde{\mathbf{x}}(k) = Q_z \tilde{\mathbf{z}}(k)\tag{9.58}$$

Substituting from Eq. (9.58) into the closed-loop Eq. (9.53) gives

$$\mathbf{x}(k+1) = (A - BK)\mathbf{x}(k) + BKQ_z \tilde{\mathbf{z}}(k) + B\mathbf{v}(k)\tag{9.59}$$

We evaluate $\tilde{\mathbf{z}}(k+1) = \mathbf{z}(k+1) - \hat{\mathbf{z}}(k+1)$ by subtracting Eq. (9.45) from Eq. (9.43) and use Eq. (9.44) to substitute $\bar{A}_2 \mathbf{z}(k)$ for \mathbf{y}_z

$$\tilde{\mathbf{z}}(k+1) = (\bar{A}_4 - L\bar{A}_2) \tilde{\mathbf{z}}(k)\tag{9.60}$$

Combining Eqs. (9.59) and (9.60), we obtain the equation

$$\begin{bmatrix} \mathbf{x}(k+1) \\ \tilde{\mathbf{z}}(k+1) \end{bmatrix} = \begin{bmatrix} A - BK & BKQ_z \\ 0_{n-l \times n} & \bar{A}_4 - L\bar{A}_2 \end{bmatrix} \begin{bmatrix} \mathbf{x}(k) \\ \tilde{\mathbf{z}}(k) \end{bmatrix} + \begin{bmatrix} B \\ 0 \end{bmatrix} \mathbf{v}(k)\tag{9.61}$$

The state matrix of Eq. (9.61) is block triangular and its characteristic polynomial is

$$\det[\lambda I - (\bar{A}_4 - L\bar{A}_2)] \det[\lambda I - (A - BK)]\tag{9.62}$$

Thus, as for the full-order observer, the closed-loop eigenvalues for the reduced-order observer state feedback can be selected separately from those of the reduced-order observer. The separation theorem therefore applies for reduced-order observers as well as for full-order observers.

In addition, combining the plant state equation and the estimator Eq. (9.47) and using the output equation, we have

$$\begin{bmatrix} \mathbf{x}(k+1) \\ \bar{\mathbf{x}}(k+1) \end{bmatrix} = \begin{bmatrix} A & 0_{n \times n-l} \\ A_y C & A_o \end{bmatrix} \begin{bmatrix} x(k) \\ \bar{x}(k) \end{bmatrix} + \begin{bmatrix} B \\ B_0 \end{bmatrix} \mathbf{u}(k) \quad (9.63)$$

We express the estimator state feedback of Eq. (9.51) as

$$\begin{aligned} \mathbf{u}(k) &= -K\hat{\mathbf{x}}(k) + \mathbf{v}(k) \\ &= -KT_o \begin{bmatrix} \mathbf{y}(k) \\ \bar{\mathbf{x}}(k) \end{bmatrix} + \mathbf{v}(k) = -K \begin{bmatrix} T_{oy} & | & T_{ox} \end{bmatrix} \begin{bmatrix} C\mathbf{x}(k) \\ \bar{\mathbf{x}}(k) \end{bmatrix} + \mathbf{v}(k) \end{aligned}$$

where T_{oy} and T_{ox} are partitions of T_o of Eq. (9.49) of order $n \times l$ and $n \times n-l$, respectively. Substituting in Eq. (9.63), we have

$$\begin{bmatrix} \mathbf{x}(k+1) \\ \bar{\mathbf{x}}(k+1) \end{bmatrix} = \begin{bmatrix} A - BKT_{oy}C & -BKT_{ox} \\ A_y C - B_o K T_{oy} C & A_o - B_o K T_{ox} \end{bmatrix} \begin{bmatrix} \mathbf{x}(k) \\ \bar{\mathbf{x}}(k) \end{bmatrix} + \begin{bmatrix} B \\ B_0 \end{bmatrix} \mathbf{v}(k) \quad (9.65)$$

Eq. (9.65) can be used to simulate the complete estimator state feedback system.

9.6.1 Choice of observer eigenvalues

In the selection of the observer poles or the associated characteristic polynomial, expression (9.55) or (9.62) must be considered. The choice of observer poles is not based on the constraints related to the control effort discussed in Section 9.2.3. However, the response of the closed-loop system must be dominated by the poles of the controller that meet the performance specifications. Therefore, as a rule of thumb, the poles of the observer should be selected from 3 to 10 times faster than the poles of the controller. An upper bound on the speed of response of the observer is imposed by the presence of the unavoidable measurement noise. Inappropriately fast observer dynamics will result in tracking the noise rather than the actual state of the system. Hence, a deadbeat observer, although appealing in theory, is avoided in practice.

The choice of observer poles is also governed by the same considerations related to the robustness of the system discussed in Section 9.2.3 for the state feedback control. Thus, the sensitivity of the eigenvalues to perturbations in the system matrices must be considered in the selection of the observer poles.

We emphasize that the selection of the observer poles does not influence the performance of the overall control system if the initial conditions are estimated perfectly. We prove this fact for the full-order observer. However, the result also holds for the reduced-order observer, but the proof is left as an exercise for the reader. To demonstrate this fact, we consider the state Eq. (9.54) with the output Eq. (9.1) modified for the augmented state vector:

$$\begin{aligned} \begin{bmatrix} \mathbf{x}(k+1) \\ \tilde{\mathbf{x}}(k+1) \end{bmatrix} &= \begin{bmatrix} A - BK & BK \\ 0 & A - LC \end{bmatrix} \begin{bmatrix} \mathbf{x}(k) \\ \tilde{\mathbf{x}}(k) \end{bmatrix} + \begin{bmatrix} B \\ 0 \end{bmatrix} \mathbf{v}(k) \\ \mathbf{y}(k) &= [C \quad 0] \begin{bmatrix} \mathbf{x}(k) \\ \tilde{\mathbf{x}}(k) \end{bmatrix} \end{aligned} \quad (9.66)$$

The zero input–output response of the system ($\mathbf{v}(k) = 0$) can be determined iteratively as

$$\begin{aligned}\mathbf{y}(0) &= C\mathbf{x}(0) \\ \mathbf{y}(1) &= [C \quad 0] \begin{bmatrix} A - BK & BK \\ 0 & A - LC \end{bmatrix} \begin{bmatrix} \mathbf{x}(0) \\ \tilde{\mathbf{x}}(0) \end{bmatrix} = C(A - BK)\mathbf{x}(0) + CBK\tilde{\mathbf{x}}(0) \\ \mathbf{y}(2) &= [C \quad 0] \begin{bmatrix} A - BK & BK \\ 0 & A - LC \end{bmatrix} \begin{bmatrix} \mathbf{x}(1) \\ \tilde{\mathbf{x}}(1) \end{bmatrix} = C(A - BK)\mathbf{x}(1) + CBK\tilde{\mathbf{x}}(1) \\ &= C(A - BK)^2\mathbf{x}(0) - C(A - BK)BK\tilde{\mathbf{x}}(0) - CBK(A - LC)\tilde{\mathbf{x}}(0) \\ &\vdots\end{aligned}$$

Clearly, the observer matrix L influences the transient response if and only if $\tilde{\mathbf{x}}(k) \neq 0$. This fact is confirmed by the determination of the z -transfer function from Eq. (9.66), which implicitly assumes zero initial conditions

$$G(z) = [C \quad 0] \left(zI - \begin{bmatrix} A - BK & BK \\ 0 & A - LC \end{bmatrix} \right)^{-1} \begin{bmatrix} B \\ 0 \end{bmatrix} = C(zI - A + BK)^{-1}B$$

where the observer gain matrix L does not appear.

Another important transfer function relates the output $\mathbf{y}(k)$ to the control $\mathbf{u}(k)$. It is obtained from the observer state Eq. (9.34) with Eq. (9.51) as the output equation and is given by

$$G_{co}(z) = -K(zI_n - A + BK + LC)^{-1}L$$

Note that any realization of this z -transfer function can be used to implement the controller–observer regardless of the realization of the system used to obtain the transfer function. The denominator of the transfer function is

$$\det(zI_n - A + BK + LC)$$

which is clearly different from Eq. (9.62). This should come as no surprise, since the transfer function represents the feedback path in Fig. 9.14 and not the closed-loop system dynamics.

Example 9.12

Consider the armature-controlled DC motor of Example 9.4. Let the true initial condition be $\mathbf{x}(0) = [1, 1, 1]$, and let its estimate be the zero vector $\hat{\mathbf{x}}(0) = [0, 0, 0]^T$. Design a full-order observer state feedback for a zero-input response with a settling time of 0.2 s.

Solution

As Example 9.4 showed, a choice for the control system eigenvalues that meets the design specification is $\{0.6, 0.4 \pm j0.33\}$. This yields the gain vector

$$K = 10^3 [1.6985 \quad -0.70088 \quad 0.01008]$$

Example 9.12—cont'd

The observer eigenvalues must be selected so that the associated modes are sufficiently faster than those of the controller. We select the eigenvalues $\{0.1, 0.1 \pm j0.1\}$. This yields the observer gain vector

$$L = 10^2 \begin{bmatrix} 0.02595 \\ 0.21663 \\ 5.35718 \end{bmatrix}$$

Using Eq. (9.66), we obtain the space-space equations

$$\begin{bmatrix} \mathbf{x}(k+1) \\ \tilde{\mathbf{x}}(k+1) \end{bmatrix} = \begin{bmatrix} 0.9972 & 0.0989 & 0 & 0.0028 & 0.0011 & 0 \\ -0.8188 & 0.6616 & 0.0046 & 0.8188 & 0.3379 & 0.0049 \\ -160.813 & -66.454 & -0.0589 & 160.813 & 66.359 & 0.9543 \\ & & & -1.5949 & 0.1 & 0 \\ 0 & & & -21.663 & 0.9995 & 0.0095 \\ & & & -535.79 & -0.0947 & 0.8954 \end{bmatrix} \begin{bmatrix} \mathbf{x}(k) \\ \tilde{\mathbf{x}}(k) \end{bmatrix}$$

$$\mathbf{y}(k) = \begin{bmatrix} 1 & 0 & 0 & | & 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{x}(k) \\ \tilde{\mathbf{x}}(k) \end{bmatrix}$$

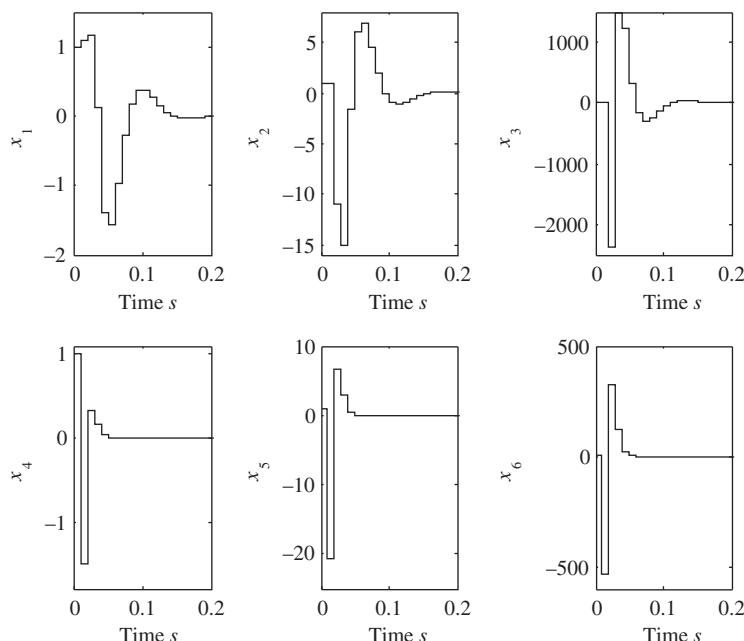


Figure 9.15
Zero-input response for Example 9.12.

Example 9.12—cont'd

The response to the initial condition $[1, 1, 1, 1, 1, 1]$ is plotted in Fig. 9.15. We compare the plant state variables x_i , $i = 1, 2, 3$ to the estimation errors \tilde{x}_i , $i = 4, 5, 6$. We observe that the estimation errors decay to zero faster than the system states and that the system has an overall settling time less than 0.2 s.

Example 9.13

Solve Example 9.12 using a reduced-order observer.

Solution

In this case, we have $l = 1$, and because the measured output corresponds to the first element of the state vector, we do not need similarity transformation, that is, $Q_0 = Q_0^{-1} = I_3$. Thus, we obtain

$$\bar{A}_1 = 1 \quad \bar{A}_2 = 0.1 \quad \bar{A}_3 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \bar{A}_4 = \begin{bmatrix} 0.9995 & 0.0095 \\ -0.0947 & 0.8954 \end{bmatrix}$$

$$B_1 = 1.622 \times 10^{-6} \quad \bar{B}_2 = \begin{bmatrix} 4.821 \times 10^{-4} \\ 9.468 \times 10^{-2} \end{bmatrix}$$

We select the reduced-order observer eigenvalues as $\{0.1 \pm j0.1\}$ and obtain the observer gain vector

$$L = 10^2 \begin{bmatrix} 0.16949 \\ 6.75538 \end{bmatrix}$$

and the associated matrices

$$A_o = \begin{bmatrix} -0.6954 & 0.0095 \\ -67.6485 & 0.8954 \end{bmatrix} \quad A_y = 10^3 \begin{bmatrix} -0.02232 \\ -1.21724 \end{bmatrix} \quad B_o = 10^{-2} \begin{bmatrix} 0.045461 \\ 9.358428 \end{bmatrix}$$

Partitioning $Q_0 = T_o = I_3$ gives

$$Q_z = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \quad T_{ox} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \quad T_{oy} = 10^2 \begin{bmatrix} 0.01 \\ 0.16949 \\ 6.75538 \end{bmatrix}$$

We have that the state-space equation of Eq. (9.61) is

$$\begin{bmatrix} x(k+1) \\ \tilde{z}(k+1) \end{bmatrix} = \begin{bmatrix} 0.99725 & 0.09886 & 0.00002 & 0.00114 & 0.00002 \\ -0.81884 & 0.66161 & 0.00464 & 0.33789 & 0.00486 \\ -160.813 & -66.4540 & -0.05885 & 66.3593 & 0.95435 \\ & & 0 & -0.6954 & 0.0095 \\ & & & -67.6485 & 0.8954 \end{bmatrix} \begin{bmatrix} x(k) \\ \tilde{z}(k) \end{bmatrix} \quad (9.67)$$

whereas the state-space equation of Eq. (9.65) is

Example 9.13—cont'd

$$\begin{bmatrix} \mathbf{x}(k+1) \\ \bar{\mathbf{x}}(k+1) \end{bmatrix} = \begin{bmatrix} 0.96693 & 0.1 & 0 & -0.00114 & -0.00002 \\ -9.82821 & 0.9995 & 0.0095 & -0.33789 & -0.00486 \\ -1930.17 & -0.0947 & 0.8954 & -66.3593 & -0.95425 \\ -31.5855 & 0 & 0 & -1.01403 & 0.00492 \\ -3125.07 & 0 & 0 & -133.240 & 0.04781 \end{bmatrix} \begin{bmatrix} \mathbf{x}(k) \\ \bar{\mathbf{x}}(k) \end{bmatrix}$$

The response of the state-space system Eq. (9.67) to the initial condition [1, 1, 1, 1, 1] is plotted in Fig. 9.16. As in Example 9.12, we observe that the estimation errors x_i , $i = 4$, and 5 decay to zero faster than the system states x_i , $i = 1$, 2, and 3, and that the system has an overall settling time less than 0.2 s.

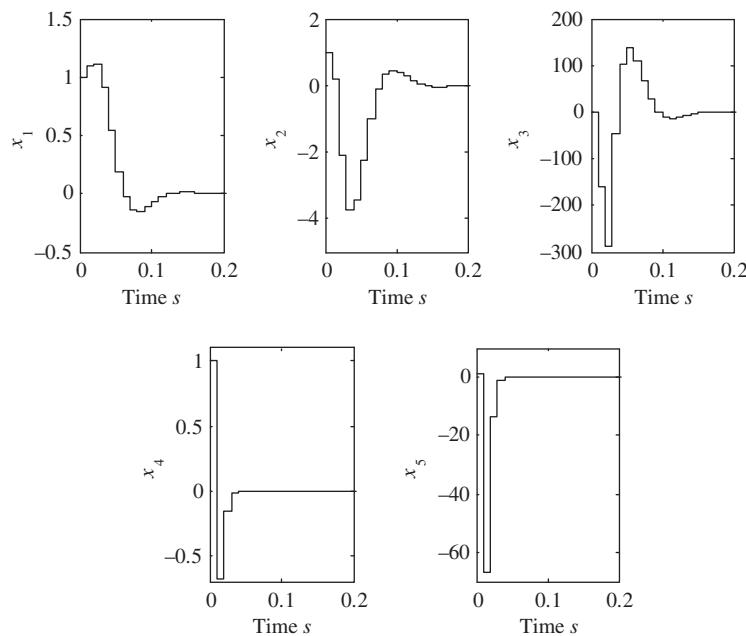


Figure 9.16
Zero-input response for Example 9.13.

Example 9.14

Design a full-order observer state feedback for the armature controller DC motor of Example 9.4 with feedforward action, a settling time of 0.2 s, and overshoot less than 10%.

Example 9.14—cont'd**Solution**

As in Example 9.12, we select the control system eigenvalues as $\{0.4, 0.6 \pm j0.33\}$ and the observer eigenvalues as $\{0.1, 0.1 \pm j0.1\}$. This yields the gain vectors

$$K = 10^3 [1.6985 \quad -0.70088 \quad 0.01008], \quad L = 10^2 \begin{bmatrix} 0.02595 \\ 0.21663 \\ 5.35718 \end{bmatrix}$$

From the matrix $A_{cl} = A - BK$, we determine the feedforward term using Eq. (9.28) as

$$F = [C(I_n - (A - BK))^{-1}B]^{-1} = 1698.49$$

Substituting $Fr(k)$ for $v(k)$ in Eq. (9.66), we obtain the closed-loop state-space equations

$$\begin{bmatrix} \mathbf{x}(k+1) \\ \tilde{\mathbf{x}}(k+1) \end{bmatrix} = \begin{bmatrix} 0.9972 & 0.0989 & 0 & | & 0.0028 & 0.0011 & 0 \\ -0.8188 & 0.6616 & 0.0046 & | & 0.8188 & 0.3379 & 0.0049 \\ -160.813 & -66.454 & -0.0589 & | & 160.813 & 66.359 & 0.9543 \\ 0 & | & -1.5949 & | & 0.1 & 0 \\ 0 & | & -21.663 & | & 0.9995 & 0.0095 \\ 0 & | & -535.79 & | & -0.0947 & 0.8954 \end{bmatrix} \begin{bmatrix} \mathbf{x}(k) \\ \tilde{\mathbf{x}}(k) \end{bmatrix} + \begin{bmatrix} 0.0028 \\ 0.8188 \\ 160.813 \\ 0 \end{bmatrix} r(k)$$

$$\mathbf{y}(k) = \begin{bmatrix} 1 & 0 & 0 & | & 0 \end{bmatrix} \begin{bmatrix} \mathbf{x}(k) \\ \tilde{\mathbf{x}}(k) \end{bmatrix}$$

The step response of the system shown in Fig. 9.17 has a settling time of about 0.1 and a percentage overshoot of about 6%. The controller meets the design specification.

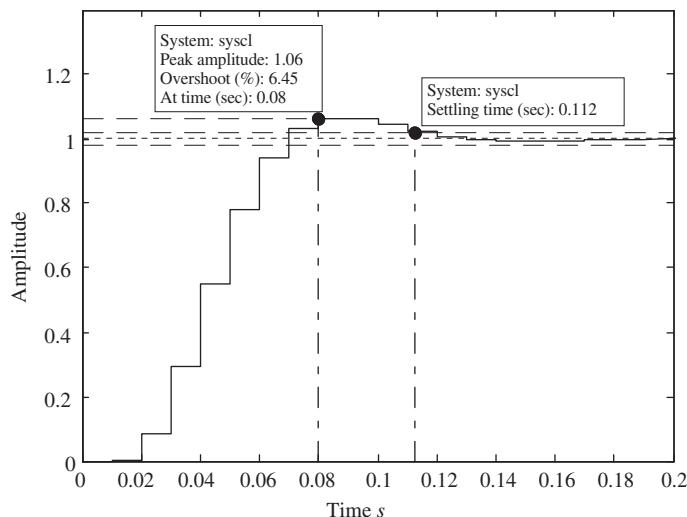


Figure 9.17
Step response for Example 9.14.

9.7 Pole assignment using transfer functions

The pole assignment problem can be solved in the framework of transfer functions. Consider the state-space equations of the two-degree-of-freedom controller shown in Fig. 9.7 with the state vector estimated using a full-order observer. For a SISO plant with observer state feedback, we have

$$\begin{aligned}\hat{\mathbf{x}}(k+1) &= A\hat{\mathbf{x}}(k) + Bu(k) + L(y(k) - C\hat{\mathbf{x}}(k)) \\ u(k) &= -K\hat{\mathbf{x}}(k) + Fr(k)\end{aligned}$$

or equivalently,

$$\begin{aligned}\hat{\mathbf{x}}(k+1) &= (A - BK - LC)\hat{\mathbf{x}}(k) + [BF \quad L] \begin{bmatrix} r(k) \\ y(k) \end{bmatrix} \\ u(k) &= -K\hat{\mathbf{x}}(k) + Fr(k)\end{aligned}$$

The corresponding z -transfer function from $[r, y]$ to $\hat{\mathbf{x}}$ is

$$\hat{\mathbf{X}}(z) = (zI - A + BK + LC)^{-1} [BF \quad L] \begin{bmatrix} R(z) \\ Y(z) \end{bmatrix}$$

The transfer function from $[r, y]$ to u is

$$U(z) = \left[-K(zI - A + BK + LC)^{-1} BF + F \right] R(z) - \left[K(zI - A + BK + LC)^{-1} L \right] Y(z)$$

Thus, the full-order observer state feedback is equivalent to the transfer function model depicted in Fig. 9.18 with

$$U(z) = G_f(z)R(z) - H(z)Y(z)$$

where

$$H(z) = \frac{S(z)}{D(z)}$$

is the feedback gain with

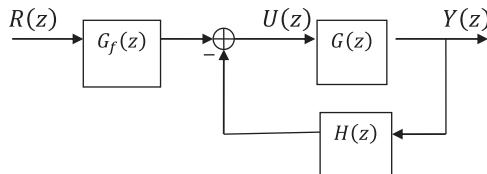


Figure 9.18

Block diagram for pole assignment with transfer functions.

$$H(z) = K(zI - A + BK + LC)^{-1}L$$

and

$$G_f(z) = \frac{N(z)}{D(z)} = -K(zI - A + BK + LC)^{-1}BF + F$$

is the prefilter gain.

The plant transfer function $G(z) = P(z)/Q(z)$ is assumed strictly realizable; that is, the degree of $P(z)$ is less than the degree of $Q(z)$. We also assume that $P(z)$ and $Q(z)$ are **coprime** (i.e., they have no common factors). Further, it is assumed that $Q(z)$ is **monic**, that is, the coefficient of the term with the highest power in z is one.

From the block diagram of Fig. 9.18, simple block diagram manipulations give the closed-loop transfer function

$$\frac{Y(z)}{R(z)} = \frac{G(z)}{1 + G(z)H(z)}G_f(z) = \frac{P(z)N(z)}{Q(z)D(z) + P(z)S(z)}$$

and the polynomial equation

$$(Q(z)D(z) + P(z)S(z))Y(z) = P(z)N(z)R(z) \quad (9.68)$$

Therefore, the closed-loop characteristic polynomial is

$$\Delta_{cl}(z) = Q(z)D(z) + P(z)S(z) \quad (9.69)$$

The pole placement problem thus reduces to finding polynomials $D(z)$ and $S(z)$ that satisfy Eq. (9.69) for given $P(z)$, $Q(z)$, and for a given desired characteristic polynomial $\Delta_{cl}(z)$.

Eq. (9.69) is called a **Diophantine equation**, and its solution can be found by first expanding its RHS terms as

$$\begin{aligned} P(z) &= p_{n-1}z^{n-1} + p_{n-2}z^{n-2} + \dots + p_1z + p_0 \\ Q(z) &= z^n + q_{n-1}z^{n-1} + \dots + q_1z + q_0 \\ D(z) &= d_mz^m + d_{m-1}z^{m-1} + \dots + d_1z + d_0 \\ S(z) &= s_mz^m + s_{m-1}z^{m-1} + \dots + s_1z + s_0 \end{aligned}$$

The closed-loop characteristic polynomial $\Delta_{cl}(z)$, which can be obtained from the desired pole locations, is of degree $n + m$ and has the form

$$\Delta_{cl}(z) = z^{n+m} + \delta_{n+m-1}z^{n+m-1} + \dots + \delta_1z + \delta_0$$

Thus, Eq. (9.69) can be rewritten as

$$\begin{aligned}
z^{n+m} + \delta_{n+m-1} z^{n+m-1} + \cdots + \delta_1 z + \delta_0 &= (z^n - q_{n-1} z^{n-1} + \cdots + q_1 z + q_0) \\
(d_m z^m + d_{m-1} z^{m-1} + \cdots + d_1 z + d_0) + (p_{n-1} z^{n-1} + p_{n-2} z^{n-2} + \cdots + p_1 z + p_0) \\
(s_m z^m + s_{m-1} z^{m-1} + \cdots + s_1 z + s_0)
\end{aligned} \tag{9.70}$$

[Eq. \(9.70\)](#) is linear in the $2m$ unknowns d_i and s_i , $i = 0, 1, 2, \dots, m-1$, and its LHS is a known polynomial with $n + m - 1$ coefficients. The solution of the Diophantine equation is unique if $n + m - 1 = 2m$ —that is, if $m = n - 1$. [Eq. \(9.70\)](#) can be written in the matrix form

$$\left[\begin{array}{ccccccccc|c|c} 1 & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 & d_m & 1 \\ q_{n-1} & 1 & 0 & \dots & 0 & p_{n-1} & 0 & \dots & 0 & d_{m-1} & \delta_{sn-2} \\ q_{n-2} & q_{n-1} & 1 & \dots & 0 & p_{n-2} & p_{n-1} & \dots & 0 & d_{m-2} & \delta_{2n-3} \\ \vdots & \vdots & \vdots & & \vdots & \vdots & \vdots & & \vdots & \vdots & \vdots \\ q_0 & q_1 & q_2 & \dots & q_{n-1} & p_0 & p_1 & \dots & p_{n-1} & d_0 & \delta_{n-1} \\ 0 & q_0 & q_1 & \dots & q_{n-2} & 0 & p_0 & \dots & p_{n-2} & s_m & \delta_{n-2} \\ 0 & 0 & q_0 & \dots & q_{n-3} & 0 & 0 & \dots & p_{n-3} & s_{m-1} & \delta_{n-3} \\ \vdots & \vdots & \vdots & & \vdots & \vdots & \vdots & & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & q_0 & 0 & 0 & \dots & p_0 & s_0 & \delta_0 \end{array} \right] = \left[\begin{array}{c} \delta_{sn-2} \\ \delta_{2n-3} \\ \vdots \\ \delta_{n-1} \\ \delta_{n-2} \\ \delta_{n-3} \\ \vdots \\ \delta_0 \end{array} \right] \tag{9.71}$$

It can be shown that the matrix on the LHS is nonsingular if and only if the polynomials $P(z)$ and $Q(z)$ are coprime, which we assume. As discussed in [Section 9.2.3](#), the matrix must have a small condition number for the system to be robust with respect to errors in the known parameters. The condition number becomes larger as the matrix becomes almost singular.

The structure of the matrix shows that it will be almost singular if the coefficients of the numerator polynomial $P(z)$ and denominator polynomial $Q(z)$ are almost identical. We therefore require that the roots of the polynomials $P(z)$ and $Q(z)$ be sufficiently different to avoid an ill-conditioned matrix. From the discussion of pole-zero matching of [Section 6.3.2](#), it can be deduced that the poles of the discretized plant approach the zeros as the sampling interval is reduced (see also [Section 12.2.2](#)). Thus, when the controller is designed by pole assignment, the sampling interval must not be excessively short to avoid an ill-conditioned matrix in [Eq. \(9.71\)](#).

We now discuss the choice of the desired characteristic polynomial. From the equivalence of the transfer function design to the state–space design described in [Section 9.6](#), the separation principle implies that $\Delta_{cl}^d(z)$ can be written as the product

$$\Delta_{cl}^d(z) = \Delta_c^d(z) \Delta_o^d(z)$$

where $\Delta_c^d(z)$ is the controller characteristic polynomial and $\Delta_o^d(z)$ is the observer characteristic polynomial. We select the polynomial $N(z)$ as

$$N(z) = k_{ff} \Delta_o^d(z) \quad (9.72)$$

so that the observer polynomial $\Delta_o^d(z)$ cancels in the transfer function from the reference input to the system output. The scalar constant k_{ff} is selected so that the steady-state output is equal to the constant reference input

$$\frac{Y(1)}{R(1)} = \frac{P(1)N(1)}{\Delta_c^d(1)\Delta_o^d(1)} = \frac{P(1)k_{ff}\Delta_o^d(1)}{\Delta_c^d(1)\Delta_o^d(1)} = 1$$

The condition for zero steady-state error is

$$k_{ff} = \frac{\Delta_c^d(1)}{P(1)} \quad (9.73)$$

Example 9.15

Solve Example 9.14 using the transfer function approach.

Solution

The plant transfer function is

$$G(z) = \frac{P(z)}{Q(z)} = 10^{-6} \frac{1.622z^2 + 45.14z + 48.23}{z^3 - 2.895z^2 + 2.791z - 0.8959}$$

Thus, we have the polynomials

$$P(z) = 1.622 \times 10^{-6} z^2 + 45.14 \times 10^{-6} z + 48.23 \times 10^{-6}$$

with coefficients

$$p_2 = 1.622 \times 10^{-6}, \quad p_1 = 45.14 \times 10^{-6}, \quad p_0 = 48.23 \times 10^{-6}$$

$$Q(z) = z^3 - 2.895z^2 + 2.791z - 0.8959$$

with coefficients

$$q_2 = -2.895, \quad q_1 = 2.791, \quad q_0 = -0.8959$$

The plant is third order—that is, $n = 3$ —and the solvability condition of the Diophantine equation is $m = n - 1 = 2$. The order of the desired closed-loop characteristic polynomial is $m + n = 5$. We can therefore select the controller poles as $\{0.6, 0.4 \pm j0.33\}$ and the observer poles as $\{0.1, 0.2\}$ with the corresponding polynomials

$$\Delta_c^d(z) = z^3 - 1.6z^2 + 0.9489z - 0.18756$$

Example 9.15—cont'd

$$\Delta_o^d(z) = z^2 - 0.3z + 0.02$$

$$\Delta_{cl}^d(z) = \Delta_c^d(z)\Delta_o^d(z) = z^5 - 1.9z^4 + 1.4489z^3 - 0.50423z^2 + 0.075246z - 0.0037512$$

In other words, $\delta_4 = -1.9$, $\delta_3 = 1.4489$, $\delta_2 = -0.50423$, $\delta_1 = 0.075246$, and $\delta_0 = -0.0037512$.

Using the matrix Eq. (9.71) gives

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ -2.8949 & 1 & 0 & 1.622 \times 10^{-6} & 0 & 0 \\ 2.790752 & -2.8949 & 1 & 45.14 \times 10^{-6} & 1.622 \times 10^{-6} & 0 \\ -0.895852 & 2.790752 & -2.8949 & 48.23 \times 10^{-6} & 45.14 \times 10^{-6} & 1.622 \times 10^{-6} \\ 0 & -0.895852 & 2.790752 & 0 & 48.23 \times 10^{-6} & 45.14 \times 10^{-6} \\ 0 & 0 & -0.895852 & 0 & 0 & 48.23 \times 10^{-6} \end{bmatrix} \begin{bmatrix} d_2 \\ d_1 \\ d_0 \\ s_2 \\ s_1 \\ s_0 \end{bmatrix} = \begin{bmatrix} 1 \\ -1.9 \\ 1.4489 \\ -0.50423 \\ 0.075246 \\ -0.00375 \end{bmatrix}$$

The MATLAB command **linsolve** gives the solution

$$\begin{aligned} d_2 &= 1, \quad d_1 = 0.9645, \quad d_0 = 0.6526, \quad s_2 = 1.8735 \cdot 10^4, \quad s_1 = -2.9556 \cdot 10^4, \quad s_0 \\ &= 1.2044 \cdot 10^4 \end{aligned}$$

and the polynomials

$$D(z) = z^2 + 0.9645z + 0.6526$$

$$S(z) = 1.8735 \cdot 10^4 z^2 - 2.9556 \cdot 10^4 z + 1.2044 \cdot 10^4$$

Then, from Eq. (9.72), we have $k_{ff} = 1698.489$ and the numerator polynomial

$$N(z) = 1698.489(z^2 - 0.3z + 0.02)$$

The step response of the control system of Fig. 9.19 has a settling time of 0.1 s and a percentage overshoot of less than 7%. The response meets all the design specifications.

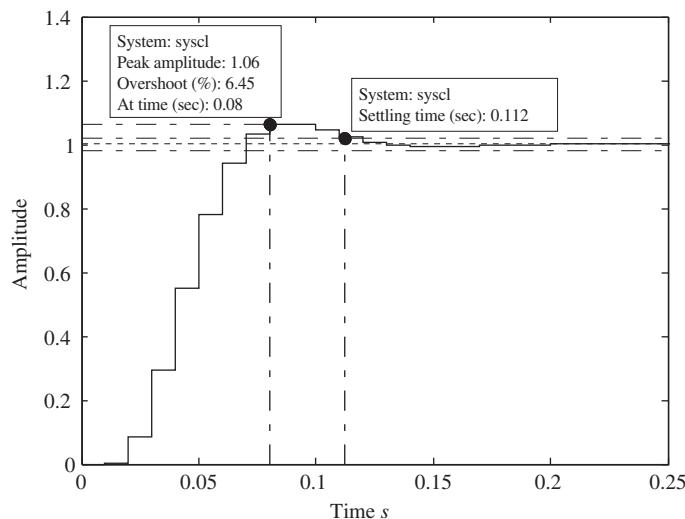
Example 9.15—cont'd

Figure 9.19
Step response for Example 9.15.

Further reading

- Bass, R.W., Gura, I., October 1983. High-order System Design via State-Space considerations Proceedings of JACC. Troy, New York.
- Chen, C.T., 1984. Linear System Theory and Design. HRW, New York.
- Delchamps, D.F., 1988. State-Space and Input-Output Linear Systems. Springer-Verlag.
- D'Azzo, J.J., Houpis, C.H., 1988. Linear Control System Analysis and Design. McGraw-Hill, New York.
- Kailath, T., 1980. Linear Systems. Prentice Hall, Englewood Cliffs, NJ.
- Kautsky, J.N., Nichols, K., Van Dooren, P., 1985. Robust pole assignment in linear state feedback. Int. J. Control 41 (5), 1129–1155.
- Mayne, D.Q., Murdoch, P., 1970. Modal control of linear time invariant systems. Int. J. Control 11 (2), 223–227.
- Patel, R.V., Munro, N., 1982. Multivariable System Theory and Design. Pergamon Press, Oxford, England.

Problems

- 9.1 Show that the closed-loop quadruple for (A, B, C, D) with the state feedback $\mathbf{u}(k) = -K\mathbf{x}(k) + \mathbf{v}(k)$ is $(A - BK, B, C - DK, D)$.
- 9.2 Show that for the pair (A, B) with state feedback gain matrix K to have the closed-loop state matrix $A_{cl} = A - BK$, a necessary condition is that for any vector \mathbf{w}^T satisfying $\mathbf{w}^T B = \mathbf{0}^T$ and $\mathbf{w}^T A = \lambda \mathbf{w}^T$, A_{cl} must satisfy $\mathbf{w}^T A_{cl} = \lambda \mathbf{w}^T$. Explain the significance of this necessary condition. (Note that the condition is also sufficient.)
- 9.3 Show that for the pair (A, B) with $m \times n$ state feedback gain matrix K to have the closed-loop state matrix $A_{cl} = A - BK$, a sufficient condition is

$$\text{rank}\{B\} = \text{rank}\{[A - A_{cl}|B]\} = m$$

Is the matrix K unique for given matrices A and A_{cl} ? Explain.

- 9.4 Using the results of Problem 9.3, determine if the closed-loop matrix can be obtained using state feedback for the pair

$$A = \begin{bmatrix} 1.0 & 0.1 & 0 \\ 0 & 1 & 0.01 \\ 0 & -0.1 & 0.9 \end{bmatrix} \quad B = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

a. $A_{cl} = \begin{bmatrix} 1.0 & 0.1 & 0 \\ -12.2 & -1.2 & 0 \\ 0.01 & 0.01 & 0 \end{bmatrix}$

b. $A_{cl} = \begin{bmatrix} 0 & 0.1 & 0 \\ -12.2 & -1.2 & 0 \\ 0.01 & 0.01 & 0 \end{bmatrix}$

- 9.5 Show that for the pair (A, C) with observer gain matrix L to have the observer matrix $A_o = A - LC$, a necessary condition is that for any vector \mathbf{v} satisfying $C\mathbf{v} = \mathbf{0}$ and $A\mathbf{v} = \lambda\mathbf{v}$, $A_o\mathbf{v} = \lambda\mathbf{v}$. Explain the significance of this necessary condition. (Note that the condition is also sufficient.)
- 9.6 Show that for the pair (A, C) with $n \times l$ observer gain matrix L to have the observer matrix $A_o = A - LC$, a sufficient condition is

$$\text{rank}\{C\} = \text{rank}\left\{\begin{bmatrix} C \\ \hline A - A_o \end{bmatrix}\right\} = l$$

Is the matrix L unique for given matrices A and A_o ? Explain.

- 9.7 Design a state feedback control law to assign the eigenvalues to the set $\{0, 0.1, 0.2\}$ for the systems with

$$\text{a. } A = \begin{bmatrix} 0.1 & 0.5 & 0 \\ 2 & 0 & 0.2 \\ 0.2 & 1 & 0.4 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0.01 \\ 0 \\ 0.005 \end{bmatrix}$$

$$\text{b. } A = \begin{bmatrix} -0.2 & -0.2 & 0.4 \\ 0.5 & 0 & 1 \\ 0 & -0.4 & -0.4 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0.01 \\ 0 \\ 0 \end{bmatrix}$$

- 9.8 Show that the eigenvalues $\{\lambda_i, i = 1, \dots, n\}$ of the closed-loop system for the pair (A, B) with state feedback are scaled by α if we use the state feedback gains from the pole placement with the pair $(A/\alpha, B/\alpha)$. Verify this result using MATLAB for the pair

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -0.005 & -0.11 & -0.7 \end{bmatrix} \quad B = \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$$

and the eigenvalues $\{0.1, 0.2, 0.3\}$ with $\alpha = 0.5$.

- 9.9 Using eigenvalues that are two to three times as fast as those of the plant, design a state estimator for the system

$$\text{a. } A = \begin{bmatrix} 0.2 & 0.3 & 0.2 \\ 0 & 0 & 0.3 \\ 0.3 & 0 & 0.3 \end{bmatrix} \quad C = [1 \ 1 \ 0]$$

$$\text{b. } A = \begin{bmatrix} 0.2 & 0.3 & 0.2 \\ 0 & 0 & 0.3 \\ 0.3 & 0 & 0.3 \end{bmatrix} \quad C = [1 \ 0 \ 0]$$

- 9.10 Consider the system

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -0.005 & -0.11 & -0.7 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$C = [0.5 \ 1 \ 0] \quad d = 0$$

- a. Design a controller that assigns the eigenvalues $\{-0.8, -0.3 \pm j0.3\}$. Why is the controller guaranteed to exist?

- b. Why can we design an observer for the system with the eigenvalues $\{-0.5, -0.1 \pm j0.1\}$? Explain why the value (-0.5) must be assigned. (*Hint:* $(s+0.1)^2(s+0.5) = s^3 + 0.7 s^2 + 0.11 s + 0.005$.)
- c. Obtain a similar system with a second-order observable subsystem, for which an observer can be easily designed, as in [Section 8.3.3](#). Design an observer for the transformed system with two eigenvalues shifted as in (b), and check your design using the MATLAB command **place** or **acker**. Use the result to obtain the observer for the original system. (*Hint:* Obtain an observer gain \mathbf{I}_r for the similar third-order system from your design by setting the first element equal to zero. Then obtain the observer gain for the original system using $\mathbf{I} = T_r \mathbf{I}_r$, where T_r is the similarity transformation matrix.)
- d. Design an observer-based feedback controller for the system with the controller and observer eigenvalues selected as in (a) and (b), respectively.

9.11 Design a reduced-order estimator state feedback controller for the discretized system

$$A = \begin{bmatrix} 0.1 & 0 & 0.1 \\ 0 & 0.5 & 0.2 \\ 0.2 & 0 & 0.4 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0.01 \\ 0 \\ 0.005 \end{bmatrix} \quad \mathbf{c}^T = [1, 1, 0]$$

to obtain the eigenvalues $\{0.1, 0.4 \pm j0.4\}$.

9.12 Consider the following model of an armature-controlled DC motor, which is slightly different from that described in [Example 7.7](#):

$$\begin{aligned} \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \end{bmatrix} &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & -11 & -11.1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 10 \end{bmatrix} u \\ y &= [1 \ 0 \ 0] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \end{aligned}$$

For digital control with $T = 0.02$, apply the state feedback controllers determined in [Example 9.4](#) in order to verify their robustness.

9.13 Consider the following model of a DC motor speed control system, which is slightly different from that described in [Example 6.9](#):

$$G(s) = \frac{1}{(1.2s + 1)(s + 10)}$$

For a sampling period $T = 0.02$, obtain a state-space representation corresponding to the discrete-time system with DAC and ADC, then use it to verify the robustness of the state controller described in [Example 9.6](#).

- 9.14 Verify the robustness of the state controller determined in Example 9.7 by applying it to the model shown in Problem 9.13.
- 9.15 Prove that the pair (\tilde{A}, \tilde{B}) of Eq. (9.31) for the system with integral control is controllable if and only if the pair (A, B) is controllable.
- 9.16 Consider the DC motor position control system described in Example 3.6, where the (type 1) analog plant has the transfer function

$$G(s) = \frac{1}{s(s+1)(s+10)}$$

For the digital control system with $T = 0.02$, design a state feedback controller to obtain a step response with zero steady-state error, zero overshoot, and a settling time of less than 0.5 s.

- 9.17 Design a digital state feedback controller for the analog system

$$G(s) = \frac{-s+1}{(5s+1)(10s+1)}$$

with $T = 0.1$ to place the closed-loop poles at $\{0.4, 0.6\}$. Show that the zero of the closed-loop system is the same as the zero of the open-loop system.

- 9.18 Design a state feedback law for the pair

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -0.005 & -0.11 & -0.7 & 2 \end{bmatrix}, B = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$

to assign the eigenvalues $\{0, 0.25, 0.3, 0.4\}$.

- 9.19 Write the closed-loop system state-space equations of a full-order observer state feedback system with integral action.
- 9.20 Consider the continuous-time model of the overhead crane proposed in Problem 7.11 with $m_c = 1000$ kg, $m_l = 1500$ kg, and $l = 8$ m. Design a discrete full-order observer state feedback control to provide motion of the load without sway.
- 9.21 We present a model for the female population of a species with a maximum age of 3 years based on an example from¹ The population is divided into three groups according to age; namely the age groups $(0, 1)$, $[1, 2]$, $[2, 3]$. These three populations are state variables of the model $\{x_{1k}, x_{2k}, x_{3k}\}$, respectively. The second and third populations can produce offspring but the second population is more fertile.

¹ Mazen Shahin, *Explorations of Mathematical Models in Biology with MATLAB*, J. Wiley, Hoboken, NJ, 2014.

$$\mathbf{x}_{k+1} = \begin{bmatrix} 0 & 6 & 10/3 \\ 0.6 & 0 & 0 \\ 0 & 0.4 & 0 \end{bmatrix} \mathbf{x}_k$$

Because of the difficulty of estimating the age of members of the species, the only measurement available is the total population and the measurement equation is

$$y_k = [1 \quad 1 \quad 1] \mathbf{x}_k$$

- (a) Check the observability of the system

Design an observer for the system with observer eigenvalues $\{0, -0.1 \pm j0.1\}$

- 9.22 Consider the continuous-time model of the overhead crane proposed in Problem 7.11 with $m_c = 1000$ kg, $m_l = 1500$ kg, and $l = 8$ m. Design a control system based on pole assignment using transfer functions in order to provide motion of the load without sway.

Computer exercises

- 9.23 Write a MATLAB script to evaluate the feedback gains using Ackermann's **formula for any pair** (A, B) and any desired poles $\{\lambda_1, \dots, \lambda_n\}$.
- 9.24 Write a MATLAB program that, for a multi-input system with known state and input matrices, determines (i) a set of feasible eigenvectors and (ii) a state feedback matrix that assign a given set of eigenvalues. Test your program using Example 9.5.
- 9.25 Write a MATLAB function that, given the system state-space matrices A , B , and C , the desired closed-loop poles, and the observer poles, determines the closed-loop system state-space matrices of a full-observer state feedback system with integral action.
- 9.26 Write a MATLAB function that uses the transfer function approach to determine the closed-loop system transfer function for a given plant transfer function $G(z)$, desired closed-loop system poles, and observer poles.

Optimal control

Objectives

After completing this chapter, the reader will be able to do the following:

1. Find the unconstrained optimum values of a function by minimization or maximization.
2. Find the constrained optimum values of a function by minimization or maximization.
3. Design an optimal digital control system.
4. Design a digital linear quadratic regulator.
5. Design a digital steady-state regulator.
6. Design a digital output regulator.
7. Design a regulator to track a nonzero constant input.
8. Use the Hamiltonian system to obtain a solution of the Riccati equation.
9. Characterize the eigenstructure of the system with optimal control in terms of the eigenstructure of the Hamiltonian system.
10. Obtain the return difference equality and characterize the robustness of the linear quadratic regulator.
11. Design model predictive control for a linear time-invariant system.

In this chapter, we introduce optimal control theory for discrete-time systems. We begin with unconstrained optimization of a cost function and then generalize to optimization with equality constraints. We then cover the optimization or optimal control of discrete-time systems. We specialize to the linear quadratic regulator and obtain the optimality conditions for a finite and for an infinite planning horizon. We also address the regulator problem where the system is required to track a nonzero constant signal. We examine the Hamiltonian system, its use in solving the Riccati equation, and its relation to the eigenstructure of the system with optimal control. We also derive the return difference equality and obtain expression for the stability margins of the linear quadratic regulator. Finally, we provide a brief introduction to a control design strategy known as model predictive control (MPC).

Chapter Outline**10.1 Optimization 442**

- 10.1.1 Unconstrained optimization 442
- 10.1.2 Constrained optimization 445

10.2 Optimal control 447**10.3 The linear quadratic regulator 453**

- 10.3.1 Free final state 455

10.4 Steady-state quadratic regulator 466

- 10.4.1 Output quadratic regulator 467
- 10.4.2 MATLAB solution of the steady-state regulator problem 468
- 10.4.3 Linear quadratic tracking controller 470

10.5 Hamiltonian system 473

- 10.5.1 Eigenstructure of the Hamiltonian matrix 477

10.6 Return difference equality and stability margins 481**10.7 Model predictive control 488**

- 10.7.1 Model 489
- 10.7.2 Cost function 489
- 10.7.3 Computation of the control law 490
- 10.7.4 Constraints 490
- 10.7.5 MATLAB commands 491

10.8 Modification of the reference signal 491

- 10.8.1 Dynamic Matrix Control 492

Further reading 498**Problems 498****Computer exercises 502**

10.1 Optimization

Many problems in engineering can be solved by minimizing a measure of **cost** or maximizing a measure of **performance**. The designer must select a suitable performance measure based on his or her understanding of the problem to include the most important performance criteria and reflect their relative importance. The designer must also select a mathematical form of the function that makes solving the optimization problem tractable.

We observe that any maximization problem can be recast as a minimization, or vice versa. This is because the location of the maximum of a function $f(\mathbf{x})$ is the same as that of the minimum $-f(\mathbf{x})$, as demonstrated in Fig. 10.1 for a scalar x . We therefore consider minimization only throughout this chapter. We first consider the problem of minimizing a cost function or performance measure, then we extend our solution to problems with equality constraints.

10.1.1 Unconstrained optimization

Consider the problem of minimizing a **cost function** or **performance measure** of the form

$$J(\mathbf{x}) \quad (10.1)$$

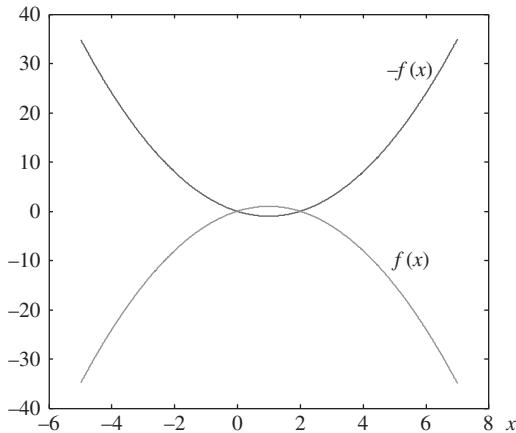


Figure 10.1
Minimization and maximization of a function of a single variable.

where \mathbf{x} is an $n \times 1$ vector of parameters to be selected. Let the optimal parameter vector be \mathbf{x}^* and expand the function $J(\mathbf{x})$ in the vicinity of the optimum as

$$J(\mathbf{x}) = J(\mathbf{x}^*) + \left[\frac{\partial J}{\partial \mathbf{x}} \right]_{\mathbf{x}^*}^T \Delta \mathbf{x} + \frac{1}{2!} \Delta \mathbf{x}^T \left[\frac{\partial^2 J}{\partial \mathbf{x}^2} \right]_{\mathbf{x}^*} \Delta \mathbf{x} + O(\|\Delta \mathbf{x}\|^3) \quad (10.2)$$

$$\Delta \mathbf{x} = \mathbf{x} - \mathbf{x}^*$$

where

$$\left[\frac{\partial J}{\partial \mathbf{x}} \right]_{\mathbf{x}^*}^T = \left[\frac{\partial J}{\partial x_1} \frac{\partial J}{\partial x_2} \cdots \frac{\partial J}{\partial x_n} \right]_{\mathbf{x}^*} \quad (10.3)$$

$$\left[\frac{\partial^2 J}{\partial \mathbf{x}^2} \right]_{\mathbf{x}^*} = \left[\frac{\partial^2 J}{\partial x_i \partial x_j} \right]_{\mathbf{x}^*} \quad (10.4)$$

and the subscript denotes that the matrix is evaluated at \mathbf{x}^* . At a **local minimum** or **local maximum**, the first-order terms of the expansion that appear in the gradient vector are zero.

To guarantee that the point \mathbf{x}^* is a minimum, any perturbation vector $\Delta \mathbf{x}$ away from \mathbf{x}^* must result in an increase of the value of $J(\mathbf{x})$. Thus, the second-order term of the expansion must at least be positive or zero for \mathbf{x}^* to have any chance of being a minimum. If the second-order term is positive, then we can guarantee that \mathbf{x}^* is indeed a minimum. The sign of the second term is determined by the characteristics of the second derivative matrix, or **Hessian**. The second term is positive for any perturbation if the Hessian matrix is positive definite, and it is positive or zero if the Hessian matrix is positive semidefinite. We summarize this discussion in Theorem 10.1.

Theorem 10.1

If \mathbf{x}^* is a local minimum of $J(\mathbf{x})$, then

$$\left[\frac{\partial J}{\partial \mathbf{x}} \right]_{\mathbf{x}^*}^T = \mathbf{0}_{1 \times n} \quad (10.5)$$

$$\left[\frac{\partial^2 J}{\partial \mathbf{x}^2} \right]_{\mathbf{x}^*} \geq 0 \quad (10.6)$$

A sufficient condition for \mathbf{x}^* to be a minimum is

$$\left[\frac{\partial^2 J}{\partial \mathbf{x}^2} \right]_{\mathbf{x}^*} > 0 \quad (10.7)$$

Example 10.1

Obtain the least-squares estimates of the linear resistance

$$v = iR$$

using N noisy measurements

$$z(k) = i(k)R + v(k), \quad k = 1, \dots, N$$

Solution

We begin by stacking the measurements to obtain the matrix equation

$$\mathbf{z} = \mathbf{i}R + \mathbf{v}$$

in terms of the vectors

$$\begin{aligned} \mathbf{z} &= [z(1) \quad \cdots \quad z(N)]^T \\ \mathbf{i} &= [i(1) \quad \cdots \quad i(N)]^T \\ \mathbf{v} &= [v(1) \quad \cdots \quad v(N)]^T \end{aligned}$$

We minimize the sum of the squares of the errors

$$\begin{aligned} J &= \sum_{k=1}^N e^2(k) \\ e(k) &= z(k) - i(k)\hat{R} \end{aligned}$$

where the caret ($\hat{\cdot}$) denotes the estimate. We rewrite the performance measure in terms of the vectors as

Example 10.1—cont'd

$$J = \mathbf{e}^T \mathbf{e} = \mathbf{z}^T \mathbf{z} - 2\mathbf{i}^T \mathbf{z} \hat{R} + \mathbf{i}^T \mathbf{i} \hat{R}^2$$

$$\mathbf{e} = [e(1) \quad \cdots \quad e(N)]^T$$

The necessary condition for a minimum gives

$$\frac{\partial J}{\partial \hat{R}} = -2\mathbf{i}^T \mathbf{z} + 2\mathbf{i}^T \mathbf{i} \hat{R} = 0$$

We now have the least-squares estimate

$$\hat{R}_{LS} = \frac{\mathbf{i}^T \mathbf{z}}{\mathbf{i}^T \mathbf{i}}$$

The solution is indeed a minimum because the second derivative is the positive sum of the squares of the currents

$$\frac{\partial^2 J}{\partial \hat{R}^2} = 2\mathbf{i}^T \mathbf{i} > 0$$

10.1.2 Constrained optimization

In most practical applications, the entries of the parameter vector \mathbf{x} are subject to physical and economic constraints. Assume that our vector of parameters is subject to the equality constraint

$$\mathbf{m}(\mathbf{x}) = \mathbf{0}_{m \times 1} \tag{10.8}$$

In the simplest cases only, we can use the constraints to eliminate m parameters and then solve an optimization problem for the remaining $n-m$ parameters. Alternatively, we can use Lagrange multipliers to include the constraints in the optimization problem. We add the constraint to the performance measure weighted by the **Lagrange multipliers** to obtain the **Lagrangian**

$$L(\mathbf{x}) = J(\mathbf{x}) + \lambda^T \mathbf{m}(\mathbf{x}) \tag{10.9}$$

We then solve for the vectors \mathbf{x} and λ that minimize the Lagrangian as in unconstrained optimization. Example 10.2 demonstrates the use of Lagrange multipliers in constrained optimization. Note that the example is simplified to allow us to solve the problem with and without Lagrange multipliers.

Example 10.2

A manufacturer decides the production level of two products based on maximizing profit subject to constraints on production. The manufacturer estimates profit using the simplified measure

$$J(\mathbf{x}) = x_1^\alpha x_2^\beta$$

where x_i is the quantity produced for product i , $i = 1, 2$, and the parameters (α, β) are determined from sales data. The sum of the quantities of the two products produced cannot exceed a fixed level b . Determine the optimum production level for the two products subject to the production constraint

$$x_1 + x_2 = b$$

Solution

To convert the maximization problem into minimization, we use the negative of the profit. We obtain the optimum first without the use of Lagrange multipliers. We solve for x_2 using the constraint

$$x_2 = b - x_1$$

then substitute in the negative of the profit function to obtain

$$J(x_1) = -x_1^\alpha (b - x_1)^\beta$$

The necessary condition for a minimum gives

$$\frac{\partial J}{\partial x_1} = -\alpha x_1^{\alpha-1} (b - x_1)^\beta + \beta x_1^\alpha (b - x_1)^{\beta-1} = 0$$

which simplifies to

$$x_1 = \frac{\alpha b}{\alpha + \beta}$$

From the production constraint we solve for production level

$$x_2 = \frac{\beta b}{\alpha + \beta}$$

We now show that the same answer can be obtained using a Lagrange multiplier. We add the constraint multiplied by the Lagrange multiplier to obtain the Lagrangian

$$L(\mathbf{x}) = -x_1^\alpha x_2^\beta + \lambda(x_1 + x_2 - b)$$

The necessary conditions for the minimum are

$$\frac{\partial L}{\partial x_1} = -\alpha x_1^{\alpha-1} x_2^\beta + \lambda = 0$$

$$\frac{\partial L}{\partial x_2} = -\beta x_1^\alpha x_2^{\beta-1} + \lambda = 0$$

$$\frac{\partial L}{\partial \lambda} = x_1 + x_2 - b = 0$$

Example 10.2—cont'd

The first two conditions give

$$\lambda = \alpha x_1^{\alpha-1} x_2^\beta = \beta x_1^\alpha x_2^{\beta-1}$$

which readily simplifies to

$$x_2 = \frac{\beta}{\alpha} x_1$$

Substituting in the third necessary condition (i.e., in the constraint) gives the solution

$$x_1 = \frac{\alpha b}{\alpha + \beta}$$

$$x_2 = \frac{\beta b}{\alpha + \beta}$$

10.2 Optimal control

To optimize the performance of a discrete-time dynamic system, we minimize the performance measure

$$J = J_f(\mathbf{x}(k_f), k_f) + \sum_{k=k_0}^{k_f-1} L(\mathbf{x}(k), \mathbf{u}(k), k) \quad (10.10)$$

subject to the constraint

$$\mathbf{x}(k+1) = A\mathbf{x}(k) + B\mathbf{u}(k), \quad k = k_0, \dots, k_f - 1 \quad (10.11)$$

We assume that the pair (A, B) is stabilizable; otherwise, there is no point in control system design, optimal or otherwise. If the system is not stabilizable, then its structure must first be changed by selecting different control variables that allow its stabilization.

The first term of the performance measure is a **terminal penalty**, and each of the remaining terms represents a cost or penalty at time k . We change the problem into unconstrained minimization using Lagrange multipliers. We have the new performance measure

$$\bar{J} = J_f(\mathbf{x}(k_f), k_f) + \sum_{k=k_0}^{k_f-1} \left\{ L(\mathbf{x}(k), \mathbf{u}(k), k) + \lambda^T(k+1)[A\mathbf{x}(k) + B\mathbf{u}(k) - \mathbf{x}(k+1)] \right\} \quad (10.12)$$

We define the Hamiltonian function as

$$H(\mathbf{x}(k), \mathbf{u}(k), k) = L(\mathbf{x}(k), \mathbf{u}(k), k) + \lambda^T(k+1)[A\mathbf{x}(k) + B\mathbf{u}(k)],$$

$$k = k_0, \dots, k_f - 1 \quad (10.13)$$

and expand the performance measure as

$$\begin{aligned}\bar{J} &= J_f(\mathbf{x}(k_f), k_f) + \{L(\mathbf{x}(k_0), \mathbf{u}(k_0), k_0) + \lambda^T(k_0 + 1)[A\mathbf{x}(k_0) + B\mathbf{u}(k_0) - \mathbf{x}(k_0 + 1)]\} \\ &\quad + \{L(\mathbf{x}(k_0 + 1), \mathbf{u}(k_0 + 1), k_0 + 1) + \lambda^T(k_0 + 2)[A\mathbf{x}(k_0 + 1) + B\mathbf{u}(k_0 + 1) - \mathbf{x}(k_0 + 2)]\} \\ &\quad \vdots \\ &\quad \{L(\mathbf{x}(k_f - 1), \mathbf{u}(k_f - 1), k_f - 1) + \lambda^T(k_f)[A\mathbf{x}(k_f - 1) + B\mathbf{u}(k_f - 1) - \mathbf{x}(k_f)]\}\end{aligned}$$

Regrouping terms allows us to rewrite the performance measure in terms of the Hamiltonian as

$$\begin{aligned}\bar{J} &= J_f(\mathbf{x}(k_f), k_f) - \lambda^T(k_f)\mathbf{x}(k_f) + \\ &\quad \sum_{k=k_0+1}^{k_f-1} \{H(\mathbf{x}(k), \mathbf{u}(k), k) - \lambda^T(k)\mathbf{x}(k)\} + H(\mathbf{x}(k_0), \mathbf{u}(k_0), k_0)\end{aligned}\tag{10.14}$$

Each term of the preceding expression can be expanded as a truncated Taylor series in the vicinity of the optimum point of the form

$$\bar{J} = \bar{J}^* + \delta_1 \bar{J} + \delta_2 \bar{J} + \dots\tag{10.15}$$

where the * denotes the optimum, δ denotes a variation from the optimal, and the subscripts denote the order of the variation in the expansion. From basic calculus, we know the necessary condition for a minimum or maximum is that the first-order term must be zero. For a minimum, the second-order term must be positive or zero, and a sufficient condition for a minimum is a positive second term. We therefore need to evaluate the terms of the expansion to determine necessary and sufficient conditions for a minimum.

Each term of the expansion includes derivatives with respect to the arguments of the performance measure and can be obtained by expanding each term separately. For example, the Hamiltonian can be expanded as

$$\begin{aligned}H(\mathbf{x}(k), \mathbf{u}(k), k) &= H(\mathbf{x}^*(k), \mathbf{u}^*(k), k) + \frac{\partial H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{x}(k)} \Big|_*^T \delta \mathbf{x}(k) + \frac{\partial H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{u}(k)} \Big|_*^T \delta \mathbf{u}(k) \\ &\quad + [\mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k)]_*^T \delta \lambda(k + 1) \\ &\quad + \delta \mathbf{x}^T(k) \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{x}^2(k)} \Big|_* \delta \mathbf{x}(k) + \delta \mathbf{u}^T(k) \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{u}^2(k)} \Big|_* \delta \mathbf{u}(k) \\ &\quad + \delta \mathbf{x}^T(k) \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{x}(k) \partial \mathbf{u}(k)} \Big|_* \delta \mathbf{u}(k) + \delta \mathbf{u}^T(k) \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{u}(k) \partial \mathbf{x}(k)} \Big|_* \delta \mathbf{x}(k) \\ &\quad + \text{higher-order terms}\end{aligned}\tag{10.16}$$

where δ denotes a variation from the optimal, the superscript $(*)$ denotes the optimum value, and the subscript $(*)$ denotes evaluation at the optimum point. The second-order terms can be written more compactly in terms of the Hessian matrix of second derivatives in the form

$$[\delta \mathbf{x}^T(k) \quad \delta \mathbf{u}^T(k)] \begin{bmatrix} \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{x}^2(k)} & \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{x}(k) \partial \mathbf{u}(k)} \\ \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{u}(k) \partial \mathbf{x}(k)} & \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{u}^2(k)} \end{bmatrix}_* \begin{bmatrix} \delta \mathbf{x}(k) \\ \delta \mathbf{u}(k) \end{bmatrix} \quad (10.17)$$

We expand the linear terms of the performance measure as

$$\lambda^T(k) \mathbf{x}(k) = \lambda^{*T}(k) \mathbf{x}^*(k) + \delta \lambda^T(k) \mathbf{x}^*(k) + \lambda^{*T}(k) \delta \mathbf{x}(k) \quad (10.18)$$

The terminal penalty can be expanded as

$$\begin{aligned} J_f(\mathbf{x}(k_f), k_f) &= J_f(\mathbf{x}^*(k_f), k_f) + \left. \frac{\partial J_f(\mathbf{x}(k_f), k_f)}{\partial \mathbf{x}(k_f)} \right|_*^T \delta \mathbf{x}(k_f) \\ &\quad + \delta \mathbf{x}^T(k_f) \left[\frac{\partial^2 J_f(\mathbf{x}^*(k_f), k_f)}{\partial \mathbf{x}^2(k_f)} \right]_* \delta \mathbf{x}(k_f) \end{aligned} \quad (10.19)$$

We now combine the first-order terms to obtain the first variation

$$\begin{aligned} \delta_1 \bar{J} &= J_f(\mathbf{x}(k_f), k_f) \lambda^T(k_f) \mathbf{x}(k_f) + \sum_{k=k_0+1}^{k_f-1} \left\{ \left[\left. \frac{\partial H(\mathbf{x}(k), u(k), k)}{\partial \mathbf{x}(k)} \right|_* \right]^T - \lambda^*(k) \right. \\ &\quad \left. \delta \mathbf{x}(k) + [\mathbf{A} \mathbf{x}^*(k) + \mathbf{B} u^*(k) - \mathbf{x}^*(k+1)]^T \delta \lambda(k+1) \right\} \\ &\quad + \left. \frac{\partial H(\mathbf{x}(k), u(k), k)}{\partial u(k)} \right|_*^T \delta u(k) + \left[\frac{\partial J_f(\mathbf{x}(k_f), k_f)}{\partial \mathbf{x}(k_f)} \right]_*^T - \lambda^*(k_f) \delta \mathbf{x}(k_f) \end{aligned}$$

From Eq. (10.15) and the discussion following it, we know that a necessary condition for a minimum of the performance measure is that the first variation must be equal to zero for any combination of variations in its arguments. A zero for any combination of variations occurs only if the coefficient of each increment is equal to zero. Equating each coefficient to zero, we have the necessary conditions

$$\mathbf{x}^*(k+1) = \mathbf{A} \mathbf{x}^*(k) + \mathbf{B}^* \mathbf{u}(k), \quad k = k_0, \dots, k_f - 1 \quad (10.21)$$

$$\lambda^*(k) = \left. \frac{\partial H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{x}(k)} \right|_* \quad k = k_0, \dots, k_f - 1 \quad (10.22)$$

$$\frac{\partial H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{u}(k)} \Big|_* = \mathbf{0} \quad (10.23)$$

$$0 = \frac{\partial J_f(\mathbf{x}(k_f), k_f)}{\partial \mathbf{x}(k_f)} - \lambda(k_f) \Big|_*^T \delta \mathbf{x}(k_f) \quad (10.24)$$

If the terminal point is fixed, then its perturbation is zero and the terminal optimality condition is satisfied. Otherwise, we have the terminal conditions

$$\lambda(k_f) = \frac{\partial J_f(\mathbf{x}(k_f), k_f)}{\partial \mathbf{x}(k_f)} \quad (10.25)$$

Similarly, conditions can be developed for the case of a free initial state with an initial cost added to the performance measure.

The necessary conditions for a minimum of Eqs. (10.21) through (10.23) are known, respectively, as the **state equation**, the **costate equation**, and the **minimum principle of Pontryagin**. The minimum principle tells us that the cost is minimized by choosing the control that minimizes the Hamiltonian. This condition can be generalized to problems where the control is constrained, in which case the solution can be at the boundary of the allowable control region. We summarize the necessary conditions for an optimum in [Table 10.1](#).

The second-order term is

$$\delta_2 \bar{J} = \sum_{k=k_0}^{k_f-1} [\delta \mathbf{x}^T(k) \quad \delta \mathbf{u}^T(k)] \begin{bmatrix} \frac{\partial H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{x}^2(k)} & \frac{\partial H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{x}(k) \partial \mathbf{u}(k)} \\ \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{u}(k) \partial \mathbf{x}(k)} & \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{u}^2(k)} \end{bmatrix}_* \begin{bmatrix} \delta \mathbf{x}(k) \\ \delta \mathbf{u}(k) \end{bmatrix} \quad (10.26)$$

$$+ \delta \mathbf{x}^T(k_f) \frac{\partial^2 J_f(\mathbf{x}(k_f), k_f)}{\partial \mathbf{x}^2(k_f)} \Big|_* \delta \mathbf{x}(k_f)$$

Table 10.1: Optimality conditions.

Condition	Equation
State equation	$\mathbf{x}^*(k+1) = A\mathbf{x}^*(k) + B^*\mathbf{u}(k) \quad k = k_0, \dots, k_f - 1$
Costate equation	$\lambda^*(k) = \frac{\partial H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{x}(k)} \Big _* \quad k = k_0, \dots, k_f - 1$
Minimum principle	$\frac{\partial H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{u}(k)} \Big _* = \mathbf{0}, \quad k = k_0, \dots, k_f - 1$

For the second-order term to be positive or at least zero for any perturbation, we need the Hessian matrix to be positive definite, or at least positive semidefinite. We have the sufficient condition

$$\begin{bmatrix} \frac{\partial H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{x}^2(k)} & \frac{\partial H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{x}(k) \partial \mathbf{u}(k)} \\ \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{u}(k) \partial \mathbf{x}(k)} & \frac{\partial^2 H(\mathbf{x}(k), \mathbf{u}(k), k)}{\partial \mathbf{u}^2(k)} \end{bmatrix}_* > 0 \quad (10.27)$$

If the Hessian matrix is positive semidefinite, then the condition is necessary for a minimum but not sufficient because there will be perturbations for which the second-order terms are zero. Higher-order terms are then needed to determine if these perturbations result in an increase in the performance measure.

For a fixed terminal state, the corresponding perturbations are zero and no additional conditions are required. For a free terminal state, we have the additional condition

$$\left. \frac{\partial^2 J(\mathbf{x}(k_f), k_f)}{\partial \mathbf{x}^2(k_f)} \right|_* > 0 \quad (10.28)$$

Example 10.3

Because of their preference for simple models, engineers often use first-order systems for design. These are known as scalar systems. We consider the scalar system

$$x(k+1) = ax(k) + bu(k), x(0) = x_0, b > 0$$

If the system is required to reach the zero state, find the control that minimizes the control effort to reach the desired final state.

Solution

For minimum control effort, we have the performance measure

$$J = \frac{1}{2} \sum_{k=0}^{k_f} u(k)^2$$

The Hamiltonian is given by

$$H(\mathbf{x}(k), \mathbf{u}(k), k) = \frac{1}{2} u(k)^2 + \lambda(k+1)[a\mathbf{x}(k) + b\mathbf{u}(k)], \quad k = 0, \dots, k_f - 1$$

We minimize the Hamiltonian by selecting the control

$$u^*(k) = -b\lambda^*(k+1)$$

The costate equation is

$$\lambda^*(k) = a\lambda^*(k+1) \quad k = 0, \dots, k_f - 1$$

Example 10.3—cont'd

and its solution is given by

$$\lambda^*(k) = a^{k_f - k} \lambda^*(k_f), \quad k = 0, \dots, k_f - 1$$

The optimal control can now be written as

$$u^*(k) = -ba^{k_f - k - 1} \lambda^*(k_f), \quad k = 0, \dots, k_f - 1$$

We next consider iterations of the state equation with the optimal control until the terminal zero state is reached

$$\begin{aligned} x^*(1) &= ax(0) + bu^*(0) = ax_0 - b^2 a^{k_f - 1} \lambda^*(k_f) \\ x^*(2) &= ax(1) + bu^*(1) = a^2 x_0 - b^2 \{a^{k_f} + a^{k_f - 2}\} \lambda^*(k_f) \\ x^*(3) &= ax(2) + bu^*(2) = a^3 x_0 - b^2 \{a^{k_f + 1} + a^{k_f - 1} + a^{k_f - 3}\} \lambda^*(k_f) \\ &\vdots \\ x^*(k_f) &= ax(k_f - 1) + bu^*(k_f - 1) = a^{k_f} x_0 - b^2 \{(a^2)^{k_f - 1} + \dots + a^2 + 1\} \lambda^*(k_f) = 0 \end{aligned}$$

We solve the last equation for the terminal Lagrange multiplier

$$\lambda^*(k_f) = \frac{(a^{k_f}/b^2)x_0}{(a^2)^{k_f - 1} + \dots + a^2 + 1}$$

Substituting for the terminal Lagrange multiplier gives the optimal control

$$u^*(k) = -\frac{(a^{2k_f - k - 1}/b)x_0}{(a^2)^{k_f - 1} + \dots + a^2 + 1} \quad k = 0, \dots, k_f - 1$$

For an open-loop stable system $|a| < 1$, and we can write the control in the form

$$u^*(k) = -\frac{a^2 - 1}{(a^2)^{k_f}} (a^{2k_f - k - 1}/b)x_0 \quad k = 0, \dots, k_f - 1, |a| < 1$$

To obtain a state feedback control law, we write the state at time k as

$$\begin{aligned} x^*(k) &= ax(k - 1) + bu^*(k - 1) \\ &= a^k x_0 - b^2 \{a^{k_f + 1} + a^{k_f - 1} + \dots + a^{k_f - k}\} \lambda^*(k_f) \\ &= a^k x_0 + \frac{b}{a^{k_f - k - 1}} \{a^{k_f + 1} + a^{k_f - 1} + \dots + a^{k_f - k}\} u^*(k) \end{aligned}$$

We solve for the control law

$$u^*(k) = \frac{x(k) - a^k x_0}{ba \{a^{k+1} + a^{k-1} + \dots + 1\}}, \quad k = 0, \dots, k_f - 1$$

The closed-loop system is given by

$$x(k+1) = ax(k) + \frac{x(k) - a^k x_0}{a \{a^{k+1} + a^{k-1} + \dots + 1\}}, \quad k = 0, \dots, k_f - 1$$

Example 10.3—cont'd

$$x(0) = x_0, b > 0$$

Note that the closed-loop system dynamics of the optimal system are time varying even though the open-loop system is time invariant.

For an initial state $x_0 = 1$ and the unstable dynamics with $a = 1.5$ and $b = 0.5$, the optimal trajectories are shown in Fig. 10.2. The optimal control stabilizes the system and drives its trajectories to the origin.

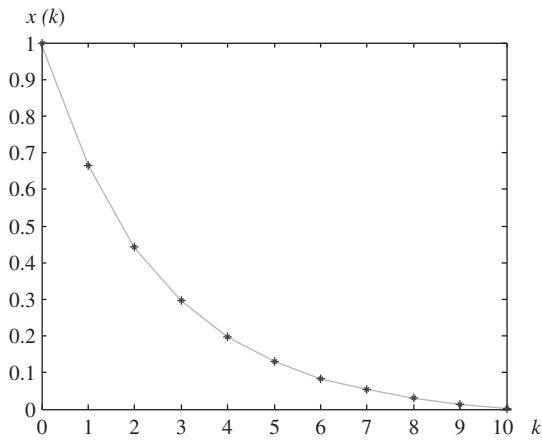


Figure 10.2
Optimal trajectory for the scalar system of Example 10.3.

10.3 The linear quadratic regulator

The choice of performance index in optimal control determines the performance of the system and the complexity of the optimal control problem. The most popular choice for the performance index is a quadratic function of the state variable and the control inputs. The performance measure is of the form

$$J = \frac{1}{2} \mathbf{x}^T(k_f) S(k_f) \mathbf{x}(k_f) + \frac{1}{2} \sum_{k=k_0}^{k_f-1} (\mathbf{x}^T(k) Q(k) \mathbf{x}(k) + \mathbf{u}^T(k) R(k) \mathbf{u}(k)) \quad (10.29)$$

where the matrices $S(k_f)$, $Q(k)$, $k = k_0, \dots, k_f-1$ are positive semidefinite symmetric $n \times n$, and the matrices $R(k)$, $k = k_0, \dots, k_f-1$ are positive definite symmetric $m \times m$. This choice is known as the **linear quadratic regulator**, and in addition to its desirable mathematical properties, it can be physically justified. In a regulator problem, the purpose is to maintain

the system close to the zero state, and the quadratic function of the state variable is a measure of error. On the other hand, the effort needed to minimize the error must also be minimized, as must the control effort represented by a quadratic function of the controls. Thus, the quadratic performance measure achieves a compromise between minimizing the regulator error and minimizing the control effort.

For the special case of diagonal weighting matrices, the performance index is the sum of the squares of the state variables and inputs. The relative importance of the different states and inputs is reflected in the choice of their weights, with the more important terms given larger weights. The larger weights ensure that the corresponding terms are kept small because they would otherwise be eliminated upon minimization.

To obtain the necessary and sufficient conditions for the linear quadratic regulator, we use the results of [Section 10.2](#). We first write an expression for the Hamiltonian:

$$\begin{aligned} H(\mathbf{x}(k), \mathbf{u}(k), k) = & \frac{1}{2} \mathbf{x}^T(k) Q(k) \mathbf{x}(k) + \frac{1}{2} \mathbf{u}^T(k) R(k) \mathbf{u}(k) \\ & + \lambda^T(k+1) [A\mathbf{x}(k) + B\mathbf{u}(k)], \quad k = k_0, \dots, k_f - 1 \end{aligned} \quad (10.30)$$

The state equation is unchanged, but the costate equation becomes

$$\lambda^*(k) = Q(k)\mathbf{x}^*(k) + A^T\lambda^*(k+1) \quad k = k_0, \dots, k_f - 1 \quad (10.31)$$

Pontryagin's minimum principle gives

$$R(k)\mathbf{u}^*(k) + B^T\lambda^*(k+1) = 0$$

which yields the optimum control expression

$$\mathbf{u}^*(k) = -R^{-1}(k)B^T\lambda^*(k+1) \quad (10.32)$$

Substituting in the state equation, we obtain the optimal dynamics

$$\mathbf{x}^*(k+1) = A\mathbf{x}^*(k) - B^*R^{-1}(k)B^T\lambda^*(k+1), \quad k = k_0, \dots, k_f - 1 \quad (10.33)$$

[Eq. \(10.27\)](#) gives the sufficient optimality condition

$$\begin{bmatrix} Q(k) & \mathbf{0}_{n \times m} \\ \mathbf{0}_{m \times n} & R(k) \end{bmatrix} > \mathbf{0} \quad (10.34)$$

Because R is positive definite, a sufficient condition for a minimum is that Q must be positive definite. A necessary condition is that Q must be positive semidefinite.

10.3.1 Free final state

If the final state of the optimal control is not fixed, we have the terminal condition based on Eq. (10.25):

$$\lambda^*(k_f) = \frac{\partial J_f(\mathbf{x}(k_f), k_f)}{\partial \mathbf{x}(k_f)} = S(k_f)\mathbf{x}^*(k_f) \quad (10.35)$$

where the matrix $S(k_f)$ is the terminal weight matrix of Eq. (10.29). The condition suggests a general relationship between the state and costate—that is,

$$\lambda^*(k) = S(k)\mathbf{x}^*(k) \quad (10.36)$$

If the proposed relation yields the correct solution, this would allow us to conclude that the relation is indeed correct. We substitute the proposed relation Eq. (10.36) in the state Eq. (10.33):

$$\mathbf{x}^*(k+1) = A\mathbf{x}^*(k) - BR^{-1}(k)B^T S(k+1)\mathbf{x}^*(k+1)$$

This yields the recursion

$$\mathbf{x}^*(k+1) = \{I + BR^{-1}(k)B^T S(k+1)\}^{-1} A\mathbf{x}^*(k) \quad (10.37)$$

Using the proposed relation Eq. (10.36) and substituting in the costate equation, we have

$$\begin{aligned} \lambda^*(k) &= Q(k)\mathbf{x}^*(k) + A^T S(k+1)\mathbf{x}^*(k+1) \\ &= \left\{ Q(k) + A^T S(k+1)[I + BR^{-1}(k)B^T S(k+1)]^{-1} A \right\} \mathbf{x}^*(k) \\ &= S(k)\mathbf{x}^*(k) \end{aligned} \quad (10.38)$$

Eq. (10.38) holds for all values of the state vector and hence we have the matrix equality

$$Q(k) + A^T S(k+1)[I + B^* R^{-1}(k)B^T S(k+1)]^{-1} A = S(k)$$

We apply the matrix inversion lemma (see Appendix III) to obtain

$$\begin{aligned} S(k) &= A^T \left\{ S(k+1) - S(k+1)B(B^T S(k+1)B + R(k))^{-1} B^T S(k+1) \right\} A \\ &\quad + Q(k), S(k_f) \end{aligned} \quad (10.39)$$

The preceding equation is known as the **matrix Riccati equation** and can be solved iteratively backward in time to obtain the matrices $S(k)$, $k = k_0, \dots, k_{f-1}$.

We return to the expression for the optimal control Eq. (10.32) and substitute for the costate to obtain

$$\begin{aligned}\mathbf{u}^*(k) &= -R^{-1}(k)B^T S(k+1)\mathbf{x}^*(k+1) \\ &= -R^{-1}(k)B^T S(k+1)[A\mathbf{x}^*(k) + B\mathbf{u}(k)]\end{aligned}$$

We solve for the control

$$\mathbf{u}^*(k) = -[I + R^{-1}(k)B^T S(k+1)B]^{-1} R^{-1}(k)B^T S(k+1)A\mathbf{x}^*(k)$$

Using the rule for the inverse of a product, we have the optimal state feedback

$$\begin{aligned}\mathbf{u}^*(k) &= -K(k)\mathbf{x}^*(k) \\ K(k) &= [R(k) + B^T S(k+1)B]^{-1} B^T S(k+1)A\end{aligned}\tag{10.40}$$

Thus, with the offline solution of the Riccati equation, we can implement the optimal state feedback.

The Riccati equation can be written in terms of the optimal gain of Eq. (10.40) in the form

$$\begin{aligned}S(k) &= A^T S(k+1) \left\{ A - B(B^T S(k+1)B + R(k))^{-1} B^T S(k+1)A \right\} + Q(k) \\ &= A^T S(k+1) \{A - BK(k)\} + Q(k)\end{aligned}$$

In terms of the closed-loop state matrix $A_{cl}(k)$, we have

$$\begin{aligned}S(k) &= A^T S(k+1)A_{cl}(k) + Q(k) \\ A_{cl}(k) &= A - BK(k)\end{aligned}\tag{10.41}$$

Note that the closed-loop matrix and the optimal feedback gain matrix are time varying even for a time-invariant pair (A, B) .

A more useful form of the Riccati is obtained by adding and subtracting terms to obtain

$$S(k) = [A - BK(k)]^T S(k+1)A_{cl}^T(k) + Q(k) + K^T(k)B^T(k)S(k+1)[A - BK(k)]$$

We expand the added term and use Eq. (10.40) to write it as

$$\begin{aligned}K^T(k)B^T(k)S(k+1)[A - BK(k)] &= K^T(k)B^T(k)S(k+1)A \\ &\quad + K^T(k)\{R(k) - R(k) - B^T S(k+1)B\}K(k) \\ &= K^T(k)B^T(k)S(k+1)A + K^T(k)R(k)K(k) \\ - K^T(k)[R(k) + B^T S(k+1)B] &[R(k) + B^T S(k+1)B]^{-1} B^T S(k+1)A \\ &= K^T(k)R(k)K(k)\end{aligned}$$

We now have the **Joseph form** of the Riccati equation:

$$S(k) = A_{cl}^T(k)S(k+1)A_{cl}(k) + K^T(k)R(k)K(k) + Q(k) \quad (10.42)$$

Because of its symmetry, this form performs better in iterative numerical computation.

If the Riccati equation in Joseph form is rewritten with the optimal gain replaced by a suboptimal gain, then the equation becomes a **Lyapunov difference equation** (see Chapter 11). For example, if the optimal gain is replaced by a constant gain K to simplify implementation, then the performance of the system will be suboptimal and is governed by a Lyapunov equation.

To find the optimal cost, we first consider the sum of quadratic forms

$$\begin{aligned} & \frac{1}{2} \sum_{k=k_0}^{k_f-1} \mathbf{x}^T(k) \{ A_{cl}^T S(k+1) A_{cl}(k) - S(k) \} \mathbf{x}(k) \\ &= \frac{1}{2} \sum_{k=k_0}^{k_f-1} \{ \mathbf{x}^T(k+1) S(k+1) \mathbf{x}(k+1) - \mathbf{x}^T(k) S(k) \mathbf{x}(k) \} \end{aligned}$$

where we used the state equation $\mathbf{x}(k+1) = A_{cl}(k)\mathbf{x}(k)$. We expand the right hand side as

$$\begin{aligned} & \sum_{k=k_0}^{k_f-1} \{ \mathbf{x}^T(k+1) S(k+1) \mathbf{x}(k+1) - \mathbf{x}^T(k) S(k) \mathbf{x}(k) \} \\ &= \mathbf{x}^T(k_f) S(k_f) \mathbf{x}(k_f) - \mathbf{x}^T(k_f-1) S(k_f-1) \mathbf{x}(k_f-1) \\ &+ \mathbf{x}^T(k_f-1) S(k_f-1) \mathbf{x}(k_f-1) - \mathbf{x}^T(k_f-2) S(k_f-2) \mathbf{x}(k_f-2) \\ &\quad \vdots \\ &+ \mathbf{x}^T(k_0+1) S(k_0+1) \mathbf{x}(k_0+1) - \mathbf{x}^T(k_0) S(k_0) \mathbf{x}(k_0) \\ &= \mathbf{x}^T(k_f) S(k_f) \mathbf{x}(k_f) - \mathbf{x}^T(k_0) S(k_0) \mathbf{x}(k_0) \end{aligned}$$

The terminal cost can therefore be written as

$$\begin{aligned} & \mathbf{x}^T(k_f) S(k_f) \mathbf{x}(k_f) \\ &= \sum_{k=k_0}^{k_f-1} \{ \mathbf{x}^T(k+1) S(k+1) \mathbf{x}(k+1) - \mathbf{x}^T(k) S(k) \mathbf{x}(k) \} + \mathbf{x}^T(k_0) S(k_0) \mathbf{x}(k_0) \end{aligned}$$

We substitute for the terminal cost in the performance index

$$\begin{aligned} J &= \frac{1}{2} \mathbf{x}^T(k_0) S(k_0) \mathbf{x}(k_0) \\ &+ \frac{1}{2} \sum_{k=k_0}^{k_f-1} \left\{ \mathbf{x}^T(k+1) S(k+1) \mathbf{x}(k+1) + \mathbf{x}^T(k) [Q(k) - S(k)] \mathbf{x}(k) + u^T(k) R(k) u(k) \right\} \end{aligned}$$

Substituting the optimal control and using the closed-loop state equation, we have the optimal cost

$$\begin{aligned} J^* &= \frac{1}{2} \mathbf{x}^T(k_0) S(k_0) \mathbf{x}(k_0) \\ &+ \frac{1}{2} \mathbf{x}^T(k) \left\{ \sum_{k=k_0}^{k_f-1} \left\{ A_{cl}^T S(k+1) A_{cl}(k) + Q(k) - S(k) + K^T(k) R(k) K(k) \right\} \right\} \mathbf{x}(k) \end{aligned}$$

From the Joseph form of Eq. (10.42), we observe that each term in the summation is zero. The optimal cost is

$$J^* = \frac{1}{2} \mathbf{x}^T(k_0) S(k_0) \mathbf{x}(k_0) \quad (10.43)$$

Example 10.4

A variety of mechanical systems can be approximately modeled by the double integrator

$$\ddot{x}(t) = u(t)$$

where $u(t)$ is the applied force. With digital control and a sampling period $T = 0.02$ s, the double integrator has the discrete state-space equation

$$\begin{aligned} \mathbf{x}(k+1) &= \begin{bmatrix} 1 & T \\ 0 & 1 \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} T^2/2 \\ T \end{bmatrix} u(k) \\ y(k) &= [1 \ 0] \mathbf{x}(k) \end{aligned}$$

1. Design a linear quadratic regulator for the system with terminal weight $S(100) = \text{diag}\{10, 1\}$, $Q = \text{diag}\{10, 1\}$, and control weight $r = 0.1$, then simulate the system with initial condition vector $\mathbf{x}(0) = [1, 0]^T$.
2. Repeat the design of part 1 with $S(100) = \text{diag}\{10, 1\}$, $Q = \text{diag}\{10, 1\}$, and control weight $r = 100$. Compare the results to part 1, and discuss the effect of changing the value of r .

Example 10.4—cont'd**Solution**

The Riccati equation can be expanded with the help of a symbolic manipulation package to obtain

$$\begin{aligned} & \begin{bmatrix} s_{11} & s_{12} \\ s_{12} & s_{22} \end{bmatrix}_k \\ &= \frac{\begin{bmatrix} (s_{11}s_{22} - s_{12}^2)T^2 + s_{11}r & 2(s_{11}s_{22} - s_{12}^2)T^3 + 4r(s_{11}T + 2s_{12}) \\ 2(s_{11}s_{22} - s_{12}^2)T^3 + 4r(s_{11}T + 2s_{12}) & (s_{11}s_{22} - s_{12}^2)T^4 + 4(s_{11}T^2 + 2Ts_{12} + s_{22})r \end{bmatrix}}{T^4s_{11} + 4T^3s_{12} + 4T^2s_{22} + 4r} \Big|_{k+1} \\ &+ Q \end{aligned}$$

where the subscript denotes the time at which the matrix S is evaluated by backward-in-time iteration starting with the given terminal value. We use the simple MATLAB script:

```
% simlqr: simulate a scalar optimal control DT system
t(1) = 0;
x{1} = [1;0]; % Initial state
T = 0.02; % Sampling period
N = 150; % Number of steps
S = diag([10,1]);r = 0.1; % Weights
Q = diag([10,1]);
A = [1,T;0,1];B = [T/2;T]; % System matrices
for i = N:1:1
K{i} = (r + B'*S*A*B)\B'*S*A; % Calculate the optimal feedback % gains
% Note that K(1) is really K(0)
kp(i) = K{i}(1); % Position gain
kv(i) = K{i}(2); % Velocity gain
Acl = A-B*K{i};
S = Acl'*S*Acl + K{i}'*r*K{i}+Q; % Iterate backward (Riccati-% Joseph form)
end
for i = 1:N
t(i+1) = t(i)+T;
u(i) = -K{i}*x{i};
x{i+1} = A*x{i}+B*u(i); % State equation
```

Example 10.4—cont'd

```

end
xmat = cell2mat(x); % Change cell to mat to extract data
xx = xmat(1,:); % Position
xv = xmat(2,:); % Velocity
plot(t,xx,'.') % Plot Position
figure % New figure
plot(t,xv,'.') % Plot Velocity
figure % New figure
plot(xx,xv,'.')% Plot phase plane trajectory
figure % New figure
plot(t(1:N),u,'.') % Plot control input
figure % New figure
plot(t(1:N),kp,'.',t(1:N),kv,'.') % Plot position gain and velocity
%gain

```

Five plots are obtained using the script for each part.

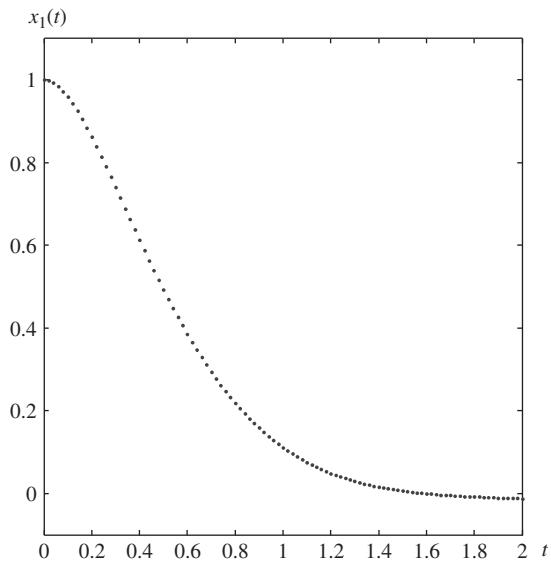
1. The first three plots are the position trajectory of Fig. 10.3, the velocity trajectory of Fig. 10.4, and the phase-plane trajectory of Fig. 10.5. The three plots show that the optimal control drives the system toward the origin. However, because the final state is free, the origin is not reached. From the gain plot shown in Fig. 10.6, we observe that the controller gains, which are obtained by iteration backward in time, approach a fixed level. Consequently, the optimal control shown in Fig. 10.7 also approaches a fixed level.
2. The first three plots are the position trajectory of Fig. 10.8, the velocity trajectory of Fig. 10.9, and the phase-plane trajectory of Fig. 10.10. The three plots show that the optimal control drives the system toward the origin. However, the final state reached is farther from the origin than that of part 1 because the error weight matrix is now smaller relative to the value of r . The larger control weight r results in smaller control gain as shown in Fig. 10.11. This corresponds to a reduction in the control effort as shown in Fig. 10.12. Note that the controller gains approach a fixed level at a slower rate than in part 1.

Calculating the matrices backward in time gives the matrix

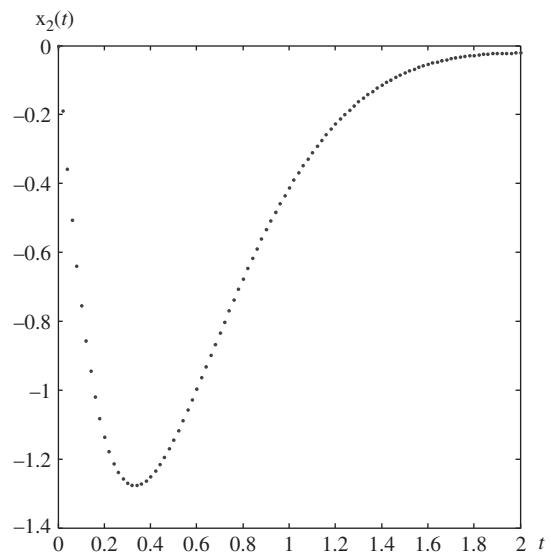
$$S(0) = \begin{bmatrix} 278.9351 & 50.0204 \\ 50.0204 & 27.9074 \end{bmatrix}$$

The optimal cost is

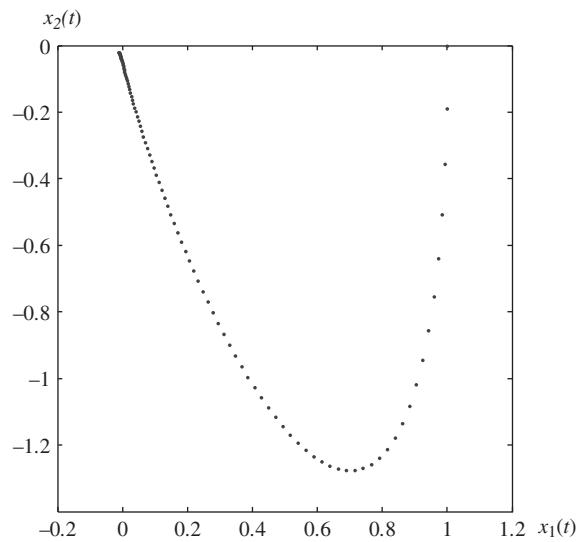
$$J^* = \frac{1}{2} x^T(0) S(k_0) x(0) = \frac{1}{2} [1 \ 0]^T \begin{bmatrix} 278.9351 & 50.0204 \\ 50.0204 & 27.9074 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 139.4675$$

Example 10.4—cont'd**Figure 10.3**

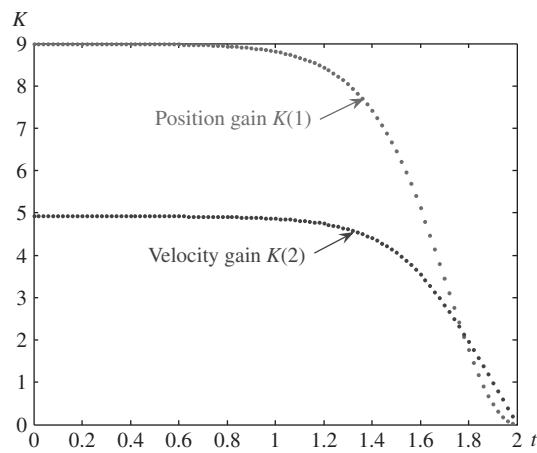
Position trajectory for the inertial control system described in Example 10.4 (1).

**Figure 10.4**

Velocity trajectory for the inertial control system described in Example 10.4 (1).

Example 10.4—cont'd**Figure 10.5**

Phase plane trajectory for the inertial control system described in Example 10.4 (1).

**Figure 10.6**

Plot of feedback gains versus time for Example 10.4 (1).

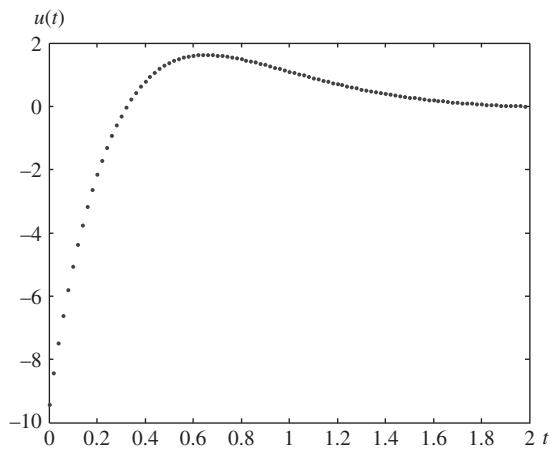
Example 10.4—cont'd

Figure 10.7
Optimal control input for Example 10.4 (1).

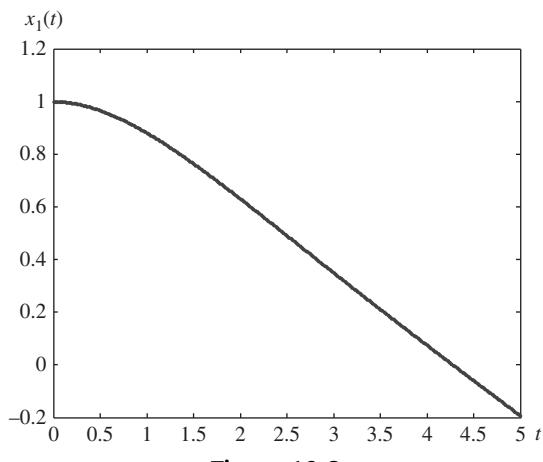
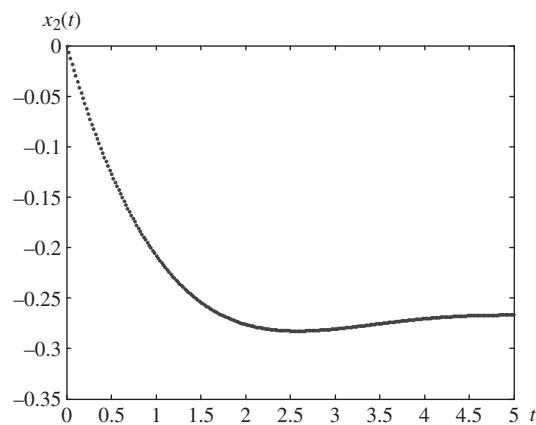
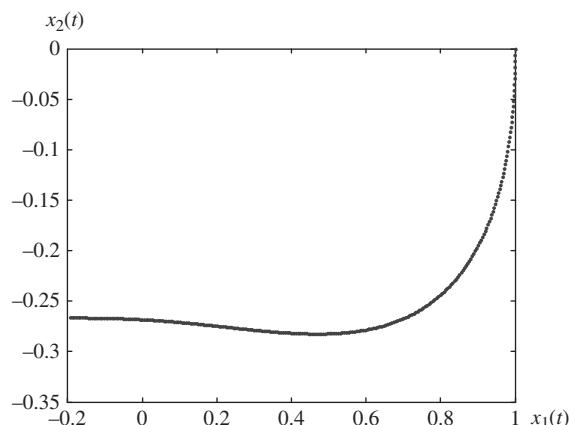


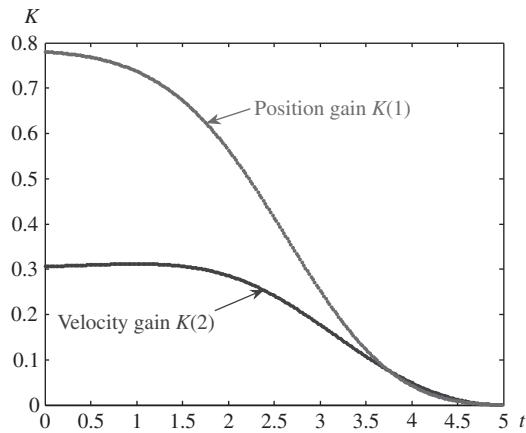
Figure 10.8
Position trajectory for the inertial control system of Example 10.4 (2).

Example 10.4—cont'd**Figure 10.9**

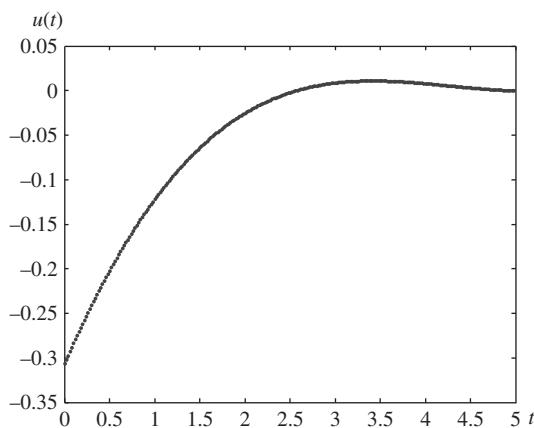
Velocity trajectory for the inertial control system of Example 10.4 (2).

**Figure 10.10**

Phase plane trajectory for the inertial control system of Example 10.4 (2).

Example 10.4—cont'd**Figure 10.11**

Plot of feedback gains versus time for Example 10.4 (2).

**Figure 10.12**

Optimal control input for Example 10.4 (2).

10.4 Steady-state quadratic regulator

Implementing the linear quadratic regulator is rather complicated because of the need to calculate and store gain matrices. From Example 10.4 we observe that the gain values converge to fixed values. This occurs in general with the fulfillment of some simple requirements that are discussed in this section.

For many applications, it is possible to simplify implementation considerably by using the steady-state gain exclusively in place of the optimal gains. This solution is only optimal if the summation interval in the performance measure, known as the **planning horizon**, is infinite. For a finite planning horizon, the simplified solution is suboptimal (i.e., gives a higher value of the performance measure) but often performs almost as well as the optimal control. Thus, it is possible to retain the performance of optimal control without the burden of implementing it if we solve the **steady-state regulator problem** with the performance measure of the form

$$J = \frac{1}{2} \sum_{k=k_0}^{\infty} (\mathbf{x}^T(k) Q \mathbf{x}(k) + \mathbf{u}^T(k) R \mathbf{u}(k)) \quad (10.44)$$

We assume that the weighting matrices Q and R are constant, with Q positive semidefinite and R positive definite. We can therefore decompose the matrix Q as

$$Q = Q_a^T Q_a \quad (10.45)$$

where Q_a is the square root matrix. This allows us to write the state error terms of “measurements” in terms of an equivalent measurement vector

$$\begin{aligned} \mathbf{y}(k) &= Q_a \mathbf{x}(k) \\ \mathbf{x}^T(k) Q \mathbf{x}(k) &= \mathbf{y}^T(k) \mathbf{y}(k) \end{aligned} \quad (10.46)$$

The matrix Q_a is positive definite for Q positive definite and positive semidefinite for Q positive semidefinite. In the latter case, large state vectors can be mapped by Q_a to zero $\mathbf{y}(k)$, and a small performance measure cannot guarantee small errors or even stability. We think of the vector $\mathbf{y}(k)$ as a measurement of the state and recall the detectability condition from Chapter 8.

We recall that if the pair (A, Q_a) is detectable, then $\mathbf{x}(k)$ decays asymptotically to zero with $\mathbf{y}(k)$. Hence, this detectability condition is required for acceptable steady-state regulator behavior. If the pair is observable, then it is also detectable and the system behavior is acceptable. For example, if the matrix Q_a is positive definite, there is a one-one mapping between the states \mathbf{x} and the measurements \mathbf{y} , and the pair (A, Q_a) is always observable. Hence, a positive definite Q (and Q_a) is sufficient for acceptable system behavior.

If the detectability (observability) condition is satisfied and the system is stabilizable, the Riccati equation does not diverge and a steady-state condition is reached. The resulting algebraic Riccati equation is in the form

$$S = A^T \left\{ S - SB(B^T SB + R)^{-1} B^T S \right\} A + Q \quad (10.47)$$

It can be shown that, under these conditions, the algebraic Riccati equation has a unique positive definite solution. However, the equation is clearly difficult to solve in general and is typically solved numerically. The MATLAB solution of the algebraic Riccati equation is discussed in [Section 10.4.2](#).

The optimal state feedback corresponding to the steady-state regulator is

$$\begin{aligned} \mathbf{u}^* &= -K\mathbf{x}^*(k) \\ K &= [R + B^T SB]^{-1} B^T SA \end{aligned} \quad (10.48)$$

Using the gain expression, we can write the algebraic Riccati equation in the Joseph form:

$$\begin{aligned} S &= A_{cl}^T S A_{cl} + K^T R K + Q \\ A_{cl} &= A - BK \end{aligned} \quad (10.49)$$

If the optimal gain K is replaced by a suboptimal gain, then the algebraic Riccati equation becomes an **algebraic Lyapunov equation** (see Chapter 11). The Lyapunov equation is clearly linear and its solution is simpler than that of the Riccati equation.

10.4.1 Output quadratic regulator

In most practical applications, the designer is interested in optimally controlling the output $\mathbf{y}(k)$ rather than the state $\mathbf{x}(k)$. To optimally control the output, we need to consider a performance index of the form

$$J = \frac{1}{2} \sum_{k=k_0}^{\infty} (\mathbf{y}^T(k) Q_y \mathbf{y}(k) + \mathbf{u}^T(k) R \mathbf{u}(k)) \quad (10.50)$$

From [Eq. \(10.46\)](#), we observe that this is equivalent to the original performance measure of [Eq. \(10.44\)](#) with the state weight matrix

$$\begin{aligned} Q &= C^T Q_y C \\ &= C^T Q_{ya}^T Q_{ya} C \\ &= Q_a^T Q_a \\ Q_a &= Q_{ya} C \end{aligned} \quad (10.51)$$

where Q_{ya} is the square root of the output weight matrix. As in the state quadratic regulator, the Riccati equation for the output regulator can be solved using the MATLAB commands discussed in Section 10.4.2. For a stabilizable pair (A, B) , the solution exists provided that the pair (A, Q_a) is detectable with Q_a as in Eq. (10.51).

10.4.2 MATLAB solution of the steady-state regulator problem

MATLAB has several commands that allow us to conveniently design steady-state regulators. The first is **dare**, which solves the discrete algebraic Riccati Eq. (10.39). The command is

$$\gg [\mathbf{S}, \mathbf{E}, \mathbf{K}] = \text{dare}(\mathbf{A}, \mathbf{B}, \mathbf{Q}, \mathbf{R})$$

The input arguments are the state matrix A , the input matrix B , and the weighting matrices Q and R . The output arguments are the solution of the discrete algebraic Riccati equation S , the feedback gain matrix K , and the eigenvalues E of the closed-loop optimal system $A - BK$.

The second command for discrete optimal control is **dlqr**, which solves the steady-state regulator problem. The command has the form

$$\gg [\mathbf{K}, \mathbf{S}, \mathbf{E}] = \text{dlqr}(\mathbf{A}, \mathbf{B}, \mathbf{Q}, \mathbf{R})$$

where the input arguments are the same as the command **dare**, and the output arguments, also the same, are the gain K , the solution of the discrete algebraic Riccati equation S , and the eigenvalues E of the closed-loop optimal system $A - BK$.

For the output regulator problem, we can use the commands **dare** and **dlqr** with the matrix Q replaced by $C^T Q_y C$. Alternatively, MATLAB provides the command

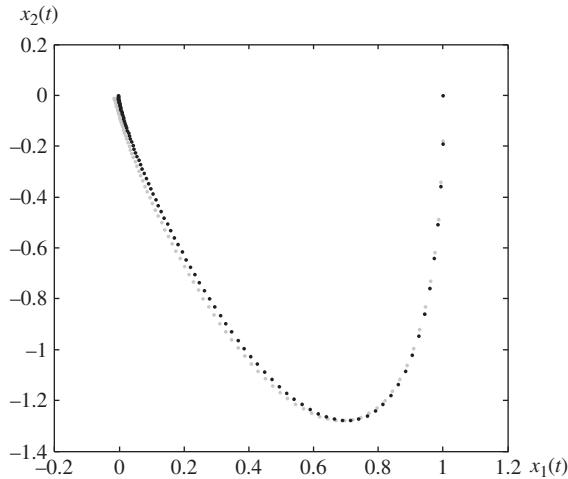
$$\gg [\mathbf{K}_y, \mathbf{S}, \mathbf{E}] = \text{dlqry}(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \mathbf{Q}_y, \mathbf{R})$$

Example 10.5

Design a steady-state regulator for the inertial system of Example 10.4, and compare its performance to the optimal control by plotting the phase plane trajectories of the two systems. Explain the advantages and disadvantages of the two controllers.

The pair (A, B) is observable. The state-error weighting matrix is positive definite, and the pair (A, Q_a) is observable. We are therefore assured that the Riccati equation will have a steady-state solution. For the inertial system presented in Example 10.4, we have the MATLAB output

```
>> [K, S, E] = dlqr(A, B, Q, r)
K = 9.4671  5.2817
S = 278.9525  50.0250  50.0250  27.9089
E = 0.9462 + 0.0299i
    0.9462 - 0.0299i
```

Example 10.5—cont'd**Figure 10.13**

Phase plane trajectories for the inertial system of Example 10.5 with optimal control (dark gray) and suboptimal steady-state regulator (light gray).

If we simulate the system with the optimal control of Example 10.4 and superimpose the trajectories for the suboptimal steady-state regulator, we obtain Fig. 10.13. The figure shows that the suboptimal control, although much simpler to implement, provides an almost identical trajectory to that of the optimal control. For practical implementation, the suboptimal control is often preferable because it is far cheaper to implement and provides almost the same performance as the optimal control. However, there may be situations where the higher accuracy of the optimal control justifies the additional cost of its implementation.

Example 10.6

Design a digital output regulator for the double integrator system

$$\begin{aligned}\mathbf{x}(k+1) &= \begin{bmatrix} 1 & T \\ 0 & 1 \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} T^2/2 \\ T \end{bmatrix} u(k) \\ y(k) &= [1 \quad 0] \mathbf{x}(k)\end{aligned}$$

with sampling period $T = 0.02$ s, output weight $Q_y = 1$, and control weight $r = 0.1$, and plot the output response for the initial condition vector $\mathbf{x}(0) = [1, 0]^T$.

Example 10.6—cont'd**Solution**

We use the MATLAB command

```
>> [Ky, S, E] = dlqry(A, B, C, 0, 1, 0.1)
Ky = 3.0837 2.4834
S = 40.2667 15.8114 15.8114 12.5753
E = 0.9749 + 0.0245i 0.9749 - 0.0245i
```

The response of the system to initial conditions $\mathbf{x}(0) = [1, 0]^T$ is shown in Fig. 10.14. The output quickly converges to the target zero state.

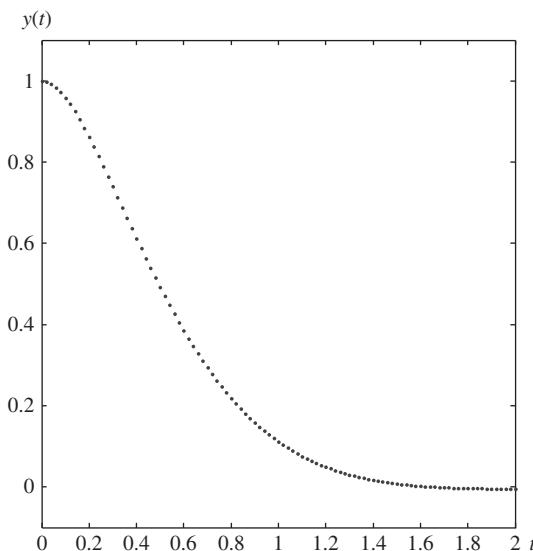


Figure 10.14

Time response of the output regulator discussed in Example 10.6.

10.4.3 Linear quadratic tracking controller

The minimization of the performance index Eq. (10.44) yields an optimal state feedback matrix K that minimizes integral square error and control effort (see Eq. (10.48)). The error is defined relative to the zero state. In many practical applications, the system is required to follow a specified function of time. The design of a controller to achieve this objective is known as the **tracking problem**. If the control task is to track a nonzero constant reference input, we can exploit the techniques described in Section 9.3 to solve the new optimal regulator problem.

In Section 9.3, we showed that the tracking or servo problem can be solved using an additional gain matrix for the reference input as shown in Fig. 9.7, thus providing an additional degree of freedom in our design. For a square system (equal number of inputs and outputs), we can implement the two degree-of-freedom control of Fig. 9.7 with the reference gain matrix

$$F = \left[C(I_n - A_{cl})^{-1}B \right]^{-1} \quad (10.52)$$

where

$$A_{cl} = A - BK \quad (10.53)$$

All other equations for the linear quadratic regulator are unchanged.

Alternatively, to improve robustness, but at the expense of a higher-order controller, we can introduce integral control, as in Fig. 9.9. In particular, as in Eq. (9.29) we consider the state-space equations

$$\begin{aligned} \mathbf{x}(k+1) &= A\mathbf{x}(k) + B\mathbf{u}(k) \\ \bar{\mathbf{x}}(k+1) &= \bar{\mathbf{x}}(k) + \mathbf{r}(k) - \mathbf{y}(k) \\ \mathbf{y}(k) &= C\mathbf{x}(k) \\ \mathbf{u}(k) &= -K\mathbf{x}(k) - \bar{K}\bar{\mathbf{x}}(k) \end{aligned} \quad (10.54)$$

with integral control gain \bar{K} . This yields the closed-loop state-space equations Eq. (9.30)

$$\begin{aligned} \tilde{\mathbf{x}}(k+1) &= (\tilde{A} - \tilde{B}\tilde{K})\tilde{\mathbf{x}}(k) + \begin{bmatrix} \mathbf{0} \\ I_l \end{bmatrix} \mathbf{r}(k) \\ \mathbf{y}(k) &= [C \quad \mathbf{0}] \tilde{\mathbf{x}}(k) \end{aligned} \quad (10.55)$$

where $\tilde{\mathbf{x}}(k) = [\mathbf{x}(k) \quad \bar{\mathbf{x}}(k)]^T$ and the matrices are given by

$$\begin{aligned} \tilde{A} &= \begin{bmatrix} A & \mathbf{0} \\ -C & I_l \end{bmatrix} \\ \tilde{B} &= \begin{bmatrix} B \\ \mathbf{0} \end{bmatrix} \quad \tilde{C} = [C \quad \mathbf{0}] \quad \tilde{K} = [K \quad \bar{K}] \end{aligned} \quad (10.56)$$

The state feedback gain matrix can be computed as

$$\tilde{K} = [R + \tilde{B}^T S \tilde{B}]^{-1} \tilde{B}^T S \tilde{A} \quad (10.57)$$

Example 10.7

Design an optimal linear quadratic state-space tracking controller for the inertial system of Example 10.4 to obtain zero steady-state error due to a unit step input.

Example 10.7—cont'd**Solution**

The state-space matrices are

$$A = \begin{bmatrix} 1 & 0.02 \\ 0 & 1 \end{bmatrix} \quad B = \begin{bmatrix} 0.0002 \\ 0.02 \end{bmatrix}$$

$$C = [1 \quad 0]$$

Adding integral control, we calculate the augmented matrices

$$\tilde{A} = \begin{bmatrix} A & 0 \\ -C & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0.02 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \quad \tilde{B} = \begin{bmatrix} B \\ 0 \end{bmatrix} = \begin{bmatrix} 0.0002 \\ 0.02 \\ 0 \end{bmatrix}$$

We select the weight matrices with larger error weighting as

$$Q = \begin{bmatrix} 10 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad r = 0.1$$

The MATLAB command **dlqr** yields the desired solution:

```
>> [Ktilde, S, E] = dlqr(Atilde, Btilde, Q, r)
```

Ktilde = 57.363418836015583 10.818259776359207 -2.818765217602360

S = 1.0e+003 * 2.922350387967343 0.314021626641202 -0.191897141854919

0.314021626641202 0.058073322228013 -0.015819292019556

-0.191897141854919 -0.015819292019556 0.020350548700473

E = 0.939332687303670 + 0.083102739623114i

0.939332687303670 -0.083102739623114i

0.893496746098272

The closed-loop system state-space equation is

$$\tilde{\mathbf{x}}(k+1) = (\tilde{A} - \tilde{B}\tilde{K})\tilde{\mathbf{x}}(k) + \begin{bmatrix} \mathbf{0} \\ I_I \end{bmatrix} \mathbf{r}(k)$$

$$\mathbf{y}(k) = [C \quad \mathbf{0}] \tilde{\mathbf{x}}(k)$$

The closed-loop step response shown in Fig. 10.15 is satisfactory with a small overshoot and zero steady-state error due to step.

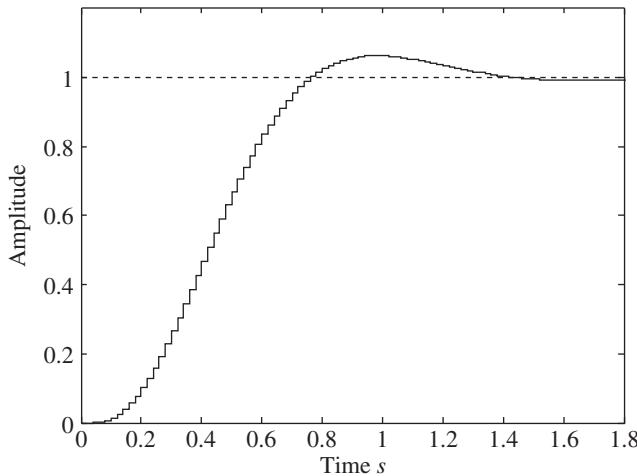
Example 10.7—cont'd

Figure 10.15
Step response for Example 10.7.

10.5 Hamiltonian system

In this section, we obtain a solution of the steady-state regulator problem using linear dynamics without directly solving the Riccati Eq. (10.47). We show that we can obtain the solution of the Riccati equation from the state-transition matrix of a linear system. We consider the case of constant weighting matrices. We can combine the state and costate equations to obtain the $2n$ -dimensional **Hamiltonian system**:

$$\begin{bmatrix} \mathbf{x}^*(k+1) \\ \lambda^*(k) \end{bmatrix} = \begin{bmatrix} A & -B^* R^{-1} B^T \\ Q & A^T \end{bmatrix} \begin{bmatrix} \mathbf{x}^*(k) \\ \lambda^*(k+1) \end{bmatrix}, \quad k = k_0, \dots, k_f - 1 \quad (10.58)$$

If the state matrix A is obtained by discretizing a continuous-time system, then it is a state transition matrix and is always invertible. We can then rewrite the state equation in the same form as the equation of the costate

$$\begin{aligned} \mathbf{x}^*(k) &= A^{-1} \left[\mathbf{x}^*(k+1) + B R(k)^{-1} B^T \lambda^*(k+1) \right] \\ \lambda^*(k) &= Q A^{-1} \left[\mathbf{x}^*(k+1) + B R(k)^{-1} B^T \lambda^*(k+1) \right] + A^T \lambda^*(k+1) \end{aligned}$$

We obtain the Hamiltonian system

$$\begin{aligned} \begin{bmatrix} \mathbf{x}^*(k) \\ \lambda^*(k) \end{bmatrix} &= H \begin{bmatrix} \mathbf{x}^*(k+1) \\ \lambda^*(k+1) \end{bmatrix}, \quad k = k_0, \dots, k_f - 1 \\ H &= \begin{bmatrix} A^{-1} & A^{-1}B^*R^{-1}B^T \\ QA^{-1} & A^T + QA^{-1}B^*R^{-1}B^T \end{bmatrix} \\ \lambda^*(k_f) &= S(k_f)\mathbf{x}(k_f) \end{aligned} \quad (10.59)$$

The Hamiltonian system of Eq. (10.59) describes the state and costate evolution backward in time because it provides an expression for evaluating the vector at time k from the vector at time $k+1$. We can solve the linear equation and write its solution in terms of the state-transition matrix for the Hamiltonian

$$\begin{aligned} \begin{bmatrix} \mathbf{x}^*(k) \\ \lambda^*(k) \end{bmatrix} &= H^{k_f-k} \begin{bmatrix} \mathbf{x}^*(k_f) \\ \lambda^*(k_f) \end{bmatrix}, \quad k = k_0, \dots, k_f \\ H^{k_f-k} &= \begin{bmatrix} \phi_{11}(k_f-k) & \phi_{12}(k_f-k) \\ \phi_{21}(k_f-k) & \phi_{22}(k_f-k) \end{bmatrix} \end{aligned} \quad (10.60)$$

The solution of the state equation yields

$$\begin{aligned} \mathbf{x}^*(k) &= \{\phi_{11}(k_f-k) + \phi_{12}(k_f-k)S(k_f)\}\mathbf{x}^*(k_f) \\ \lambda^*(k) &= \{\phi_{21}(k_f-k) + \phi_{22}(k_f-k)S(k_f)\}\mathbf{x}^*(k_f) \end{aligned} \quad (10.61)$$

We solve for the costate in terms of the state vector

$$\begin{aligned} \lambda^*(k) &= M(k)\mathbf{x}^*(k) \\ M(k) &= \{\phi_{21}(k_f-k) + \phi_{22}(k_f-k)S(k_f)\}\{\phi_{11}(k_f-k) \\ &\quad + \phi_{12}(k_f-k)S(k_f)\}^{-1} \end{aligned} \quad (10.62)$$

Substituting in Eq. (10.32), we now have the control

$$\begin{aligned} \mathbf{u}^*(k) &= -R^{-1}B^TM(k+1)\mathbf{x}^*(k+1) \\ &= -R^{-1}B^TM(k+1)\{A\mathbf{x}^*(k) + B\mathbf{u}^*(k)\} \\ k &= k_0, \dots, k_f - 1 \end{aligned} \quad (10.63)$$

We multiply both sides by the matrix R and then solve for the control

$$\begin{aligned} \mathbf{u}^*(k) &= -K(k)\mathbf{x}^*(k) \\ K(k) &= \{R + B^TM(k+1)B\}^{-1}B^TM(k+1)A, \quad k = k_0, \dots, k_f - 1 \end{aligned} \quad (10.64)$$

We compare the gain expression of Eq. (10.64) with that obtained by solving the Riccati equation in Eq. (10.39). Because the two expressions are solutions of the same problem

and must hold for any state vector, we conclude that they must be equal and that $M(k+1)$ is identical to $S(k+1)$. Hence, the solution of the Riccati equation is given in terms of the state-transition matrix of the Hamiltonian system by Eq. (10.59). We can now write

$$S(k) = \{\phi_{21}(k_f - k) + \phi_{22}(k_f - k)S(k_f)\}\{\phi_{11}(k_f - k) + \phi_{12}(k_f - k)S(k_f)\}^{-1} \quad (10.65)$$

It may seem that the expression of Eq. (10.65) is always preferable to direct solution of the Riccati equation because it eliminates the nonlinearity. However, because the dimension of the Hamiltonian system is $2n$, the equation doubles the dimension of the problem. In many cases, it is preferable to solve the Riccati equation in spite of its nonlinearity.

Example 10.8

Consider the double integrator described in Example 10.4 with a sampling period $T = 0.02$ s, terminal weight $S(10) = \text{diag}\{10, 1\}$, $Q = \text{diag}\{10, 1\}$, and control weight $r = 0.1$. Obtain the Hamiltonian system and use it to obtain the optimal control.

Solution

The Hamiltonian system is

$$\begin{aligned} \begin{bmatrix} \mathbf{x}^*(k) \\ \lambda^*(k) \end{bmatrix} &= H \begin{bmatrix} \mathbf{x}^*(k+1) \\ \lambda^*(k+1) \end{bmatrix}, \quad k = k_0, \dots, k_f - 1 \\ H &= \begin{bmatrix} A^{-1} & A^{-1}B^*R^{-1}B^T \\ QA^{-1} & A^T + QA^{-1}B^*R^{-1}B^T \end{bmatrix} = \left[\begin{array}{cc|cc} 1 & -0.02 & 0 & 0 \\ 0 & 1 & 0 & 0.004 \\ \hline 10 & -0.2 & 1 & -0.0004 \\ 0 & 1 & 0.02 & 1.004 \end{array} \right] \\ \lambda^*(10) &= \text{diag}\{10, 1\}\mathbf{x}(10) \end{aligned}$$

Because the dynamics describe the evolution backward in time, each multiplication of the current state by the matrix H yields the state at an earlier time. We define a backward-time variable to use in computing as

$$k_b = k - (k_f - 1)$$

Proceeding backward in time, k_b starts with the value zero at $k = k_{f-1}$, then increases to $k_b = k_{f-1}$ with $k = 0$. We solve for the transition matrices and the gains using the following MATLAB program:

```
% hamilton
% Form the Hamiltonian system for backward dynamics
[n,n] = size(a); % Order of state matrix
n2 = 2*n; % Order of Hamiltonian matrix
kf = 11; % Final time
q = diag([10,1]); r = 0 = .1; % Weight matrices
```

Example 10.8—cont'd

```

sf = q; % Final state weight matrix

% Calculate the backward in time Hamiltonian matrix

a1 = inv(a);

He = b/r*b';

H3 = q/a;

H = [a1,a1'*He; H3,a'+H3*He]; % Hamiltonian matrix

fi = H; % Initialize the state-transition matrix

% i is the backward time variable kb = k-(kf-1), k = discrete time

for i = 1:kf-1

    fi11 = fi(1:n,1:n); % Partition the state-transition matrix

    fi12 = fi(1:n,n+1:n2);

    fi21 = fi(n+1:n2,1:n);

    fi22 = fi(n+1:n2,n+1:n2);

    s = (fi21 + fi22*s)/(fi11 + fi12*s); % Compute the Riccati equation. % solution

    K = (r + b'*s*b)\b'*s*a % Calculate the gains

    fi = H*fi; % Update the state-transition matrix

end

```

The optimal control gains are given in [Table 10.2](#). Almost identical results are obtained using the program described in Example 10.4. Small computational errors result in small differences between the results of the two programs.

Table 10.2: Optimal Gains for the Integrator with a Planning Horizon $k_f = 10$.

Time	Optimal gain	Vector K
0	2.0950	2.1507
1	1.7717	1.9632
2	1.4656	1.7730
3	1.1803	1.5804
4	0.9192	1.3861
5	0.6855	1.1903
6	0.4823	0.9935
7	0.3121	0.7958
8	0.1771	0.5976
9	0.0793	0.3988

10.5.1 Eigenstructure of the Hamiltonian matrix

The costate Eq. (10.31) can be rewritten as

$$\lambda^*(k+1) = -A^{-T}Q\mathbf{x}^*(k) + A^{-T}\lambda^*(k)$$

Substituting in the state Eq. (10.33), we write it in terms of the variables at time k in the form

$$\mathbf{x}^*(k+1) = [A + BR^{-1}B^TA^{-T}Q]\mathbf{x}^*(k) - BR^{-1}B^TA^{-T}\lambda^*(k)$$

We now write the forward expression for the state and Lagrange multiplier vectors

$$\begin{bmatrix} \mathbf{x}^*(k+1) \\ \lambda^*(k+1) \end{bmatrix} = \begin{bmatrix} A + BR^{-1}B^TA^{-T}Q & -BR^{-1}B^TA^{-T} \\ -A^{-T}Q & A^{-T} \end{bmatrix} \begin{bmatrix} \mathbf{x}^*(k) \\ \lambda^*(k) \end{bmatrix} \quad (10.66)$$

Clearly, based on the backward recursion of Eq. (10.59), the inverse of the Hamiltonian matrix is

$$H^{-1} = \begin{bmatrix} A + BR^{-1}B^TA^{-T}Q & -BR^{-1}B^TA^{-T} \\ -A^{-T}Q & A^{-T} \end{bmatrix} \quad (10.67)$$

In fact, the same matrix is obtained more laboriously by directly inverting the Hamiltonian matrix using the formula for the inversion of a partitioned matrix given in Appendix III.

Transposing the matrix gives

$$H^{-T} = \begin{bmatrix} A^T + QA^{-1}BR^{-1}B^T & -QA^{-1} \\ -A^{-1}BR^{-1}B^T & A^{-1} \end{bmatrix}$$

We observe that the entries of the matrix are the same as those of the Hamiltonian matrix with some rearrangement and sign changes. Hence, we can rearrange terms in the Hamiltonian matrix and change signs to obtain the equality

$$JH = H^{-T}J \quad (10.68)$$

where J is the matrix

$$J = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix} \quad (10.69)$$

By multiplying J and its transpose to obtain the identity matrix, we can verify that the inverse of the matrix J is

$$J^{-1} = J^T = -J \quad (10.70)$$

This allows us to write the inverse of the Hamiltonian matrix as

$$H^{-1} = (JHJ^{-1})^T = JHJ^T$$

[Eq. \(10.68\)](#) allows us to examine the properties of the eigenvalues and eigenvectors of the Hamiltonian matrix. Multiplying by any right eigenvector gives

$$JH\mathbf{v}_i = \mu_i J\mathbf{v}_i = H^{-T}J\mathbf{v}_i, i = 1, 2, \dots, 2n \quad (10.71)$$

where (μ_i, \mathbf{v}_i) are an eigenpair of the Hamiltonian matrix. Clearly, [Eq. \(10.71\)](#) shows that $J\mathbf{v}$ is a right eigenvector of H^{-T} corresponding to the same eigenvalue. Since the eigenvalues of any matrix are the same as those of its transpose, we conclude that the Hamiltonian matrix H and its inverse H^{-1} have the same eigenvalues. Because the eigenvalues of a matrix are the reciprocals of those of its inverse, we conclude that the eigenvalues of H can be written as

$$\left\{ \mu_i, \frac{1}{\mu_i}, i = 1, \dots, n \right\} \quad (10.72)$$

If we arrange the eigenvalues of the Hamiltonian matrix so that the first n are inside the unit circle, then the remaining n eigenvalues will be outside the unit circle. It turns out that the stable eigenvalues of the Hamiltonian matrix are also the closed-loop eigenvalues of the system with optimal control.

From [Eq. \(10.71\)](#), we also observe that $J\mathbf{v}_i, i = 1, 2, \dots, 2n$ are the eigenvectors of H^{-T} . Premultiplying [Eq. \(10.71\)](#) by H^T then transposing gives

$$\mu_i \mathbf{v}_i^T J^T H = \mathbf{v}_i^T J^T, i = 1, 2, \dots, 2n$$

Using [Eq. \(10.70\)](#), this can be rewritten as

$$\mathbf{v}_i^T J H = \frac{1}{\mu_i} \mathbf{v}_i^T J, i = 1, 2, \dots, 2n$$

Thus, the left eigenvectors of H are

$$\{\mathbf{v}_i^T J, i = 1, \dots, 2n\}$$

Note that if \mathbf{v}_i is a right eigenvector corresponding to a stable eigenvalue, then $\mathbf{v}_i^T J$ is a left eigenvector corresponding to an unstable eigenvalue, and vice versa.

The eigenstructure of the inverse Hamiltonian matrix governs the evolution of the state and costate vectors. Because the closed-loop optimal control system is stable in the steady-state provided that the stabilizability and detectability conditions are met, we expect the stable eigenvalues of the matrix to be related to the dynamics of the closed-loop optimal control system.

Theorem 10.2

The n stable eigenvalues of H^{-1} are the eigenvalues of the closed-loop optimal control state matrix $A_{cl} = A - BK$ with the optimal control of Eq. (10.48).

Proof

Let the initial state of the control system be given by the i th eigenvector \mathbf{v}_{xi} . Then the response of the system will remain on the eigenvector and will be of the form (see Chapter 7)

$$\mathbf{x}^*(k) = \mu_i^k \mathbf{v}_{xi},$$

and for time $k+1$

$$\mathbf{x}^*(k+1) = \mu_i^{k+1} \mathbf{v}_{xi}$$

This shows that $\mathbf{x}^*(k)$, and therefore \mathbf{v}_{xi} is an eigenvector of A_{cl} with eigenvalue μ_i for $i = 1, \dots, 2n$.

Using Eq. (10.36), we can write the costate as

$$\lambda^*(k) = S\mathbf{x}^*(k) = \mu_i^k S\mathbf{v}_{xi} = \mu_i^k \mathbf{v}_{\lambda i} \quad (10.73)$$

where $\mathbf{v}_{\lambda i} = S\mathbf{v}_{xi}$.

For time $k+1$, the costate is

$$\lambda^*(k+1) = \mu_i^{k+1} \mathbf{v}_{\lambda i}$$

The two equations at time $k+1$ combine to give

$$\begin{bmatrix} \mathbf{x}^*(k+1) \\ \lambda^*(k+1) \end{bmatrix} = \mu_i^{k+1} \begin{bmatrix} \mathbf{v}_{xi} \\ \mathbf{v}_{\lambda i} \end{bmatrix}$$

From Eq. (10.66), we also have

$$\begin{aligned} \begin{bmatrix} \mathbf{x}^*(k+1) \\ \lambda^*(k+1) \end{bmatrix} &= H^{-1} \begin{bmatrix} \mathbf{x}^*(k) \\ \lambda^*(k) \end{bmatrix} \\ &= H^{-1} \mu_i^k \begin{bmatrix} \mathbf{v}_{xi} \\ \mathbf{v}_{\lambda i} \end{bmatrix} = \mu_i^{k+1} \begin{bmatrix} \mathbf{v}_{xi} \\ \mathbf{v}_{\lambda i} \end{bmatrix} \end{aligned}$$

Because this expression holds for any Hamiltonian system, we now have

$$H^{-1} \begin{bmatrix} \mathbf{v}_{xi} \\ \mathbf{v}_{\lambda i} \end{bmatrix} = \mu_i \begin{bmatrix} \mathbf{v}_{xi} \\ \mathbf{v}_{\lambda i} \end{bmatrix}$$

Proof—cont'd

In addition, multiplying by H gives

$$H \begin{bmatrix} \mathbf{v}_{xi} \\ \mathbf{v}_{\lambda i} \end{bmatrix} = \frac{1}{\mu_i} \begin{bmatrix} \mathbf{v}_{xi} \\ \mathbf{v}_{\lambda i} \end{bmatrix}$$

so that the eigenvectors of the eigenvalue μ_i are also those of the eigenvalue $1/\mu_i$, $i = 1, 2, \dots, n$, that is, each stable eigenvalue has the same eigenvector as the reciprocal unstable eigenvalue of the Hamiltonian matrix.

The eigenvector relations Eq. (10.73) of the proof allow us to write the following expression for the solution of the algebraic Riccati equation

$$V_\lambda = SV_x$$

where the eigenvector matrices are given by

$$\begin{aligned} V_x &= [\mathbf{v}_{x1} \quad \cdots \quad \mathbf{v}_{xn}] \\ V_\lambda &= [\mathbf{v}_{\lambda 1} \quad \cdots \quad \mathbf{v}_{\lambda n}] \end{aligned}$$

Thus, the solution of the algebraic Riccati equation is given in terms of the stable eigenvectors of the inverse Hamiltonian matrix as

$$S = V_\lambda V_x^{-1} \quad (10.74)$$

We return to the optimal feedback gain of Eq. (10.48) and multiply both sides by $[R + B^T S B]$

$$[R + B^T S B] K = B^T S A$$

Rearranging terms gives

$$RK = B^T S (A - BK)$$

Thus, we can rewrite the gain in terms of the closed-loop state matrix

$$K = R^{-1} B^T S A_{cl}$$

We then substitute from Eq. (10.74) for S and for A_{cl} in terms of its eigenstructure

$$\begin{aligned} K &= R^{-1} B^T V_\lambda V_x^{-1} \times V_x M V_x^{-1} \\ A_{cl} &= V_x M V_x^{-1}, M = \text{diag}\{\mu_1, \dots, \mu_n\} \end{aligned}$$

With the stabilizability and detectability assumptions that guarantee that the optimal control stabilizes the system, we now have an expression for the optimal gain in terms of the stable eigenpairs of the inverse Hamiltonian matrix

$$K = R^{-1} B^T V_\lambda M V_x^{-1} \quad (10.75)$$

From Chapter 9, we know that the feedback gain that assigns a given set of eigenvalues is unique for a single-input system but not multiinput systems. In the latter case, there is some freedom in choosing the closed-loop eigenvectors. The single-input case allows us to examine the effect of the weighting matrices on the locations of the closed-loop poles.

Example 10.9

Obtain the eigenstructure of the inverse Hamiltonian system and use it to obtain the optimal control for the steady-state regulator for the double integrator system of Example 10.4.

The discretized double integrator system has the state equation

$$\mathbf{x}(k+1) = \begin{bmatrix} 1 & T \\ 0 & 1 \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} T^2/2 \\ T \end{bmatrix} u(k)$$

The inverse of the state matrix is

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

The inverse Hamiltonian matrix is

$$H^{-1} = \begin{bmatrix} A + BR^{-1}B^TA^{-T}Q & -BR^{-1}B^TA^{-T} \\ -A^{-T}Q & A^{-T} \end{bmatrix} = \begin{bmatrix} 1 & 0.02 & -4 \times 10^{-7} & -3.999 \times 10^{-5} \\ 4 \times 10^{-4} & 1.004 & -4 \times 10^{-5} & -4 \times 10^{-3} \\ -10 & 0 & 1 & 0 \\ 0.2 & -1 & -0.02 & 1 \end{bmatrix}$$

Using MATLAB, we obtain the eigenvalues $\{1.0558 \pm j0.03217, 0.9462 \pm 0.0310\}$, of which the latter two are inside the unit circle. The corresponding eigenvectors form the matrix

$$V = \begin{bmatrix} V_x \\ V_\lambda \end{bmatrix} = \begin{bmatrix} 5.1155 \times 10^{-3} - j2.9492 \times 10^{-3} & 5.1155 \times 10^{-3} + j2.9492 \times 10^{-3} \\ -9.1786 \times 10^{-3} + j16.4508 \times 10^{-3} & -9.1786 \times 10^{-3} - j16.4508 \times 10^{-3} \\ \hline 0.95087 & 0.95087 \\ -0.014069 + j0.3086 & -0.014069 - j0.3086 \end{bmatrix}$$

Finally, we have the optimal feedback gain matrix

$$K = R^{-1}B^TV_\lambda MV_x^{-1} = [8.5862 \quad 5.088]$$

10.6 Return difference equality and stability margins

We consider a linear time-invariant system and a linear quadratic performance measure with constant weighting matrices. In the steady state, the Riccati equation in Joseph's form becomes

$$S = (A - BK)^T S (A - BK) + K^T R K + Q \quad (10.76)$$

We expand the equation and regroup terms to write it as

$$S - A^T S A + K^T (B^T S B + R) K - K^T B^T S A - A^T S B K = Q \quad (10.77)$$

We first show that

i) $S - A^T S A = (z^{-1} I_n - A)^T S (zI_n - A) + (z^{-1} I_n - A)^T S A + A^T S (zI_n - A)$

Expand the RHS to obtain

$$S - zA^T S - z^{-1} S A + A^T A + z^{-1} S A - A^T S A + zA^T S - A^T S A = RHS$$

ii) $K^T (B^T S B + R) K = K^T B^T S A = A^T S B K$

Substitute using the optimal gain

$$(B^T S B + R) K = B^T S A$$

then premultiply by K^T to obtain the first equality.

Transpose to obtain

$$K^T (B^T S B + R) = A^T S B$$

then postmultiply by K to obtain the second equality.

We substitute the two equalities in the Riccati equation to rewrite it as

$$(z^{-1} I_n - A)^T S (zI_n - A) + (z^{-1} I_n - A)^T S A + A^T S (zI_n - A) + K^T (B^T S B + R) K = Q$$

Premultiplying by $B^T (z^{-1} I_n - A)^{-T}$ and postmultiplying by $(zI_n - A)^{-1} B$ then adding R , we rewrite the LHS as

$$\begin{aligned} & B^T S B + B^T S A (zI_n - A)^{-1} B + B^T (z^{-1} I_n - A)^{-T} A^T S B \\ & + B^T (z^{-1} I_n - A)^{-T} K^T (B^T S B + R) K (zI_n - A)^{-1} B + R \\ & = (B^T S B + R) + (B^T S B + R) K (zI_n - A)^{-1} B + B^T (z^{-1} I_n - A)^{-T} K^T (B^T S B + R) \\ & + B^T (z^{-1} I_n - A)^{-T} K^T (B^T S B + R) K (zI_n - A)^{-1} B \\ & = \left(I + K (zI_n - A)^{-1} B \right)^T (B^T S B + R) \left(I + K (zI_n - A)^{-1} B \right) \end{aligned}$$

We now have the return difference equality

$$\begin{aligned} & \left(I + K (z^{-1} I_n - A)^{-1} B \right)^T (B^T S B + R) \left(I + K (zI_n - A)^{-1} B \right) \\ & = R + B^T (z^{-1} I_n - A)^{-T} Q (zI_n - A)^{-1} B \end{aligned} \quad (10.78)$$

The name of the inequality refers to the return difference $I_n - \left[-K(zI_n - A)^{-1}B \right]$, where the negative sign before the return difference transfer function $K(zI_n - A)^{-1}B$ from the plant input to control is due to negative feedback.

Next, we relate the return difference to the closed-loop characteristic polynomial with LQR control. The closed-loop characteristic polynomial is

$$\begin{aligned} p_{cl}(z) &= \det[zI_n - A + BK] = \det\left[I_n + BK(zI_n - A)^{-1}\right]\det[zI_n - A] \\ &= \det\left[I_m + K(zI_n - A)^{-1}B\right]\det[zI_n - A] \end{aligned} \quad (10.79)$$

while the open-loop characteristic polynomial is

$$p_{ol}(z) = \det[zI_n - A] \quad (10.80)$$

This shows that the two characteristic polynomials are related by

$$p_{cl}(z) = \det\left[I_m + K(zI_n - A)^{-1}B\right]p_{ol}(z) \quad (10.81)$$

We now consider the determinant of the return difference equality using the determinant properties

$$\begin{aligned} \det[A] &= \det[A^T] = \frac{1}{\det[A^{-1}]} = \frac{1}{\det[A^{-T}]} \\ \det[AB] &= \det[A]\det[B] \end{aligned} \quad (10.82)$$

The determinant of the return difference equality yields

$$\begin{aligned} &\det\left[\left(I_n + K(zI_n - A)^{-1}B\right)^T\right]\det[B^T S B + R]\det\left[I + K(zI_n - A)^{-1}B\right] \\ &= \frac{p_{cl}(z)p_{cl}(z^{-1})}{p_{ol}(z)p_{ol}(z^{-1})}\det[B^T S B + R] \\ &= \det\left[R + B^T(z^{-1}I_n - A)^{-T}Q(zI_n - A)^{-1}B\right] \end{aligned}$$

The closed-loop characteristic polynomial satisfies

$$p_{cl}(z)p_{cl}(z^{-1}) = p_{ol}(z)p_{ol}(z^{-1}) \frac{\det\left[R + B^T(z^{-1}I_n - A)^{-T}Q(zI_n - A)^{-1}B\right]}{\det[B^T S B + R]} \quad (10.83)$$

Because the denominator does not include z terms, the closed-loop poles can be determined without solving the algebraic Riccati equation for S by finding the root of the numerator. Thus, we only need to consider the polynomial

$$\begin{aligned}
& p_{ol}(z)p_{ol}(z^{-1}) \det \left[R + (z^{-1}I_n - A)^{-T} Q(zI_n - A)^{-1} \right] \\
&= \det \left[(z^{-1}I_n - A)^T \right] \det \left[R + B^T (z^{-1}I_n - A)^{-T} Q(zI_n - A)^{-1} B \right] \det[zI_n - A] \\
&= \det \left[(z^{-1}I_n - A)^T \right] \left\{ \det[R] + \det \left[B^T (z^{-1}I_n - A)^{-T} Q(zI_n - A)^{-1} B \right] \right\} \det[zI_n - A]
\end{aligned}$$

Note that the second term between braces is

$$\det \left[(z^{-1}I_n - A)^T \right] \det[B^T] \det \left[(z^{-1}I_n - A)^{-T} Q(zI_n - A)^{-1} \right] \det[B] \det[zI_n - A]$$

Cancellation then gives the equality

$$p_{cl}(z)p_{cl}(z^{-1}) = \det \left[(z^{-1}I_n - A)^T R(zI_n - A) + Q \right] \quad (10.84)$$

The roots of the polynomial are symmetric with respect to the unit circle such that if z_0 is a root then so is its reciprocal z_0^{-1} . The roots inside the unit circle are the closed-loop poles by virtue of the stability of the LQR solution.

While it is clear that LQR yields a stable solution, its relative stability should also be considered. To do so, we need to examine its gain margin and phase margin. We limit this to SISO case where the return difference in the frequency domain is

$$1 + \mathbf{k}^T (e^{j\omega T} I_n - A)^{-1} \mathbf{b}$$

and the return difference equality is

$$\left| 1 + \mathbf{k}^T (e^{j\omega T} I_n - A)^{-1} \mathbf{b} \right|^2 = \frac{r + \mathbf{b}^T (e^{-j\omega T} I_n - A)^{-T} Q(e^{j\omega T} I_n - A)^{-1} \mathbf{b}}{\mathbf{b}^T S \mathbf{b} + r} \geq \frac{r}{\mathbf{b}^T S \mathbf{b} + r}$$

The square root gives the inequality

$$\left| 1 + \mathbf{k}^T (e^{j\omega T} I_n - A)^{-1} \mathbf{b} \right| \geq \sqrt{\frac{r}{\mathbf{b}^T S \mathbf{b} + r}} = k_s \quad (10.85)$$

We can also express k_s more succinctly in terms of the ratio $\alpha = \mathbf{b}^T S \mathbf{b} / r$ as

$$k_s = 1 / \sqrt{1 + \alpha}$$

Because S is positive definite, $\mathbf{b}^T S \mathbf{b} > 0$, $\alpha > 0$, and $k_s < 1$ and it goes to zero as $r \rightarrow 0$. Hence, it lies in the range

$$0 < k_s < 1 \quad (10.86)$$

The Nyquist criterion tells us that $\mathbf{k}^T (e^{j\omega T} I_n - A)^{-1} \mathbf{b}$ is stable if its Nyquist plot does not encircle the origin, or equivalently, the plot of $1 + \mathbf{k}^T (e^{j\omega T} I_n - A)^{-1} \mathbf{b}$ does not encircle the point $(-1, 0)$. The latter plot is given by the tips of the vectors $\mathbf{k}^T (e^{j\omega T} I_n - A)^{-1} \mathbf{b}$ with

their tails at the point $(-1, 0)$. However, the lower bound on $|1 + \mathbf{k}^T (e^{j\omega T} I_n - A)^{-1} \mathbf{b}|$ indicates that the vectors remain outside a circle of radius k_s . Therefore, the plot of $\mathbf{k}^T (e^{j\omega T} I_n - A)^{-1} \mathbf{b}$ lies outside a circle of radius $k_s < 1$ centered at the origin as shown in Fig. 10.16. The Nyquist plot hits the point $(-1, 0)$ if scaled by a factor outside the range

$$(GM_L, GM_H) = \left(\frac{1}{1+k_s}, \frac{1}{1-k_s} \right) \quad (10.87)$$

To express the gain margins in terms of α , we observe that

$$\frac{1}{1+k_s} = \frac{\sqrt{1+\alpha}}{1+\sqrt{1+\alpha}} = \frac{\sqrt{1+\alpha}(\sqrt{1+\alpha}-1)}{\alpha} = \frac{1+\alpha-\sqrt{1+\alpha}}{\alpha}$$

In terms of α , the gain margins are

$$(GM_L, GM_H) = \left(\frac{1+\alpha-\sqrt{1+\alpha}}{\alpha}, \frac{1+\alpha+\sqrt{1+\alpha}}{\alpha} \right) \quad (10.88)$$

This range defines an upper and lower gain margin for the system. Thus, the system becomes unstable if subjected to again perturbation of magnitude greater than GM_H or magnitude less than GM_L .

For the phase margin, we need to find the coordinates of the intersection of the two circles shown in Fig. 10.17. The intersection is at the point whose real axis coordinate is $k_s^2/2 - 1$

The phase margin PM is the phase lag that makes the plot intersect the point $(-1, 0)$, that is the angle (the proof is left as an exercise)

$$PM = \pm \cos^{-1}(|x|) = \pm \cos^{-1}\left(1 - k_s^2/2\right) \quad (10.89)$$

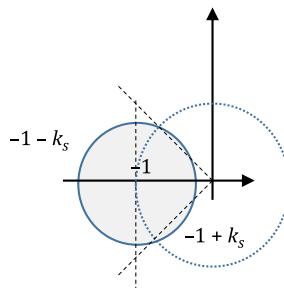
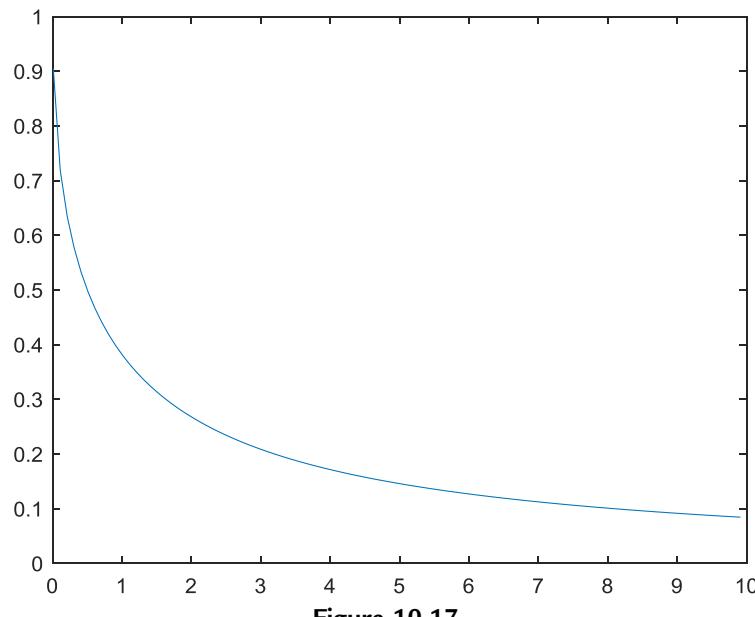


Figure 10.16
Stability region for the frequency response $k^T (e^{j\omega T} I_n - A)^{-1} b$.

**Figure 10.17**

Plot of closed-loop pole location for optimal control with $\alpha = 1$, $r = 1$.

In terms of α , the phase margin is

$$PM = \pm \cos^{-1} \left(1 - \frac{1}{2(1 + \alpha)} \right) = \pm \cos^{-1} \left(\frac{1}{2} \times \frac{1 + 2\alpha}{1 + \alpha} \right) \quad (10.90)$$

Using trigonometric identities, it can be shown that another formula for the phase margin is

$$PM = \pm 2 \sin^{-1} \left(\frac{k_s}{2} \right) = \pm 2 \sin^{-1} \left(\frac{1}{2\sqrt{1 + \alpha}} \right) \quad (10.91)$$

We leave the proof as an exercise.

The stability margins show the robustness of LQR when using state feedback. Unfortunately, the robustness of LQR deteriorates if state estimator feedback is used.

Example 10.10

Consider the scalar system

$$x(k+1) = ax(k) + bu(k)$$

Example 10.10—cont'd

(i) Solve the algebraic Riccati equation for the system to show that

$$\alpha = \frac{b^2 s}{r} = \frac{1}{2} \left\{ a^2 + \beta - 1 + \sqrt{\left((a-1)^2 + \beta \right) \left((a+1)^2 + \beta \right)} \right\}$$

$$\beta = b^2 q / r$$

- (ii) Obtain an expression for the return difference equality and the closed-loop poles. For $a = 1$, $r = 1$, plot the stable pole location versus q as q varies in the range [0.01, 10] and comment on the results.
- (iii) Plot the gain margins and the phase margins for the system with LQR control as a function of $\alpha = b^2 s / r$ and discuss the effect of the relative values of q and r on relative stability.

Solution

(i) The algebraic Riccati equation for this simple system is

$$s = a^2 \left\{ s - \frac{s^2 b^2}{b^2 s + r} \right\} + q$$

The positive solution of the quadratic equation is

$$s = \frac{a^2 r + b^2 q - r + \sqrt{\left((a-1)^2 r + b^2 q \right) \left((a+1)^2 r + b^2 q \right)}}{2b^2}$$

$$\alpha = \frac{b^2 s}{r} = \frac{1}{2} \left\{ a^2 + \beta - 1 + \sqrt{\left((a-1)^2 + \beta \right) \left((a+1)^2 + \beta \right)} \right\}$$

$$\beta = b^2 q / r$$

(ii) The return difference equality for the scalar system is

$$(b^2 s + r) \left(1 + \frac{kb}{z^{-1} - a} \right) \left(1 + \frac{kb}{z - a} \right) = r + \frac{b^2 q}{(z^{-1} - a)(z - a)}$$

The closed-loop equality is

$$p_{cl}(z)p_{cl}(z^{-1}) = r(z^{-1} - a)(z - a) + q$$

This gives the quadratic polynomial

$$az^2 - [a^2 + 1 + q/r]z + a$$

The roots of the polynomial are

$$p_{1,2} = \frac{1}{2a} \left\{ a^2 + 1 + \frac{q}{r} \pm \sqrt{\left(a^2 + 1 + \frac{q}{r} \right)^2 - 4a^2} \right\}$$

Example 10.10—cont'd

The product of the two roots is

$$\frac{1}{4a^2} \left\{ \left(a^2 + 1 + \frac{q}{r} \right)^2 - \left[\left(a^2 + 1 + \frac{q}{r} \right)^2 - 4a^2 \right] \right\} = 1$$

The closed-loop pole of the optimal control system is the stable pole

$$p_{1,2} = \frac{1}{2a} \left\{ a^2 + 1 + \frac{q}{r} - \sqrt{\left(a^2 + 1 + \frac{q}{r} \right)^2 - 4a^2} \right\}$$

Using MATLAB, we obtain the plot of Fig. 10.17. The plot shows that at the ratio q/r approaches zero the control action is severely limited and the system approaches its open-loop behavior. As the ratio becomes large, the control action is free to place the pole arbitrarily close to the origin and the pole approaches zero in the limit as $q/r \rightarrow \infty$.

- (iii) Using MATLAB, we obtain plots of the gain and phase margins Figs. 10.18 and 10.19, respectively). The plots show that the margins decrease as the value of α increases. The phase margins approach $\pm 60^\circ$ and the gain margins approach $(0.5, \infty)$ as $\alpha \rightarrow 0$. As α becomes very large, the system approaches the stability boundary. Part (i) shows that α increases with $\beta = b^2 q/r$, that is with the ratio of the weights q/r .

10.7 Model predictive control

MPC is a discrete-time optimal control strategy that exploits knowledge of a system model to calculate the control law to minimize a given cost function. Different MPC algorithms can be implemented by considering different classes of models and different cost functions. The main components of a model predictive controller are as follows:

- A model of the system to predict its output over a future time interval known as the **prediction horizon**.
- An optimizer to calculate the future control sequence by minimizing an objective function.
- A receding strategy to only use the first element of the calculated control sequence and shift the prediction horizon by one step at each sampling point.

The application of an MPC is based on the following steps:

1. Minimize a cost function over a suitable prediction horizon using a mathematical model of the system to obtain the optimal sequence of control inputs. The cost function usually involves the weighted sum of the tracking errors and of the control efforts and may be subject to constraints on the control variable and/or on the system output.

-
2. Use the first value of the control sequence and discard the rest, then shift the prediction horizon forward by one step.
 3. Go to Step 1.

MPC exploits the model of the system and the reference input to optimize performance by explicitly considering the system constraints and compensating for model estimation errors and disturbances by means of the receding horizon strategy.

10.7.1 Model

MPC requires a predictive model. The model can be in different forms; it can be based on impulse response coefficients (see Eq. (2.27)), on step response coefficients (see Eq. (2.28)), which are more easily obtained experimentally especially in process control applications, on a transfer function Eq. (2.31) or on a state space description Eq. (7.52). A step response model is easier to obtain than a transfer function or a state-space model but has the disadvantage that a larger number of parameters must be stored in the controller memory. Further, impulse response or step response models cannot be used for unstable processes.

Consider a single-input-single-output model based on step response coefficients. The system output at time $k + i$ predicted at time k is written as

$$\hat{y}(k+i|k) = \sum_{j=1}^N g(j)\Delta u(k+i-j|k) \quad (10.92)$$

where $\Delta u(k) = u(k) - u(k-1)$ and N is the number of coefficients of the step response. Although the step response has an infinite number of coefficients, it is sufficient to truncate the response and only use N coefficients in practice provided that $g(N) - g(N-1) \approx 0$.

10.7.2 Cost function

MPC optimizes the trade-off between tracking performance and control effort. Hence, the cost function to be minimized includes two elements, one related to the tracking error and one related to the control variable. A general representation is

$$J(N_1, N_2, N_u) = \sum_{j=N_1}^{N_2} \delta(j) \{\hat{y}(k+j|k) - \mathbf{w}(k+j)\}^2 + \sum_{j=1}^{N_u} \lambda(j) \Delta \mathbf{u}(k+j-1)^2 \quad (10.93)$$

where

- N_1 and N_2 are the minimum and maximum prediction horizons and they represent the time interval in which the tracking performance is relevant (for example, if the system has a time delay d , it is reasonable to set $N_1 = d$).
- N_u is the control horizon, that is, the number of future samples of the control signal to be calculated in the optimization procedure
- $\delta(j)$ and $\lambda(j)$ are weight sequences for the two parts of the cost function. These weights are usually constant or decrease exponentially.
- $w(k)$ is the reference signal. We again observe how the future reference signal is considered in the controller, unlike standard feedback control.

10.7.3 Computation of the control law

At each sampling point, the control law is computed by calculating the predicted output as a function of past input and output values and future control values using the prediction model. The values of $\hat{y}(k+i|k)$ are then substituted in the cost function, which is minimized (by means of a suitable optimization procedure) with respect to Δu . Finally, the first value of the vector Δu is applied and the procedure is repeated at the next sampling point. Note that the predicted output is given by the superposition of the free response (i.e., the response obtained when the future control actions are zero) and of the forced response (i.e., the response obtained from the future control actions and zero initial conditions). In fact, the optimization procedure modifies the future control actions to obtain the desired response. For example, if the free response is satisfactory, then the future control actions are equal to zero.

10.7.4 Constraints

A major advantage of MPC is that it allows us to include constraints in the optimization problem to be solved to determine future control inputs. These constraints are usually related to the control input u , to the control input variation (slew rate) Δu , or to the system output y . For example, constraints may be added to limit the magnitudes of overshoots or undershoots. To obtain the solution of the optimization problem, all these constraints must be expressed as a function of the control inputs. For example, if the system output is constrained between y_{min} and y_{max} ,

$$y_{min} \leq y \leq y_{max}$$

we substitute for the output

$$y_{min} \leq Gu + f \leq y_{max}$$

and write the constraints as a function of \mathbf{u}

$$\mathbf{y}_{min} - \mathbf{f} \leq G\mathbf{u} \leq \mathbf{y}_{max} - \mathbf{f}$$

where G is the model of the system and \mathbf{f} is the free response vector.

The constraints can, in general, be expressed in the form

$$R\mathbf{u} \leq \mathbf{c}$$

For example, the above constraints on the system output can be written as

$$R = \begin{bmatrix} G \\ G \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} \mathbf{y}_{max} - \mathbf{f} \\ \mathbf{f} - \mathbf{y}_{min} \end{bmatrix}$$

We now have an optimization problem where we minimize a quadratic function subject to linear constraints. This is known as a **quadratic programming** problem.

10.7.5 MATLAB commands

Consider the quadratic programming problem

$$\min_{\mathbf{u}} \frac{1}{2} \mathbf{u}^T H \mathbf{u} + \mathbf{b}^T \mathbf{u}$$

subject to

$$R\mathbf{u} \leq \mathbf{c}$$

The MATLAB command to solve the problem is.

`>> u = quadprog(H,b,R,c)`

10.8 Modification of the reference signal

For a control system to track set-point step signals $\mathbf{r}(k)$, the reference signal $\mathbf{w}(k)$ can be chosen to be equal to $\mathbf{r}(k)$, but it can also be conveniently recomputed at each sampling point to provide a feasible reference trajectory to the system. The reference trajectory is equal to the measured process output at the current time and then attains the set-point value as a first-order system response. For simplicity, we only consider a single set-point signal r associated with a scalar output y . Formally, the reference trajectory is

$$\begin{aligned} w(k) &= y(k) \\ w(k+i) &= \alpha w(k+i-1) + (1-\alpha)r(k+i) \end{aligned} \tag{10.94}$$

where $\alpha \in [0, 1]$ is a parameter that is selected to specify the speed of the convergence of the trajectory toward the set-point value, and therefore the speed of the response. An example of the calculation of $w(k)$ starting from a process output $y(5)$ for different values of α is shown in Fig. 10.18.

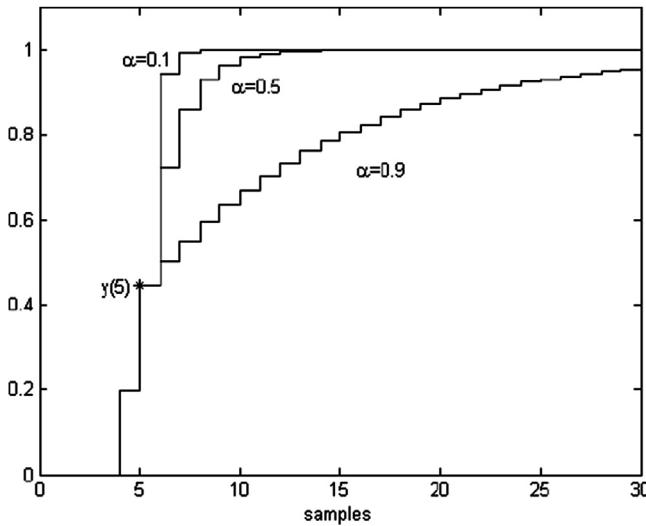


Figure 10.18
Reference signals $w_i(t)$ for different values of α (with $r_i = 1$).

10.8.1 Dynamic Matrix Control

As an illustrative example of MPC, we consider the application of Dynamic Matrix Control (DMC), which has been often applied in industry, to a single-input–single-output process. The process is assumed to be asymptotically stable and the prediction model is determined based on the step response coefficients. A disturbance term \hat{n} is added resulting in

$$\hat{y}(k+i|k) = \sum_{j=1}^{\infty} g(j)\Delta u(k+i-j|k) + \hat{n}(k+i|k) \quad (10.95)$$

This can be rewritten as

$$\hat{y}(k+i|k) = \sum_{j=1}^i g(i)\Delta u(k+i-j|k) + \sum_{j=i+1}^{\infty} g(i)\Delta u(k+i-j|k) + \hat{n}(k+i|k) \quad (10.96)$$

If the disturbance is constant along the prediction horizon, we have

$$\hat{n}(k+i|k) = \hat{n}(k|k) = y_m(k|k) - \hat{y}(k|k) \quad (10.97)$$

where $y_m(k)$ is the measured output.

By substituting (10.97) into (10.96) we obtain

$$\begin{aligned} \hat{y}(k+i|k) &= \sum_{j=1}^i g(j)\Delta u(k+i-j|k) + \sum_{j=i+1}^{\infty} g(j)\Delta u(k+i-j|k) + y_m(k) - \sum_{j=1}^{\infty} g(j)\Delta u(k-j|k) \\ &\quad \sum_{j=1}^i g(j)\Delta u(k+i-j|k) + f(k+i) \end{aligned} \quad (10.98)$$

where

$$f(k+i) \approx y(k) + \sum_{j=1}^{\infty} [g(i+j) - g(j)] \Delta u(k-j)$$

is the free response of the system, as this term depends only on past control actions. As expected, the predicted output is given by the superposition of the forced and free responses of the system.

For an asymptotically stable system, the coefficients of the step response become almost constant after N sampling periods, that is,

$$g(i+j) - g(j) \approx 0, \quad i > N$$

Thus, the free response can be written as

$$f(k+k) \approx y(k) + \sum_{j=1}^N [g(i+j) - g(j)] \Delta u(k-j)$$

The process outputs determined along the prediction horizon $k = 1, \dots, N_2$ (note that $N_1 = 1$) by using N_u control actions are therefore determined as

$$\begin{aligned} \hat{y}(k+1|k) &= g(1)\Delta u(k) + f(k+1) \\ \hat{y}(k+2|k) &= g(2)\Delta u(k) + g(1)\Delta u(k+1) + f(k+2) \\ &\vdots \\ \hat{y}(k+N_2|k) &= \sum_{j=N_2-N_u+1}^{N_2} g(j)\Delta u(k+N_2-j) + f(k+N_2) \end{aligned}$$

We define the **dynamic matrix** as

$$G = \begin{bmatrix} g(1) & 0 & \cdots & 0 \\ g(2) & g(1) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ g(N_u) & g(N_u-1) & \cdots & g(1) \\ g(N_2) & g(N_2-1) & \cdots & g(N_2-N_u+1) \end{bmatrix}$$

The vector of predicted outputs can now be written in compact form as

$$\mathbf{Y} = G\mathbf{u} + \mathbf{f} \quad (10.99)$$

where

$$\mathbf{Y} = \begin{bmatrix} \hat{y}(k+1|k) \\ \hat{y}(k+2|k) \\ \vdots \\ \hat{y}(k+N_2|k) \end{bmatrix}$$

$$\mathbf{u} = \begin{bmatrix} \Delta u(k) \\ \Delta u(k+1) \\ \vdots \\ \Delta u(k+N_u - 1) \end{bmatrix}$$

$$\mathbf{f} = \begin{bmatrix} f(k) \\ f(k+1) \\ \vdots \\ f(k+N_u - 1) \end{bmatrix}$$

The cost function is selected as Eq. (10.93) with $\delta = 1$ and a constant λ along the whole horizon yields

$$J = \mathbf{Y}^T \mathbf{Y} + \lambda \mathbf{u}^T \mathbf{u} \quad (10.100)$$

In the simple unconstrained case, the solution of the optimization problem (that is, the array \mathbf{u} that minimizes the cost function Eq. (10.100) is obtained by computing the derivative of J and equating it to 0. We first write an expression for the cost function in terms of the input

$$\begin{aligned} J &= (\hat{\mathbf{y}} - \mathbf{w})^T (\hat{\mathbf{y}} - \mathbf{w}) + \lambda \mathbf{u}^T \mathbf{u} \\ &= (\mathbf{G}\mathbf{u} + \mathbf{f} - \mathbf{w})^T (\mathbf{G}\mathbf{u} + \mathbf{f} - \mathbf{w}) + \lambda \mathbf{u}^T \mathbf{u} \\ &= \frac{1}{2} \mathbf{u}^T H \mathbf{u} + \mathbf{b}^T \mathbf{u} + f_0 \end{aligned}$$

where

$$\begin{aligned} H &= 2(G^T G + \lambda I) \\ \mathbf{b}^T &= 2(\mathbf{f} - \mathbf{w})^T G \\ f_0 &= (\mathbf{f} - \mathbf{w})^T (\mathbf{f} - \mathbf{w}) \end{aligned}$$

We differentiate to solve for the control

$$\frac{\partial J}{\partial \mathbf{u}} = 0 \Rightarrow \mathbf{u} = -H^{-1} \mathbf{b} = -(G^T G + \lambda I) G^T (\mathbf{f} - \mathbf{w})$$

which can be rewritten as

$$\mathbf{u} = -K(\mathbf{f} - \mathbf{w})$$

where

$$K = (G^T G + \lambda I) G^T$$

Because of the receding horizon concept, only the first element of the vector \mathbf{u} is applied as input to the process. The procedure is repeated to calculate the next control input.

Example 10.11

Apply DMC to the process of Example 6.21 with no constraints to obtain the closed-loop unit step response with $\alpha = 0$ and with $\lambda = 0.1, 1, 2$, and $N_1 = 1, N_2 = 10, N_u = 5$. Discuss the effect of varying λ on the results.

Solution

By examining the open-loop step response of the system of Example 6.21, we observe that we can let $N = 60$. Implementing the DMC method in a MATLAB program (see Problem 10.25), we obtain the results shown in Figs. 10.19–10.21. As expected, increasing λ reduces the control effort and increases the rise time. Thus, λ can be selected by the user to handle the trade-off between performance and control effort.

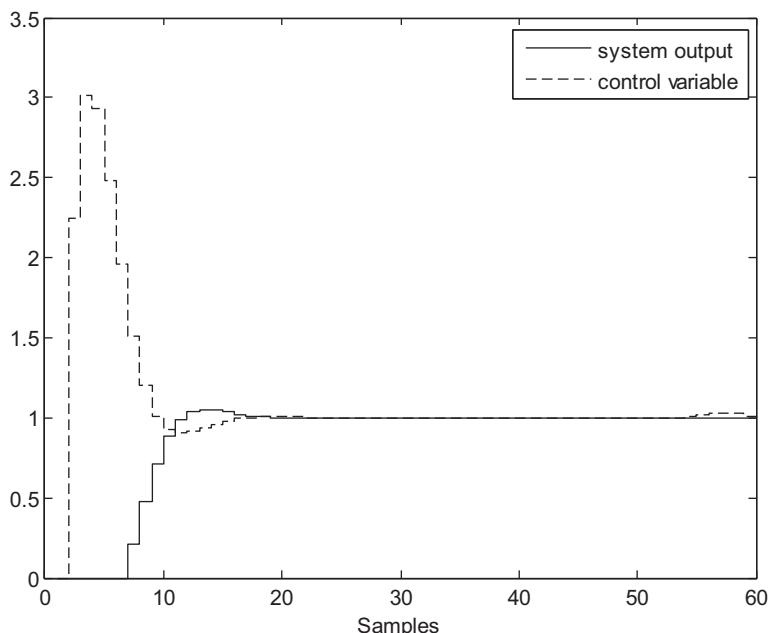


Figure 10.19
Closed-loop unit step response for $\lambda = 0.1$.

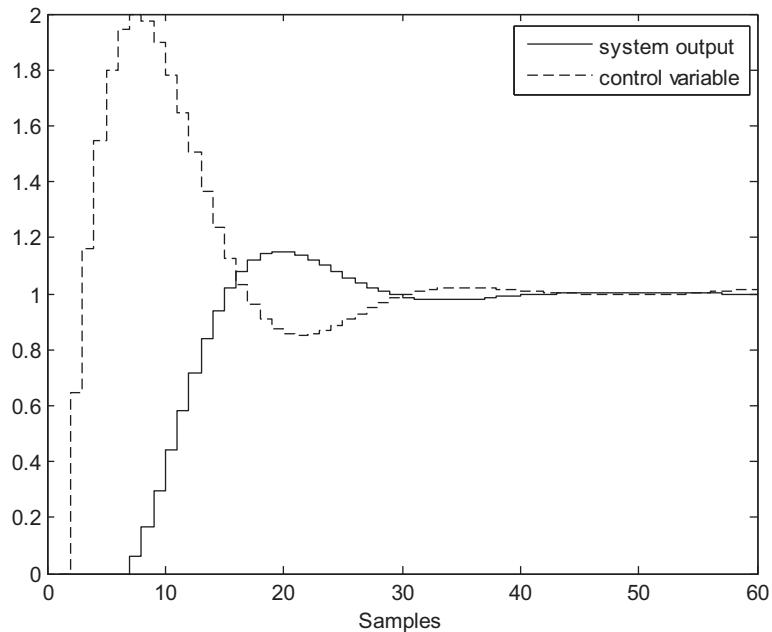
Example 10.11—cont'd

Figure 10.20
Closed-loop unit step response for $\lambda = 1$.

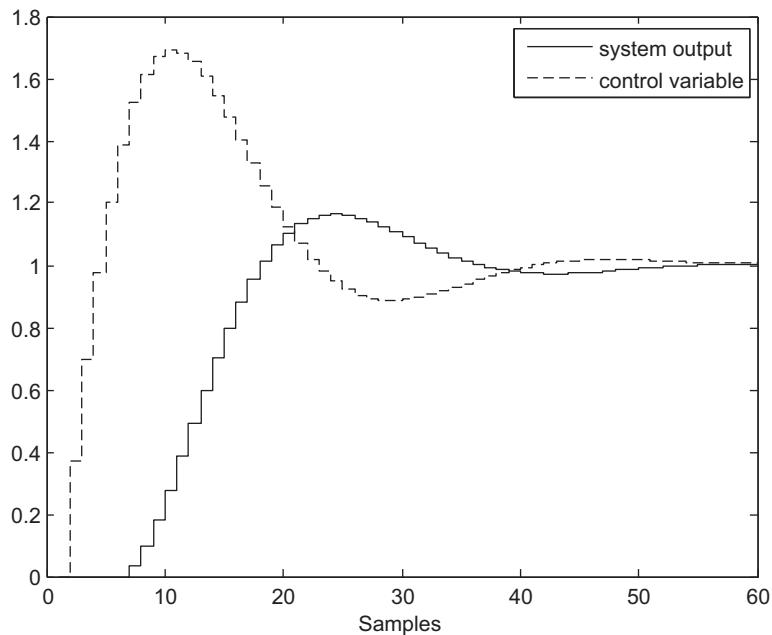


Figure 10.21
Closed-loop unit step response for $\lambda = 2$.

Example 10.12

Apply DMC to the process of Example 6.21 with no constraints and analyze the different results obtained for a closed-loop unit step response with $\lambda = 1$ and with different values of α , in particular $\alpha = 0.1$ and $\alpha = 0.5$. Consider $N_1 = 1$, $N_2 = 10$, $N_u = 5$.

Solution

Using the MATLAB program used for Example 10.10 with the two values of α , we obtain the results shown in Figs. 10.22 and 10.23. It can be seen that α allows the trade-off between performance and control effort. In fact, by increasing α the control effort is reduced, the rise time increases and the overshoot decreases.

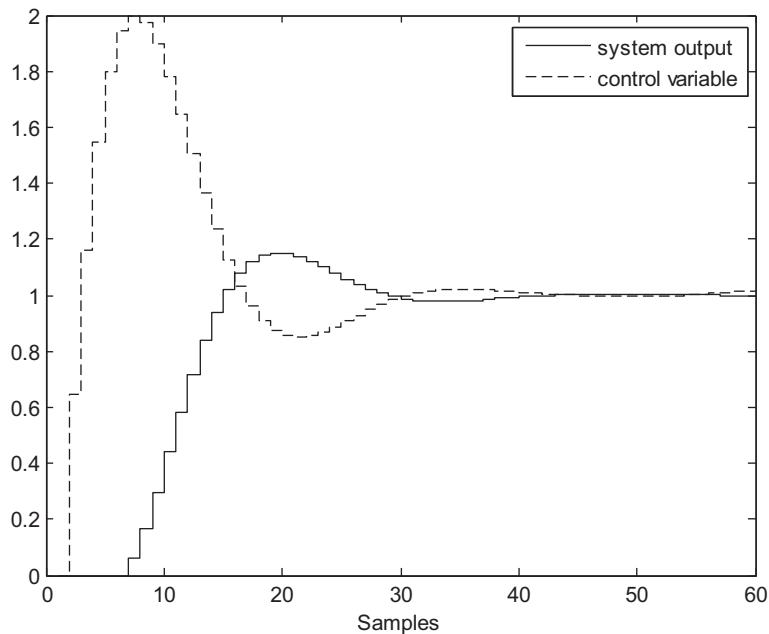


Figure 10.22
Closed-loop unit step response for $\alpha = 0.1$.

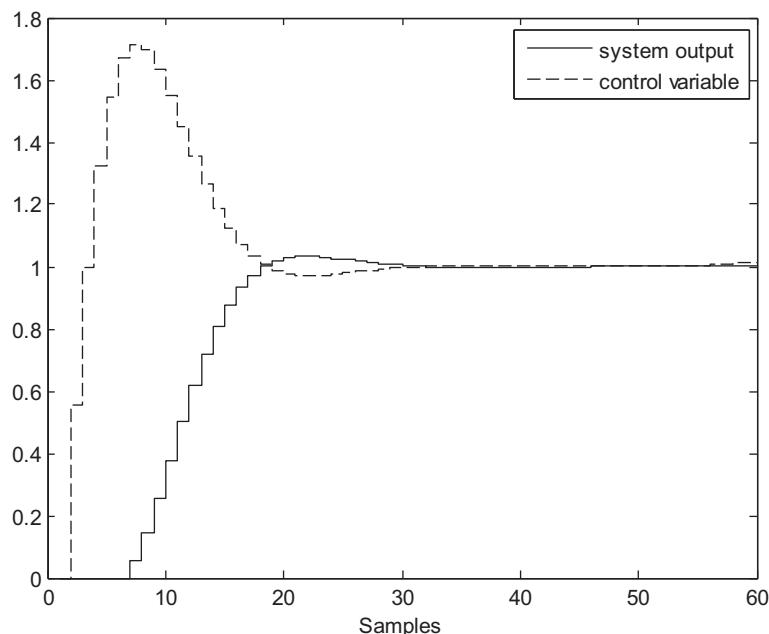
Example 10.12—cont'd

Figure 10.23
Closed-loop unit step response for $\alpha = 0.1$.

Further reading

- Anderson, B.D.O., Moore, J.B., 1990. Optimal Control: Linear Quadratic Methods. Prentice Hall, Englewood Cliffs, NJ.
- Camacho, E.F., Bordons, C., 2007. Model Predictive Control, second ed. Springer, London.
- Chong, E.D., Zak, S.H., 1996. An Introduction to Optimization. Wiley-Interscience, New York.
- Jacquot, R.G., 1981. Modern Digital Control Systems. Marcel Dekker, New York.
- Kwakernaak, H., Sivan, R., 1972. Linear Optimal Control Systems. Wiley-Interscience, New York.
- Lewis, F.L., Syrmos, V.L., 1995. Optimal Control. Wiley-Interscience, New York.
- Naidu, D.S., 2002. Optimal Control Systems. CRC Press, Boca Raton, FL.
- Rossiter, J.A., 2017. Model-Based Predictive Control: A Practical Approach. CRC Press, Boca Raton.
- Sage, A.P., White, C.C., 1977. Optimum Systems Control. Prentice-Hall, Englewood Cliffs, NJ.

Problems

- 10.1 Show that for a voltage source versus with source resistance R_s connected to a resistive load R_L , the maximum power transfer to the load occurs when $R_L = R_s$.

- 10.2 Let \mathbf{x} be an $n \times 1$ vector whose entries are the quantities produced by a manufacturer. The profit of the manufacturer is given by the quadratic form

$$J(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T P \mathbf{x} + \mathbf{q}^T \mathbf{x} + \mathbf{r}$$

where P is a negative definite symmetric matrix and \mathbf{q} and \mathbf{r} are constant vectors. Find the vector \mathbf{x} to maximize profit.

- With no constraints on the quantity produced
- If the quantity produced is constrained by

$$B\mathbf{x} = \mathbf{c}$$

where B is an $m \times n$ matrix, $m < n$, and \mathbf{c} is a constant vector.

- 10.3 Prove that the rectangle of the largest area that fits inside a circle of diameter D is a square of diagonal D .
- 10.4 With $q = 1$ and $r = 2$, $S(k_f) = 1$, write the design equations for a digital optimal quadratic regulator for the integrator

$$\dot{x} = u$$

- 10.5 The discretized state-space model of the INFANTE AUV of Problem 7.15 is given by

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} 0.9932 & -0.03434 & 0 \\ -0.009456 & 0.9978 & 0 \\ -0.0002368 & 0.04994 & 1 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 0.002988 \\ -0.0115 \\ -0.0002875 \end{bmatrix} u(k)$$

$$\begin{bmatrix} y_1(k) \\ y_2(k) \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}$$

Design a steady-state linear quadratic regulator for the system using the weight matrices $Q = I_3$ and $r = 2$.

- 10.6 A simplified linearized model of a drug delivery system to maintain blood glucose and insulin levels at prescribed values is given by

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} -0.04 & -4.4 & 0 \\ 0 & -0.025 & 1.3 \times 10^{-5} \\ 0 & 0.09 & 0 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 0.1 \end{bmatrix} \begin{bmatrix} u_1(k) \\ u_2(k) \end{bmatrix}$$

$$\begin{bmatrix} y_1(k) \\ y_2(k) \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}$$

where all variables are perturbations from the desired steady-state levels.¹ The state variables are the blood glucose concentration x_1 in mg/dL, the blood insulin concentration x_3 in mg/dL, and a variable describing the accumulation of blood insulin x_2 . The controls are the rate of glucose infusion u_1 and the rate of insulin infusion u_2 , both in mg/dL/min. Discretize the system with a sampling period $T = 5$ min and design a steady-state regulator for the system with weight matrices $Q = I_3$ and $R = 2I_2$. Simulate the system with the initial state $\mathbf{x}(0) = [6 \ 0 \ -1]^T$, and plot the trajectory in the x_1-x_3 plane as well as the time evolution of the glucose concentration.

- 10.7 Optimal control problems can be solved as unconstrained optimization problems without using Lagrange multipliers. This exercise shows the advantages of using Lagrange multipliers and demonstrates that discrete-time optimal control is equivalent to an optimization problem over the control history.
- Substitute the solution of the state equation

$$\begin{aligned} x(k) &= A^k \mathbf{x}(0) + \sum_{i=0}^{k-1} A^{k-i-1} B \mathbf{u}(i) \\ &= A^k \mathbf{x}(0) + \mathcal{C}(k) \mathcal{U}(k) \\ \mathcal{C}(k) &= [B \ AB \ \cdots \ A^{k-1}B] \\ \mathcal{U}(k) &= \text{col}\{\mathbf{u}(k-1), \dots, \mathbf{u}(0)\} \end{aligned}$$

In the performance measure

$$J = \frac{1}{2} \mathbf{x}^T(k_f) S(k_f) \mathbf{x}(k_f) + \frac{1}{2} \sum_{k=k_0}^{k_f-1} (\mathbf{x}^T(k) Q(k) \mathbf{x}(k) + \mathbf{u}^T(k) R(k) \mathbf{u}(k))$$

to eliminate the state vector and obtain

$$J = \frac{1}{2} \sum_{k=0}^{k_f} \left(\mathbf{u}^T(k) \bar{R}(k) \mathbf{u}(k) + 2^T \mathbf{x}(0) (A^T)^k Q(k) \mathcal{C}(k) \mathbf{u}(k) \right)$$

with $Q(k_f) = S(k_f)$, $R(k_f) = 0_{m \times m}$.

- Without tediously evaluating the matrix R_{eq} and the vector \mathbf{l} , explain why it is possible to rewrite the performance measure in the equivalent form

$$J_{eq} = \frac{1}{2} \mathbf{u}^T(k_f) R_{eq} \mathbf{u}(k_f) + \mathbf{u}^T(k_f) \mathbf{l}$$

¹ Chee, F., Savkin, A.V., Fernando, T.L., Nahavandi, S., 2005. Optimal H^∞ insulin injection control for blood glucose regulation in diabetic patients. IEEE Trans. Biomed. Eng. 52 (10), 1625–31.

- c. Show that the solution of the optimal control problem is given by

$$\mathbf{u}(k_f) = -R_{eq}^{-1}\mathbf{1}$$

- 10.8 For (A, B) stabilizable and $(A, Q^{1/2})$ detectable, the linear quadratic regulator yields a closed-loop stable system. To guarantee that the eigenvalues of the closed-loop system will lie inside a circle of radius $1/\alpha$, we solve the regulator problem for the scaled state and control

$$\bar{\mathbf{x}}(k) = \alpha^k \mathbf{x}(k) \quad \bar{\mathbf{u}}(k) = \alpha^{k+1} \mathbf{u}(k)$$

- a. Obtain the state equation for the scaled state vector.
 - b. Show that if the scaled closed-loop state matrix with the optimal control $\bar{\mathbf{u}}(k) = -\bar{K}\bar{\mathbf{x}}(k)$ has eigenvalues inside the unit circle, then the eigenvalues of the original state matrix with the control $\mathbf{u}(k) = -K\mathbf{x}(k)$, $K = \bar{K}/\alpha$ are inside a circle of radius $1/\alpha$.
- 10.9 Repeat Problem 10.5 with a design that guarantees that the eigenvalues of the closed-loop system are inside a circle of radius equal to one-half.
- 10.10 Using the fact that for (A, B) stabilizable and $(A, Q^{1/2})$ detectable the linear quadratic regulator yields a closed-loop stable system, show that the eigenvalues of the closed-loop system are placed inside a circle of radius $\rho < 1$ using the LQR gain for the pair $(A/\rho, B/\rho)$.
- 10.11 Show that the linear quadratic regulator with cross-product term of the form

$$J = \mathbf{x}^T(k_f)S(k_f)\mathbf{x}(k_f) + \sum_{k=k_0}^{k_f-1} (\mathbf{x}^T(k)Q(k)\mathbf{x}(k) + 2\mathbf{x}^T S \mathbf{u} \\ + \mathbf{u}^T(k)R(k)\mathbf{u}(k))$$

is equivalent to a linear quadratic regulator with no cross-product term with the cost

$$J = \mathbf{x}^T(k_f)S(k_f)\mathbf{x}(k_f) + \sum_{k=k_0}^{k_f-1} (\mathbf{x}^T(k)\bar{Q}(k)\mathbf{x}(k) \\ + \bar{\mathbf{u}}^T(k)R(k)\bar{\mathbf{u}}(k)) \\ \bar{Q} = Q - SR^{-1}S^T \\ \bar{\mathbf{u}}(k) = \mathbf{u}(k) + R^{-1}S^T\mathbf{x}(k)$$

and the plant dynamics

$$\mathbf{x}(k+1) = \bar{A}\mathbf{x}(k) + B\bar{\mathbf{u}}(k), \quad k = k_0, \dots, k_f - 1 \\ \bar{A} = A - BR^{-1}S^T$$

- 10.12 Using the multiplication of partitioned matrices, verify that the inverse of the Hamiltonian matrix

$$H = \begin{bmatrix} A^{-1} & A^{-1}BR^{-1}B^T \\ QA^{-1} & A^T + QA^{-1}BR^{-1}B^T \end{bmatrix}$$

is given by

$$H^{-1} = \begin{bmatrix} A + BR^{-1}B^TA^{-T}Q & -BR^{-1}B^TA^{-T} \\ -A^{-T}Q & A^{-T} \end{bmatrix}$$

- 10.13 Rewrite the performance measure shown in Problem 10.11 in terms of a combined input and state vector $\text{col}\{\mathbf{x}(k), \mathbf{u}(k)\}$. Then use the Hamiltonian to show that for the linear quadratic regulator with cross-product term, a sufficient condition for a minimum is that the matrix

$$\left[\begin{array}{c|c} Q & N \\ \hline N^T & R \end{array} \right]$$

must be positive definite.

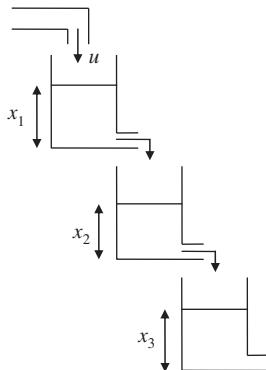
- 10.14 Obtain a predictive model for the state-space representation of a linear time-invariant system.
- 10.15 For a prediction horizon N_2 and a control horizon N_u , write the expression of the predicted outputs of a linear time-invariant system in state-space form in a compact matrix form, such as in expression (10.99).

Computer exercises

- 10.16 Write a MATLAB script to place the eigenvalues of the closed-loop steady-state regulator inside a circle of radius ρ . Use the program redesign the steady-state regulator of Example 10.5 to place its closed-loop eigenvalues inside a circle of radius $\rho = 0.5$
- 10.17 Design a steady-state regulator for the INFANTE AUV presented in Problem 10.5 with the performance measure modified to include a cross-product term with

$$S = [1 \quad 0.2 \quad 0.1]^T$$

- a. Using the MATLAB command **dlqr** and the equivalent problem with no cross-product as in Problem 10.10
 - b. Using the MATLAB command **dlqr** with the cross-product term
- 10.18 Design an output quadratic regulator for the INFANTE UAV presented in Problem 10.5 with the weights $Q_y = 1$ and $r = 100$. Plot the output response for the initial condition vector $\mathbf{x}(0) = [1, 0, 1]^T$.
- 10.19 Design an optimal LQ state-space tracking controller for the drug delivery system presented in Problem 10.6 to obtain zero steady-state error due to a unit step input.

**Figure P10.24**

Schematic of the three tanks system.

- 10.20 Write a MATLAB script that determines the steady-state quadratic regulator for the inertial system of Example 10.4 for $r = 1$ and different values of Q . Use the following three matrices, and discuss the effect of changing Q on the results of your computer simulation:
- $Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$
 - $Q = \begin{bmatrix} 10 & 0 \\ 0 & 10 \end{bmatrix}$
 - $Q = \begin{bmatrix} 100 & 0 \\ 0 & 100 \end{bmatrix}$
- 10.21 The linearized analog state-space model of the three tanks system shown in Fig. P10.24 can be written as²

$$\begin{bmatrix} \dot{x}_1(t) \\ \dot{x}_2(t) \\ \dot{x}_3(t) \end{bmatrix} = \begin{bmatrix} -1 & 0 & 0 \\ 1 & -1 & 0 \\ 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} x_1(t) \\ x_2(t) \\ x_3(t) \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} u(t)$$

$$y(t) = [0 \ 0 \ 1] \begin{bmatrix} x_1(t) \\ x_2(t) \\ x_3(t) \end{bmatrix}$$

where $x_i(t)$ is the level of the i th tank and $u(t)$ is the inlet flow of the first tank. A digital controller is required to control the fluid levels in the tanks. The controller must provide a fast step response without excessive overshoot that could lead to

² Koenig, D.M., 2009. Practical Control Engineering. McGraw-Hill, New York.

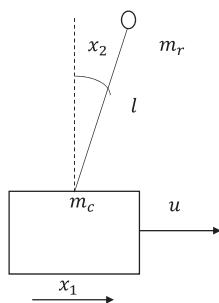


Figure P10.25
Schematic of the inverted pendulum.

fluid overflow. Using the appropriate MATLAB commands, design a digital LQ tracking controller for the system, and then simulate the closed-loop dynamics using SIMULINK. Simulate the closed-loop system using SIMULINK. Use a sampling period of 0.1 s, a state-weighting matrix $Q = 0.1 I_4$, and an input-weighting matrix $R = 1$.

- 10.22 The linearized analog state-space model of the inverted pendulum shown in Fig. P10.25 can be written as

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \\ \dot{x}_4 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & -m_r g / m_c & 0 & 0 \\ 0 & (m_r + m_c)g / (m_c l) & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1/m_c \\ -1/(m_c l) \end{bmatrix} u(t)$$

$$y(t) = [0 \quad 1 \quad 0 \quad 0] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

where $x_1(t)$ is the position of the cart, $x_2(t)$ is the angle of the pendulum, $x_3(t)$ is the velocity of the cart, $x_4(t)$ is the angular velocity of the pendulum, and $u(t)$ is the force applied to the cart. The system parameters are $m_c = 1$ kg is the mass of the cart, $m_r = 0.2$ kg is the mass of the rod, and $l = 0.5$ is the center of gravity of the rod. A digital controller is required to keep the pendulum in the upright position starting from a nonzero initial condition. Write a MATLAB program to design a digital LQ tracking controller for the system then simulate the closed-loop system to obtain its response from the initial condition $[0, 1, 0, 0]$. Use a sampling period of 0.01 s, a state-weighting matrix $Q = I_4$, and an input-weighting factor $r = 0.1$.

- 10.23 For the system in Problem 10.21, write a MATLAB program that plots the locus of the eigenvalues of the closed-loop system for $R = 1$ and $Q = q I_4$ with q in the range $0.1 \leq q \leq 10$.
- 10.24 Show that the phase margin is given by the formula

$$PM = \pm \cos^{-1}(|x|) = \pm \cos^{-1}\left(1 - \frac{k_s^2}{2}\right)$$

where (x, y) is the intersection of the two circles of Fig. 10.16 with one centered at the origin and of unity radius and the other centered at $(-1, 0)$ with radius k_S .

Show that the formula for the phase margin is equivalent to

$$PM = 2 \sin^{-1}\left(\frac{1}{2\sqrt{1+\alpha}}\right)$$

- 10.25 Consider a water heater system whose transfer function is³

$$G(z) = \frac{0.2713}{z^2(z - 0.8351)}$$

Write a MATLAB program that implements an unconstrained DMC algorithm with $N_1 = 1$, $N_2 = 10$, $N_u = 5$, $\lambda = 1$, $\alpha = 0$ and obtain the unitary set-point step response. Then, change the values of λ and α in order to understand their role in the algorithm.

- 10.26 Consider again the water heater system of the Problem 10.25. Write a MATLAB program that implements a constrained DMC algorithm with $N_1 = 1$, $N_2 = 10$, $N_u = 5$, $\lambda = 1$, $\alpha = 0$, $u_{min} = -10$, $u_{max} = 0.8$, $\Delta u_{min} = -10$, $\Delta u_{min} = -10$, $\Delta u_{max} = 10$, $u_{max} = 0.8$, $y_{min} = -1$, $y_{max} = 1.1$, and obtain the unitary set-point step response.

³ E. F. Camacho, C. Bordons, Model Predictive Control—second edition, Springer, 2007.

Elements of nonlinear digital control systems

Objectives

After completing this chapter, the reader will be able to do the following:

1. Discretize special types of nonlinear systems.
2. Determine the equilibrium point of a nonlinear discrete-time system.
3. Classify equilibrium points of discrete-time systems based on their state–plane trajectories.
4. Determine the stability or instability of nonlinear discrete-time systems.
5. Design controllers for nonlinear discrete-time systems.

Most physical systems do not obey the principle of superposition and can therefore be classified as nonlinear. By limiting the operating range of physical systems, it is possible to approximately model their behavior as linear. This chapter examines the behavior of nonlinear discrete systems without this limitation. We begin by examining the behavior of nonlinear continuous-time systems with piecewise constant inputs. We discuss Lyapunov stability theory and input–output stability both for nonlinear and linear systems. We provide a simple controller design based on Lyapunov stability theory.

Chapter Outline

11.1 Discretization of nonlinear systems	508
11.1.1 Extended linearization by input redefinition	509
11.1.2 Extended linearization by input and state redefinition	511
11.1.3 Extended linearization by output differentiation	512
11.1.4 Extended linearization using matching conditions	514
11.2 Nonlinear difference equations	517
11.2.1 Logarithmic transformation	517
11.3 Equilibrium of nonlinear discrete-time systems	518
11.4 Lyapunov stability theory	522
11.4.1 Lyapunov functions	522
11.4.2 Stability theorems	524
11.4.3 Rate of convergence	526

11.4.4 Lyapunov stability of linear systems 527

11.4.5 MATLAB 531

11.4.6 Lyapunov's linearization method 532

11.4.7 Instability theorems 533

11.4.8 Estimation of the domain of attraction 534

11.5 Stability of analog systems with digital control 537

11.6 State-plane analysis 539

11.7 Discrete-time nonlinear controller design 543

11.7.1 Controller design using extended linearization 543

11.7.2 Controller design based on Lyapunov stability theory 546

11.8 Input-output stability and the small gain theorem 548

11.8.1 Absolute stability 556

Further reading 560

Problems 561

Computer exercises 565

11.1 Discretization of nonlinear systems

Discrete-time models are easily obtained for linear continuous-time systems from their transfer functions or the solution of the state equations. For nonlinear systems, transfer functions are not defined, and the state equations are analytically solvable in only a few special cases. Thus, it is no easy task to obtain discrete-time models for nonlinear systems.

We examine the nonlinear differential equation

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + B(\mathbf{x})\mathbf{u} \quad (11.1)$$

where \mathbf{x} and \mathbf{f} are $n \times 1$ vectors, \mathbf{u} is an $m \times 1$ vector, $B(\mathbf{x})$ is an $n \times m$ matrix, and $B(\cdot)$ and $\mathbf{f}(\cdot)$ are continuous functions of all their arguments. We assume that the input is defined by

$$\mathbf{u}(t) = \mathbf{u}(kT) = \mathbf{u}(k), \quad t \in [kT, (k+1)T] \quad (11.2)$$

For each sampling period, we have the model

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + B(\mathbf{x})\mathbf{u}(k) \quad (11.3)$$

where $k = 0, 1, 2, \dots$. The behavior of the discrete-time system can, in theory, be obtained from the solution of equations of the form (11.3) with the appropriate initial conditions. In practice, only a numerical solution is possible except in a few special cases. One way to obtain a solution for nonlinear systems is to transform the dynamics of the system to obtain equivalent linear dynamics. The general theory governing such transformations, known as **global** or **extended linearization**, is beyond the scope of this text. However, in some special cases, it is possible to obtain equivalent linear models and

use them for discretization quite easily without resorting to the general theory. These include models that occur frequently in some applications. We present four cases where extended linearization is quite simple.

11.1.1 Extended linearization by input redefinition

Consider the nonlinear model

$$M(\mathbf{q})\ddot{\mathbf{q}} + \mathbf{m}(\mathbf{q}, \dot{\mathbf{q}}) = \mathbf{v}(t) \quad (11.4)$$

where $M(\mathbf{q})$ is an invertible $m \times m$ matrix, and \mathbf{m} , \mathbf{q} , and \mathbf{v} are $m \times 1$ vectors. A natural choice of state variables for the system is the $2m \times 1$ vector

$$\mathbf{x} = \text{col}\{\mathbf{x}_1, \mathbf{x}_2\} = \text{col}\{\mathbf{q}, \dot{\mathbf{q}}\} \quad (11.5)$$

To obtain equivalent linear dynamics for the system, we redefine the input vector as

$$\mathbf{u}(t) = \ddot{\mathbf{q}} \quad (11.6)$$

This yields the set of double integrators with state equations

$$\begin{aligned} \dot{\mathbf{x}}_1 &= \mathbf{x}_2 \\ \dot{\mathbf{x}}_2 &= \mathbf{u}(t) \\ \mathbf{u}(t) &= M^{-1}(\mathbf{q})[\mathbf{v}(t) - \mathbf{m}(\mathbf{q}, \dot{\mathbf{q}})] \end{aligned} \quad (11.7)$$

The solution of the state equations can be obtained by Laplace transformation, and it yields the discrete-time state equation

$$\mathbf{x}(k+1) = A_d \mathbf{x}(k) + B_d \mathbf{u}(k) \quad (11.8)$$

with

$$A_d = \begin{bmatrix} I_m & TI_m \\ 0 & I_m \end{bmatrix} \quad B_d = \begin{bmatrix} (T^2/2)I_m \\ TI_m \end{bmatrix} \quad (11.9)$$

With digital control, the input is piecewise constant and is obtained using the system model

$$\mathbf{v}(k) = M(\mathbf{x}_1(k))\mathbf{u}(k) + \mathbf{m}(\mathbf{x}_1(k), \mathbf{x}_2(k)) \quad (11.10)$$

The expression for the input \mathbf{v} is approximate because it assumes that the state does not change appreciably over a sampling period with a fixed acceleration input. The approximation is only acceptable for sufficiently slow dynamics.

The model shown in (11.4) includes many important physical systems. In particular, a large class of mechanical systems, including the m -D.O.F. (degree-of-freedom)

manipulator presented in Example 7.4, are in the form (11.4). Recall that the manipulator is governed by the equation of motion

$$M(\mathbf{q})\ddot{\mathbf{q}} + V(\mathbf{q}, \dot{\mathbf{q}})\dot{\mathbf{q}} + g(\mathbf{q}) = \boldsymbol{\tau} \quad (11.11)$$

where

\mathbf{q} = vector of generalized coordinates

$M(\mathbf{q})$ = $m \times m$ positive definite inertia matrix

$V(\mathbf{q}, \dot{\mathbf{q}})$ = $m \times m$ matrix of velocity-related terms

$\mathbf{g}(\mathbf{q})$ = $m \times 1$ vector of gravitational terms

$\boldsymbol{\tau}$ = $m \times 1$ vector of generalized forces

Clearly, the dynamics of the manipulator are a special case of (11.4). As applied to robotic manipulators, extended linearization by input redefinition is known as the **computed torque method**. This is because the torque is computed from the acceleration if measurements of positions and velocities are available.

Example 11.1

Find an approximate discrete-time model for the 2-D.O.F. robotic manipulator described in Example 7.5 and find an expression for calculating the torque vector.

Solution

The manipulator dynamics are governed by

$$\begin{aligned} M(\theta) &= \begin{bmatrix} (m_1 + m_2)l_1^2 + m_2l_2^2 + 2m_2l_1l_2 \cos(\theta_2) & m_2l_2^2 + m_2l_1l_2 \cos(\theta_2) \\ m_2l_2^2 + m_2l_1l_2 \cos(\theta_2) & m_2l_2^2 \end{bmatrix} \\ V(\theta, \dot{\theta})\dot{\theta} &= \begin{bmatrix} -m_2l_1l_2 \sin(\theta_2)\dot{\theta}_2(2\dot{\theta}_1 + \dot{\theta}_2) \\ m_2l_1l_2 \sin(\theta_2)\dot{\theta}_1^2 \end{bmatrix} \\ \mathbf{g}(\theta) &= \begin{bmatrix} (m_1 + m_2)gl_1 \sin(\theta_1) + m_2gl_2 \sin(\theta_1 + \theta_2) \\ m_2gl_2 \sin(\theta_1 + \theta_2) \end{bmatrix} \end{aligned}$$

For this system, the coordinate vector is $\mathbf{q} = [\theta_1 \ \theta_2]^T$. We have a fourth-order linear system of the form

$$\begin{aligned} \dot{\mathbf{x}}_1(t) &= \mathbf{x}_2(t) \\ \dot{\mathbf{x}}_2(t) &= \mathbf{u}(t) \\ \mathbf{x}_1(t) &= \mathbf{q} = [x_{11}(t) \ x_{12}(t)]^T \\ \mathbf{x}_2(t) &= \dot{\mathbf{q}} = [x_{21}(t) \ x_{22}(t)]^T \\ \mathbf{u}(t) &= [\tau_1(t) \ \tau_2(t)]^T \end{aligned}$$

As in (11.10), the torque is calculated using the equation

$$\boldsymbol{\tau}(k) = M(\mathbf{x}_1(k))\mathbf{u}(k) + V(\mathbf{x}_1(k), \mathbf{x}_2(k))\mathbf{x}_2(k) + \mathbf{g}(\mathbf{x}_1(k))$$

11.1.2 Extended linearization by input and state redefinition

Consider the nonlinear state equations

$$\begin{aligned}\dot{\mathbf{z}}_1 &= \mathbf{f}_1(\mathbf{z}_1, \mathbf{z}_2) \\ \dot{\mathbf{z}}_2 &= \mathbf{f}_2(\mathbf{z}_1, \mathbf{z}_2) + G(\mathbf{z}_1, \mathbf{z}_2)\mathbf{v}(t)\end{aligned}\quad (11.12)$$

where $\mathbf{f}_i(\cdot)$, $i = 1, 2$, \mathbf{v} are $m \times 1$, and G is $m \times m$. We redefine the state variables and input as

$$\begin{aligned}\mathbf{x}_1 &= \mathbf{z}_1 \\ \mathbf{x}_2 &= \dot{\mathbf{z}}_1 \\ \mathbf{u}(t) &= \ddot{\mathbf{z}}_1\end{aligned}\quad (11.13)$$

The new variables have the linear dynamics of (11.7) and the discrete-time model of (11.8–11.9).

Using the state equations, we can rewrite the new input as

$$\mathbf{u}(t) = \frac{\partial \mathbf{f}_1(\mathbf{z}_1, \mathbf{z}_2)}{\partial \mathbf{z}_1} \mathbf{f}_1(\mathbf{z}_1, \mathbf{z}_2) + \frac{\partial \mathbf{f}_1(\mathbf{z}_1, \mathbf{z}_2)}{\partial \mathbf{z}_2} [\mathbf{f}_2(\mathbf{z}_1, \mathbf{z}_2) + G(\mathbf{z}_1, \mathbf{z}_2)\mathbf{v}(t)] \quad (11.14)$$

We solve for the nonlinear system input at time k to obtain

$$\begin{aligned}\mathbf{v}(k) &= \left[\frac{\partial \mathbf{f}_1(\mathbf{z}_1, \mathbf{z}_2)}{\partial \mathbf{z}_2} G(\mathbf{z}_1, \mathbf{z}_2) \right]^{-1} \\ &\quad \times \left\{ \mathbf{u}(k) - \frac{\partial \mathbf{f}_1(\mathbf{z}_1, \mathbf{z}_2)}{\partial \mathbf{z}_1} \mathbf{f}_2(\mathbf{z}_1, \mathbf{z}_2) - \frac{\partial \mathbf{f}_1(\mathbf{z}_1, \mathbf{z}_2)}{\partial \mathbf{z}_2} \mathbf{f}_2(\mathbf{z}_1, \mathbf{z}_2) \right\}\end{aligned}\quad (11.15)$$

The expression for the input $\mathbf{v}(k)$ is approximate because the state of the system changes over a sampling period for a fixed input $\mathbf{u}(k)$.

Example 11.2

Discretize the nonlinear differential equation

$$\begin{aligned}\dot{z}_1 &= -\frac{1}{2}(z_1 z_2)^2 \\ \dot{z}_2 &= 4z_1 z_2^3 + 1.5z_2 - v, \quad z_1 \neq 0, z_2 \neq 0\end{aligned}$$

Solution

We differentiate the first state equation to obtain

$$\begin{aligned}\ddot{z}_1 &= -z_1 z_2^2 \dot{z}_1 - z_1^2 z_2 \dot{z}_2 \\ &= -3.5z_1^3 z_2^4 - 1.5z_1^2 z_2^2 + z_1^2 z_2 v = u(t), \quad z_1 \neq 0, z_2 \neq 0\end{aligned}$$

Example 11.2—cont'd

We have the linear model

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= u(t)\end{aligned}$$

This gives the equivalent discrete-time model

$$\mathbf{x}(k+1) = \begin{bmatrix} 1 & T \\ 0 & 1 \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} (T^2/2) \\ T \end{bmatrix} \mathbf{u}(k)$$

with $\mathbf{x}(k) = [x_1(k) \ x_2(k)]^T = [z_1(k) \ z_1(k+1)]^T$. If measurements of the actual state of the system are available, then the input is computed using

$$v(k) = 3.5z_1(k)z_2^2(k) + 1.5z_2(k) + \frac{u(k)}{z_1^2(k)z_2(k)}, \quad z_1 \neq 0, z_2 \neq 0$$

11.1.3 Extended linearization by output differentiation

Consider the nonlinear state-space model

$$\begin{aligned}\dot{\mathbf{z}} &= \mathbf{f}(\mathbf{z}) + G(\mathbf{z})\mathbf{v}(t) \\ \mathbf{y} &= \mathbf{c}(\mathbf{z})\end{aligned}\tag{11.16}$$

where \mathbf{z} and \mathbf{f} are $n \times 1$, \mathbf{v} is $m \times 1$, \mathbf{y} and \mathbf{c} are $l \times 1$, and G is $m \times m$. If we differentiate each scalar output equation and substitute from the state equation, we have

$$\begin{aligned}\frac{dy_i}{dt} &= \frac{\partial c_i^T(\mathbf{z})}{\partial \mathbf{z}} \dot{\mathbf{z}} \\ &= \frac{\partial c_i^T(\mathbf{z})}{\partial \mathbf{z}} [\mathbf{f}(\mathbf{z}) + G(\mathbf{z})\mathbf{v}(t)] \quad i = 1, \dots, l\end{aligned}\tag{11.17}$$

The gradient of a scalar function y multiplied by another vector f is known as the Lie derivative of y with respect to f and is written as

$$L_f y = \frac{\partial y(z)}{\partial z} f\tag{11.18}$$

If the coefficient of the input vector \mathbf{v} gives

$$\frac{\partial c_i^T(\mathbf{z})}{\partial \mathbf{z}} G(\mathbf{z}) = 0^T\tag{11.19}$$

then the time derivative of an output y_i is equivalent to the operation $L_f y_i$. If the coefficient of the input vector is nonzero, we stop to avoid input differentiation. If the

coefficient is zero, we repeat the process until the coefficient becomes nonzero. We denote r_i repetitions of this operation by the symbol $L_z^{r_i}y_i$. We define a new input equal to the derivative where the input appears and obtain a linear model in the form

$$\begin{aligned}\frac{d^{r_i}y_i}{dt^{r_i}} &= u_i(t) \\ u_i(t) &= L_z^{r_i}y_i, \quad i = 1, 2, \dots, l\end{aligned}\quad (11.20)$$

where the number of derivatives r_i required for the input to appear in the expression is known as the i^{th} **relative degree**. The linear system is in the form of l sets of integrators, and its order is

$$n_l = \sum_{i=1}^l r_i \leq n \quad (11.21)$$

Because the order of the linear model can be less than the original order of the nonlinear system n , the equivalent linear model of the nonlinear system can have unobservable dynamics. The unobservable dynamics are known as the **internal dynamics** or **zero dynamics**. For the linear model to yield an acceptable discrete-time representation, we require the internal dynamics to be stable and sufficiently fast so as not to significantly alter the time response of the system. Under these conditions, the linear model provides an adequate though incomplete description of the system dynamics, and we can discretize each of the l linear subsystems to obtain the overall discrete-time model.

Example 11.3

Consider a model of drug receptor binding in a drug delivery system. The drug is assumed to be divided between four compartments in the body. The concentrations in the four compartments are given by

$$\begin{aligned}\dot{z}_1 &= a_{11}z_1 + a_{12}z_1z_2 + b_1d \\ \dot{z}_2 &= a_{21}z_2 + a_{22}z_1z_2 + b_2d \\ \dot{z}_3 &= a_{31}z_1 + a_{32}z_3 \\ \dot{z}_4 &= a_{41}z_2\end{aligned}$$

where d is the drug dose, z_i is the concentration in the i^{th} compartment, and a_{ij} are time-varying coefficients. The second compartment represents the drug receptor complex. Hence, the output equation is given by

$$y = z_2$$

Obtain an equivalent linear model by output differentiation then use it to obtain an approximate discrete-time model of the system.

Solution

We differentiate the output equation once

$$\dot{y} = \dot{z}_2 = a_{21}z_2 + a_{22}z_1z_2 + b_2d$$

The derivative includes the input, and no further differentiation is possible because it would lead to input derivative terms.

Example 11.3—cont'd

The model of the linearized system is

$$\begin{aligned}\dot{x} &= u(t) \\ u(t) &= a_{21}z_2 + a_{22}z_1z_2 + b_2d\end{aligned}$$

The model is first order, even though we have a third-order plant, and it is only equivalent to part of the dynamics of the original system. We assume the system is detectable with fast unobservable dynamics, which is reasonable for the drug delivery system. Hence, the linear model provides an adequate description of the system, and we obtain the discrete-time model

$$x(k+1) = x(k) + Tu(k)$$

The drug dose is given by

$$d(k) = \frac{1}{b_2}[u(k) - a_{21}z_2(k) + a_{22}z_1(k)z_2(k)]$$

11.1.4 Extended linearization using matching conditions

Theorem 11.1 gives the analytical solution in a special case where an equivalent linear model can be obtained for a nonlinear system by a simple state transformation.

Theorem 11.1

For the system of (11.1), let $B(\mathbf{x})$ and $f(\mathbf{x})$ satisfy the matching conditions

$$\begin{aligned}B(\mathbf{x}) &= B_1(\mathbf{x})B_2 \\ \mathbf{f}(\mathbf{x}) &= B_1(\mathbf{x})A\mathbf{h}(\mathbf{x})\end{aligned}\tag{11.22}$$

where $B_1(\mathbf{x})$ is an $n \times n$ matrix invertible in some region D , B_2 is an $n \times m$ constant vector, A is an $n \times n$ matrix, and $\mathbf{h}(\mathbf{x})$ has the derivative given by

$$\frac{\partial \mathbf{h}(\mathbf{x})}{\partial \mathbf{x}} = [B_1(\mathbf{x})]^{-1}\tag{11.23}$$

with $\mathbf{h}(\cdot)$ invertible in the region D . Then the solution of (11.1) over D is

$$\mathbf{x}(t) = \mathbf{h}^{-1}(\mathbf{w}(t))\tag{11.24}$$

where \mathbf{w} is the solution of the linear equation

$$\dot{\mathbf{w}}(t) = A\mathbf{w}(t) + B_2\mathbf{u}(t)\tag{11.25}$$

Proof

From (11.24), $\mathbf{w} = \mathbf{h}(\mathbf{x})$. Differentiating with respect to time and substituting from (11.1) and (11.22) gives

$$\begin{aligned}\dot{\mathbf{w}} &= \frac{\partial \mathbf{h}(\mathbf{x})}{\partial \mathbf{x}} \dot{\mathbf{x}} = [B_1(\mathbf{x})]^{-1} \{ \mathbf{f}(\mathbf{x}) + B_1(\mathbf{x})B_2 \mathbf{u} \} \\ &= A\mathbf{h}(\mathbf{x}) + B_2 \mathbf{u} \\ &= A\mathbf{w} + B_2 \mathbf{u}\end{aligned}$$

Note that, if the decomposition of Theorem 11.1 is possible, we do not need to solve a partial differential equation to obtain the vector function $\mathbf{h}(\mathbf{x})$, but we do need to find its inverse function to obtain the vector \mathbf{x} . The solution of a partial differential equation is a common requirement for obtaining transformations of nonlinear systems to linear dynamics that can be avoided in special cases.

The solution of the state equation for \mathbf{w} over any sampling period T is

$$\mathbf{w}(k+1) = A_w \mathbf{w}(k) + B_w \mathbf{u}(k)$$

where

$$A_w = e^{AT}, \quad B_w = \int_0^T e^{A\tau} B_2 d\tau$$

This is a discrete state equation for the original nonlinear system.

Example 11.4

Discretize the nonlinear differential equation

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} -(x_1/x_2)^2 \\ (8x_2^3/x_1) + 3x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ -x_2^3 \end{bmatrix} u, \quad x_1 \neq 0, \quad x_2 \neq 0$$

Solution

Assuming a piecewise constant input u , we rewrite the nonlinear equation as

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_1^2 & 0 \\ 0 & -x_2^3 \end{bmatrix} \left\{ \begin{bmatrix} 0 & -1 \\ -4 & 3 \end{bmatrix} \begin{bmatrix} -x_1^{-1} \\ x_2^{-2}/2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u(k) \right\}$$

Example 11.4—cont'd

To apply Theorem 11.1, we define the terms

$$B_1(\mathbf{x}) = \begin{bmatrix} x_1^2 & 0 \\ 0 & -x_2^3 \end{bmatrix} \quad B_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & -1 \\ -4 & 3 \end{bmatrix} \quad h(\mathbf{x}) = \begin{bmatrix} -x_1^{-1} \\ x_2^{-2}/2 \end{bmatrix}$$

We verify that the Jacobian of the vector satisfies

$$\frac{\partial \mathbf{h}(\mathbf{x})}{\partial \mathbf{x}} = \begin{bmatrix} \partial h_1 / \partial x_1 & \partial h_1 / \partial x_2 \\ \partial h_2 / \partial x_1 & \partial h_2 / \partial x_2 \end{bmatrix} = [B_1(\mathbf{x})]^{-1} = \begin{bmatrix} x_1^{-2} & 0 \\ 0 & -x_2^{-3} \end{bmatrix}$$

and that Theorem 11.1 applies.

We first solve the linear state equation

$$\begin{bmatrix} \dot{w}_1 \\ \dot{w}_2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ -4 & 3 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u(k)$$

to obtain

$$\begin{bmatrix} w_1(k+1) \\ w_2(k+1) \end{bmatrix} = \frac{1}{5} \left\{ \begin{bmatrix} 1 & -1 \\ -4 & 4 \end{bmatrix} e^{4T} + \begin{bmatrix} 4 & 1 \\ 4 & 1 \end{bmatrix} e^{-T} \right\} \begin{bmatrix} w_1(k) \\ w_2(k) \end{bmatrix} + \frac{1}{5} \left\{ \begin{bmatrix} -1 \\ 4 \end{bmatrix} \frac{e^{4T} - 1}{4} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} (1 - e^{-T}) \right\} u(k)$$

where $\mathbf{w}(k) = \mathbf{h}[\mathbf{x}(k)]$, $k = 0, 1, 2, \dots$. We now have the recursion

$$\begin{bmatrix} w_1(k+1) \\ w_2(k+1) \end{bmatrix} = \frac{1}{5} \left\{ \begin{bmatrix} 1 & -1 \\ -4 & 4 \end{bmatrix} e^{4T} + \begin{bmatrix} 4 & 1 \\ 4 & 1 \end{bmatrix} e^{-T} \right\} \begin{bmatrix} w_1(k) \\ w_2(k) \end{bmatrix} + \frac{1}{5} \left\{ \begin{bmatrix} -1 \\ 4 \end{bmatrix} \frac{e^{4T} - 1}{4} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} (1 - e^{-T}) \right\} u(k)$$

Using the inverse of the function $\mathbf{h}[\mathbf{x}(k)]$, we solve for the state

$$x_1(k+1) = -1/w_1(k+1), \quad x_2(k+1) = \pm (2w_2(k+1))^{1/2}$$

which is clearly nonunique. However, one would be able to select the appropriate solution because changes in the state variable should not be large over a short sampling period. For a given sampling period T , we can obtain a discrete nonlinear recursion for $\mathbf{x}(k)$.

It is obvious from Example 11.4 that the analysis of discrete-time systems based on nonlinear analog systems is only possible in special cases. In addition, the transformation is sensitive to errors in the system model. We conclude that this is often not a practical approach to controlling system design. However, most nonlinear control systems are today digitally implemented even if based on an analog design. We therefore return to a discussion of system (11.1) with the control (11.2) in Section 11.5.

11.2 Nonlinear difference equations

For a nonlinear system with a DAC and ADC, system identification can yield a discrete-time model directly. A discrete-time model can also be obtained analytically in a few special cases, as discussed in Section 11.1. We thus have a nonlinear difference equation to describe a nonlinear discrete-time system. Unfortunately, nonlinear difference equations can be solved analytically in only a few special cases. We discuss one special case where such solutions are available.

11.2.1 Logarithmic transformation

Consider the nonlinear difference equation

$$[y(k+n)]^{\alpha_n} [y(k+n-1)]^{\alpha_{n-1}} \dots [y(k)]^{\alpha_0} = u(k) \quad (11.26)$$

where α_i , $i = 0, 1, \dots, n$ are constant. Then, taking the natural log of (11.26) gives

$$\alpha_n x(k+n) + \alpha_{n-1} x(k+n-1) \dots + \alpha_0 x(k) = v(k) \quad (11.27)$$

with $x(k+i) = \ln[y(k+i)]$, $i = 0, 1, 2, \dots, n$, and $v(k) = \ln[u(k)]$. This is a linear difference equation that can be easily solved by z -transformation. Finally, we obtain

$$y(k) = e^{x(k)} \quad (11.28)$$

Example 11.5

Solve the nonlinear difference equation

$$[y(k+2)][y(k+1)]^5[y(k)]^4 = u(k)$$

with zero initial conditions and the input

$$u(k) = e^{-10k}$$

Solution

Taking the natural log of the equation, we obtain

$$x(k+2) + 5x(k+1) + 4x(k) = -10k$$

Example 11.5—cont'd

The z-transform of the equation with zero initial conditions

$$[z^2 + 5z + 4]X(z) = \frac{-10z}{(z - 1)^2}$$

yields $X(z)$ as

$$X(z) = \frac{-10z}{(z - 1)^2(z + 1)(z + 4)}$$

Inverse z-transforming gives the discrete-time function

$$x(k) = -k + 0.7 - 0.8333(-1)^k + 0.1333(-4)^k, \quad k \geq 0$$

Hence, the solution of the nonlinear difference equation is

$$y(k) = e^{-k+0.7-0.8333(-1)^k+0.1333(-4)^k}, \quad k \geq 0$$

11.3 Equilibrium of nonlinear discrete-time systems

Equilibrium is defined as a condition in which a system remains indefinitely unless it is disturbed. If a discrete-time system is described by the nonlinear difference equation

$$\mathbf{x}(k+1) = \mathbf{f}[\mathbf{x}(k)] + B[\mathbf{x}(k)]\mathbf{u}(k) \quad (11.29)$$

then at an equilibrium, it is governed by the identity

$$\mathbf{x}_e = \mathbf{f}[\mathbf{x}_e] + B[\mathbf{x}_e]\mathbf{u}(k) \quad (11.30)$$

where \mathbf{x}_e denotes the equilibrium state. We are typically interested in the equilibrium for an unforced system, and we therefore rewrite the equilibrium condition (11.30) as

$$\mathbf{x}_e = \mathbf{f}[\mathbf{x}_e] \quad (11.31)$$

In mathematics, such an equilibrium point is known as a **fixed point**.

Note that the behavior of the system in the vicinity of its equilibrium point determines whether we classify the equilibrium as stable or unstable. For a stable equilibrium, we expect the trajectories of the system to remain arbitrarily close to or to converge to the equilibrium. Unlike continuous-time systems, convergence to an equilibrium point can occur after a finite time period.

Clearly, both (11.30) and (11.31), in general, have more than one solution so that nonlinear systems often have several equilibrium points. This is demonstrated in Example 11.6.

Example 11.6

Find the equilibrium points of the nonlinear discrete-time system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} x_2(k) \\ x_1^3(k) \end{bmatrix}$$

Solution

The equilibrium points are determined from the condition

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

or equivalently,

$$x_2(k) = x_1(k) = x_1^3(k)$$

Thus, the system has the three equilibrium points: $(0, 0)$, $(1, 1)$, and $(-1, -1)$.

Nonlinear systems often have multiple equilibrium points, some of which are stable, while others are unstable. We therefore talk about the stability of an equilibrium point rather than the stability of the system. We also need to be more careful in our definitions of stability. To characterize equilibrium points, we need the following definitions.

Definition 11.1**Lyapunov Stability**

The equilibrium \mathbf{x}_e is Lyapunov stable if

$\forall \varepsilon > 0, \exists \delta(\varepsilon) > 0$ such that

$$\|\mathbf{x}(k_0) - \mathbf{x}_e\|_e < \delta(\varepsilon) \Rightarrow \|\mathbf{x}(k) - \mathbf{x}_e\| < \varepsilon, \forall k \geq k_0 \quad (11.32)$$

Otherwise, \mathbf{x}_e is unstable.

Lyapunov stability requires the trajectories of the system to remain arbitrarily close to the equilibrium by appropriate choice of the initial state. In other words, if we start at a point that is sufficiently close to the equilibrium, we will stay as close to the equilibrium as we desire as depicted in Fig. 11.1. However, Lyapunov stability does not require asymptotic convergence to the equilibrium. The following definition does.

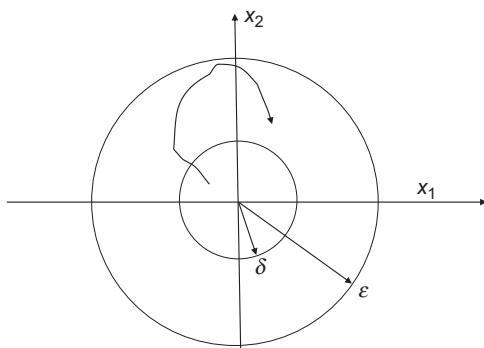


Figure 11.1
Lyapunov stability.

Definition 11.2

Asymptotic Stability

The equilibrium \mathbf{x}_e is **locally asymptotically stable** if it is both

1. Stable in the sense of Lyapunov and
2. $\exists \delta_1 > 0$ s.t.

$$\|\mathbf{x}(k_0) - \mathbf{x}_e\| < \delta_1 \Rightarrow \lim_{k \rightarrow \infty} \mathbf{x}(k) = \mathbf{x}_e. \quad (11.33)$$

If the second property holds globally, i.e., for any real $\mathbf{x}(k_0)$, then the equilibrium is **globally asymptotically stable**.

For the equilibrium to be globally asymptotically stable, it must have only one equilibrium point. This is because the system will never converge to it if it starts at another equilibrium point and will remain there indefinitely. For a system with one stable equilibrium point, we can unambiguously say that the system is stable.

The following property allows us to conclude asymptotic stability for some nonlinear systems.

Definition 11.3

Contraction: A function $f(\mathbf{x})$ is known as a contraction if it satisfies

$$\|\mathbf{f}(\mathbf{x} - \mathbf{y})\| \leq \alpha \|\mathbf{x} - \mathbf{y}\|, \quad |\alpha| < 1 \quad (11.34)$$

where α is known as a **contraction constant** and $\|\cdot\|$ is any vector norm.

Theorem 11.2 provides conditions for the existence of an equilibrium point for a discrete-time system. The proof is left as an exercise.

Theorem 11.2

A contraction $\mathbf{f}(\mathbf{x})$ has a unique fixed point.

Example 11.7

Determine if the following nonlinear discrete-time system converges to the origin

$$\begin{aligned}\mathbf{x}(k+1) &= \mathbf{f}[\mathbf{x}(k)], \quad \|\mathbf{f}[\mathbf{x}(k)]\| \leq \alpha \|\mathbf{x}(k)\|, \quad 0 < \alpha < 1 \\ k &= 0, 1, 2, \dots\end{aligned}$$

Solution

We have the inequality

$$\|\mathbf{x}(k+1)\| = \|\mathbf{f}[\mathbf{x}(k)]\| \leq \alpha \|\mathbf{x}(k)\|, \quad k = 0, 1, 2, \dots$$

For any initial state $\mathbf{x}(0)$, the solution of the difference equation satisfies

$$\|\mathbf{x}(k)\| \leq \alpha^k \|\mathbf{x}(0)\|, \quad k = 0, 1, 2, \dots$$

Hence, the system converges to a fixed point at the origin as $k \rightarrow \infty$.

If a continuous-time system is discretized, then all its equilibrium points will be equilibrium points of the discrete-time system. This is because discretization corresponds to the solution of the differential equation governing the original system at the sampling points. Because the analog system at an equilibrium remains there, the discretized system will have the same behavior and therefore the same equilibrium.

The definition of asymptotic stability does not address the issue of rate of convergence to the equilibrium. The following definition includes the rate of convergence.

Definition 11.4**Exponential Stability**

The equilibrium \mathbf{x}_e is exponentially stable if $\exists \alpha, \lambda > 0$ s.t.

$$\mathbf{x}(k) - \mathbf{x}_e \leq \alpha \mathbf{x}(k_0) - \mathbf{x}_e \lambda^k, \quad \forall k \geq k_0, \quad \forall \mathbf{x}(k_0) - \mathbf{x}_e < \delta \quad (11.35)$$

The equilibrium \mathbf{x}_e is globally exponentially stable if the condition holds $\forall \mathbf{x} \in \mathcal{R}^n$

Clearly, exponentially stable implies asymptotically stability, which in turn implies stability in the sense of Lyapunov.

11.4 Lyapunov stability theory

Lyapunov stability theory is based on the idea that at a stable equilibrium, the energy of the system has a local minimum, whereas at an unstable equilibrium, it is at a maximum. This property is not restricted to energy and is in fact shared by a class of function that depends on the dynamics of the system. We call such functions **Lyapunov functions**.

11.4.1 Lyapunov functions

If a Lyapunov function can be found for an equilibrium point, then it can be used to determine its stability or instability. This is particularly simple for linear systems but can be complicated for a nonlinear system.

We begin by examining the properties of energy functions that we need to generalize and retain for a Lyapunov function. We note the following:

- Energy is a nonnegative quantity.
- Energy changes continuously with its arguments.

We call functions that are positive except at the origin **positive definite**. We provide a formal definition of this property.

Definition 11.5

A scalar continuous function $V(x)$ is said to be positive definite if

- $V(\mathbf{0}) = 0$.
- $V(\mathbf{x}) > 0$ for any nonzero \mathbf{x} .

Similar definitions are also useful where the greater-than sign in the second condition is replaced by other inequalities with all other properties unchanged. Thus, we can define **positive semidefinite** (\geq), **negative definite** ($<$), and **negative semidefinite** (\leq) functions. If none of these definitions apply, the function is said to be **indefinite**. Definition 11.5 may hold locally in some region in the vicinity of the origin, and then the function is called **locally positive definite**, or globally, in which case it is **globally positive definite**. Similarly, we characterize other properties, such as negative definiteness, as local or global.

A common choice of definite function is the **quadratic form** $\mathbf{x}^T P \mathbf{x}$. The sign of the quadratic form is determined by the eigenvalues of the matrix P (see Appendix III). The quadratic form is positive definite if the matrix P is positive definite, in which case its

eigenvalues are all positive. Similarly, we can characterize the quadratic form as negative definite if the eigenvalues of P are all negative, positive semidefinite if the eigenvalues of P are positive or zero, negative semidefinite if negative or zero, and indefinite if P has positive and negative eigenvalues. Quadratic forms are a common choice of Lyapunov function because of their simple mathematical properties. However, it is often preferable to use other Lyapunov functions with properties tailored to suit the particular problem. We now list the mathematical properties of Lyapunov functions.

Definition 11.6

A scalar function $V(\mathbf{x})$ is a Lyapunov function in a region D if it satisfies the following conditions in D :

- It is positive definite.
- It decreases along the trajectories of the system—that is,

$$\Delta V(k) = V(\mathbf{x}(k+1)) - V(\mathbf{x}(k)) < 0, \quad k = 0, 1, 2, \dots \quad (11.36)$$

Definition 11.6 is used in local stability theorems. To prove global stability, we need an additional condition in addition to extending the two listed earlier.

Definition 11.7

A scalar function $V(\mathbf{x})$ is a Lyapunov function if it satisfies the following conditions:

- It is positive definite.
- It decreases along the trajectories of the system.
- It is radially unbounded, i.e., it uniformly satisfies

$$V(\mathbf{x}(k)) \rightarrow \infty, \quad \text{as} \quad \|\mathbf{x}(k)\| \rightarrow \infty \quad (11.37)$$

This last condition ensures that whenever the function V remains bounded, the state vector will also remain bounded. By “uniformly,” we mean that the condition must be satisfied regardless of how the norm of the state vector grows unbounded so that a finite value of V can always be associated with a finite state. Note that unlike Definition 11.6, Definition 11.7 requires that the first and second conditions be satisfied globally.

11.4.2 Stability theorems

We provide some sufficient conditions for the stability of nonlinear discrete-time systems, then we specialize our results to linear time-invariant systems.

Theorem 11.3

The equilibrium point $\mathbf{x} = \mathbf{0}$ of the nonlinear discrete-time system

$$\mathbf{x}(k+1) = \mathbf{f}[\mathbf{x}(k)], \quad k = 0, 1, 2, \dots \quad (11.38)$$

is asymptotically stable if there exists a locally positive definite Lyapunov function for the system satisfying Definition 11.6.

Proof

For any motion along the trajectories of the system, we have

$$\begin{aligned}\Delta V(k) &= V(\mathbf{x}(k+1)) - V(\mathbf{x}(k)) \\ &= V(\mathbf{f}[\mathbf{x}(k)]) - V(\mathbf{f}[\mathbf{x}(k-1)]) < 0, \quad k = 0, 1, 2, \dots\end{aligned}$$

This implies that as the motion progresses, the value of V decreases continuously. However, because V is bounded below by zero, it converges to zero. Because V is only zero for zero argument, the state of the system must converge to zero.

Example 11.8

Investigate the stability of the system using the Lyapunov stability approach:

$$\begin{aligned}x_1(k+1) &= 0.2x_1(k) - 0.08x_2^2(k) \\ x_2(k+1) &= -0.3x_1(k)x_2(k), \quad k = 0, 1, 2, \dots\end{aligned}$$

Solution

We first observe that the system has an equilibrium point at the origin. We select the quadratic Lyapunov function

$$V(\mathbf{x}) = \mathbf{x}^T P \mathbf{x}, \quad P = \text{diag}\{p_1, 1\}$$

Note that the unity entry simplifies the notation without affecting the results because the form of the function is unchanged if it is multiplied by any positive scalar p_2 .

Example 11.8—cont'd

The corresponding difference is

$$\begin{aligned}\Delta V(k) &= p_1 \left\{ [0.2x_1(k) - 0.08x_2^2(k)]^2 - x_1^2(k) \right\} + \left\{ [-0.3x_1(k)x_2(k)]^2 - x_2^2(k) \right\} \\ &= -0.96p_1x_1^2(k) + \{0.09x_1^2(k) - 0.032p_1x_1(k) + 0.0064p_1x_2^2(k) - 1\}x_2^2(k) \\ k &= 0, 1, 2, \dots\end{aligned}$$

The difference remains negative provided that the term between the braces is negative. We restrict the magnitude of $x_2(k)$ to less than 12.5 (the square root of the reciprocal of 0.0064) and then simplify the condition to

$$9x_1^2(k) - 3.2p_1x_1(k) + 100(p_1 - 1) < 0, k = 0, 1, 2, \dots$$

Choosing a very small but positive value for p_1 makes the two middle terms of $\Delta V(k)$ negligible. This leaves two terms that are negative for values of x_1 of magnitude smaller than 3.33. For example, a plot of the LHS of the inequality for p_1 verifies that it is negative in the selected x_1 range. Thus, the difference remains negative for initial conditions that are inside a circle of radius approximately equal to 3.33 centered at the origin. By Theorem 11.3, we conclude that the system is asymptotically stable in the region $\|\mathbf{x}\| < 3.33$.

The system is actually stable outside this region, but our choice of Lyapunov function cannot provide a better estimate of the stability region than the one we obtained.

For global asymptotic stability, the system must have a unique equilibrium point. Otherwise, starting at an equilibrium point prevents the system from converging to another equilibrium state. We can therefore say that a system is globally asymptotically stable and not just its equilibrium point. The following theorem gives a sufficient condition for global asymptotic stability.

Theorem 11.4

The nonlinear discrete-time system

$$\mathbf{x}(k+1) = \mathbf{f}[\mathbf{x}(k)], \quad k = 0, 1, 2, \dots \quad (11.39)$$

with equilibrium $\mathbf{x}(0) = \mathbf{0}$ is globally asymptotically stable if there exists a globally positive definite, radially unbounded Lyapunov function for the system satisfying Definition 11.7.

Proof

The proof is similar to that of Theorem 11.3 and the details are omitted. The radial unboundedness condition guarantees that the state of the system will converge to zero with the Lyapunov function.

The results of this section require the difference ΔV of the Lyapunov function to be negative definite. In some cases, this condition can be relaxed to negative semidefinite. In particular, if the nonzero values of the vector \mathbf{x} for which ΔV is zero are ones that are never reached by the system, then they have no impact on the stability analysis. This leads to the following result.

Corollary 11.1

The equilibrium point $\mathbf{x} = \mathbf{0}$ of the nonlinear discrete-time system of (11.39) is asymptotically stable if there exists a locally positive definite Lyapunov function with negative semidefinite difference $\Delta V(k)$ for all k for the system and with $\Delta V(k)$ zero only for $\mathbf{x} = \mathbf{0}$.

Note that the preceding theorems only provide a sufficient stability condition. Thus, failure of the stability test does not prove instability. In the linear time-invariant case, a much stronger result is available.

11.4.3 Rate of convergence

In some cases, we can use a Lyapunov function to determine the rate of convergence of the system to the origin. In particular, if we can rewrite the difference of the Lyapunov function in the form

$$\Delta V(k) \leq -\alpha V(\mathbf{x}(k)), \quad 0 < \alpha < 1 \quad (11.40)$$

substituting for the difference gives the recursion

$$V(\mathbf{x}(k+1)) \leq (1 - \alpha)V(\mathbf{x}(k)) \quad (11.41)$$

The upper bound of the constant α guarantees that the coefficient of V is positive. The solution of the difference equation is

$$V(\mathbf{x}(k)) \leq (1 - \alpha)^k V(\mathbf{x}(0)) \quad (11.42)$$

If V is a Lyapunov function, then it converges to zero with the state of the system and its rate of convergence allows us to estimate the rate of convergence to the equilibrium point. If in addition V is a quadratic form, then the rate of convergence of the state is the square root of that of V .

Example 11.9

If the quadratic form $\mathbf{x}^T \mathbf{x}$ is a Lyapunov function for a discrete-time system with difference

$$\Delta V(k) = -0.25\|\mathbf{x}(k)\|^2 - 0.5\|\mathbf{x}(k)\|^4 < 0$$

Characterize the convergence of the system trajectories to the origin.

Solution

$$\begin{aligned}\Delta V(k) &= -0.25\|\mathbf{x}(k)\|^2 - 0.5\|\mathbf{x}(k)\|^4 \leq -0.25\|\mathbf{x}(k)\|^2 \\ &= -0.25V(\mathbf{x}(k))\end{aligned}$$

$$\begin{aligned}V(\mathbf{x}(k)) &= \|\mathbf{x}(k)\|^2 \\ &\leq (1 - 0.25)^k V(\mathbf{x}(0)) = 0.75^k \|\mathbf{x}(0)\|^2 \\ \|\mathbf{x}(k)\| &\leq 0.866^k \|\mathbf{x}(0)\|, \quad k = 0, 1, 2, \dots\end{aligned}$$

The trajectories of the system converge to the origin exponentially with convergence rate faster than 0.86.

11.4.4 Lyapunov stability of linear systems

Lyapunov stability results typically provide us with sufficient conditions. Failure to meet the conditions of a Lyapunov test leaves us with no conclusion and with the need to repeat the test using a different Lyapunov function or to try a different test. For linear systems, Lyapunov stability can provide us with necessary and sufficient stability conditions.

Theorem 11.5

The linear time-invariant discrete-time system

$$\mathbf{x}(k+1) = A_d \mathbf{x}(k), \quad k = 0, 1, 2, \dots \quad (11.43)$$

is asymptotically stable if and only if for any positive definite matrix Q , there exists a unique positive definite solution P to the discrete Lyapunov equation

$$A_d^T P A_d - P = -Q \quad (11.44)$$

Proof

We drop the subscript d for brevity.

Sufficiency

Consider the Lyapunov function

$$V(\mathbf{x}(k)) = \mathbf{x}^T(k)P\mathbf{x}(k)$$

with P a positive definite matrix. The change in the Lyapunov function along the trajectories of the system is

$$\begin{aligned}\Delta V(k) &= V(\mathbf{x}(k+1)) - V(\mathbf{x}(k)) \\ &= \mathbf{x}^T(k)[A^T P A - P]\mathbf{x}(k) = -\mathbf{x}^T(k)Q\mathbf{x}(k) < 0\end{aligned}$$

Hence, the system is stable by Theorem 11.3.

Necessity

We first show that the solution of the Lyapunov equation is given by

$$P = \sum_{k=0}^{\infty} (A^T)^k Q A^k$$

This is easily verified by substitution in the Lyapunov equation and then changing the index of summation as follows:

$$A^T \left[\sum_{k=0}^{\infty} (A^T)^k Q A^k \right] A - \sum_{k=0}^{\infty} (A^T)^k Q A^k = \sum_{j=1}^{\infty} (A^T)^j Q A^j - \sum_{k=0}^{\infty} (A^T)^k Q A^k = -Q$$

To show that for any positive definite Q , P is positive definite, consider the quadratic form

$$\mathbf{x}^T(0)Px(0) = \sum_{k=0}^{\infty} \mathbf{x}^T(0)(A^T)^k Q A^k \mathbf{x}(0) = \sum_{k=1}^{\infty} \mathbf{x}^T(k)Q\mathbf{x}(k)$$

For positive definite Q , each term on the right hand side is positive for any nonzero $\mathbf{x}(0)$. It follows that P is positive definite.

Let the system be stable but for some positive definite Q there is no finite solution P to the Lyapunov equation. We show that this leads to a contradiction.

Recall that the state-transition matrix of the discrete system can be written as

$$A^k = \sum_{i=1}^n Z_i \lambda_i^k \quad (11.45)$$

Let λ_{\max} be the spectral radius of the state matrix and let the largest norm of its constituent matrices be $\|Z\|_{\max}$. We use any matrix norm to obtain the inequality

$$\|A^k\| = \left\| (A^T)^k \right\| \leq n \|Z\|_{\max} \lambda_{\max}^k = \alpha \lambda_{\max}^k, \quad k \geq 0$$

where $\alpha = n\|Z\|$. For a stable system, $|\lambda_{\max}| < 1$ and we have

Proof—cont'd

$$\begin{aligned}\|P\| &= \left\| \sum_{k=0}^{\infty} (A^T)^k Q A^k \right\| \leq \sum_{k=0}^{\infty} \| (A^T)^k \| \|Q\| \|A^k\| \\ &\leq \sum_{k=0}^{\infty} \alpha^2 \lambda_{\max}^{2k} \|Q\| = \frac{\alpha^2}{1 - \lambda_{\max}^{2k}} \|Q\|\end{aligned}$$

This contradicts the assumption and establishes necessity.

Uniqueness

Let P_1 be a second solution of the Lyapunov equation. Then we can write it in the form of infinite summation including Q , and then substitute for Q in terms of P using the Lyapunov equation. This yields the equation

$$\begin{aligned}P_1 &= \sum_{k=0}^{\infty} (A^T)^k Q A^k = - \sum_{k=0}^{\infty} (A^T)^k [A^T P A - P] A^k \\ &= - \sum_{j=1}^{\infty} (A^T)^j P A^j + \sum_{k=0}^{\infty} (A^T)^k P A^k = P\end{aligned}$$

As in the nonlinear case, it is possible to relax the stability condition as follows.

Corollary 11.2

The linear time-invariant discrete-time system of (11.43) is asymptotically stable if and only if for any detectable pair (A_d, C) there exists a unique positive definite solution P to the discrete Lyapunov equation

$$A_d^T P A_d - P = -C^T C \quad (11.46)$$

Proof

The proof follows the same steps as the theorem. We only show that zero values of the difference do not impact stability. We first obtain

$$\begin{aligned}\Delta V(k) &= V(\mathbf{x}(k+1)) - V(\mathbf{x}(k)) \\ &= \mathbf{x}^T(k) [A^T P A - P] \mathbf{x}(k) \\ &= -\mathbf{x}^T(k) C^T C \mathbf{x}(k) = -\mathbf{y}^T(k) \mathbf{y}(k) \leq 0\end{aligned}$$

If (A_d, C) is observable, $\mathbf{y}(k) = C\mathbf{x}(k)$ is zero only if $\mathbf{x}(k)$ is zero. From the eigenvector test for observability, if the pair is only detectable, then the only nonzero values of $\mathbf{x}(k)$ for which $\mathbf{y}(k)$

Proof—cont'd

is zero are $\mathbf{x}(k) = \mathbf{v}$ where \mathbf{v} is an eigenvector corresponding to a stable mode. A trajectory starting on such an eigenvector will decay to zero asymptotically while remaining on the eigenvector and can be ignored in our stability analysis. We can then conclude that the system is asymptotically stable even though it can reach a point on \mathbf{v} for which $\Delta V \leq 0$.

Although using a semidefinite matrix in the stability test may simplify computation, checking the pair (A_d, C) for detectability (guaranteed by observability) eliminates the gain from the simplification. However, the corollary is of theoretical interest and helps clarify the properties of the discrete Lyapunov equation.

The discrete Lyapunov equation is clearly a linear equation in the matrix P , and by rearranging terms we can write it as a linear system of equations. The equivalent linear system involves an $n^2 \times 1$ vector of unknown entries of P obtained using the operation

$$\begin{aligned} st(P) &= \text{col}\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n\} \\ P &= [\mathbf{p}_1 \quad \mathbf{p}_2 \quad \dots \quad \mathbf{p}_n] \end{aligned}$$

The $n^2 \times n^2$ coefficient matrix of the linear system is obtained using the Kronecker product of matrices. In this operation, each entry of the first matrix is replaced by the second matrix scaled by the original entry. The Kronecker product is thus defined by

$$A \otimes B = [a_{ij}B] \quad (11.47)$$

It can be shown that the Lyapunov Eq. (11.44) is equivalent to the linear system

$$\begin{aligned} L\mathbf{p} &= -\mathbf{q} \\ \mathbf{p} &= st(P) \\ L &= A^T \otimes A^T - I \otimes I \\ \mathbf{q} &= st(Q) \end{aligned} \quad (11.48)$$

We only need to solve for the upper half of the matrix P because of its symmetry. In other words, we need to solve $\sum_{i=1}^n i = n(n+1)/2$ rather than n^2 equations for a $n \times n$. However, extracting the linearly independent equations from the linear system is not simple.

Because the eigenvalues of any matrix are identical to those of its transpose, the Lyapunov equation can be written in the form

$$A_d P A_d^T - P = -Q \quad (11.49)$$

The first form of the Lyapunov equation of (11.44) is the **controller form**, whereas the second form of (11.49) is known as the **observer form**.

11.4.5 MATLAB

To solve the discrete Lyapunov equation using MATLAB, we use the command **dlyap**. The command solves the observer form of the Lyapunov equation.

$$\gg \mathbf{P} = \text{dlyap}(\mathbf{A}, \mathbf{Q})$$

To solve the controller form, we simply replace the matrix A with its transpose A' . To solve the equivalent linear system, we use the Kronecker product and the command

$$\gg \text{kron}(\mathbf{A}', \mathbf{A}')$$

Example 11.10

Use the Lyapunov approach with $\mathbf{Q} = I_3$ to determine the stability of the linear time-invariant system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \\ x_3(k+1) \end{bmatrix} = \begin{bmatrix} -0.2 & -0.2 & 0.4 \\ 0.5 & 0 & 1 \\ 0 & -0.4 & -0.5 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \\ x_3(k) \end{bmatrix}$$

Is it possible to investigate the stability of the system using $\mathbf{Q} = \text{diag}\{0, 0, 1\}$?

Solution

Using MATLAB, we obtain the solution as follows:

$$\gg \mathbf{P} = \text{dlyap}(\mathbf{A}, \mathbf{Q}) \quad % \text{Observer form of the Lyapunov equation}$$

$$\begin{aligned} \mathbf{P} = & \begin{bmatrix} 1.5960 & 0.5666 & 0.0022 \\ 0.5666 & 3.0273 & -0.6621 \\ 0.0022 & -0.6621 & 1.6261 \end{bmatrix} \end{aligned}$$

$$\gg \text{eig}(\mathbf{P}) \quad % \text{Check the signs of the eigenvalues of P}$$

$$\begin{aligned} \text{ans} = & \begin{bmatrix} 1.1959 \\ 1.6114 \\ 3.4421 \end{bmatrix} \end{aligned}$$

Because the eigenvalues of P are all positive, the matrix is positive definite and the system is asymptotically stable.

For $\mathbf{Q} = \text{diag}\{0, 0, 1\} = \mathbf{C}^T \mathbf{C}$, with $\mathbf{C} = [0, 0, 1]$, we check the observability of the pair (A, C) using MATLAB and the rank test

$$\begin{aligned} \gg \text{rank}(\text{obsv}(\mathbf{a}, [0, 0, 1])) \quad & % \text{obsv computes the observability matrix} \\ \text{ans} = & 3 \end{aligned}$$

The observability matrix is full rank, and the system is observable. Thus, we can use the matrix to check the stability of the system.

11.4.6 Lyapunov's linearization method

It is possible to characterize the stability of an equilibrium for a nonlinear system by examining its approximately linear behavior in the vicinity of the equilibrium. Without loss of generality, we assume that the equilibrium is at the origin and rewrite the state equation in the form

$$\mathbf{x}(k+1) = \frac{\partial \mathbf{f}[\mathbf{x}(k)]}{\partial \mathbf{x}(k)} \Big|_{\mathbf{x}(k)=0} \mathbf{x}(k) + \mathbf{f}_2[\mathbf{x}(k)], \quad k = 0, 1, 2, \dots \quad (11.50)$$

where $\mathbf{f}_2[\cdot]$ is a function including all terms of order higher than one. We then rewrite the equation in the form

$$\begin{aligned} \mathbf{x}(k+1) &= A\mathbf{x}(k) + \mathbf{f}_2[\mathbf{x}(k)], \quad k = 0, 1, 2, \dots \\ A &= \frac{\partial \mathbf{f}[\mathbf{x}(k)]}{\partial \mathbf{x}(k)} \Big|_{\mathbf{x}(k)=0} \end{aligned} \quad (11.51)$$

In the vicinity of the origin the behavior is approximately the same as that of the linear system

$$\mathbf{x}(k+1) = A\mathbf{x}(k) \quad (11.52)$$

This intuitive fact can be more rigorously justified using Lyapunov stability theory, but we omit the proof. Thus, the equilibrium is stable if the linear approximation is stable and unstable if the linear approximation is unstable. If the linear system (11.52) has an eigenvalue on the unit circle, then the stability of the nonlinear system cannot be determined from the first-order approximation alone. This is because higher-order terms can make the nonlinear system either stable or unstable. Based on our discussion, we have the following theorem.

Theorem 11.6

The equilibrium point of the nonlinear system of (11.51) with linearized model (11.52) is as follows:

- Asymptotically stable if all the eigenvalues of A are inside the unit circle.
- Unstable if one or more of the eigenvalues are outside the unit circle.
- If A has one or more eigenvalues on the unit circle, then the stability of the nonlinear system cannot be determined from the linear approximation.

Example 11.11

Show that the origin is an unstable equilibrium for the system

$$x_1(k+1) = 2x_1(k) + 0.08x_2^2(k)$$

$$x_2(k+1) = x_1(k) + 0.1x_2(k) + 0.3x_1(k)x_2(k), \quad k = 0, 1, 2, \dots$$

Solution

We first rewrite the state equations in the form

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 1 & 0.1 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 0.08x_2^2(k) \\ 0.3x_1(k)x_2(k) \end{bmatrix}, \quad k = 0, 1, 2, \dots$$

The state matrix A of the linear approximation has one eigenvalue $= 2 > 1$. Hence, the origin is an unstable equilibrium of the nonlinear system.

11.4.7 Instability theorems

The weakness of the Lyapunov approach for nonlinear stability investigation is that it only provides sufficient stability conditions. Although no necessary and sufficient conditions are available, it is possible to derive conditions for instability and use them to test the system if one is unable to establish its stability. Clearly, failure of both sufficient conditions is inconclusive, and it is only in the linear case that we have the stronger necessary and sufficient condition.

Theorem 11.7

The equilibrium point $\mathbf{x} = \mathbf{0}$ of the nonlinear discrete-time system

$$\mathbf{x}(k+1) = \mathbf{f}[\mathbf{x}(k)], \quad k = 0, 1, 2, \dots \quad (11.53)$$

is unstable if there exists a locally positive function for the system with locally uniformly positive definite changes along the trajectories of the system.

Proof

For any motion along the trajectories of the system, we have

$$\begin{aligned} \Delta V(k) &= V(\mathbf{x}(k+1)) - V(\mathbf{x}(k)) \\ &= V(\mathbf{f}[\mathbf{x}(k)]) - V(\mathbf{f}[\mathbf{x}(k-1)]) > 0, \quad k = 0, 1, 2, \dots \end{aligned}$$

Proof—cont'd

This implies that as the motion progresses, the value of V increases continuously. However, because V is only zero with argument zero, the trajectories of the system will never converge to the equilibrium at the origin and the equilibrium is unstable.

Example 11.12

Show that the origin is an unstable equilibrium for the system

$$\begin{aligned}x_1(k+1) &= -2x_1(k) + 0.08x_2^2(k) \\x_2(k+1) &= 0.3x_1(k)x_2(k) + 2x_2(k), \quad k = 0, 1, 2, \dots\end{aligned}$$

Solution

Choose the Lyapunov function

$$V(\mathbf{x}) = \mathbf{x}^T P \mathbf{x}, \quad P = \text{diag}\{p_1, 1\}$$

The corresponding difference is given by

$$\begin{aligned}\Delta V(k) &= p_1 \left\{ [-2x_1(k) + 0.08x_2^2(k)]^2 - x_1^2(k) \right\} + \left\{ [0.3x_1(k)x_2(k) + 2x_2(k)]^2 - x_2^2(k) \right\} \\&= 3p_1x_1^2(k) + 0.0064p_1x_2^4(k) + \{0.09x_1^2(k) + (1.2 - 0.32p_1)x_1(k) + 3\}x_2^2(k), \\k &= 0, 1, 2, \dots\end{aligned}$$

We complete the squares for the last term by choosing $p_1 = 3.75(1 - \sqrt{3}/2) = 0.5024$ and reduce the difference to

$$\begin{aligned}\Delta V(k) &= 1.5072x_1^2(k) + 0.003x_2^4(k) + (0.3x_1(k) + \sqrt{3})^2 x_2^2(k) \\&\geq 1.5072x_1^2(k) + 0.0032x_2^4(k), \quad k = 0, 1, 2, \dots\end{aligned}$$

The inequality follows from the fact that the last term in ΔV is positive semidefinite. We conclude that ΔV is positive definite because it is greater than the sum of even powers and that the equilibrium at $\mathbf{x} = \mathbf{0}$ is unstable.

11.4.8 Estimation of the domain of attraction

We consider a system

$$\mathbf{x}(k+1) = A\mathbf{x} + \mathbf{f}[\mathbf{x}(k)], \quad k = 0, 1, 2, \dots \quad (11.54)$$

where the matrix A is stable and $\mathbf{f}(\cdot)$ includes second-order terms and higher and satisfies the inequality

$$\|\mathbf{f}[\mathbf{x}(k)]\| \leq \alpha \|\mathbf{x}(k)\|^2, \quad k = 0, 1, 2, \dots \quad (11.55)$$

for some constant $\alpha > 0$. Because the linear approximation of the system is stable, we can solve the associated Lyapunov equation for a positive definite matrix P with any positive definite matrix Q . This yields the Lyapunov function

$$V(\mathbf{x}(k)) = \mathbf{x}^T(k) P \mathbf{x}(k)$$

and the difference

$$\begin{aligned} \Delta V(\mathbf{x}(k)) &= \mathbf{x}^T(k) [A^T P A - P] \mathbf{x}(k) + 2\mathbf{f}^T[\mathbf{x}(k)] P A \mathbf{x}(k) + \mathbf{f}^T[\mathbf{x}(k)] P \mathbf{f}[\mathbf{x}(k)] \\ &= -\mathbf{x}^T(k) Q \mathbf{x}(k) + 2\mathbf{f}^T[\mathbf{x}(k)] P A \mathbf{x}(k) + \mathbf{f}^T[\mathbf{x}(k)] P \mathbf{f}[\mathbf{x}(k)] \end{aligned}$$

To simplify this expression, we make use of the inequalities

$$\begin{aligned} \lambda_{\min}(P) \|\mathbf{x}\|^2 &\leq \mathbf{x}^T P \mathbf{x} \leq \lambda_{\max}(P) \|\mathbf{x}\|^2 \\ \|\mathbf{f}^T P A \mathbf{x}\| &\leq \|\mathbf{f}\| \|P A\| \|\mathbf{x}\| \leq \alpha \|P A\| \|\mathbf{x}\|^3 \end{aligned} \quad (11.56)$$

where the norm $\|A\|$ denotes the square root of the largest eigenvalue of the matrix $A^T A$ or its largest **singular value** (see **Appendix III**).

Using the inequality (11.56), we have

$$\Delta V(k) \leq [\alpha^2 \lambda_{\max}(P) \|\mathbf{x}(k)\|^2 + 2\alpha \|P A\| \|\mathbf{x}(k)\| - \lambda_{\min}(Q)] \|\mathbf{x}(k)\|^2$$

Note that, because of the negative sign of the quadratic term including Q , its upper bound is in terms of the minimum rather than the maximum of Q .

The difference $\Delta V(k)$ is negative if the following quadratic is negative

$$\alpha^2 \lambda_{\max}(P) \|\mathbf{x}(k)\|^2 + 2\alpha \|P A\| \|\mathbf{x}(k)\| - \lambda_{\min}(Q) \quad (11.57)$$

Because the coefficient of the quadratic is positive, its second derivative is also positive and it has negative values between its two roots. The positive root defines a bound on $\mathbf{x}(k)$ that guarantees a negative difference:

$$\|\mathbf{x}\| < \frac{1}{\alpha \lambda_{\max}(P)} \left[-\|P A\| + \sqrt{\|P A\|^2 + \lambda_{\max}(P) \lambda_{\min}(Q)} \right]$$

An estimate of the domain of attraction is given by

$$B(\|\mathbf{x}\|) = \left\{ \mathbf{x}: \|\mathbf{x}\| < \frac{1}{\alpha \lambda_{\max}(P)} \left[-\|P A\| + \sqrt{\|P A\|^2 + \lambda_{\max}(P) \lambda_{\min}(Q)} \right] \right\} \quad (11.58)$$

Example 11.13

Estimate the domain of attraction of the system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 0.1 & 0 \\ -1 & 0.5 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 0.25x_2^2(k) \\ 0.1x_1^2(k) \end{bmatrix}$$

Solution

The nonlinear vector satisfies

$$\|\mathbf{f}[\mathbf{x}(k)]\| \leq 0.27\|\mathbf{x}(k)\|, \quad k = 0, 1, 2, \dots$$

We solve the Lyapunov equation with $Q = I_2$ to obtain

$$P = \begin{bmatrix} 2.4987 & -0.7018 \\ -0.7018 & 1.3333 \end{bmatrix}$$

whose largest eigenvalue is equal to 2.8281. The norm $\|PA\| = 1.8539$ can be computed with the MATLAB command

`>> norm(P*A)`

Our estimate of the domain of attraction is

$$\begin{aligned} B(\mathbf{x}) &= \left\{ \mathbf{x}: \|\mathbf{x}\| < \frac{1}{0.25 \times \sqrt{2.8281}} \left[-1.8539 + \sqrt{(1.8539)^2 + 0.25} \right] \right\} \\ &= \{ \mathbf{x}: \|\mathbf{x}\| < 0.0937 \} \end{aligned}$$

The state portrait presented in Fig. 11.2 shows that the estimate of the domain of attraction is quite conservative and that the system is stable well outside the estimated region.

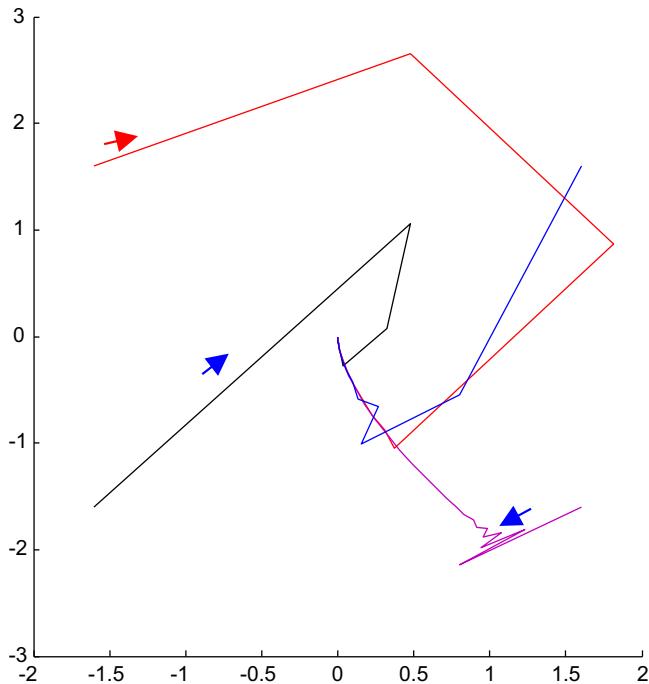


Figure 11.2
Phase portrait for the nonlinear system of Example 11.13.

11.5 Stability of analog systems with digital control

Although most nonlinear design methodologies are analog, they are often digitally implemented. The designer typically completes an analog controller design and then obtains a discrete approximation of the analog controller. For example, a discrete approximation can be obtained using a differencing approach to approximate the derivatives (see Chapter 6). The discrete-time approximation is then used with the analog plant, assuming that the time response will be approximately the same as that of the analog design. This section examines the question “Is the stability of the digital control system guaranteed by the stability of the original analog design?” As in the linear case discussed in Chapter 6, the resulting digital control system may be unstable or may have unacceptable intersample oscillations.

We first examine the system of (11.1) with the digital control (11.2). Substituting (11.2) in (11.1) gives the equation

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + B(\mathbf{x})\mathbf{u}(k) = \mathbf{f}_k(\mathbf{x}), \quad t \in [kT, (k+1)T] \quad (11.59)$$

Let the solution of (11.59) be

$$\mathbf{x}(t) = \mathbf{g}_k(\mathbf{x}(kT), t) \quad (11.60)$$

where \mathbf{g}_k is a continuous function of all its arguments. Then at the sampling points, we have the discrete model

$$\mathbf{x}(k+1) = \mathbf{g}_k(\mathbf{x}(k)), \quad k = 0, 1, \dots \quad (11.61)$$

Theorem 11.8 shows that the stability of the discrete-time system of (11.61) is a necessary condition for the stability of the analog system with digital control.

Theorem 11.8

Analog System with Digital Control

The continuous-time system of (11.59) with piecewise constant control is exponentially stable only if the discrete-time system of (11.61) is exponentially stable.

Proof

We show that the discrete-time system of (11.61) is stable for any stable system in (11.59). For an exponentially stable continuous-time system, we have

$$\|\mathbf{x}(t)\| \leq \|\mathbf{x}(k)\|e^{-\alpha_k t}, \quad t \in [kT, (k+1)T] \quad (11.62)$$

for some positive constant α_k . Hence, at the sampling points we have the inequality

$$\|\mathbf{x}(k)\| \leq \|\mathbf{x}(k-1)\|e^{-\alpha_{k-1} T}, \quad k = 0, 1, \dots$$

Proof—cont'd

Applying this inequality repeatedly gives

$$\begin{aligned}\|\mathbf{x}(k)\| &\leq \|\mathbf{x}(0)\| \prod_{i=0}^{k-1} e^{-\alpha_i T} \\ &\leq \|\mathbf{x}(0)\| e^{-\alpha_m kT}, \quad k = 0, 1, 2, \dots\end{aligned}\tag{11.63}$$

with $\alpha_m = \min_k \alpha_k$. Thus, the discrete-time system is exponentially stable.

Theorem 11.8 only provides a necessary stability condition for the discretization of an analog plant with a digital controller. If a stable analog controller is implemented digitally and used with the analog plant, its stability cannot be guaranteed. The resulting system may be unstable, as Example 11.14 demonstrates.

Example 11.14

Consider the system

$$\dot{x}(t) = x^2(t) + x(t)u(t)$$

with the continuous control

$$u(t) = -x(t) - \alpha x^2(t), \quad \alpha > 0$$

Then the closed-loop system is

$$\dot{x}(t) = -\alpha x^3(t)$$

which is asymptotically stable.¹

Now consider the digital implementation of the analog control

$$u(t) = -x(k) - \alpha x^2(k) \quad \alpha > 0, \quad t \in [kT, (k+1)T)$$

and the corresponding closed-loop system

$$\dot{x}(t) = x(t) \{x(t) - [\alpha x(k) + 1]x(k)\}, \quad t \in [kT, (k+1)T]$$

We solve the equation by separation of variables:

$$\begin{aligned}\frac{dx(t)}{x(t)(x(t) - \beta)} &= \frac{dx(t)}{\beta} \left[\frac{1}{x - \beta} - \frac{1}{x} \right] = dt, \quad t \in [kT, (k+1)T) \\ \beta &= [\alpha x(k) + 1]x(k)\end{aligned}$$

Integrating from kT to $(k+1)T$, we obtain

$$\begin{aligned}\ln \left(1 - \frac{\beta}{x(t)} \right) \Big|_{t=kT}^{t=(k+1)T} &= dt \Big|_{t=kT}^{t=(k+1)T} \\ \left(1 - \frac{\beta}{x(k+1)} \right) &= \left(1 - \frac{\beta}{x(k)} \right) e^T\end{aligned}$$

Example 11.14—cont'd

$$\begin{aligned}x(k+1) &= \frac{\beta}{1 - \left[1 - \frac{\beta}{x(k)}\right]e^T} \\&= \frac{[\alpha x(k) + 1]x(k)}{1 + \alpha x(k)e^T}\end{aligned}$$

For stability, we need the condition $|x(k+1)| < |x(k)|$ —that is,

$$|1 + \alpha x(k)e^T| > |1 + \alpha x(k)|$$

This condition is satisfied for positive $x(k)$ but not for all negative $x(k)$! For example, with $T = 0.01$ s and $\alpha x(k) = -0.5$, the LHS is 0.495 and the RHS is 0.5.

¹ Slotine (1991, p. 66): $\dot{x} + c(x) = 0$ is asymptotically stable if $c(\cdot)$ is continuous and satisfies $c(x)x > 0$, $\forall x \neq 0$. Here $c(x) = x^3$.

11.6 State–plane analysis

As shown in [Section 11.4.6](#), the stability of an equilibrium of a nonlinear system can often be determined from the linearized model of the system in the vicinity of the equilibrium. Moreover, the behavior of system trajectories of linear discrete-time systems in the vicinity of an equilibrium point can be visualized for second-order systems in the state plane. State–plane trajectories can be plotted based on the solutions of the state equations for discrete-time equations similar to those for continuous-time systems (see Chapter 7). Consider the unforced second-order difference equation

$$y(k+2) + a_1 y(k+1) + a_0 y(k) = 0, \quad k = 0, 1, \dots \quad (11.64)$$

The associated characteristic equation is

$$z^2 + a_1 z + a_0 = (z - \lambda_1)(z - \lambda_2) = 0 \quad (11.65)$$

We can characterize the behavior of the system based on the location of the characteristic values of the system λ_i , $i = 1, 2$ in the complex plane. If the system is represented in state–space form, then a similar characterization is possible using the eigenvalues of the state matrix. [Table 11.1](#) gives the names and characteristics of equilibrium points based on the locations of the eigenvalues.

Table 11.1: Equilibrium point classification.

Equilibrium type	Eigenvalue location
Stable node	Real positive inside unit circle
Unstable node	Real positive outside unit circle
Saddle point	Real eigenvalues with one inside and one outside unit circle
Stable focus	Complex conjugate or both real negative inside unit circle
Unstable focus	Complex conjugate or both real negative outside unit circle
Vortex or center	Complex conjugate on unit circle

Trajectories corresponding to different types of equilibrium points are shown in [Figs. 11.3–11.10](#). We do not include some special cases such as two real eigenvalues equal to unity, in which case the system always remains at the initial state. The reader is invited to explore other pairs of eigenvalues through MATLAB simulation.

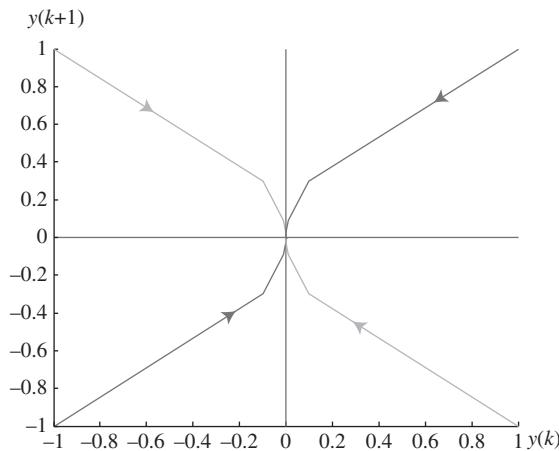


Figure 11.3
Stable node (eigenvalues 0.1, 0.3).

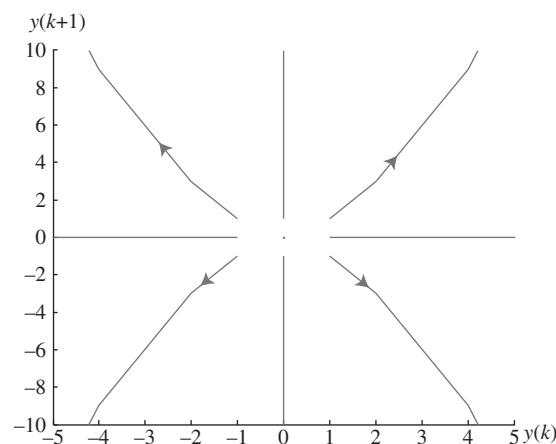


Figure 11.4
Unstable node (eigenvalues 2, 3).

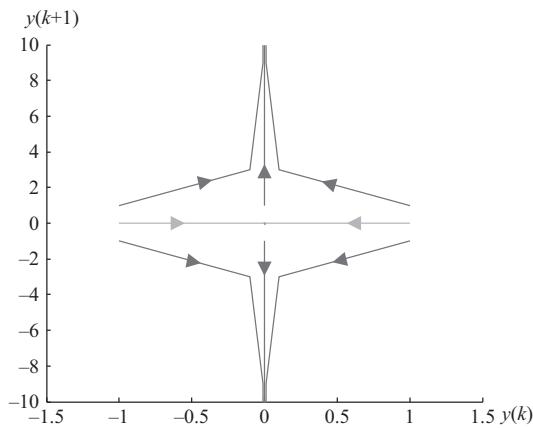


Figure 11.5
Saddle point (eigenvalues 0.1, 3).

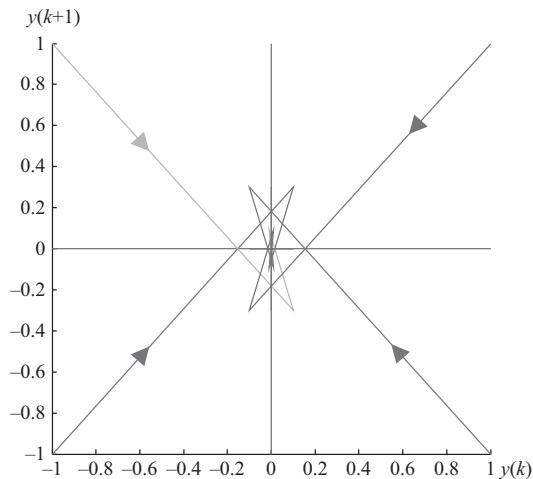


Figure 11.6
Stable focus (eigenvalues $-0.1, -0.3$).

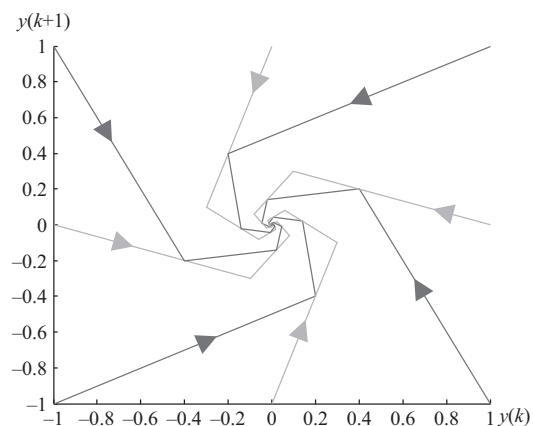


Figure 11.7
Stable focus (eigenvalues $0.1 \pm j0.3$).

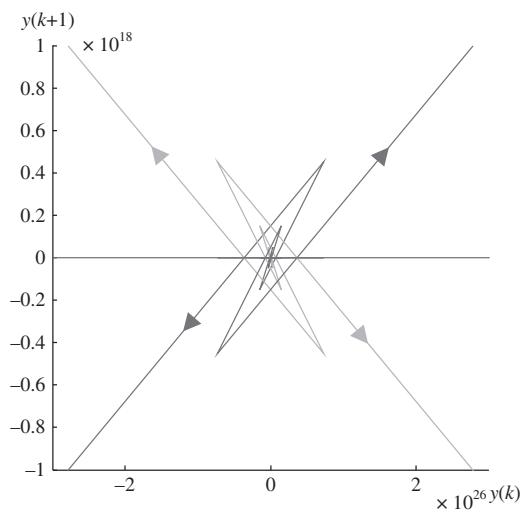


Figure 11.8
Unstable focus (eigenvalues $-5, -3$).

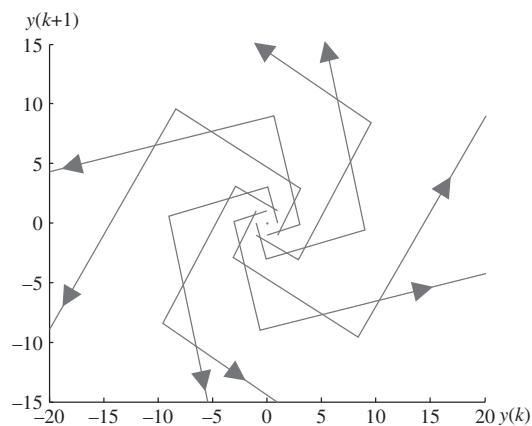


Figure 11.9
Unstable focus (eigenvalues $0.1 \pm j3$).

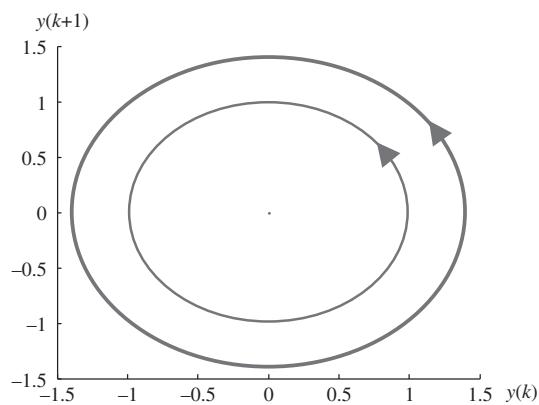


Figure 11.10
Center (eigenvalues $\cos(45^\circ) \pm j \sin(45^\circ)$).

11.7 Discrete-time nonlinear controller design

Nonlinear control system design for discrete-time systems is far more difficult than linear design. It is also often more difficult than the design of analog nonlinear systems. For example, one of the most powerful approaches to nonlinear system design is to select a control that results in a closed-loop system for which a suitable Lyapunov function can be constructed. It is usually easier to construct a Lyapunov function for an analog nonlinear system than it is for a discrete-time system.

We discuss some simple approaches to nonlinear design for discrete-time systems that are possible for special classes of systems.

11.7.1 Controller design using extended linearization

If one of the extended linearization approaches presented in Section 11.1 is applicable, then we can obtain a linear discrete-time model for the system and use it in linear control system design. The nonlinear control can then be recovered from the linear design. We demonstrate this approach with Example 11.15.

Example 11.15

Consider the mechanical system

$$\ddot{x} + b(\dot{x}) + c(x) = f$$

with $b(0) = 0$, $c(0) = 0$. For the nonlinear damping $b(\dot{x}) = 0.1\dot{x}^3$ and the nonlinear spring $c(x) = 0.1x + 0.01x^3$, design a digital controller for the system by redefining the input and discretization. Use a sampling period $T = 0.02$ s.

Solution

The state equations of the system are

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= -b(x_2) - c(x_1) + f = u\end{aligned}$$

With input u , the system reduces to a double integrator. Using the results for the discretization of the equivalent linear model obtained in Example 11.2 with $T = 0.02$, we have

$$\mathbf{x}(k+1) = \begin{bmatrix} 1 & 0.02 \\ 0 & 1 \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} 2 \times 10^{-4} \\ 0.02 \end{bmatrix} \mathbf{u}(k)$$

with $\mathbf{x}(k) = [x_1(k) \ x_2(k)]^T = [x(k) \ x(k+1)]^T$.

We select the eigenvalues $\{0.1 \pm j0.1\}$ and design a state feedback controller for the system as shown in Chapter 9. Using the MATLAB command `place`, we obtain the feedback gain matrix

$$\mathbf{k}^T = [2050 \quad 69.5]$$

Example 11.15—cont'd

For a reference input $r(k)$, we have the nonlinear control

$$f(k) = u(k) + b(x_2(k)) + c(x_1(k))$$

$u(k) = r(k) - \mathbf{k}^T \mathbf{x}(k)$. The closed-loop dynamics is

$$\begin{aligned}\mathbf{x}(k+1) &= \begin{bmatrix} 1 & 0.02 \\ 0 & 1 \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} 2 \times 10^{-4} \\ 0.02 \end{bmatrix} (r(k) - [2050 \quad 69.5] \mathbf{x}(k)) \\ &= \begin{bmatrix} 0.5900 & 0.0061 \\ -41.0000 & -0.3900 \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} 2 \times 10^{-4} \\ 0.02 \end{bmatrix} r(k)\end{aligned}$$

The simulation diagram for the system is shown in Fig. 11.11, and the simulation diagram for the controller block is shown in Fig. 11.12. For a steady-state position of unity, we use a step input r and we select its amplitude using the equilibrium condition

$$\mathbf{x}(k) = \begin{bmatrix} 0.5900 & 0.0061 \\ -41.0000 & -0.3900 \end{bmatrix} \mathbf{x}(k) + \begin{bmatrix} 2 \times 10^{-4} \\ 0.02 \end{bmatrix} r = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

This simplifies to

$$\begin{bmatrix} 2 \times 10^{-4} \\ 0.02 \end{bmatrix} r = \begin{bmatrix} 0.41 \\ 41 \end{bmatrix}$$

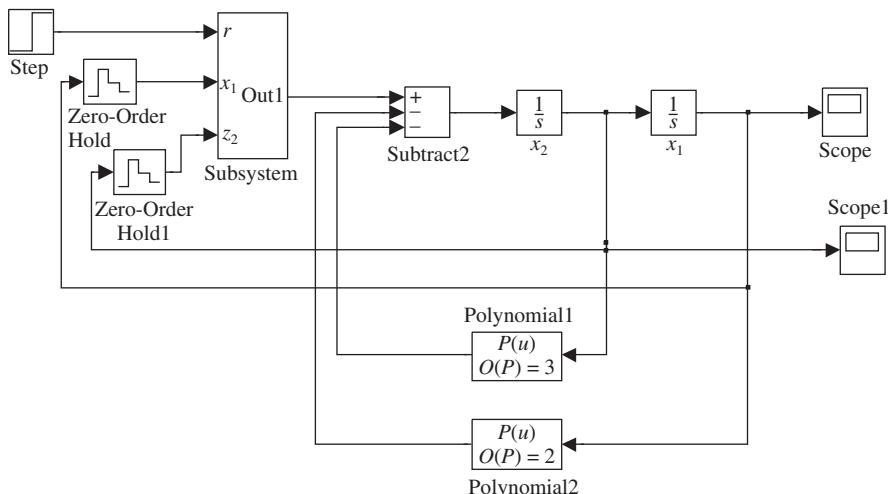


Figure 11.11
Simulation diagram for the system of Example 11.15.

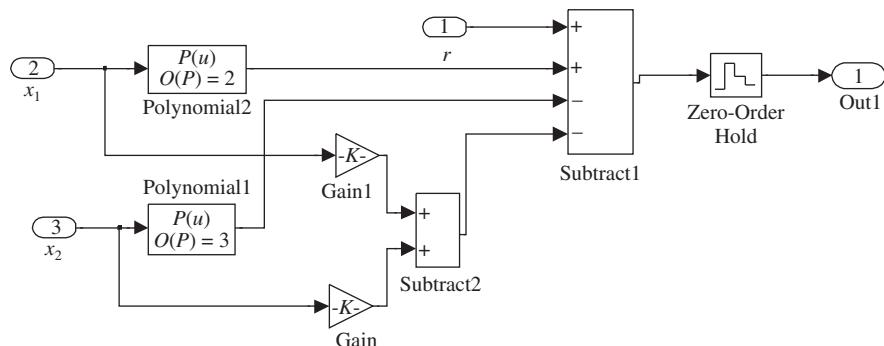
Example 11.15—cont'd

Figure 11.12
Controller simulation diagram of Example 11.15.

which gives the amplitude $r = 2050$. The step response for the nonlinear system with digital control in Fig. 11.13 shows a fast response to a step input at $t = 0.2$ s that quickly settles to the desired steady-state value of unity. Fig. 11.14 shows a plot of the velocity for the same input. The velocity increases sharply, reverses direction, then goes to zero to approach the desired steady-state position.

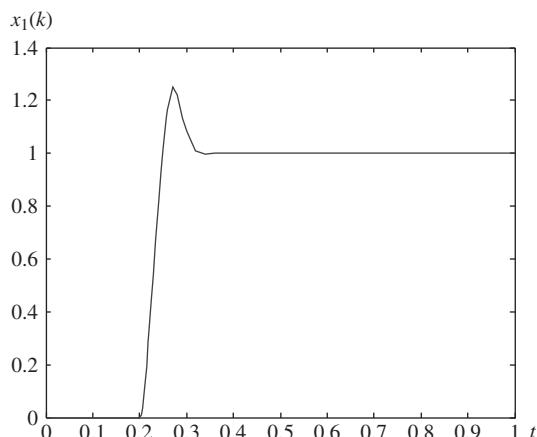
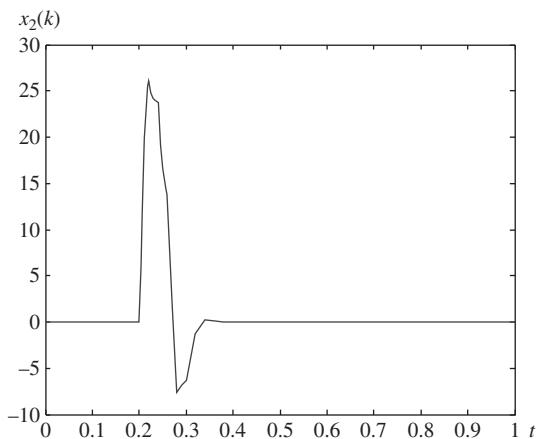


Figure 11.13
Step response for the linear design of Example 11.15.

Example 11.15—cont'd**Figure 11.14**

Velocity plot for the step response for the linear design of Example 11.15.

11.7.2 Controller design based on Lyapunov stability theory

Lyapunov stability theory provides a means of stabilizing unstable nonlinear systems using feedback control. The idea is that if one can select a suitable Lyapunov function and force it to decrease along the trajectories of the system, the resulting system will converge to its equilibrium. In addition, the control can be chosen to speed up the rate of convergence to the origin by forcing the Lyapunov function to decrease to zero faster.

To simplify our analysis, we consider systems of the form

$$\mathbf{x}(k+1) = A\mathbf{x}(k) + B(\mathbf{x}(k))\mathbf{u}(k) \quad (11.66)$$

and assume that the eigenvalues of the matrix A are all inside the unit circle. This model could approximately represent a nonlinear system in the vicinity of its stable equilibrium.

Theorem 11.9

The open-loop stable affine system with linear unforced dynamics and full-rank input matrix for all nonzero \mathbf{x} is asymptotically stable with the feedback control law

$$\mathbf{u}(k) = -[B^T(\mathbf{x}(k))PB(\mathbf{x}(k))]^{-1}B^T(\mathbf{x}(k))PA\mathbf{x}(k) \quad (11.67)$$

Theorem 11.9—cont'd

where P is the solution of the discrete Lyapunov equation

$$A^T P A - P = -Q \quad (11.68)$$

and Q is an arbitrary positive definite matrix.

Proof

For the Lyapunov function

$$V(\mathbf{x}(k)) = \mathbf{x}^T(k)P\mathbf{x}(k)$$

the difference is given by

$$\begin{aligned}\Delta V(k) &= V(\mathbf{x}(k+1)) - V(\mathbf{x}(k)) \\ &= \mathbf{x}^T(k)[A^T P A - P]\mathbf{x}(k) + 2\mathbf{u}^T(k)B^T P A \mathbf{x} + \mathbf{u}^T(k)B^T P B \mathbf{u}(k)\end{aligned}$$

where the argument of B is suppressed for brevity. We minimize the function with respect to the control to obtain

$$\frac{\partial \Delta V(k)}{\partial \mathbf{u}(k)} = 2[B^T P A \mathbf{x} + B^T P B \mathbf{u}(k)] = 0$$

and solve for the feedback control law of (11.67) using the full-rank condition for B .

By the assumption of open-loop stability, we have

$$\Delta V(k) = -\mathbf{x}^T(k)Q\mathbf{x}(k) + \mathbf{x}^T(k)A^T P B \mathbf{u}(k), \quad Q > 0$$

We substitute for the control and rewrite the equation as

$$\begin{aligned}\Delta V(k) &= -\mathbf{x}^T(k)\{Q + A^T P B [B^T P B]^{-1} B^T P A\}\mathbf{x}(k) \\ &= -\mathbf{x}^T(k)\{Q + A^T M A\}\mathbf{x}(k) \\ M &= P B [B^T P B]^{-1} B^T P\end{aligned}$$

Because the matrix P is positive definite and B is full rank, the matrix $B^T P B$, and consequently its inverse, must be positive definite. Thus, the matrix M is positive semidefinite, and the term $-\mathbf{x}^T(k)A^T M A \mathbf{x}(k)$ is negative semidefinite. We conclude that the difference of the Lyapunov function is negative definite, because it is the sum of a negative definite term and a negative semidefinite term, and that the closed-loop system is asymptotically stable.

Example 11.16

Design a controller to stabilize the origin for the system

$$\begin{aligned}x_1(k+1) &= 0.2x_1(k) + x_2(k)u(k) \\x_2(k+1) &= 0.2x_2(k) + [1 + 0.4x_1(k)]u(k), \quad k = 0, 1, 2, \dots\end{aligned}$$

Solution

We rewrite the system dynamics in the form

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 0.2 & 0 \\ 0 & 0.2 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} x_2(k) \\ 1 + 0.4x_1(k) \end{bmatrix} u(k)$$

The state matrix is diagonal with two eigenvalues equal to $0.2 < 1$. Hence, the matrix A is stable. We choose the Lyapunov function

$$V(\mathbf{x}) = \mathbf{x}^T P \mathbf{x}, \quad P = \text{diag}\{p_1, 1\}, \quad p_1 > 0$$

for which $Q = 0.96P$ is positive definite.

The corresponding stabilizing control is given by

$$\begin{aligned}u(k) &= -[B^T(\mathbf{x}(k))PB(\mathbf{x}(k))]^{-1}B^T(\mathbf{x}(k))PA\mathbf{x}(k) \\&= -\frac{0.2}{(1 + 0.4x_1(k))^2 + p_1x_2^2(k)} [p_1 x_2(k) \quad 1 + 0.4x_1(k)] \mathbf{x}(k)\end{aligned}$$

We choose $p_1 = 5$ and obtain the stabilizing control

$$u(k) = -\frac{1}{(1 + 0.4x_1(k))^2 + 5x_2^2(k)} [x_2(k) \quad 0.2 + 0.08x_1(k)] \mathbf{x}(k)$$

11.8 Input-output stability and the small gain theorem

Consider a nonlinear dynamic system as an operation that maps a sequence of $m \times 1$ input vectors $\mathbf{u}(k)$ to a sequence of $l \times 1$ output vector $\mathbf{y}(k)$. The equation representing this mapping is

$$\mathbf{y} = N(\mathbf{u}) \tag{11.69}$$

where $N(\cdot)$ represents the operation performed by the dynamic system, $\mathbf{u} = \text{col}\{\mathbf{u}(k), k = 0, 1, \dots, \infty\}$ represents the input sequence, and $\mathbf{y} = \text{col}\{\mathbf{y}(k), k = 0, 1, \dots, \infty\}$ represents the output sequence. We restrict our analysis of nonlinear systems to the class

of causal systems where the output at any given instant is caused by the history of system inputs and is unaffected by future inputs. Similarly to Chapter 4, we define input–output stability based on the magnitude of the output vector over its history for a bounded input vector. The magnitude of a vector over its history is defined using the following norm:

$$\|\mathbf{u}\|_2 = \sqrt{\sum_{k=0}^{\infty} \|\mathbf{u}(k)\|^2} \quad (11.70)$$

The norm of each vector in the summation can be any vector norm for the $m \times 1$ vectors $\mathbf{u}(k)$, $k = 0, 1, \dots, \infty$, but we restrict our discussion to the 2-norm. Thus, the norm of (11.70) is the 2-norm of an infinite-dimensional vector obtained by stacking the vectors $\mathbf{u}(k)$, $k = 0, 1, \dots, \infty$. We can define input–output stability for nonlinear systems as in Chapter 4.

Definition 11.8: Input–output Stability

A system is input–output stable if for any bounded input sequence

$$\|\mathbf{u}\|_2 < K_u \quad (11.71)$$

the output sequence is also bounded—that is,

$$\|\mathbf{y}\|_2 < K_y \quad (11.72)$$

For a system to be input–output stable, the system must have an upper bound on the ratio of the norm of its output to that of its input. We call this upper bound the **gain** of the nonlinear system and define it as

$$\gamma = \max_{\|\mathbf{u}\|_2} \frac{\|\mathbf{y}\|_2}{\|\mathbf{u}\|_2} \quad (11.73)$$

Example 11.17

Find the gain of a saturation nonlinearity (Fig. 11.15) governed by

$$y = N(u) = \begin{cases} Ku, |u| < L \\ KL, u \geq L \\ -KL, u \leq -L \end{cases}$$

$K > 0, L > 0$

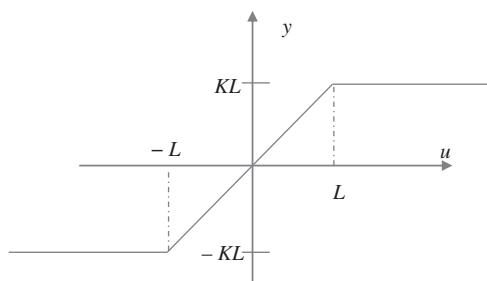
Example 11.17—cont'd

Figure 11.15
Saturation nonlinearity.

Solution

The output is bounded by

$$|y| = \begin{cases} K|u|, & |u| < L \\ KL, & |u| \geq L \end{cases}$$

The norm of the output is

$$\|\mathbf{y}\|_2^2 = \sum_{k=0}^{\infty} |y(k)|^2 \leq \sum_{k=0}^{\infty} |u(k)|^2$$

Thus, the gain of the system is the constant $\gamma_{\text{sat}} = K$.

Example 11.18

Find the gain of a causal linear system with impulse response

$$G(k), k = 0, 1, 2, \dots$$

Solution

The response of the system to any input $u(k)$, $k = 0, 1, 2, \dots$, is given by the convolution summation

$$\mathbf{y}(k) = \sum_{i=0}^{\infty} G(i)\mathbf{u}(k-i), \quad k = 0, 1, 2, \dots$$

Example 11.18—cont'd

The square of the norm of the output sequence is

$$\|\mathbf{y}\|_2^2 = \sum_{k=0}^{\infty} \|\mathbf{y}(k)\|_2^2 = \sum_{k=0}^{\infty} \left\| \sum_{i=0}^k G(i) \mathbf{u}(k-i) \right\|_2^2$$

By the triangle inequality for norms, we have

$$\|\mathbf{y}\|_2^2 = \sum_{k=0}^{\infty} \sum_{i=0}^k \|\mathbf{u}(k-i)\|_2^2 \leq \sum_{k=0}^{\infty} \sum_{i=0}^k \|G(i)\|_2^2 \|\mathbf{u}(k-i)\|_2 \|G(i)\|_2^2$$

For a causal input, $\mathbf{u}(k-i)$ is zero for $i > k$ and we can extend the summation to write

$$\|\mathbf{y}\|_2^2 \leq \sum_{i=0}^{\infty} \|G(i)\|_2^2 \sum_{k=0}^{\infty} \|\mathbf{u}(k-i)\|_2^2 = (\|G\|_2 \|\mathbf{u}\|_2)^2$$

Thus, the gain of the linear causal system is the norm of its impulse response sequence

$$\gamma_G = \|G\|_2 = \sqrt{\sum_{k=0}^{\infty} \|G(k)\|_2^2}$$

Recall from Theorem 8.2 that a linear system is said to be BIBO stable if its impulse response sequence is absolutely summable, as guaranteed if the system gain defined here is finite.

We can also express the gain of a linear system in terms of its frequency response as follows.

Theorem 11.10

The gain of a stable linear system is given by

$$\gamma_G = \|G\|_2 \leq \max_{\omega} \|G(e^{j\omega T})\| \quad (11.74)$$

Proof

Parseval's identity states that the energy of a signal in the time domain is equal to its energy in the frequency domain—that is,

Proof—cont'd

$$\|\mathbf{y}\|_2^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} \mathbf{y}(e^{-j\omega T})^T \mathbf{y}(e^{j\omega T}) d\omega \quad (11.75)$$

For a stable system, the energy of the output signal is finite, and we can use Parseval's identity to write the norm as

$$\begin{aligned} \|\mathbf{y}\|_2^2 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \mathbf{u}(e^{-j\omega T})^T G(e^{-j\omega T})^T G(e^{j\omega T}) \mathbf{u}(e^{j\omega T}) d\omega \\ &\leq \max_{\omega} \|G(e^{j\omega T})\|^2 \times \frac{1}{2\pi} \int_{-\pi}^{\pi} \mathbf{u}(e^{-j\omega T})^T \mathbf{u}(e^{j\omega T}) d\omega \\ &\leq \max_{\omega} \|G(e^{j\omega T})\|^2 \|\mathbf{u}\|_2^2 \end{aligned}$$

Using Parseval's identity on the input terms, we have

$$\|\mathbf{y}\|_2^2 \leq \max_{\omega} \|G(e^{j\omega T})\|^2 \|\mathbf{u}\|_2^2$$

Hence, the gain of the system is given by (11.74).

Example 11.19

Verify that the expression of (11.74) can be used to find the gain of the linear system with transfer function

$$G(z) = \frac{z}{z - a}, \quad |a| < 1$$

Solution

The impulse response of the linear system is

$$g(k) = a^k, \quad k = 0, 1, 2, \dots, |a| < 1$$

and its gain squared is

$$\|g\|_s^2 = \sum_{k=0}^{\infty} |g(k)|^2 = \sum_{k=0}^{\infty} |a|^{2k} = \frac{1}{1 - |a|^2}$$

Example 11.19—cont'd

The frequency response of the system is

$$G(e^{j\omega T}) = \frac{e^{j\omega T}}{e^{j\omega T} - a} = \frac{1}{1 - ae^{-j\omega T}}$$

The maximum magnitude is attained when the denominator is minimized to obtain

$$\gamma_G = \max_{\omega} G(e^{j\omega T}) = \max_{\omega} \frac{1}{|1 - e^{-j\omega T}a|} = \frac{1}{1 - |a|} \geq \frac{1}{\sqrt{1 - |a|^2}}$$

We are particularly interested in the stability of closed-loop systems of the form in Fig. 11.16, where two causal nonlinear systems are connected in a closed-loop configuration. The stability of the system is governed by Theorem 11.11.

Theorem 11.11**The Small Gain Theorem**

The closed-loop system in Fig. 11.16 is input–output stable if the product of the gains of all the systems in the loop is less than unity—that is,

$$\prod_{i=1}^2 \gamma_{N_i} < 1 \quad (11.76)$$

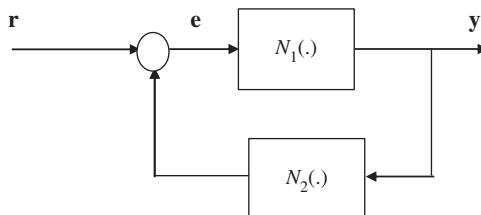


Figure 11.16
Block diagram of a nonlinear closed-loop system.

Proof

The output of the system is given by

$$\mathbf{y} = N_1(\mathbf{e})$$

Proof—cont'd

and the error is given by

$$\mathbf{e} = \mathbf{r} - N_2(\mathbf{y})$$

Using norm inequalities, we have

$$\begin{aligned}\|\mathbf{e}\| &\leq \|\mathbf{r}\| + \gamma_{N_2} \|\mathbf{y}\| \\ &\leq \|\mathbf{r}\| + \gamma_{N_2} \gamma_{N_1} \|\mathbf{e}\|\end{aligned}$$

Solving for the norm of the error, we obtain

$$\|\mathbf{e}\| \leq \frac{\|\mathbf{r}\|}{1 - \gamma_{N_2} \gamma_{N_1}}$$

which is finite if (11.76) is satisfied.

Note that Theorem 11.11 is applicable if the forward path, the feedback path, or both are replaced by a cascade of nonlinear systems so that the condition (11.76) can be generalized to n nonlinear systems in a feedback loop as

$$\prod_{i=1}^n \gamma_{Ni} < 1 \quad (11.77)$$

The results of the small gain theorem provide a conservative stability condition when applied to linear systems as compared to the Nyquist criterion. For a single-input–single-output system, the theorem implies that a closed-loop system is stable if its loop gain has magnitude less than unity. This corresponds to open-loop stable systems whose frequency response plots are completely inside the gray unit circle shown in Fig. 11.17.

These frequency response plots clearly do not encircle the -1 point and are obviously stable using the Nyquist criterion. However, the Nyquist criterion uses phase data and can show the stability for a wider class of systems, including ones for which the circle criterion fails.

Example 11.20

Simulate a feedback loop with linear subsystem transfer function ($T = 0.01$ s)

$$G(z) = \frac{z + 0.5}{z^2 + 0.5z + 0.3}$$

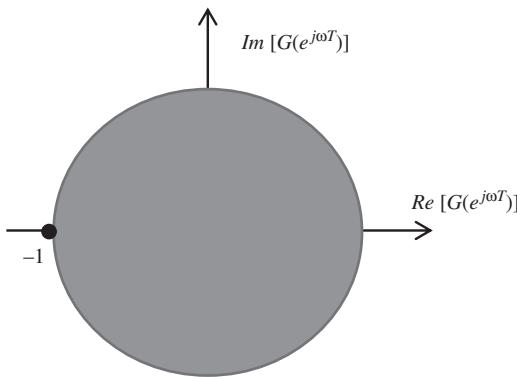


Figure 11.17
Unit circle in the complex plane.

Example 11.20—cont'd

in series with an amplifier of variable gain and a symmetric saturation nonlinearity with slope and saturation level unity. Investigate the stability of the system and discuss your results referring to the small gain theorem for two amplifier gains: 1 and 0.2.

Solution

We use SIMULINK to simulate the system, and the simulation diagram is shown in Fig. 11.18.

The maximum magnitude of the frequency response of the linear subsystem can be obtained using the command

`>> [mag, ph] = bode(g)`

The maximum magnitude is approximately 5.126. By the small gain theorem, the stability condition for the gain of the nonlinearity is $\gamma_N < 1/5.125 \approx 0.195$ (Fig. 11.18).

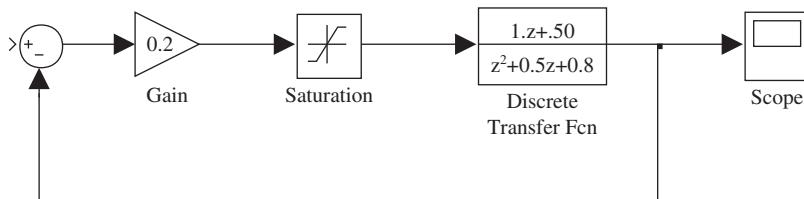
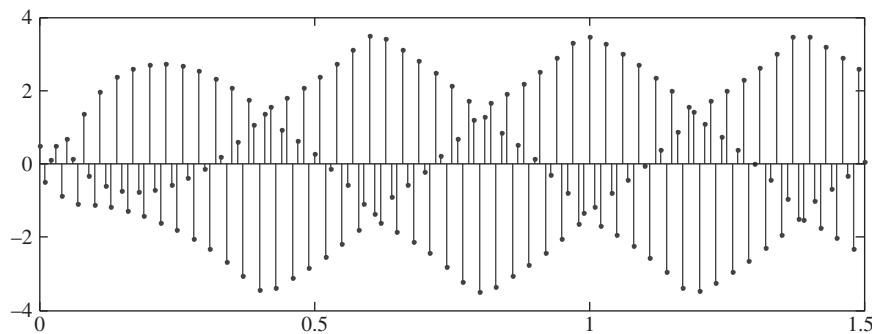


Figure 11.18
Simulation diagram of the closed-loop system with saturation nonlinearity.

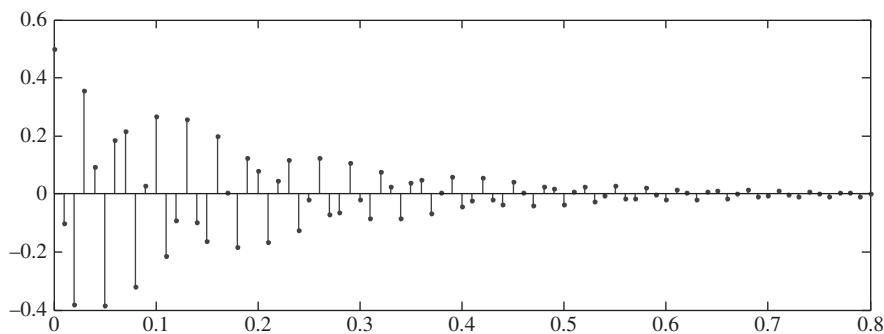
1. The gain of the nonlinear block with slope and saturation level unity is also unity. Thus, by the small gain theorem, the stability of the system cannot be guaranteed with an amplifier gain of unity. This is verified by the simulation results of Fig. 11.19. However, if the amplifier gain is reduced to 0.1, then the overall gain of the linear subsystem is reduced to about 0.513, which is also equal to the product of the gains in the loop. By the small gain theorem, the stability of the system is guaranteed.

Example 11.20—cont'd

2. The results of the small gain theorem are conservative, and the system may be stable for larger gain values. In fact, the simulation results of Fig. 11.20 show that the system is stable for a gain of 0.2, for which the small gain theorem predicts instability.

**Figure 11.19**

Response of the closed-loop system with initial conditions [1, 0] and unity gain.

**Figure 11.20**

Response of the closed-loop system with initial conditions [1, 0] and 0.2 gain.

11.8.1 Absolute stability

In this section, we consider the stability of the SISO system of Fig. 11.21 for a class of nonlinearities $N(\cdot)$. We obtain sufficient conditions for the stability of the closed-loop system if the nonlinearity belongs to the following class.

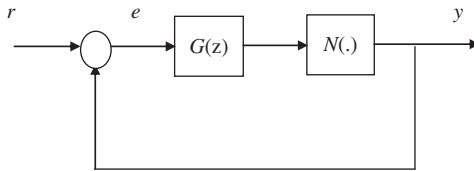


Figure 11.21
Block diagram of closed-loop system with linear and nonlinear blocks.

Definition 11.9

Sector Bound Nonlinearity

A nonlinearity (Fig. 11.22) $N(\cdot)$ belongs to the sector $[k_l, k_u]$ if it satisfies the condition

$$k_l \leq \frac{N(y)}{y} \leq k_u$$

The stability of the loop for any sector bound nonlinearity is defined as follows. We associate the sector bound nonlinearity with the disk $D(k_l, k_u)$ bounded by a circle with its center on the negative real axis with endpoints $-1/k_u$ and $-1/k_l$.

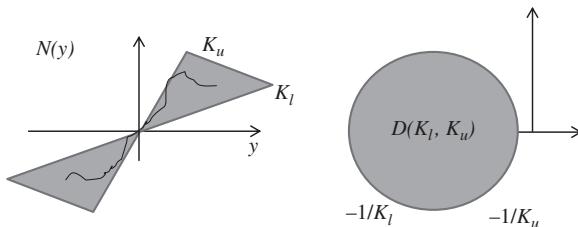


Figure 11.22
Sector bound nonlinearity and the associated disc.

Definition 11.10

Absolute Stability

A nonlinear feedback system of the form in Fig. 11.21 is absolutely stable with respect to the sector $[k_l, k_u]$ if the origin of the state space $\mathbf{x} = \mathbf{0}$ is globally asymptotically stable for all nonlinearities belonging to the sector $[k_l, k_u]$.

Theorem 11.12 gives sufficient conditions for the absolute stability of the closed system with any stable linear block.

Theorem 11.12**The Circle Criterion**

The closed-loop system of Fig. 11.23 is absolutely stable with $G(z)$ stable and nonlinearity $N(.)$ belonging to sector $[k_l, k_u]$ if one of the following conditions is satisfied:

- a. If $0 < k_l < k_u$ and the Nyquist plot of $G(e^{j\omega T})$ does not enter or encircle the disc $D(k_l, k_u)$.

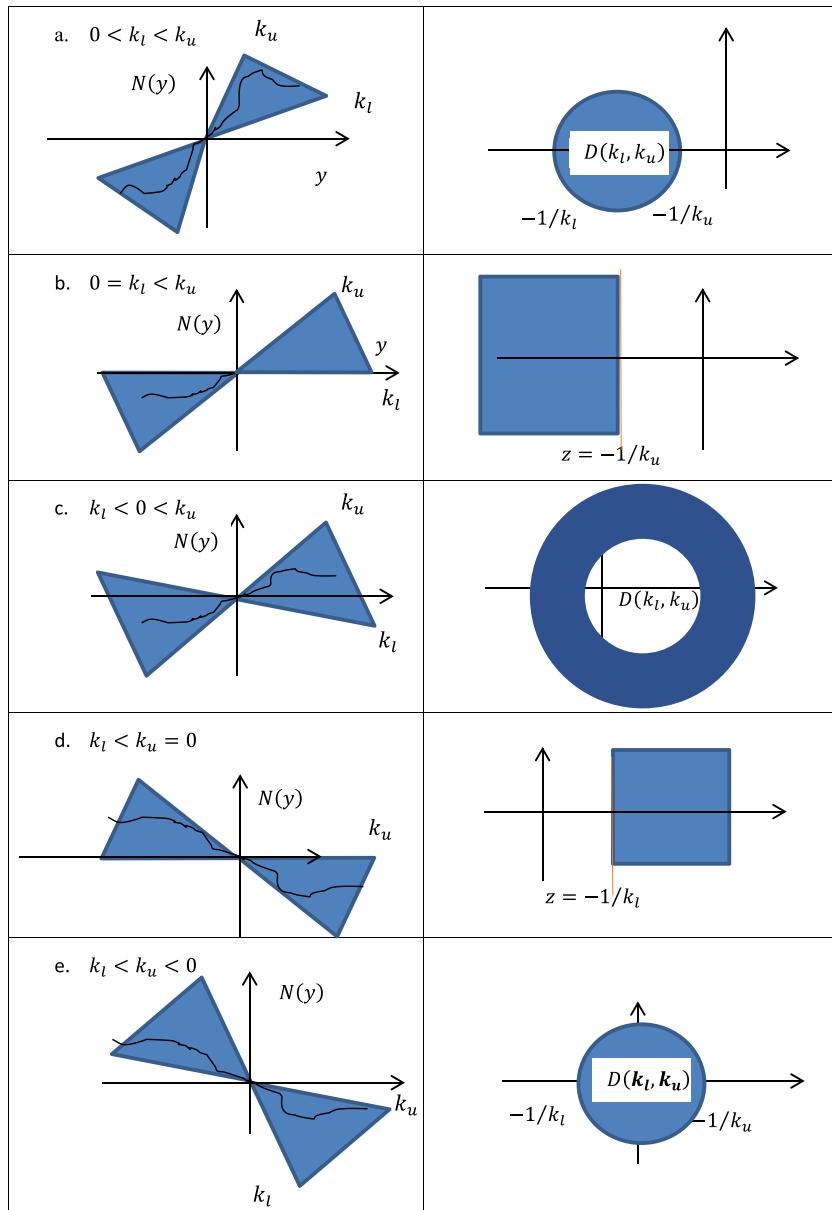


Figure 11.23
The circle criterion.

Theorem 11.12—cont'd

- b. If $0 = k_l < k_u$ and the Nyquist plot of $G(e^{j\omega T})$ lies to the right of the line $z = -1/k_u$.
- c. If $k_l < 0 < k_u$ and the Nyquist plot of $G(e^{j\omega T})$ lies in the interior of the disc $D(k_l, k_u)$.
- d. If $k_l < k_u = 0$ and the Nyquist plot of $G(e^{j\omega T})$ lies to the left of the line $z = -1/k_l$.
- e. If $k_l < k_u < 0$ and the Nyquist plot of $G(e^{j\omega T})$ does not enter or encircle the disc $D(k_l, k_u)$.

The circle criterion provides a sufficient stability condition, but its results can be conservative because it is only based on bounds on the nonlinearity without considering its shape. An overoptimistic estimate of the stability sector can be obtained using the Jury criterion with the nonlinearity replaced by a constant gain as in [Section 4.5](#). The truth will typically lie somewhere between these two estimates, as seen from Example 11.21.

Example 11.21

For the furnace and actuator in Example 4.11, determine the stability sector for the system and compare it to the stable range for a linear gain block.

Solution

The transfer function of the actuator and furnace is

$$G_a(z)G_{ZAS}(z) = 10^{-5} \frac{4.711z + 4.664}{z^3 - 2.875z^2 + 2.753z - 0.8781}$$

and the corresponding Nyquist plot is shown in [Fig. 11.24](#). We use the circle criterion to examine the stability of the closed-loop system with a sector bound nonlinearity. The plot lies inside a circle of radius $10+\epsilon$, which corresponds to a sector $(-0.1, 0.1)$. It is to the right of the line $z = -0.935$, which corresponds to a sector $(0, 1.05)$. Note that we cannot conclude that the system is stable in the sector $(-0.1, 1.05)$; that is, we cannot combine stability sectors. Nevertheless, because the circle criterion gives conservative results,

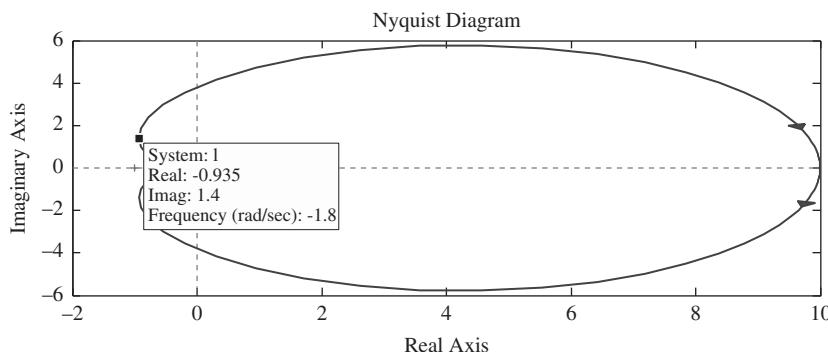
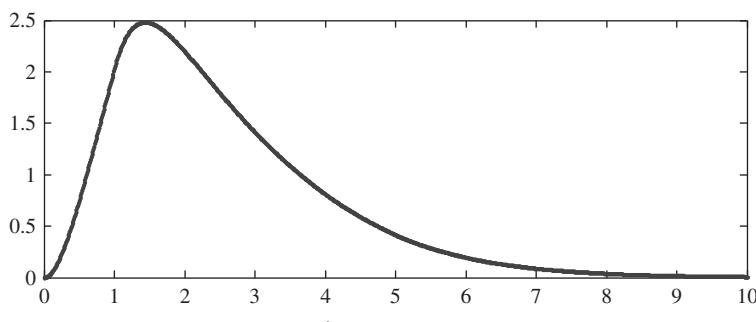


Figure 11.24
Nyquist plot of a furnace and actuator.

Example 11.21—cont'd

the actual stable sector may be wider than our estimates indicate. The Jury criterion gives the stable range $(-0.099559, 3.47398)$. The lower bound is approximately the same as predicted by the circle criterion, but the upper bound is more optimistic. Simulation results, shown in [Figure 11.25](#), verify the stability of the system with saturation nonlinearity in the sector $[-1, 1]$. This is a much larger sector than the conservative estimate of the circle criterion but is well inside the range obtained using the Jury criterion.

**Figure 11.25**

Simulation results for the nonlinear system for a saturation nonlinearity with a linear gain of unity and saturation levels $(-1, 1)$.

Further reading

- Apostol, T.M., 1975. Mathematical Analysis. Addison-Wesley, Reading, MA.
- Fadali, M.S., 1987. Continuous drug delivery system design using nonlinear decoupling: a tutorial. *IEEE Trans. Biomed. Eng.* 34 (8), 650–653.
- Goldberg, S., 1986. Introduction to Difference Equations. Dover, Mineola, NY.
- Kalman, R.E., Bertram, J.E., 1960. Control system analysis and design via the “Second Method” of Lyapunov II: discrete time systems. *J. Basic Engineering Trans. ASME.* 82 (2), 394–400.
- Khalil, H.K., 2002. Nonlinear Systems. Prentice Hall, Upper Saddle River, NJ.
- Kuo, B.C., 1992. Digital Control Systems. Saunders, Fort Worth, TX.
- LaSalle, J.P., 1986. The Stability and Control of Discrete Processes. Springer-Verlag, New York.
- Mickens, R.E., 1987. Difference Equations. Van Nostrand Reinhold, New York.
- Mutoh, Y., Shen, T., Nikiforuk, P.N., 1996. Absolute Stability of Discrete Nonlinear Systems. Academic Press.
- Oppenheim, A.V., Willsky, A.S., Nawab, S.H., 1996. Signals and Systems. Prentice Hall, Upper Saddle River, NJ.
- Rugh, W.J., 1996. Linear System Theory. Prentice Hall, Upper Saddle River, NJ.
- Slotine, J.-J., 1991. Applied Nonlinear Control. Prentice Hall, Englewood Cliffs, NJ.

Problems

11.1 Discretize the following system:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} -3x_1 - x_1^2/(2x_2^2) \\ x_2^2/x_1 + 1 \end{bmatrix} + \begin{bmatrix} x_1^2 \\ 0 \end{bmatrix}u(t)$$

11.2 The equations for rotational maneuvering² of a helicopter are given by

$$\ddot{\theta} = -\frac{1}{2} \left(\frac{I_z - I_y}{I_x} \right) \dot{\psi} \sin(2\theta) - mgA \cos(\theta) + 2 \left(\frac{I_z - I_y}{I_x} \right) \dot{\psi} \theta$$

$$\frac{\sin(2\theta)}{(I_z + I_y) + (I_z - I_y)\cos(2\theta)} T_p$$

$$\dot{\psi} = \frac{1}{(I_z + I_y) + (I_z - I_y)\cos(2\theta)} T_y$$

where

I_x , I_y , and I_z = moments of inertia about the center of gravity

m = total mass of the system

g = acceleration due to gravity

θ and ψ = pitch and yaw angles in radians

T_p and T_y = pitch and yaw input torques

Obtain an equivalent linear discrete-time model for the system and derive the equations for the torque in terms of the linear system inputs.

11.3 A single-link manipulator³ with a flexible link has the equation of motion

$$I\ddot{\theta} + MgL \sin(\theta) - mgA \cos(\theta) + k(\theta - \psi) = 0$$

$$J\ddot{\psi} + k(\psi - \theta) = \tau$$

where

L = distance from the shaft to the center of gravity of the link

M = mass of the link

I = moment of inertia of the link

J = moment of inertia of the joint

K = rotational spring constant for the flexible joint

L = distance between the center of gravity of the link and the flexible joint

θ and ψ = link and joint rotational angles in radians

τ = applied torque

² Elshafei, A.L., Karray, F., 2005. Variable structure-based fuzzy logic identification of a class of nonlinear systems. IEEE Trans. Control Systems Tech. 13 (4), 646–653.

³ Spong, M.W., Vidyasagar, M., 1989. Robot Dynamics and Control. Wiley, New York, pp. 269–273.

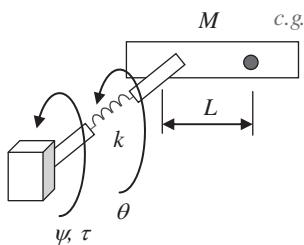


Figure P11.3
Schematic of a single-link manipulator.

Obtain a discrete-time model of the manipulator (Fig. P11.3).

- 11.4 Solve the nonlinear difference equation

$$[y(k+2)][y(k+1)]^{-2}[y(k)]^{1.25} = u(k)$$

with zero initial conditions and the input $u(k) = e^{-k}$.

- 11.5 Determine the equilibrium point for the system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} -x_1(k)/9 + 2x_2^2(k) \\ -x_2(k)/9 + 0.4x_1^2(k) \end{bmatrix} + \begin{bmatrix} 0 \\ x_1(k) \end{bmatrix} u(k)$$

- a. Unforced
 - b. With a fixed input $u_e = 1$
- 11.6 Use the Lyapunov approach to show that if the function $\mathbf{f}(\mathbf{x})$ is a contraction, then the system $\mathbf{x}(k+1) = \mathbf{f}[\mathbf{x}(k)]$ is asymptotically stable.
- 11.7 Obtain a general expression for the eigenvalues of a 2×2 matrix and use it to characterize the equilibrium points of the second-order system with the given state matrix:

- a. $\begin{bmatrix} 0.9997 & 0.0098 \\ -0.0585 & 0.9509 \end{bmatrix}$
- b. $\begin{bmatrix} 3 & 1 \\ 1 & 2 \end{bmatrix}$
- c. $\begin{bmatrix} 0.3 & -0.1 \\ 0.1 & 0.2 \end{bmatrix}$
- d. $\begin{bmatrix} 1.2 & -0.4 \\ 0.4 & 0.8 \end{bmatrix}$

- 11.8 The following model describes population growth for a single species⁴

⁴ M. Shahin, Explorations of Mathematical Models in Biology with MATLAB, J. Wiley, Hoboken, NJ, 2014.

$$p(k+1) = p(k) + ap(k) - bp(k)^2, \quad p(k) < \frac{1+a}{b}, \quad a < 1$$

where $p(k)$ is the population at the end of the k th period, a is the birth rate, and b is the death rate for the species.

- (a) Show that the constraint is required to ensure that the population remains nonnegative.
 - (b) Show that the population has an equilibrium at $p(k) = L = a/b$ (ignore the trivial equilibrium at zero).
 - (c) Use the quadratic Lyapunov function $V(k) = (p(k) - L)^2$ to investigate the stability of the equilibrium in the range of validity of the model.
- 11.9 Consider the system
- $$x(k+1) = x(k)^2$$
- (a) Show that the system has equilibrium points at $x = 0, x = 1$
 - (b) Use the properties of a contraction to show that the equilibrium at zero is stable
 - (c) Show that the equilibrium at $x = 1$ is unstable
- 11.10 Determine the stability of the origin using the linear approximation for the system

$$\begin{aligned} x_1(k+1) &= 0.2x_1(k) + 1.1x_2^3(k) \\ x_2(k+1) &= x_1(k) + 0.1x_2(k) + 2x_1(k)x_2^2(k), \quad k = 0, 1, 2, \dots \end{aligned}$$

- 11.11 Verify the stability of the origin using the Lyapunov approach and estimate the rate of convergence to the equilibrium

$$\begin{aligned} x_1(k+1) &= 0.1x_1(k)x_2(k) - 0.05x_2^2(k) \\ x_2(k+1) &= -0.5x_1(k)x_2(k) + 0.05x_2^3(k), \quad k = 0, 1, 2, \dots \end{aligned}$$

- 11.12 Show that the convergence of the trajectories of a nonlinear discrete-time system $\mathbf{x}(k+1) = \mathbf{f}[\mathbf{x}(k)]$ to a known nominal trajectory $\mathbf{x}^*(k)$ is equivalent to the stability of the dynamics of the tracking error $\mathbf{e}(k) = \mathbf{x}(k) - \mathbf{x}^*(k)$.
- 11.13 Prove that the scalar system

$$x(k+1) = -ax^3(k), \quad a > 0$$

is locally asymptotically stable in the region $|x(k)| \leq 1/a$.

- 11.14 Use Lyapunov stability theory to investigate the stability of the system

$$\begin{aligned} x_1(k+1) &= \frac{ax_1(k)}{a + bx_2^2(k)} \\ x_2(k+1) &= \frac{bx_2(k)}{b + ax_1^2(k)}, \quad a > 0, b > 0 \end{aligned}$$

- 11.15 Use the Lyapunov approach to determine the stability of the discrete-time linear time-invariant systems

a. $\begin{bmatrix} 0.3 & -0.1 \\ 0.1 & 0.22 \end{bmatrix}$

b. $\begin{bmatrix} 0.3 & -0.1 & 0 \\ 0.1 & 0.22 & 0.2 \\ 0.4 & 0.2 & 0.1 \end{bmatrix}$

- 11.16 Show that for the time-varying linear system

$$\mathbf{x}(k+1) = A(k)\mathbf{x}(k), \quad k = 0, 1, \dots$$

the system is stable if there exists a positive definite symmetric and bounded matrix $Q(k)$ that satisfies the inequality

$$A^T(k)Q(k+1)A(k) - Q(k) < 0$$

Hint: Use a quadratic Lyapunov function with a time-varying matrix $Q(k)$.

- 11.17 Use the condition of Problem 11.16 to show that the linear time invariant system with nonuniform sampling period $T(k)$, $k = 0, 1, 2, \dots$

$$\dot{\mathbf{x}}(t) = F\mathbf{x}(t), \quad \mathbf{x}(t) \in \mathcal{R}^n$$

remains stable if and only if the eigenvalues of its state matrix F are in the open LHP.

- 11.18 Show that the origin is an unstable equilibrium for the system

$$\begin{aligned} x_1(k+1) &= -1.4x_1(k) + 0.1x_2^2(k) \\ x_2(k+1) &= 1.5x_2(k)(0.1x_1(k) + 1), \quad k = 0, 1, 2, \dots \end{aligned}$$

- 11.19 Estimate the domain of attraction of the system

$$\begin{bmatrix} x_1(k+1) \\ x_2(k+1) \end{bmatrix} = \begin{bmatrix} 0.2 & 0.3 \\ -0.4 & 0.5 \end{bmatrix} \begin{bmatrix} x_1(k) \\ x_2(k) \end{bmatrix} + \begin{bmatrix} 0.3x_2^2(k) \\ 0.36x_1^2(k) \end{bmatrix}$$

- 11.20 Design a controller to stabilize the origin for the system

$$\begin{aligned} x_1(k+1) &= 0.4x_1(k) + 0.5x_2(k) + x_2^2(k)u(k) \\ x_2(k+1) &= 0.1x_1(k) + 0.2x_2(k) + [x_2(k) + x_1(k)]u(k), \quad k = 0, 1, 2, \dots \end{aligned}$$

Computer exercises

- 11.21 Write a MATLAB program to generate phase plane plots for a discrete-time second-order linear time-invariant system. The function should accept the eigenvalues of the state matrix and the initial conditions needed to generate the plots.
- 11.22 Design a controller for the nonlinear mechanical system described in Example 11.15 with the nonlinear damping $b(\dot{x}) = 0.25\dot{x}^5$, the nonlinear spring $c(x) = 0.5x + 0.02x^3$, $T = 0.02$ s, and the desired eigenvalues for the linear design equal to $\{0.2 \pm j0.1\}$. Determine the value of the reference input for a steady-state position of unity and simulate the system using Simulink.
- 11.23 Design a stabilizing digital controller with a sampling period $T = 0.01$ s for a single-link manipulator using extended linearization, then simulate the system with your design. The equation of motion of the manipulator is given by $\ddot{\theta} + 0.01 \sin(\theta) + 0.01\theta + 0.001\theta^3 = \tau$. Assign the eigenvalues of the discrete-time linear system to $\{0.6 \pm j0.3\}$. Hint: Use Simulink for your simulation, and use a **ZOH block** or a **discrete filter** block with both the numerator and denominator set to 1 for sampling.
- 11.24 Some chemical processes, such as a distillation column, can be modeled as a series of first-order processes with the same time constant. For the process with transfer function

$$G(s) = \frac{1}{(s + 1)^{10}}$$

design a digital proportional controller with a sampling period $T = 0.1$ by applying the tangent method and the Ziegler–Nichols rules (see [Sections 5.5 and 6.3.4](#)). Use the circle criterion to show that the closed-loop system with actuator saturation nonlinearity of slope unity is asymptotically stable. Obtain the step response of the closed-loop system to verify the stability of the closed-loop system with actuator saturation.

Practical issues

Objectives

After completing this chapter, the reader will be able to do the following:

1. Write pseudocode to implement a digital controller.
2. Select the sampling frequency in the presence of antialiasing filters and quantization errors.
3. Implement a proportional–integral–derivative (PID) controller effectively.
4. Design a controller that addresses changes in the sampling rate during control operation.
5. Design a controller with faster input sampling than output sampling.

Successful practical implementation of digital controllers requires careful attention to several hardware and software requirements. In this chapter, we discuss the most important of these requirements and their influence on controller design. We analyze the choice of the sampling frequency in more detail (already discussed in [Section 2.9](#)) in the presence of antialiasing filters and the effects of quantization, rounding, and truncation errors. In particular, we examine the effective implementation of a PID controller. Finally, we examine changing the sampling rate during control operation as well as output sampling at a slower rate than that of the controller.

Chapter Outline

12.1 Design of the hardware and software architecture 568

- 12.1.1 Software requirements 568
- 12.1.2 Selection of ADC and DAC 570

12.2 Choice of the sampling period 572

- 12.2.1 Antialiasing filters 572
- 12.2.2 Effects of quantization errors 574
- 12.2.3 Phase delay introduced by the zero-order hold 584

12.3 Controller structure 585

12.4 Proportional–integral–derivative control 588

- 12.4.1 Filtering the derivative action 588
- 12.4.2 Integrator windup 590
- 12.4.3 Bumpless transfer between manual and automatic mode 593
- 12.4.4 Incremental form 596

12.5 Sampling period switching 598

12.5.1 MATLAB commands 601

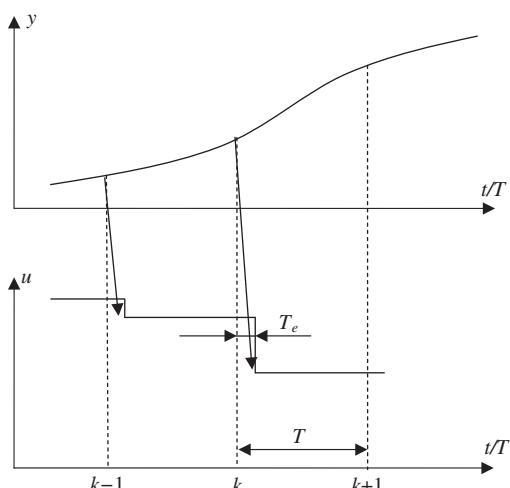
12.5.2 Dual-rate control 608

Reference 611**Further reading 611****Problems 611****Computer exercises 613****12.1 Design of the hardware and software architecture**

The designer of a digital control system must be mindful of the fact that the control algorithm is implemented as a software program that forms part of the control loop. This introduces factors that are not present in analog control loops. This section discusses several of these factors.

12.1.1 Software requirements

During the design phase, designers make several simplifying assumptions that affect the implemented controller. They usually assume uniform sampling with negligible delay due to the computation of the control variable. Thus, they assume no delay between the sampling instant and the instant at which the computed control value is applied to the actuator. This instantaneous execution assumption is not realistic because the control algorithm requires time to compute its output (Fig. 12.1). If the computational time is known and constant, we can use the modified z -transform (see Section 2.7) to obtain a

**Figure 12.1**

Execution time T_e for the computation of the control variable for a system with sampling period T .

more precise discrete model. However, the computational time of the control algorithm can vary from one sampling period to the next. The variation in the computational delay is called the **control jitter**. For example, control jitter is present when the controller implementation utilizes a switching mechanism.

Digital control systems have additional requirements such as data storage and user interface, and their proper operation depends not only on the correctness of their calculations but also on the time at which their results are available. Each task must satisfy either a start or completion timing constraint. In other words, a digital control system is a **real-time system**. To implement a real-time system, we need a real-time operating system that can provide capabilities such as multitasking, scheduling, and intertask communication, among others. In a multitasking environment, the value of the control variable must be computed and applied over each sampling interval regardless of other tasks necessary for the operations of the overall control system. Hence, the highest priority is assigned to the computation and application of the control variable.

Clearly, the implementation of a digital control system requires skills in software engineering and computer programming. There are well-known programming guidelines that help minimize execution time and control jitter for the control algorithm. For example, **if-then-else** and **case** statements must be avoided as much as possible because they can lead to paths of different lengths and, consequently, paths with different execution times. The states of the control variable must be updated after the application of the control variable. Finally, the software must be tested to ensure that no errors occur. This is known as **software verification**. In particular, the execution time and the control jitter must be measured to verify that they can be neglected relative to the sampling period, and memory usage must be analyzed to verify that it does not exceed the available memory. Fortunately, software tools are available to make such analysis possible.

Example 12.1

Write **pseudocode** that implements the following controller:

$$C(z) = \frac{U(z)}{E(z)} = \frac{10.5z - 9.5}{z - 1}$$

Then propose a possible solution to minimize the **execution time** (see Fig. 12.1).

Solution

The difference equation corresponding to the controller transfer function is

$$u(k) = u(k-1) + 10.5e(k) - 9.5e(k-1)$$

This control law can be implemented by writing the following code:

Example 12.1—cont'd

```

function controller
% This function is executed during each sampling period
% r is the value of the reference signal
% u1 and e1 are the values of the control variable and of the control
% error respectively for the previous sampling period
y = read_ADC(ch0) % Read the process output from channel 0 of the ADC
e = r-y; % Compute the tracking error
u = u1 + 10.5*e-9.5*e1; % Compute the control variable
u1 = u; % Update the control variable for the next sampling period
e1 = e; % Update the tracking error for the next sampling period
write_DAC(ch0,u); % Output the control variable to channel 0 of the %DAC

```

To decrease the execution time, two tasks are assigned different priorities (using a real-time operating system):

Task 1 (Maximum Priority)

```

y = read_ADC(ch0) % Read the process output from channel 0 of the ADC
e = r-y; % Compute the tracking error
u = u1 + 10.5*e-9.5*e1; % Compute the control variable
write_DAC(ch0,u); % Write the control variable to channel 0 of the %DAC

```

Task 2

```

u1 = u; % Update the control variable for the next sampling period
e1 = e; % Update the tracking error for the next sampling period

```

12.1.2 Selection of ADC and DAC

The ADC and DAC must be sufficiently fast for negligible conversion time relative to the sampling period. In particular, the conversion delay of the ADC creates a negative phase shift, which affects the phase margin of the system and must be minimized to preserve the stability margins of the system. In addition, the word length of the ADC affects its conversion time. With the conversion time provided by standard modern ADCs, this is not a significant issue in most applications.

The choice of the ADC and DAC word length is therefore mainly determined by the **quantization** effects. Typically, commercial ADCs and DACs are available in the range of 8–16 bits. An 8-bit ADC provides a **resolution** of a 1 in 2^8 , which corresponds to an error of 0.4%, whereas a 16-bit ADC gives an error of 0.0015%.

Clearly, the smaller the ADC resolution, the better the performance, and therefore a 16-bit ADC is preferred. However, the cost of the component increases as the word length

increases, and the presence of noise might render the presence of a high number of bits useless in practical applications. For example, if the sensor has a 5 mV noise and a 5 V range, there is no point in employing an ADC with more than 10 bits because its resolution of 1 in 2^{10} corresponds to an error of 0.1%, which is equal to the noise level. The DAC resolution is usually chosen equal to the ADC resolution, or slightly higher, to avoid introducing another source of quantization error. Once the ADC and DAC resolution have been selected, the resolution of the reference signal representation must be the same as that of the ADC and DAC. In fact, if the precision of the reference signal is higher than the ADC resolution, the control error will never go to zero, and therefore a limit cycle (periodic oscillations associated with nonlinearities) will occur.

Another important design issue, especially for multiinput–multioutput control, is the choice of **data acquisition system**. Ideally, we would use an ADC for each channel to ensure simultaneous sampling as shown in Fig. 12.2A. However, this approach can be quite expensive, especially for a large number of channels. A more economical approach is to use a **multiplexer** (MUX), with each channel sampled in sequence and the sampled value sent to a master ADC (Fig. 12.2B). If we assume that the sampled signals change

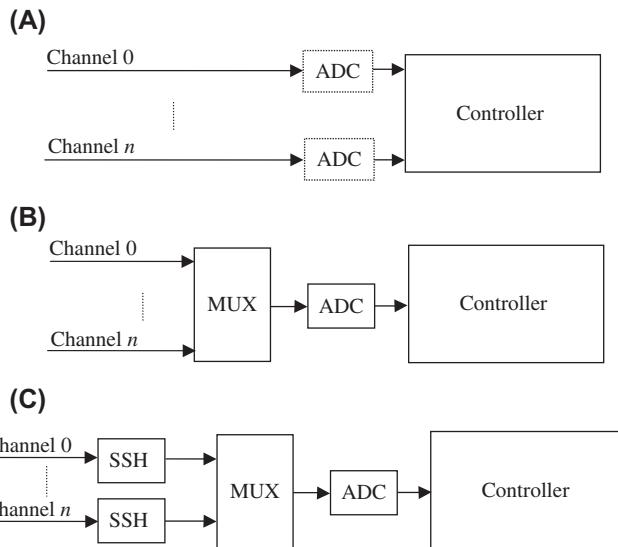


Figure 12.2

Choices for the data acquisition system. (A) Separate ADC for each channel. (B) Multiplexer with single master ADC. (C) Simultaneous sample-and-hold, multiplexer, and single master ADC.

relatively slowly, small changes in the sampling instant do not result in significant errors, and the measured variables appear to be sampled simultaneously. If this assumption is not valid, a more costly **simultaneous sample-and-hold** (SSH) system can be employed, as depicted in Fig. 12.2C. The system samples the channels simultaneously and then delivers the sampled data to the ADC through an MUX. In recent years, the cost of analog-to-digital converters has decreased significantly, and the use of a SSH system has become less popular as using an ADC for each channel has become more affordable.

We discuss other considerations related to the choice of the ADC and DAC components with respect to the sampling period in Section 12.2.2.

12.2 Choice of the sampling period

In Section 2.9, we showed that the choice of the sampling frequency must satisfy the sampling theorem and is based on the effective bandwidth ω_m of the signals in the control systems. This leads to relation (2.67), where we choose the sampling frequency in the range between 5 and 10 times the value of ω_m . We now discuss the choice of sampling frequency more thoroughly, including the effects of antialiasing filters, as well as the effects of quantization, rounding, and truncation errors.

12.2.1 Antialiasing filters

If the sampling frequency does not satisfy the sampling theorem (i.e., the sampled signal has frequency components greater than half the sampling frequency), then the sampling process creates new frequency components (see Fig. 2.10). This phenomenon is called **aliasing** and must obviously be avoided in a digital control system. Hence, the continuous signal to be sampled must not include significant frequency components greater than the Nyquist frequency $\omega_s/2$.

For this purpose, it is recommended to low-pass filter the continuous signal before sampling, especially in the presence of high-frequency noise. The analog low-pass filter used for this purpose is known as the **antialiasing filter**. The antialiasing filter is typically a simple first-order RC filter, but some applications require a higher-order filter such as a Butterworth or a Bessel filter. The overall control scheme is shown in Fig. 12.3.

Because a low-pass filter can slow down the system by attenuating high-frequency dynamics, the cutoff frequency of the low-pass filter must be higher than the bandwidth of the closed-loop system so as not to degrade the transient response. A rule of thumb is to choose the filter bandwidth equal to a constant time the bandwidth of the closed-loop

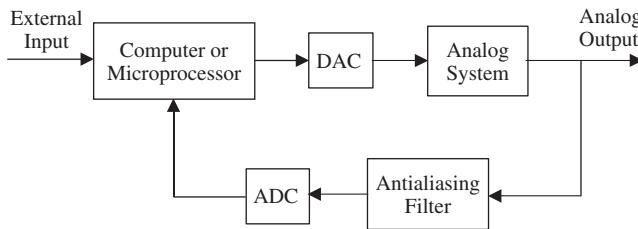


Figure 12.3
Control scheme with an antialiasing filter.

system. The value of the constant varies depending on economic and practical considerations. For a conservative but more expensive design, the cutoff frequency of the low-pass filter can be chosen as 10 times the bandwidth of the closed-loop system to minimize its effect on the control system dynamics, and then the sampling frequency can be chosen 10 times higher than the filter cutoff frequency so there is a sufficient attenuation above the Nyquist frequency. Thus, the sampling frequency is 100 times the bandwidth of the closed-loop system. To reduce the sampling frequency, and the associated hardware costs, it is possible to reduce the antialiasing filter cutoff frequency. In the extreme case, we select the cutoff frequency slightly higher than the closed-loop bandwidth. For a low-pass filter with a high roll-off (i.e., a high-order filter), the sampling frequency is chosen as five times the closed-loop bandwidth. In summary, the sampling period T can be chosen (as described in [Section 2.9](#)) in general as

$$5\omega_b \leq \frac{2\pi}{T} \leq 100\omega_b \quad (12.1)$$

where ω_b is the bandwidth of the closed-loop system.

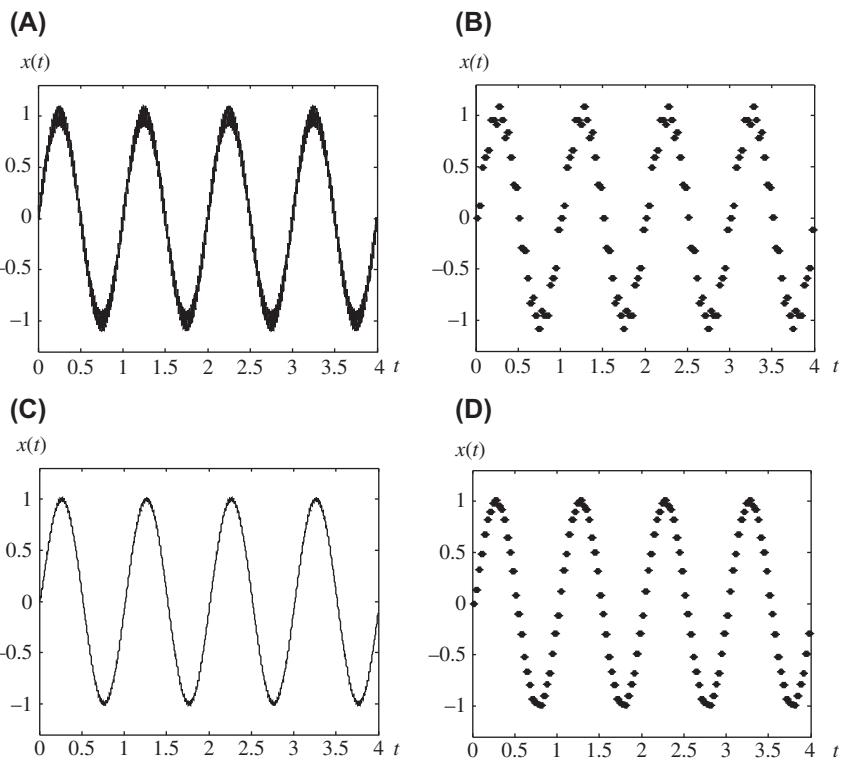
If the phase delay introduced by the antialiasing filter is significant, then [Eq. \(12.1\)](#) may not yield good results, and the filter dynamics must be considered when selecting a dynamic model for the design phase.

Example 12.2

Consider a 1 Hz sinusoidal signal of unity amplitude with an additive 50 Hz sinusoidal noise. Verify the effectiveness of an antialiasing filter if the signal is sampled at a frequency of 30 Hz.

Solution

The noisy 1 Hz analog signal to be sampled is shown in [Fig. 12.4A](#). If this signal is sampled at 30 Hz without an antialiasing filter, the result is shown in [Fig. 12.4B](#). [Fig. 12.4C](#) shows the filtered analog signal with a first-order antialiasing filter with cutoff frequency equal to 10 Hz, and the resulting sampled signal is shown in [Fig. 12.4D](#). The sampled sinusoidal signal is no

Example 12.2—cont'd**Figure 12.4**

Effect of an antialiasing filter on the analog and sampled signals of Example 12.2. (A) Noisy analog signal. (B) Signal sampled at 30 Hz with no antialiasing filter. (C) Filtered analog signal with a first-order antialiasing filter with cutoff frequency equal to 10 Hz. (D) Sampled signal with a first-order antialiasing filter with cutoff frequency equal to 10 Hz.

longer distorted because of the use of the antialiasing filter; however, a small phase delay results from the filter dynamics.

12.2.2 Effects of quantization errors

As discussed in Section 12.1, the design of the overall digital control system includes the choice of the ADC and DAC components. In this context, the effects of the **quantization** due to ADC **rounding** or **truncation** (Fig. 12.5) are considered in the selection of the sampling period. The noise due to quantization can be modeled as a uniformly distributed random process with the following mean and variance values (denoted respectively by \bar{e} and σ_e^2) in the two cases:

$$\text{Rounding: } \bar{e} = 0 \quad \sigma_e^2 = \frac{q^2}{12} \quad (12.2)$$

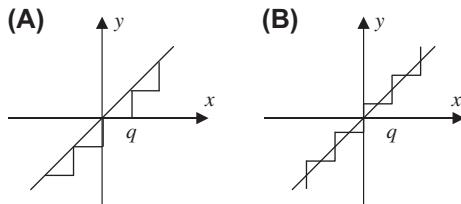


Figure 12.5

Quantization characteristics of the ADC. (A) Truncating ADC. (B) Rounding ADC.

$$\text{Truncation: } \bar{e} = \frac{q}{2} \quad \sigma_e^2 = \frac{q^2}{12} \quad (12.3)$$

where q is the quantization level—namely, the range of the ADC divided by $2n$, and n is the number of bits.

Obviously, the effects of the quantization error increase as q increases and the resolution of the ADC decreases. To evaluate the influence of q on quantization noise and on the sampling period, we consider a proportional feedback digital controller with a gain K applied to the analog first-order lag:

$$G(s) = \frac{1}{\tau s + 1}$$

The z -transfer function of the DAC (zero-order hold [ZOH]), analog subsystem, and ADC (ideal sampler) cascade has the discrete time state-space model

$$\begin{aligned} x(k+1) &= \left[e^{-T/\tau} \right] x(k) + \left[1 - e^{-T/\tau} \right] u(k) \\ y(k) &= x(k) \end{aligned}$$

For a truncating ADC as the only source of noise with zero set-point value, the control action is

$$u(k) = -K[y(k) - w(k)]$$

where w is the quantization noise governed by Eq. (12.3) and is subtracted from $y(k)$ with no loss of generality. The state-space model of the closed-loop system is

$$\begin{aligned} x(k+1) &= \left[e^{-T/\tau} - K \left(1 - e^{-T/\tau} \right) \right] x(k) + K \left[1 - e^{-T/\tau} \right] w(k) \\ y(k) &= x(k) \end{aligned}$$

For zero initial conditions, the solution of the difference equation is

$$\begin{aligned} x(k) &= \sum_{i=0}^{k-1} \left[e^{-T/\tau} - K \left(1 - e^{-T/\tau} \right) \right]^{k-i-1} K \left[1 - e^{-T/\tau} \right] w(k) \\ y(k) &= x(k) \end{aligned}$$

The mean value of the output noise is

$$\begin{aligned} m_y(k) &= E\{x(k)\} = \sum_{i=0}^{k-1} \left[e^{-T/\tau} - K(1 - e^{-T/\tau}) \right]^{k-i-1} K \left[1 - e^{-T/\tau} \right] E\{w(k)\} \\ &= K \left[1 - e^{-T/\tau} \right] \left(\frac{q}{2}\right) \sum_{i=0}^{k-1} \left[e^{-T/\tau} - K(1 - e^{-T/\tau}) \right]^{k-i-1} \end{aligned}$$

where $E\{\cdot\}$ denotes the expectation. If $T/\tau \ll 1$, we use the linear approximation of the exponential terms $e^{-T/\tau} \approx 1 - T/\tau$ to obtain

$$m_y(k) = K \left(\frac{T}{\tau}\right) \left(\frac{q}{2}\right) \sum_{i=0}^{k-1} \left[1 - \left(\frac{T}{\tau}\right)(1+K) \right]^{k-i-1}$$

We recall the relationship

$$\frac{1}{1-a} = \sum_{k=0}^{\infty} a^k, |a| < 1$$

and take the limit as $k \rightarrow \infty$ to obtain the steady-state mean

$$m_y(k) = K \left(\frac{T}{\tau}\right) \left(\frac{q}{2}\right) \frac{1}{\left(\frac{T}{\tau}\right)(1+K)} = \frac{K}{1+K} \left(\frac{q}{2}\right)$$

For small gain values, we have

$$m_y = \frac{K}{1+K} \left(\frac{q}{2}\right) \approx K \left(\frac{q}{2}\right), \quad K \ll 1$$

For large gain values, we have $m_y \approx \frac{q}{2}$.

We observe that the mean value is independent of the sampling period, is linear in the controller gain for small gains, and is almost independent of the gain for large gains. In any case, the worst mean value is half the quantization interval.

The variance of the output is

$$\sigma_y^2 = E\{x^2(k)\} - E^2\{x(k)\}$$

Using the expression for the mean and after some tedious algebraic manipulations, we can show that

$$\sigma_y^2 = \frac{K^2 (1 - e^{-T/\tau})^2}{(K+1)[(1-K)2Ke^{-T/\tau} - (K+1)e^{-2T/\tau}]} \left(\frac{q^2}{12}\right)$$

If $T/\tau \ll 1$, we use the linear approximations of the exponential terms $e^{-T/\tau} \approx 1 - T/\tau$ and $e^{-2T/\tau} \approx 1 - 2T/\tau$, and the output variance simplifies to

$$\sigma_y^2 = \frac{K^2}{K+1} \frac{T}{2\tau} \left(\frac{q^2}{12} \right) \approx \begin{cases} \frac{KT}{2\tau} \left(\frac{q^2}{12} \right), & K \gg 1 \\ K^2 \frac{\tau}{2\tau}, & K \ll 1 \end{cases}$$

Unlike the output mean, the output variance is linear in the sampling period and linear in the controller gain for large gains. Thus, the effect of the quantization noise can be reduced by decreasing the sampling period, once the ADC has been selected. We conclude that decreasing the sampling period has beneficial effects with respect to both aliasing and quantization noise. However, decreasing the sampling period requires more expensive hardware and may aggravate problems caused by the finite-word representation of parameters.

We illustrate this fact with a simple example. Consider an analog controller with poles at $s_1 = -1$ and $s_2 = -10$. For a digital implementation of the analog controller with sampling period $T = 0.001$, we have the digital controller poles as given by Eq. (6.3) as

$$z_1 = e^{-0.001} \approx 0.9990 \quad \text{and} \quad z_2 = e^{-0.01} \approx 0.9900$$

If we truncate the two values after two significant digits, we have $z_1 = z_2 = 0.99$; that is, the two poles of the digital controller become identical and correspond to two identical poles of the analog controller at

$$s_1 = s_2 = \frac{1}{T} \ln z_1 = \frac{1}{T} \ln z_2 \approx -10.05$$

For the longer sampling period $T = 0.1$, truncating after two significant digits gives the poles $z_1 = 0.90$ and $z_2 = 0.36$, which correspond to $s_1 \approx -1.05$ and $s_2 = -10.21$. This shows that a much better approximation is obtained with the longer sampling period. The next example shows the effect of a truncating ADC on the system response.

Example 12.3

Consider the process

$$G(s) = \frac{1}{s+1}$$

with a proportional digital feedback controller, a sampling period $T = 0.02$, and a gain $K = 5$.

- Determine the resolution of an 8-bit ADC with a range of 20 V and the resolution of a 14-bit ADC with the same range.

Example 12.3—cont'd

- b. Obtain a plot of the output response and the control input of the closed-loop system with and without a truncating ADC, first with the 8-bit ADC and then with the 14-bit ADC, and compare the effect of truncation and ADC resolution on the response.

Solution

The 8-bit ADC has a resolution of $20,000/(2^8) = 78.125$ mV, while the 14-bit ADC has a resolution of $20,000/(2^{14}) = 1.22$ mV. Using SIMULINK, we can simulate the closed-loop system to obtain plots of its output response and control input with and without the quantization effect. For the 8-bit ADC, the process output is shown in Fig. 12.6, and the control input is shown in Fig. 12.7. The differences between the two responses and the two control inputs are evident. In contrast, for the 14-bit ADC, the results are shown in Figs. 12.8 and 12.9, where both the two responses and the two control inputs are almost identical. With more powerful hardware, it is important to take the numerical precision of the CPU into account during the implementation of a digital controller. If, for example, the numerical representation of real numbers uses 32 bits (IEEE-754 32-bit floating point arithmetic), the result of the mathematical computation is different from the case where the numerical representation of the real numbers uses 64 bits (IEEE-754 64-bit floating point arithmetic). Although the higher accuracy associated with a 64-bit representation comes at a higher cost, the cost may be justifiable if lower accuracy impairs control performance as shown in the following example.

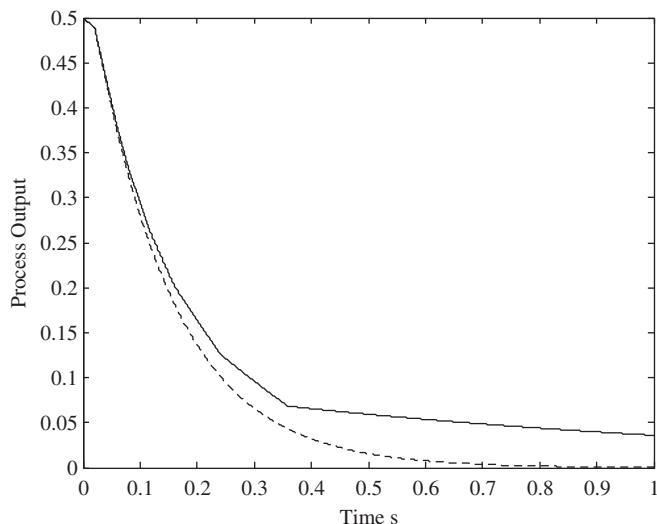
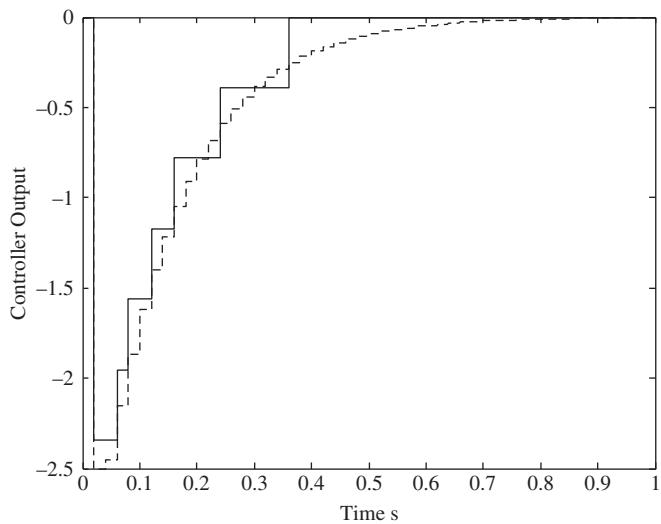
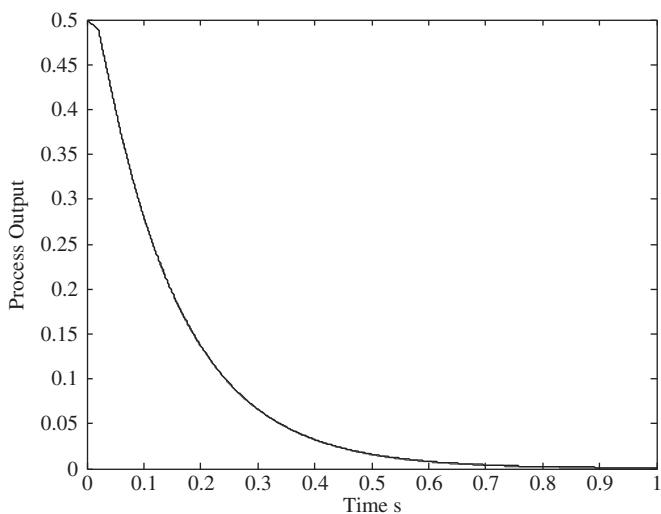


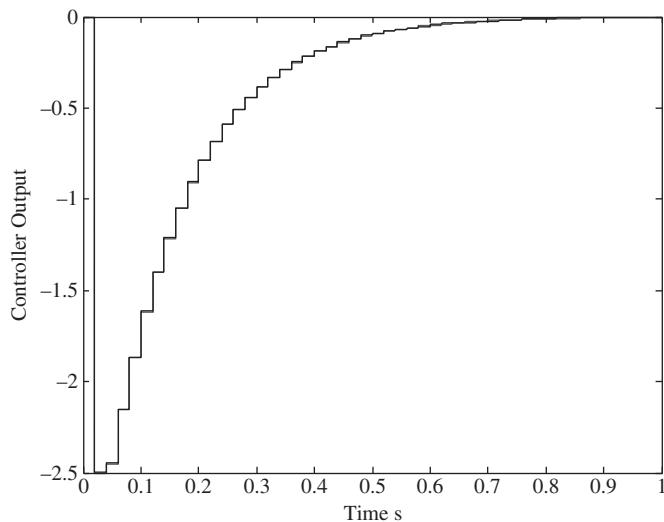
Figure 12.6
Process output with (solid line) and without (dashed line) the 8-bit ADC.

Example 12.3—cont'd**Figure 12.7**

Controller output with (*solid line*) and without (*dashed line*) the 8-bit ADC.

**Figure 12.8**

Process output with (*solid line*) and without (*dashed line*) the 8-bit ADC.

Example 12.3—cont'd**Figure 12.9**

Controller output with (solid line) and without (dashed line) the 8-bit ADC.

Example 12.4

Consider the implementation of a model-based control strategy for the overhead crane shown in Fig. 12.10 with the following parameters: cart mass $m_c = 50 \text{ kg}$, viscous friction coefficient of the cart $c_c = 5 \text{ N/m}$, viscous friction coefficient of the cable $c_p = 0.01 \text{ Nms/rad}$, payload mass $m_p = 10 \text{ kg}$, length of the cable $l = 2.5 \text{ m}$. Compare the performance of the controller with 32-bit and 64-bit hardware with the sampling period $T = 0.05 \text{ s}$, and $T = 0.2 \text{ s}$.

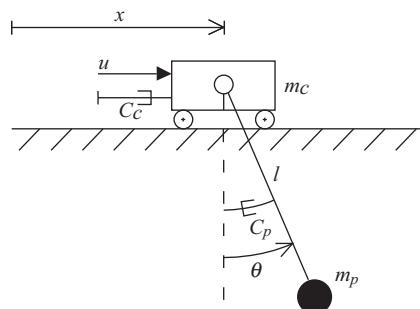


Figure 12.10
Schematic of an overhead crane.

Example 12.4—cont'd**Solution**

The state space representation of a system can be written as

$$\begin{cases} \dot{x}_{ss} = Ax_{ss} + Bu \\ y = Cx_{ss} \end{cases}$$

where x_{ss} is the state vector

$$x_{ss} = [x, \dot{x}, \theta, \dot{\theta}]^T$$

where x is the cart position, \dot{x} is the cart velocity, θ is the angular position of the payload, $\dot{\theta}$ is the angular velocity of the payload; u is the force applied to the cart, y is the output vector, and A , B , and C are, respectively,

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & -\frac{c_c}{m_c} & \frac{g m_p}{m_c} & \frac{c_p}{m_c} \\ 0 & 0 & 0 & 1 \\ 0 & \frac{c_c}{l m_c} & -\frac{g(m_p + m_c)}{l m_c} & \frac{c_p(m_p + m_c)}{l^2 m_c m_p} \end{bmatrix}$$

$$B = \begin{bmatrix} 0 & \frac{1}{m_c} & 0 & -\frac{1}{l m_c} \end{bmatrix}^T$$

$$C = \begin{bmatrix} 1 & 0 & l & 0 \\ 0 & 1 & 0 & l \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{bmatrix}.$$

Consider, for example, the possibility to estimate the load position. The transfer function between cart velocity and load position is

$$F(s) = \frac{x_l(s)}{\dot{x}_c(s)} = \frac{g m_p l^2 + c_p s}{s(m_p l^2 s^2 + c_p s + g m_p l^2)}.$$

The discretization of $F(s)$ by using the bilinear transform results in

$$F(z) = \frac{x_l(z)}{\dot{x}_c(z)} = \frac{A_1 + A_2 z^{-1} + A_3 z^{-2} + A_4 z^{-3}}{B_1 + B_2 z^{-1} + B_3 z^{-2} + B_4 z^{-3}}$$

where the coefficients of the numerator are

Example 12.4—cont'd

$$\begin{aligned}
 A_1 &= T_s^2 (T_s g m_p l^2 + 2c_p) \\
 A_2 &= 2T_s^2 (T_s g m_p l^2 + 2c_p) - T_s^2 (-T_s g m_p l^2 + 2c_p) \\
 A_3 &= T_s^2 (T_s g m_p l^2 + 2c_p) - 2T_s^2 (-T_s g m_p l^2 + 2c_p) \\
 A_4 &= -T_s^2 (-T_s g m_p l^2 + 2c_p)
 \end{aligned}$$

and the coefficients of the denominator are

$$\begin{aligned}
 B_1 &= 2g m_p T_s^2 l^2 + 4c_p T_s + 8m_p l^2 \\
 B_2 &= 2g m_p T_s^2 l^2 - 4c_p T_s + 24m_p l^2 \\
 B_3 &= -2g m_p T_s^2 l^2 - 4c_p T_s + 24m_p l^2 \\
 B_4 &= -2g m_p T_s^2 l^2 + 4c_p T_s + 8m_p l^2
 \end{aligned}$$

The numerical values of the coefficients of $F(z)$ computed are shown in [Table 12.1](#) with $T = 0.05$ s and in [Table 12.2](#) for $T = 0.2$ s. The poles of the continuous time transfer function

Table 12.1: Coefficients of $F(z)$ computed with a sampling period $T = 0.05$ s.

Precision	32-bit precision	64-bit precision
A_1	0.07669063657522015	0.076690625000000026
A_2	0.229971900582313540	0.229971875000000080
A_3	0.229871898889541630	0.229871875000000060
A_4	0.076590634882450104	0.076590625000000023
B_1	503.067626953125000000	503.067625000000020000
B_2	-1496.936401367187500000	-1496.936375000000000000
B_3	1496.932373046875000000	1496.93237500000100000
B_4	-503.063629150390620000	-503.063625000000000000

Table 12.2: Coefficients of $F(z)$ computed with $T = 0.2$ s.

	32-bit precision	64-bit precision
A_1	4.905800819396972700	4.905800000000001000
A_2	14.715802192687988000	14.715800000000005000
A_3	14.714201927185059000	14.714200000000005000
A_4	4.904200553894043000	4.904200000000001200
B_1	549.057983398437500000	549.05799999999990000
B_2	-1450.958007812500000000	-1450.958000000000100000
B_3	1450.942016601562500000	1450.942000000000000000
B_4	-549.041992187500000000	-549.042000000000030000

Example 12.4—cont'd

$F(s)$ and of the discrete time transfer function $F(z)$ are shown in [Table 12.3](#) with the sampling time $T = 0.05$ s and in [Table 12.4](#) with $T = 0.2$ s.

Table 12.3: Poles of the continuous time transfer function $F(s)$ and of the discrete time transfer function $F(z)$ with sampling period $T = 0.05$ s.

Poles of $F(s)$	p_1	$0.000000000000000 + 0.000000000000000i$
	p_2	$-0.000080000000000 + 3.132091951651483i$
	p_3	$-0.000080000000000 - 3.132091951651483i$
Poles of $F(z)$, 32 bits	p_1	$1.000002488695106 + 0.000000000000000i$
	p_2	$0.987807075215285 + 0.155649435064696i$
	p_3	$0.987807075215285 - 0.155649435064696i$
Poles of $F(z)$, 64 bits	p_1	$0.999999999999946 + 0.000000000000000i$
	p_2	$0.987808299132771 + 0.155649648079199i$
	p_3	$0.987808299132771 - 0.155649648079199i$

Table 12.4: Poles of the continuous time transfer function $F(s)$ and of the discrete time transfer function $F(z)$ with sampling period $T = 0.2$ s.

Poles of $F(s)$	p_1	$0.000000000000000 + 0.000000000000000i$
	p_2	$-0.000080000000000 + 3.132091951651483i$
	p_3	$-0.000080000000000 - 3.132091951651483i$
Poles of $F(z)$, 32 bits	p_1	$0.999999999999997 + 0.000000000000000i$
	p_2	$0.821315827912822 + 0.570448232539258i$
	p_3	$0.821315827912822 - 0.570448232539258i$
Poles of $F(z)$, 64 bits	p_1	$0.999999999999999 + 0.000000000000000i$
	p_2	$0.821315780846470 + 0.570448286274216i$
	p_3	$0.821315780846470 - 0.570448286274216i$

The continuous time transfer function has one pole at the origin and two stable poles (poles in the open LHP) so that the overall system is marginally stable. After the discretization process, stability should be maintained the poles of the discrete time transfer function are on or inside the unit circle.

For the sampling period $T = 0.2$ s ([Table 12.4](#)), stability is maintained both with 32- and 64-bit computation; all three poles lie inside the unit circle. For the sampling period $T = 0.05$ s ([Table 12.3](#)), stability is maintained with 64-bit computation; all three poles are inside the unit circle. In the case of 32-bit computation, one of the three poles (p_1) is outside of the unit circle. Thus, the system discretized by computing the coefficients with 32 bits is unstable.

12.2.3 Phase delay introduced by the zero-order hold

As shown in Section 3.3, the frequency response of the ZOH can be approximated as

$$G_{ZOH}(j\omega) = \frac{1 - e^{-j\omega T}}{j\omega} \approx e^{-j\omega T/2}$$

This introduces an additional delay in the control loop approximately equal to half of the sampling period. The additional delay reduces the stability margins of the control system, and the reduction is worse as the sampling period is increased. This imposes an upper bound on the value of the sampling period T .

Example 12.5

Let ω_c be the gain crossover frequency of an analog control system. Determine the maximum value of the sampling period for a digital implementation of the controller that decreases the phase margin by no more than 5 degrees.

Solution

Because of the presence of the ZOH, the phase margin decreases by $\omega_c T/2$, which yields the constraint

$$\omega_c \frac{T}{2} \leq 5 \frac{\pi}{180}$$

or equivalently, $T \leq 0.1745 \omega_c$.

Example 12.6

Consider the tank control system described in Example 2.1 with the transfer function

$$G(s) = \frac{1.2}{20s + 1} e^{-15x}$$

and the proportional-integral (PI) controller

$$C(s) = 7 \frac{20s + 1}{20s}$$

Let the actuator and the sensor signals be in the range 0–5 V with a sensor gain of 0.169 V/cm. Select a suitable sampling period, antialiasing filter, DAC, and ADC for the system.

Solution

The gain crossover frequency of the analog control system as obtained using MATLAB is $\omega_c = 0.42$ rad/s, and the phase margin is $\phi_m = 54^\circ$. We select a sampling period $T = 0.2$ s and use a second-order Butterworth antialiasing filter with cutoff frequency of 4 rad/s. The transfer function of the Butterworth filter is

Example 12.6—cont'd

$$F(s) = \frac{64}{s^2 + 11.31s + 64}$$

The antialiasing filter does not change the gain crossover frequency significantly. The phase margin is reduced to $\phi_m = 49.6^\circ$, which is acceptable. At the Nyquist frequency of $\pi/T = 10.47$ rad/s, the antialiasing filter decreases the magnitude of the noise by more than 40 dB. The phase delay introduced by the ZOH is $(\omega_c T/2) \times 180/\pi = 3.6^\circ$, which is also acceptable. We select a 12-bit ADC with a quantization level of 1.2 mV, which corresponds to a quantization error in the fluid level of 0.07 mm. We also select a 12-bit DAC. Because the conversion time is on the order of microseconds, this does not influence the overall design.

12.3 Controller structure

Section 12.2 demonstrates how numerical errors can affect the performance of a digital controller. To reduce numerical errors and mitigate their effects, we must select an appropriate **controller structure** for implementation. To examine the effect of controller structure on errors, consider the controller

$$\begin{aligned} C(z) &= \frac{N(z, \mathbf{q})}{D(z, \mathbf{q})} = \frac{\mathbf{b}^T \mathbf{z}_m}{\mathbf{a}^T \mathbf{z}_n} a \\ \mathbf{a} &= [a_0(\mathbf{q}) \quad a_1(\mathbf{q}) \quad \dots \quad a_n(\mathbf{q})]^T \\ \mathbf{b} &= [b_0(\mathbf{q}) \quad b_1(\mathbf{q}) \quad \dots \quad b_m(\mathbf{q})]^T \\ \mathbf{z}_l &= [1 \quad z \quad \dots \quad z^l] \end{aligned} \tag{12.4}$$

where \mathbf{q} is an $l \times 1$ vector of controller parameters. If the nominal parameter vector is \mathbf{q}^* and the corresponding poles are p_i^* , $i = 1, 2, \dots, n$, for an n^{th} -order controller, then the nominal characteristic equation of the controller is

$$D(p_i^*, q^*) = 0, \quad i = 1, 2, \dots, n \tag{12.5}$$

In practice, the parameter values are only approximately implemented and the characteristic equation of the system is

$$\begin{aligned} D(z, \mathbf{q}) &\approx D(p_i^*, \mathbf{q}^*) + \left[\frac{\partial D}{\partial z} \right]_{z=p_i^*} \delta p_i^* + \left[\frac{\partial D}{\partial \mathbf{a}} \right]_{z=p_i^*} \delta \mathbf{a} \\ &= \left[\frac{\partial D}{\partial z} \right]_{z=p_i^*} \delta p_i^* + \left[\frac{\partial D}{\partial \mathbf{a}} \right]_{z=p_i^*} \delta \mathbf{a} \approx 0 \\ \frac{\partial D}{\partial \mathbf{a}} &= \left[\frac{\partial D}{\partial a_0} \frac{\partial D}{\partial a_1} \dots \frac{\partial D}{\partial a_n} \right]^T \end{aligned} \tag{12.6}$$

In terms of the controller parameters, the perturbed characteristic equation is

$$D(z, \mathbf{q}) \approx \frac{\partial D}{\partial z} \Big|_{z=p_i^*} \delta p_i^* + \frac{\partial D^T}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{q}} \Big|_{\mathbf{q}=\mathbf{q}^*} \delta \mathbf{q} \approx 0 \quad (12.7)$$

$$\frac{\partial \mathbf{a}}{\partial \mathbf{q}} = \begin{bmatrix} \frac{\partial a_i}{\partial q_j} \end{bmatrix}$$

We solve for parameter perturbations in the location of the i^{th} pole

$$\delta p_i^* = - \frac{\partial D^T}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{q}} \Big|_{\mathbf{q}=\mathbf{q}^*} \frac{\delta \mathbf{q}}{\frac{\partial D}{\partial z} \Big|_{z=p_i^*}} \quad (12.8)$$

To characterize the effect of a particular parameter on the i^{th} pole, we can set the perturbations in all other parameters to zero to obtain

$$\delta p_i^* = - \frac{\partial D^T}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial q_i} \Big|_{\mathbf{q}=\mathbf{q}^*} \frac{\delta q_i}{\frac{\partial D}{\partial z} \Big|_{z=p_i^*}} \quad (12.9)$$

$$\frac{\partial \mathbf{a}}{\partial q_i} = \left[\frac{\partial a_0}{\partial q_i} \frac{\partial a_1}{\partial q_i} \dots \frac{\partial a_n}{\partial q_i} \right]^T$$

This concept is explained by the following example. Consider the following second-order general controller:

$$C(z) = \frac{(a+b)z - (ap_2 + bp_1)}{z^2 - (p_1 + p_z)z + p_1p_2}$$

Let $a_1 = -(p_1 + p_2)$ and $a_0 = p_1p_2$ denote the nominal coefficients of the characteristic equation of the controller, and write the characteristic equation as

$$D(z, a_1, a_0) = 0$$

When a coefficient λ_i is changed (due to numerical errors) to $\lambda_i + \delta\lambda_i$, then the position of a pole is changed according to the following equation (where second- and higher-order terms are neglected):

$$D(p_i + \delta p_i, \lambda_i + \delta\lambda_i) = D(p_i, \lambda_i) + \frac{\partial D}{\partial z} \Big|_{z=p_i} \delta p_i + \frac{\partial D}{\partial \lambda_i} \delta\lambda_i$$

That is,

$$\frac{\delta p_i}{\delta \lambda_i} = - \frac{\frac{\partial D}{\partial \lambda_i}}{\left. \frac{\partial D}{\partial z} \right|_{z=p_i}}$$

Now we have the partial derivatives

$$\frac{\partial D}{\partial z} = 2z + a_1 \quad \frac{\partial D}{\partial a_1} = z \quad \frac{\partial D}{\partial a_0} = 1$$

and therefore

$$\begin{aligned} \frac{\delta p_1}{\delta a_1} &= - \frac{p_1}{2p_1 - (p_1 + p_2)} = \frac{p_1}{p_1 - p_2}, & \frac{\delta p_2}{\delta a_1} &= - \frac{p_2}{2p_2 - (p_1 + p_2)} = \frac{p_2}{p_2 - p_1} \\ \frac{\delta p_1}{\delta a_0} &= - \frac{1}{2p_1 - (p_1 + p_2)} = \frac{1}{p_1 - p_2}, & \frac{\delta p_2}{\delta a_0} &= - \frac{1}{2p_2 - (p_1 + p_2)} = \frac{1}{p_2 - p_1} \end{aligned}$$

Thus, the controller is most sensitive to changes in the last coefficient of the characteristic equation, and its sensitivity increases when the poles are close. This concept can be generalized to high-order controllers. Note that decreasing the sampling period draws the poles closer when we start from an analog design. In fact, for $T \rightarrow 0$ we have that

$$p_z = e^{p_s T} \rightarrow 1$$

independently on the value of the analog pole p_s .

These problems can be avoided by writing the controller in an equivalent **parallel form**:

$$C(z) = \frac{a}{z - p_1} + \frac{b}{z - p_2}$$

We can now analyze the sensitivity of the two terms of the controller separately to show that the sensitivity is equal to one, which is less than the previous case if the poles are close. Thus, the parallel form is preferred. Similarly, the parallel form is also found to be superior to the **cascade form**:

$$C(z) = \frac{(a+b)z - (ap_2 + bp_1)}{z - p_1} \times \frac{1}{z - p_2}$$

Example 12.7

Write the difference equations in direct, parallel, and cascade forms for the system

$$C(z) = \frac{U(z)}{E(z)} = \frac{z - 0.4}{z^2 - 0.3z + 0.02}$$

Solution

The difference equation corresponding to the direct form of the controller is

$$u(k) = 0.3u(k-1) - 0.02u(k-2) + e(k-1) - 0.4e(k-2)$$

For the parallel form, we obtain the partial fraction expansion of the transfer function

$$C(z) = \frac{U(z)}{E(z)} = \frac{-2}{z - 0.2} + \frac{3}{z - 0.1}$$

This is implemented using the following difference equations:

$$\begin{aligned} u_1(k) &= 0.2u(k-1) - 2e(k-1) \\ u_2(k) &= 0.1u(k-1) + 3e(k-1) \\ u(k) &= u_1(k) + u_2(k) \end{aligned}$$

Finally, for the cascade form we have

$$C(z) = \frac{U(z)}{E(z)} = \frac{z - 0.4}{z - 0.2} \frac{1}{z - 0.1} = \frac{X(z)U(z)}{E(z)X(z)}$$

which is implemented using the difference equations

$$\begin{aligned} x(k) &= 0.2x(k-1) + e(k) - 0.4e(k-1) \\ u(k) &= 0.1u(k-1) + x(k-1) \end{aligned}$$

12.4 Proportional–integral–derivative control

In this section, we discuss several critical issues related to the implementation of PID controllers. Rather than providing an exhaustive discussion, we highlight a few problems and solutions directly related to digital implementation.

12.4.1 Filtering the derivative action

The main problem with derivative action is that it amplifies the high-frequency noise and may lead to a noisy control signal that can eventually cause serious damage to the actuator. It is therefore recommended that one filter the overall control action with a low-pass filter or, alternatively, filter the derivative action. In this case, the controller transfer function in the analog case can be written as (see Eq. (5.20))

$$C(s) = K_p \left(1 + \frac{1}{T_i s} + \frac{T_d s}{1 + \frac{T_d}{N} s} \right)$$

where N is a constant in the interval [1, 33], K_p is the proportional gain, T_i is the integral time constant, and T_d is the derivative time constant.

In the majority of cases encountered in practice, the value of N is in the smaller interval [8, 16]. The controller transfer function can be discretized as discussed in Chapter 6. A useful approach in practice is to use the forward differencing approximation for the integral part and the backward differencing approximation for the derivative part. This gives the discretized controller transfer function

$$C(z) = K_p \left(1 + \frac{T}{T_i(z-1)} + \frac{T_d}{T + \frac{T_d}{N}} \cdot \frac{z-1}{z - \frac{T_d}{NT+T_d}} \right) \quad (12.10)$$

which can be simplified to

$$C(z) = \frac{K_0 + K_1 z + K_2 z^2}{(z-1)(z-\gamma)} \quad (12.11)$$

where

$$\begin{aligned} K_0 &= K_p \left(\frac{T_d}{NT+T_d} - \frac{T}{T_i} \frac{T_d}{NT+T_d} + \frac{NT_d}{NT+T_d} \right) \\ K_1 &= -K_p \left(1 + \frac{T_d}{NT+T_d} - \frac{T}{T_i} + 2 \frac{NT_d}{NT+T_d} \right) \\ K_2 &= K_p \left(1 + \frac{NT_d}{NT+T_d} \right) \\ \gamma &= \frac{T_d}{NT+T_d} = \frac{1}{N(T/T_d) + 1} \end{aligned} \quad (12.12)$$

Example 12.8

Select a suitable derivative filter parameter value N for the PID controller described in Example 5.9 with a sampling period $T = 0.01$, and obtain the corresponding discretized transfer function of Eq. (12.11).

Solution

The analog PID parameters are $K_p = 2.32$, $T_i = 3.1$, and $T_d = 0.775$. We select the filter parameter $N = 20$ and use Eq. (12.12) to obtain $K_0 = 38.72$, $K_1 = -77.92$, $K_2 = 39.20$, and $\gamma = 0.79$.

Example 12.8—cont'd

The discretized PID controller expression is therefore

$$C(z) = \frac{K_0 + K_1 z + K_2 z^2}{(z - 1)(z - \gamma)} = \frac{38.72 - 77.92z + 39.20z^2}{(z - 1)(z - 0.79)}$$

12.4.2 Integrator windup

Most control systems are based on linear models and design methodologies. However, every actuator has a saturation nonlinearity, as in the control loop shown in Fig. 12.11, which affects both the analog and digital control. The designer must consider the nonlinearity at the design stage to avoid performance degradation. A common phenomenon related to the presence of actuator saturation is known as **integrator windup**. If not properly handled, it may result in a step response with a large overshoot and settling time.

In fact, if the control variable attains its saturation limit when a step input is applied, the control variable becomes independent of the feedback and the system behaves as in the open-loop case. The control error decreases more slowly than in the absence of saturation, and the integral term becomes large or **winds up**. The large integral term causes saturation of the control variable even after the process output attains its reference value and a large overshoot occurs.

Many solutions have been devised to compensate for integrator windup and retain linear behavior. The rationale of these antiwindup techniques is to design the control law, disregarding the actuator nonlinearity, and then compensate for the detrimental effects of integrator windup.

One of the main antiwindup techniques is the so-called **conditional integration**, which keeps the integral control term constant when a specified condition is met. For example, the integral control term is kept constant if the integral component of the computed

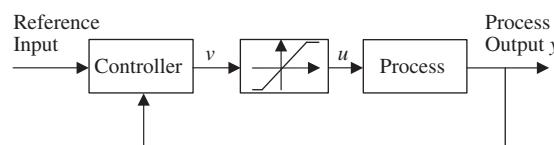


Figure 12.11
Control loop with actuator saturation.

control exceeds a given threshold specified by the designer. Alternatively, the integral controller is kept constant if the actuator saturates with the control variable and the control error having the same sign (i.e., if $u \cdot e > 0$). The condition $u \cdot e > 0$ implies that the control increases rather than corrects the error due to windup and that integral action should not increase.

On the other hand, a positive saturation with $u \cdot e < 0$ means that the error is negative and therefore the integral action is decreasing and there is no point in keeping it constant. The same reasoning can be easily applied in case of a negative saturation. Thus, the condition avoids inhibiting the integration when it helps to push the control variable away from saturation.

An alternative technique is **back-calculation**, which reduces (increases) the integral control when the maximum (minimum) saturation limit is attained by adding to the integrator a term proportional to the difference between the computed value of the control signal v and its saturated value u . In other words, the integral value $I(k)$ is determined by

$$I(k) = I(k-1) + \frac{K_p}{T_i} e(k) - \frac{1}{T_i} (v(k) - u(k)) \quad (12.13)$$

where T_i is the **tracking time constant**. This is a tuning parameter that determines the rate at which the integral term is reset.

Example 12.9

Consider the digital PI controller transfer function

$$C(z) = \frac{1.2z - 1.185}{z - 1}$$

with $K_p = 1.2$, $T_i = 8$, sampling period $T = 0.1$, and the process transfer function

$$G(s) = \frac{1}{10s + 1} e^{-5s}$$

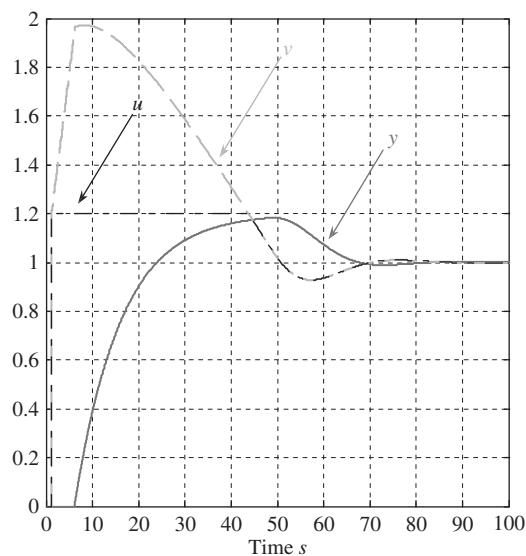
Obtain the step response of the system (1) if the saturation limits of the actuator are $u_{\min} = -1.2$ and $u_{\max} = 1.2$, and (2) with no actuator saturation. Compare and discuss the two responses, and then use back-calculation to reduce the effect of windup on the step response.

Solution

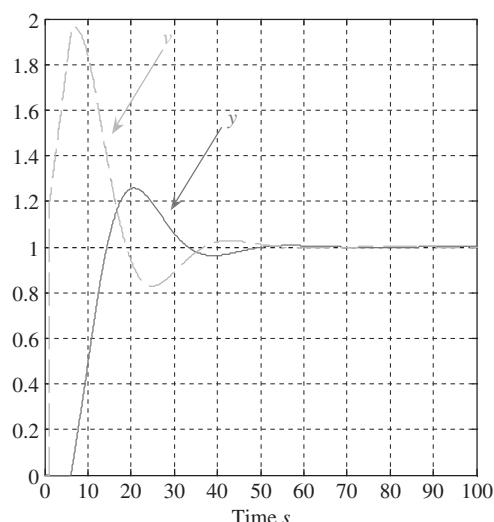
For the system with actuator saturation and no antiwindup strategy, the process output y , controller output v , and process input u are shown in Fig. 12.12. We observe that the actuator output exceeds the saturation level even when the process output attains its reference value, which leads to a large overshoot and settling time. The response after the removal of saturation nonlinearity is shown in Fig. 12.13. The absence of saturation results in a faster response with fast settling to the desired steady-state level.

Example 12.9—cont'd

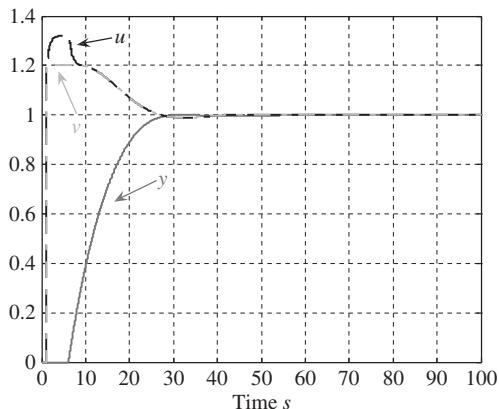
The results obtained by applying back-calculation with $T_t = T_i = 8$ are shown in Fig. 12.14. The control variable is kept at a much lower level, and this helps avoid the overshoot almost entirely. The response is slower than the response with no saturation but is significantly faster than the response with no antiwindup strategy.

**Figure 12.12**

Process input u , controller output v , and process output y with actuator saturation and no antiwindup strategy.

**Figure 12.13**

Controller output v and process output y with no actuator saturation.

Example 12.9—cont'd**Figure 12.14**

Process input u , controller output v , and process output y with actuator saturation and a back-calculation antiwindup strategy.

12.4.3 Bumpless transfer between manual and automatic mode

When the controller can operate in either **manual mode** or **automatic mode**, switching between the two modes of operation must be handled carefully to avoid a bump in the process output at the switching instant. During manual mode, the operator provides feedback control and the automatic feedback is disconnected. The integral term in the feedback controller can assume a value different from the one selected by the operator. Simply switching from automatic to manual, or vice versa, as in Fig. 12.15, leads to a bump in the control signal, even if the control error is zero. This results in an undesirable bump in the output of the system.

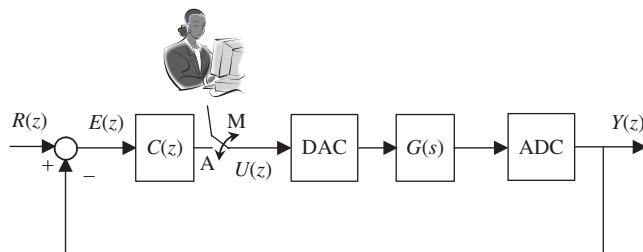


Figure 12.15
Block diagram for bumpy manual (M)/automatic (A) transfer.

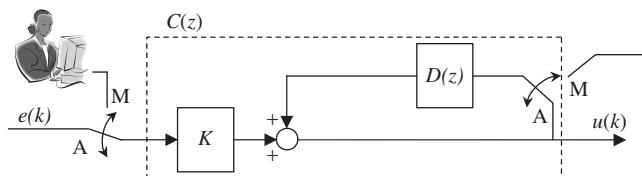


Figure 12.16
Block diagram for bumpless manual (M)/automatic (A) transfer.

For a smooth or **bumpless** transfer between manual control and the automatic digital controller $C(z)$, we use the digital scheme shown in Fig. 12.16. We write the automatic controller transfer function in terms of an asymptotically stable controller $D(z)$ with unity DC gain—that is, $D(1) = 1$ —as

$$C(z) = \frac{K}{1 - D(z)} \quad (12.14)$$

We then solve for $D(z)$ in terms of $C(z)$ to obtain

$$D(z) = \frac{C(z) - K}{C(z)}$$

If $C(z)$ is the PID controller transfer function of Eq. (12.11), we have

$$D(z) = \frac{(K_2 - K)z^2 + (K_1 + K + K_\gamma)z + K_0 - K_\gamma}{K_2 z^2 + K_1 z + K_0} \quad (12.15)$$

If the coefficient of the term z^2 in the numerator is nonzero, then the controller has the form

$$D(z) = \frac{U(z)}{E(z)} = (K_2 - K) + D_a(z)$$

where $D_a(z)$ has a first-order numerator polynomial. The controller output $u(k+2)$ is equal to the sum of two terms, one of which is the controller output $u(k+2)$ itself. Thus, the solution of the related difference equation cannot be computed by a simple recursion. This undesirable controller structure is known as an **algebraic loop**. To avoid an algebraic loop, we impose the condition

$$K = K_2 \quad (12.16)$$

to eliminate the z^2 term in the numerator. The following example illustrates the effectiveness of the bumpless manual/automatic mode scheme.

Example 12.10

Verify that a bump occurs if switching between manual and automatic operation uses the configuration shown in Fig. 12.15 for the process

$$G(s) = \frac{1}{10s + 1} e^{-2s}$$

and the PID controller ($T = 0.1$)

$$C(z) = \frac{44z^2 - 85.37z + 41.43}{z^2 - 1.368z + 0.368}$$

Design a scheme that provides a bumpless transfer between manual and automatic modes.

Solution

The unit step response for the system shown in Fig. 12.15 is shown in Fig. 12.17, together with the automatic controller output. The transfer between manual mode, where a step input

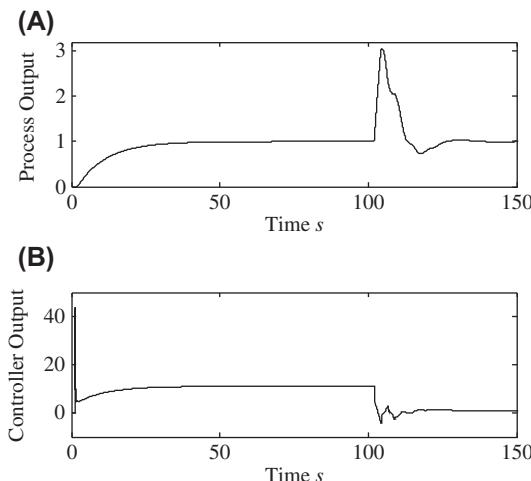


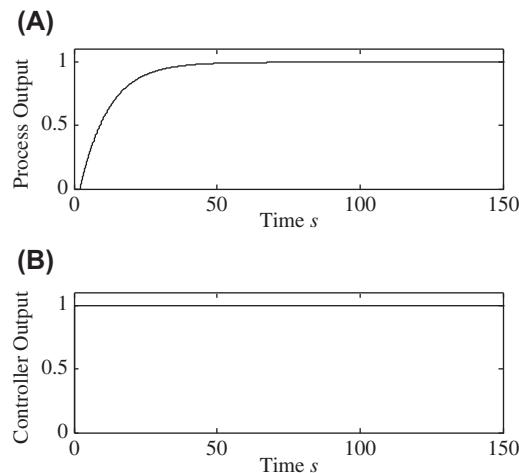
Figure 12.17

Process output (A) and controller output (B) for the system without bumpless transfer.

$u = 1$ is selected, and automatic mode occurs at time $t = 100$. The output of the PID controller is about 11, which is far from the reference value of unity at $t = 100$. This leads to a significant bump in the control variable on switching from manual to automatic control that acts as a disturbance, causing a bump in the process output. To eliminate the bump, we use the bumpless transfer configuration shown in Fig. 12.16 with the PID controller parameters

$$K_2 = 44, \quad K_1 = -85.37, \quad K_0 = 41.43, \quad \gamma = 0.368$$

Using Eqs. (12.15) and (12.16), we have

Example 12.10—cont'd**Figure 12.18**

Process output (A) and controller output (B) for the bumpless manual/automatic transfer shown in Fig. 12.16.

$$D(z) = \frac{-0.572z + 0.574}{z^2 - 1.94z + 0.942}$$

and $K = 44$. The results of Fig. 12.18 show a bumpless transfer, with the PID output at $t = 100$ equal to one.

12.4.4 Incremental form

Integrator windup and bumpless transfer issues are solved by implementing the PID controller in **incremental form**. We determine the increments in the control signal at each sampling period instead of determining the actual values of the control signal. This moves the integral action outside the control algorithm. To better understand this process, we consider the difference equation of a PID controller (the filter on the derivative action is not considered for simplicity):

$$u(k) = K_p \left(e(k) + \frac{T}{T_i} \sum_{i=0}^k e(i) + \frac{T_d}{T} (e(k) - e(k-1)) \right)$$

Subtracting the expression for $u(k-1)$ from that of $u(k)$, we obtain the increment

$$u(k) - u(k-1) = K_p \left(\left(1 + \frac{T}{T_i} + \frac{T_d}{T} \right) e(k) + \left(-1 - \frac{2T_d}{T} \right) e(k-1) + \frac{T_d}{T} e(k-2) \right)$$

which can be rewritten more compactly as

$$u(k) - u(k-1) = K_2 e(k) + K_1 e(k-1) + K_0 e(k-2) \quad (12.17)$$

where

$$K_2 = K_p \left(1 + \frac{T}{T_i} + \frac{T_d}{T} \right) \quad (12.18)$$

$$K_1 = -K_p \left(1 + \frac{2T_d}{T} \right) \quad (12.19)$$

$$K_0 = K_p \frac{T_d}{T} \quad (12.20)$$

From the difference Eq. (12.17), we can determine the increments in the control signal at each sampling period. We observe that in Eq. (12.17) there is no error accumulation (in this case the integral action can be considered “outside” the controller), and the integrator windup problem does not occur. In practice, it is sufficient that the control signal is not incremented when the actuator saturates—namely, we have $u(k) = u(k-1)$ when $v(k) = u(k)$ in Fig. 12.11. Further, the transfer between manual mode and automatic mode is bumpless as long as the operator provides the increments in the control variable rather than their total value.

The z -transform of the difference Eq. (12.17) gives the PID controller z -transfer function in incremental form as

$$C(z) = \frac{\Delta U(z)}{E(z)} = \frac{K_2 z^2 + K_1 z + K_0}{z^2} \quad (12.21)$$

where $\Delta U(z)$ is the z -transform of the increment $u(k) - u(k-1)$.

Example 12.11

For the process described in Example 12.9 and an analog PI controller with $K_p = 1.2$ and $T_i = 8$, verify that windup is avoided with the digital PI controller in incremental form ($T = 0.1$) if the saturation limits of the actuator are $u_{\min} = -1.2$ and $u_{\max} = 1.2$.

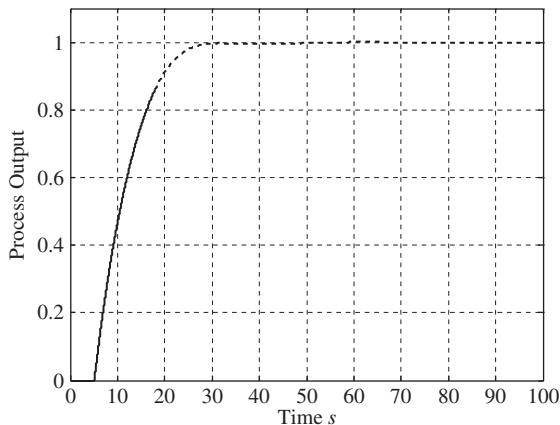
Solution

Using Eqs. (12.18) and (12.19), we obtain $K_2 = 1.215$ and $K_1 = -1.2$ and the digital controller transfer function

$$C(z) = \frac{1.215z - 1.2}{z - 1}$$

Example 12.11—cont'd

The output obtained with the PI controller in incremental form (avoiding controller updating when the actuator saturates) is shown in Fig. 12.19. Comparing the results to those shown in Fig. 12.12, we note that the response is much faster both in rising to the reference value and in settling. In Fig. 12.19, the effects of saturation are no longer noticeable.

**Figure 12.19**

Process output with the proportional–integral–derivative (PID) in incremental form.

12.5 Sampling period switching

In many control applications, it is necessary to change the sampling period during operation to achieve the optimal usage of the computational resources. In fact, a single CPU usually performs many activities such as data storage, user interface, and communication, and possibly implements more than one controller. It is therefore necessary to optimize CPU utilization by changing the sampling period. For a given digital control law, on the one hand, it is desirable to decrease the sampling period to avoid performance degradation, but on the other hand, decreasing it can overload the CPU and violate real-time constraints.

The problem of changing the control law when the sampling frequency changes can be solved by switching between controllers working in parallel, each with a different sampling period. This is a simple task provided that bumpless switching is implemented (see Section 12.4.3). However, using multiple controllers in parallel is computationally inefficient and is unacceptable if the purpose is to optimize CPU utilization. Thus, it is necessary to shift from one controller to another when the sampling period changes rather than operate controllers in parallel.

This requires computing the new controller parameters as well as past values of error and control variables that it needs to compute the control before switching. If the original sampling interval is T' and the new sampling interval is T , then we must switch from the controller

$$\begin{aligned} C'(z) &= \frac{U(z)}{E(z)} = \frac{\mathbf{b}^T \mathbf{z}_m}{\mathbf{a}^T \mathbf{z}_n} \\ \mathbf{a} &= [a_0(T') \quad a_1(T') \quad \dots \quad 1]^T \\ \mathbf{b} &= [b_0(T') \quad b_1(T') \quad \dots \quad b_m(T')]^T \\ \mathbf{z}_l &= [1 \quad z \quad \dots \quad z^l] \end{aligned} \quad (12.22)$$

to the controller

$$\begin{aligned} C(z) &= \frac{U(z)}{E(z)} = \frac{\mathbf{b}^T \mathbf{z}_m}{\mathbf{a}^T \mathbf{z}_n} \\ \mathbf{a} &= [a_0(T) \quad a_1(T) \quad \dots \quad 1]^T \\ \mathbf{b} &= [b_0(T) \quad b_1(T) \quad \dots \quad b_m(T)]^T \\ \mathbf{z}_l &= [1 \quad z \quad \dots \quad z^l] \end{aligned}$$

Equivalently, we switch from the difference equation

$$\begin{aligned} u(kT') &= -a_{n-1}(T')u((k-1)T') - \dots - a_0(T')u((k-n)T') \\ &\quad + b_0(T')e((k-n+m)T') + \dots + b_m(T')e((k-n)T') \end{aligned}$$

to the difference equation

$$\begin{aligned} u(kT) &= -a_{n-1}(T)u((k-1)T) - \dots - a_0(T)u((k-n)T) \\ &\quad + b_0(T)e((k-n+m)T) + \dots + b_m(T)e((k-n)T) \end{aligned}$$

Thus, at the switching time instant, we must recompute the values of the parameter vectors \mathbf{a} and \mathbf{b} , as well as the corresponding past $m+1$ values of the tracking error e and the past n values of the control variable u .

We compute the new controller parameters using the controller transfer function, which explicitly depends on the sampling period. For example, if the PID controller of Eq. (12.11) is used, the new controller parameters can be easily computed using Eq. (12.12) with the new value of the sampling period T , or equivalently, using Eq. (12.10) where the sampling period T appears explicitly. To compute the past values of the tracking error e and of the control variable u , different techniques can be applied depending on whether the sampling period increases or decreases. For simplicity, we only consider the case

where one sampling period is a divisor or multiple of the other. The case where the ratio between the previous and the new sampling periods (or vice versa) is not an integer is a simple extension, which is not considered here.

If the new sampling period T is a fraction of the previous sampling period T' (i.e., $T' = \lambda T$), the previous n values of the control variable $[u((k-1)T), u((k-2)T), \dots, u((k-n)T)]$ are determined with the control variable kept constant during the past λ periods. The $m + 1$ previous error values are computed using an **interpolator** such as a cubic polynomial. In particular, the coefficients c_3, c_2, c_1 , and c_0 of a third-order polynomial $\tilde{e}(t) = c_3 t^3 + c_2 t^2 + c_1 t + c_0$ can be determined by considering the past three samples and the current value of the control error. The data yield the following linear system:

$$\begin{bmatrix} ((k-3)T')^3 & ((k-3)T')^2 & (k-3)T' & 1 \\ ((k-2)T')^3 & ((k-2)T')^2 & (k-2)T' & 1 \\ ((k-1)T')^3 & ((k-1)T')^2 & (k-1)T' & 1 \\ (kT')^3 & (kT')^2 & kT' & 1 \end{bmatrix} \begin{bmatrix} c_3 \\ c_2 \\ c_1 \\ c_0 \end{bmatrix} = \begin{bmatrix} e((k-3)T') \\ e((k-2)T') \\ e((k-1)T') \\ e(kT') \end{bmatrix} \quad (12.23)$$

Once the coefficients of the polynomial function have been determined, the previous values of the control error with the new sampling period are the values of the polynomial functions at the required sampling instants. The procedure is illustrated in Fig. 12.20, where the dashed line connecting the control error values between $(k-3)T$ and kT is the polynomial function $\tilde{e}(t)$.

If the new sampling period T is a multiple of the previous sampling period T' (i.e., $T = \lambda T'$), the previous m error samples are known. However, the memory buffer must be large enough to store them. If the pole-zero difference of the process is equal to one, the equivalent n past control actions are approximately computed as the outputs estimated

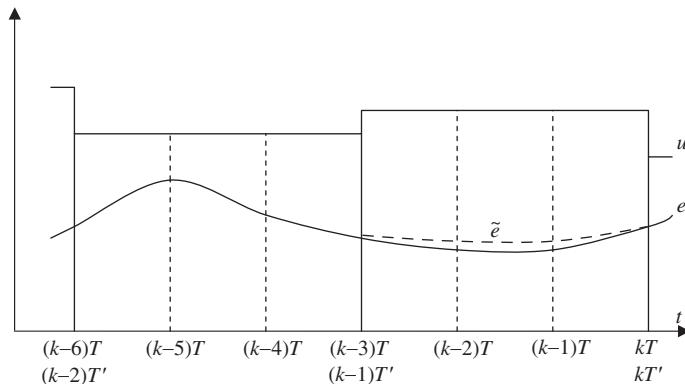


Figure 12.20

Switching controller sampling from a sampling period T' to a faster sampling rate with sampling period $T = T'/3$.

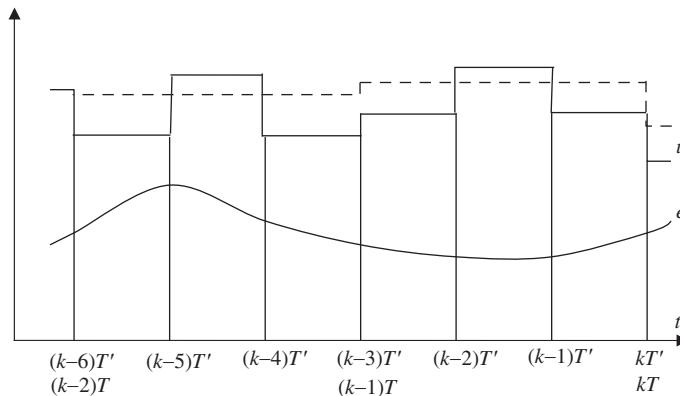


Figure 12.21

Switching controller sampling from a sampling period T' to a slower sampling rate with sampling period $T = 3T'$.

using the model of the control system with sampling period T . Specifically, let the process model obtained by discretizing the analog process with sampling period T be

$$G_{ZAS}(z) = \frac{Y(z)}{U(z)} = \frac{\beta_{h-1}z^{h-1} + \beta_{h-2}z^{h-2} + \dots + \beta_0}{z^h + \alpha_{h-1}z^{h-1} + \dots + \alpha_0}$$

where $h \geq n$. Then the equivalent past n control actions $u((k-1)T)$, $u((k-2)T)$, ..., $u((k-n)T)$ are determined by minimizing the difference between the measured output and that estimated by the model at the switching time:

$$\begin{aligned} & \min |y(kT) - (-\alpha_{h-1}y((k-1)T) - \dots - \alpha_0y((k-h)T) \\ & + \beta_{h-1}u((k-1)T) + \dots + \beta_0u((k-h)T))| \end{aligned} \quad (12.24)$$

We solve the optimization problem numerically using an appropriate approach such as the **simplex algorithm** (see Section 12.5.1). To initiate the search, initial conditions must be provided. The initial conditions can be selected as the values of the control signal at the same sampling instants determined earlier with the faster controller. The situation is depicted in Fig. 12.21.

12.5.1 MATLAB commands

When the sampling frequency is increased, the array **en** of the $m + 1$ previous error values can be computed using the following MATLAB command:

```
>> en = interp1(ts, es, tn', 'cubic')
```

where **ts** is a vector containing the last four sampling instants $[(k-3)T', (k-2)T', (k-1)T', kT']$ of the slower controller and **es** is a vector containing the corresponding control errors

$[e((k-3)T'), e((k-2)T'), e((k-1)T'), e(kT')]$. The array **tn** contains the $m + 1$ sampling instants for which the past control errors for the new controller must be determined.

Alternatively, the vector of the coefficients $\mathbf{c} = [c_3, c_2, c_1, c_0]^T$ of the cubic polynomial can be obtained by solving the linear system (12.23) using the command

```
>> c = linsolve(M, e)
```

where **M** is the matrix containing the time instants and **e** is the column vector of error samples such that Eq. (12.23) is written as $\mathbf{M}^* \mathbf{c} = \mathbf{e}$. Once the coefficients of the polynomial are computed, the values of the control error at previous sampling instants can be easily determined by interpolation.

To find the past control values using the simplex algorithm, use the command

```
>> un = fminsearch(@(u)(abs(yk - ah1 * yk1 - ...
- a0*ykh+bh1 * u(1)+...+b0 * u(n))), initcond)
```

where **initcond** is the vector of n elements containing the initial conditions, **un** is the vector of control variables $[u((k-1)T), u((k-2)T), \dots, u((k-n)T)]$, and the terms in the **abs** function are the values corresponding to Eq. (12.24), with the exception of **u(1), ..., u(n)** that must be written explicitly. Details are provided in the following examples.

Example 12.12

Discuss the effect of changing the sampling rate on the step response of the process described in Example 12.9 and an analog PI controller with $K_p = 1.2$ and $T_i = 8$. Initially, use a digital PI controller in incremental form with $T' = 0.4$ s; then switch at time $t = 22$ s to a faster controller with $T = 0.1$ s.

Solution

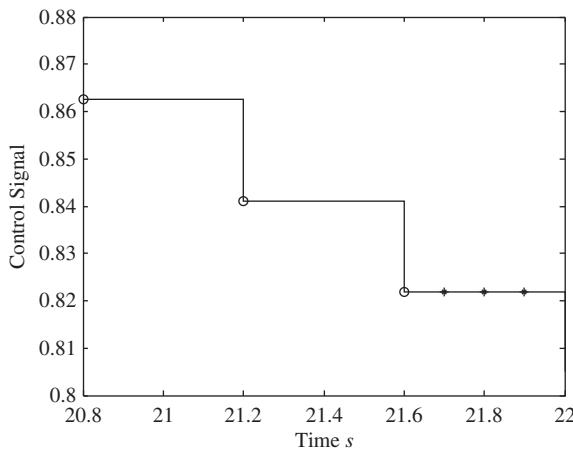
For PI control, we set T_d to zero in Eqs. (12.18)–(12.20) to obtain the parameters

$$\begin{aligned} K_2 &= K_p \left(1 + \frac{T'}{T_i} \right) = 1.2(1 + 0.4/8) = 1.26 \\ K_1 &= -K_p = -1.2 \\ K_0 &= 0 \end{aligned}$$

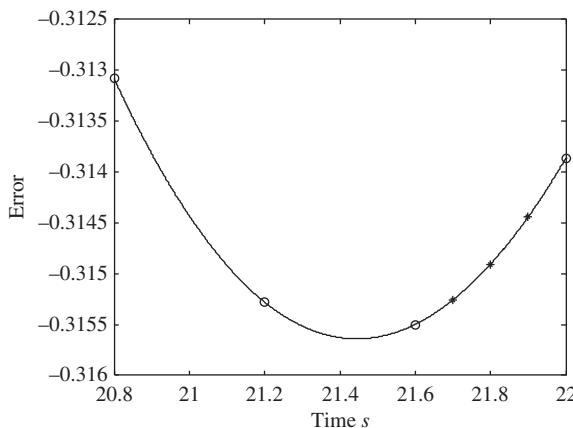
Using Eq. (12.21), the initial digital controller transfer function is

$$C(z) = \frac{U(z)}{E(z)} = \frac{1.26z - 1.2}{z - 1}$$

We simulate the system using MATLAB to compute the control and error signals. The control signal at $t = 20.8, 21.2$, and 21.6 s is plotted in Fig. 12.22, whereas the error at time $t = 20.8, 21.2, 21.6$, and 22 s is plotted as circles in Fig. 12.23. At time $t = 22$ s, the sampling

Example 12.12—cont'd**Figure 12.22**

Control signal (*solid line*) obtained by simulating the system of Example 12.12 with sampling period $T' = 0.4$ s. Circles denote the control values at the sampling instants and stars denote the control values for $T = 0.1$ s.

**Figure 12.23**

Control error interpolation $\tilde{e}(t)$ (*solid line*) obtained by simulating the system of Example 12.12 with sampling period $T' = 0.4$ s. Circles denote the error values at the sampling instants and stars denote the error values for $T = 0.1$ s.

period switches to $T = 0.1$ s and the control signal must be recomputed. The new PI controller parameters are

$$K_2 = K_p \left(1 + \frac{T'}{T_i} \right) = 1.2 \left(1 + 0.1/8 \right) = 1.215$$

$$K_1 = -K_p = -1.2$$

$$K_0 = 0$$

yielding the following transfer function:

$$C(z) = \frac{1.215z - 1.2}{z - 1}$$

Example 12.12—cont'd

The associated difference equation is therefore

$$u(k) = u(k-1) + 1.215e(k) - 1.2e(k-1) \quad (12.25)$$

Whereas at $t = 22$ s the value of $e(k) = -0.3139$ is unaffected by switching, the values of $u(k-1)$ and $e(k-1)$ must in general be recomputed for the new sampling period. For a first-order controller, the value of $u(k-1)$ is the value of the control signal at time $t = 22-T = 22-0.1 = 21.9$ s. This value is the same as that of the control signal for sampling period T' over the interval 21.6–22 s. From the simulation results shown in Fig. 12.23, we obtain $u(k-1) = 0.8219$, where the values of the control signal at $t = 21.7, 21.8$, and 21.9 s with the new controller are denoted by stars.

To calculate the previous value of the control error $e(k-1)$, we use a cubic polynomial to interpolate the last four error values with the slower controller (the interpolating function is the solid line in Fig. 12.23). Using the control errors at $t = 21.7, 21.8, 21.9$ s, $e(21.6) = -0.3155$, $e(21.2) = -0.3153$, and $e(20.8) = -0.3131$, the linear system Eq. (12.23) is solved for the coefficients $c_3 = -0.0003$, $c_2 = 0.0257$, $c_1 = -0.6784$, and $c_0 = 5.4458$. To solve the linear system, we use the MATLAB commands

```
>> M = [20.8^3 20.8^2 20.8 1; 21.2^3 21.2^2 21.2 1; 21.6^3 21.6^2 21.6 1; 22.0^3 22.0^2 22.0 1];
>> e = [-0.3131 -0.3153 -0.3155 -0.3139]';
>> c = linsolve(M,e); % Solve the linear system M c = e
>> en = c(1)*21.9^3 + c(2)*21.9^2 + c(3)*21.9 + c(4);
```

The value of the control error at time $t = 21.9$ s is therefore

$$e(k-1) = -0.0003 \cdot 21.9^3 + 0.0257 \cdot 21.9^2 - 0.6784 \cdot 21.9 + 5.4458 = -0.3144$$

Alternatively, we can use the MATLAB command for cubic interpolation

```
>> en = interp1([20.8 21.2 21.6 22], [-0.3131 -0.3153 -0.3155 -0.3139], 21.9, 'cubic');
```

The control error for the faster controller at $t = 21.7, 21.8$, and 21.9 s is denoted by stars in Fig. 12.23. From (12.25), we observe that only the error at $t = 21.9$ is needed to calculate the control at $t = 22$ s after switching to the faster controller. We compute the control value:

$$\begin{aligned} u(k) &= u(k-1) + 1.215e(k) - 1.2e(k-1) \\ &= 0.8219 + 1.215 \times (-0.3139) - 1.2 \times (-0.3144) = 0.8178 \end{aligned}$$

We compute the control at time $t = 22.1$ s and the subsequent sampling instants using the same expression without the need for further interpolation. Note that the overall performance is not significantly affected by changing the sampling period, and the resulting process output is virtually the same as the one shown in Fig. 12.13.

Example 12.13

Design a digital controller for the DC motor speed control system described in Example 6.17 with transfer function

$$G(s) = \frac{1}{(s+1)(s+10)}$$

to implement the analog PI controller

$$C(s) = 47.2 \frac{s+1}{s}$$

with the sampling period switched from $T' = 0.01$ s to $T = 0.04$ s at time $t = 0.5$ s. Obtain the step response of the closed-loop system, and discuss your results.

Solution

Applying the bilinear transformation with $T' = 0.01$ to the controller transfer function $C(s)$, we obtain the initial controller transfer function

$$C(z) = \frac{47.44z - 46.96}{z - 1}$$

We simulate the system using MATLAB to compute the error and process output values. The error values at $t = 0.42, \dots, 0.5$ s are shown in Fig. 12.24. The process output for a unit step reference input at the same sampling instants is shown in Fig. 12.25.

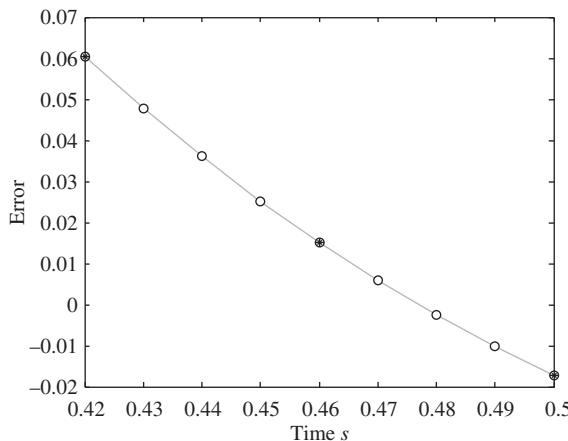
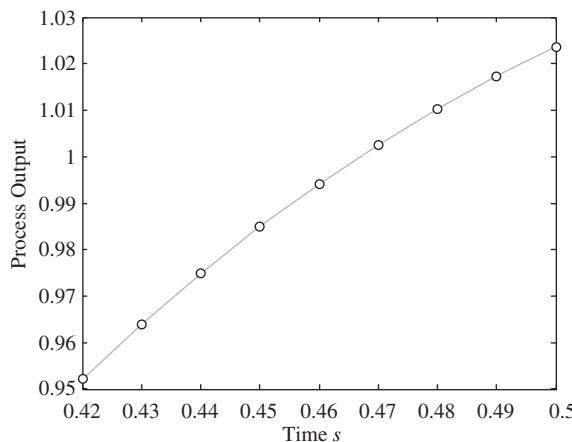


Figure 12.24

Control error for the two controllers described in Example 12.13. Circles denote the error values for the faster controller; the dark circles denote the error values for the slower controller.

Example 12.13—cont'd**Figure 12.25**

Process output for the faster controller described in Example 12.13.

Starting at $t = 0.5$ s, the controller transfer function obtained by bilinearly transforming the analog controller with sampling period $T = 0.04$ s becomes

$$C(z) = \frac{U(z)}{E(z)} = \frac{48.14z - 46.26}{z - 1}$$

The corresponding difference equation

$$u(k) = u(k-1) + 48.14e(k) - 46.26e(k-1)$$

is used to calculate the control variable starting at $t = 0.5$ s. From the MATLAB simulation results, the error values needed to compute the control at $t = 0.5$ s are $e(k) = -0.0172$ at $t = 0.5$ s and $e(k-1) = 0.0151$ at $t = 0.5 - 0.04 = 0.46$ s. The control $u(k-1)$ at $t = 0.46$ s must be determined by solving the optimization problem (12.24). The z-transfer function of the plant, ADC and DAC, with $T = 0.04$ is

$$\begin{aligned} G_{ZAS}(z) &= \frac{Y(z)}{U(z)} = 10^{-4} \frac{6.936z + 5.991}{z^2 - 1.631z + 0.644} \\ &= \frac{6.936 \times 10^{-4}z^{-1} + 5.991 \times 10^{-4}z^{-2}}{1 - 1.631z^{-1} + 0.644z^{-2}} \end{aligned}$$

and the corresponding difference equation is

$$y(k) = 1.631y(k-1) - 0.644y(k-2) + 6.936 \cdot 10^{-4}u(k-1) + 5.991 \cdot 10^{-4}u(k-2)$$

Example 12.13—cont'd

Therefore, the optimization problem is

$$\min |y(0.5) - (1.631y(0.46) - 0.644y(0.42) + 6.936 \times 10^{-4}u(0.46) + 5.991 \times 10^{-4}u(0.42))|$$

Using the output values $y(0.5) = 1.0235$, $y(0.46) = 0.9941$, and $y(0.42) = 0.9522$, the values of $u(0.46)$ and $u(0.42)$ are computed by solving the optimization problem using the following MATLAB command:

```
>> u = fminsearch(@(un)(abs(1.0235-(1.631*0.9941-0.644*0.9522 + 6.936e-4
*un(1) + 5.991e-4*un(2)))),[13.6 11.5]);
```

The initial conditions $u(k-1) = 13.6$ at $t = 0.46$ and $u(k-2) = 11.5$ at $t = 0.42$ are obtained from the values of the control variable at time $t = 0.42$ and $t = 0.46$ with the initial sampling period $T' = 0.01$ s (Fig. 12.26). The optimization yields the control values $u(0.46) = 11.3995$ and $u(0.42) = 12.4069$. The resulting value of the objective function is zero (i.e., the measured output is equal to the model estimate at the switching time). Thus, the value of the control variable at the switching time $t = 0.5$ s is

$$u(0.5) = 11.3995 + 48.14 \times -0.0172 - 46.26 \times 0.0151 = 9.8730$$

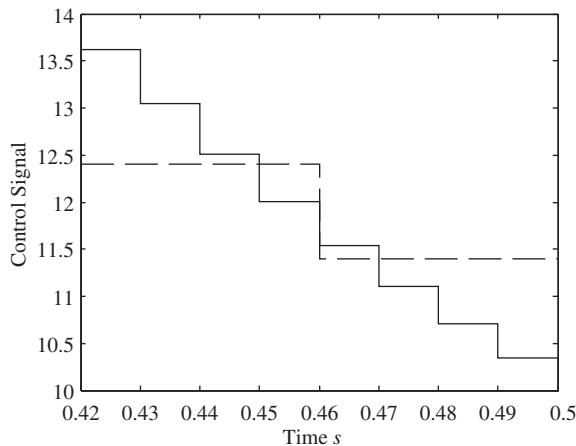
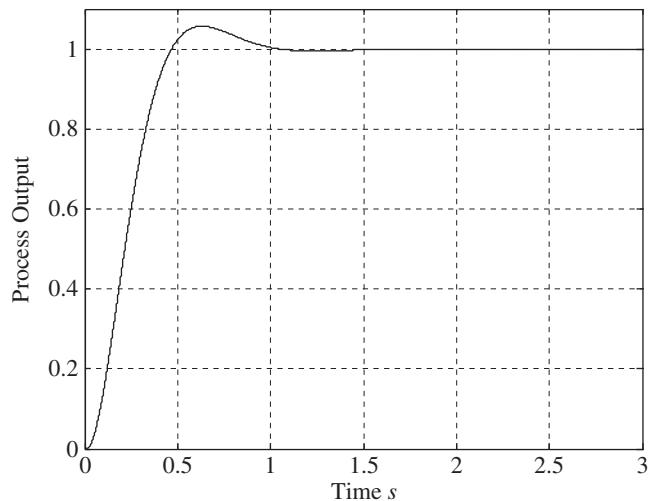


Figure 12.26

Control values for the two controllers described in Example 12.13. The solid line represents the control variable with the faster controller. The dashed line represents the equivalent control variable with the slower controller.

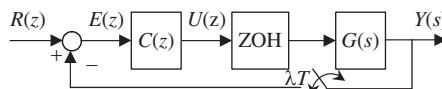
Example 12.13—cont'd**Figure 12.27**

Step response described in Example 12.13 with sampling period switching.

The step response of the system, shown in Fig. 12.27, has a small overshoot and time to first peak and a short settling time. The response is smooth and does not have a discontinuity at the switching point.

12.5.2 Dual-rate control

In some industrial applications, samples of the process output are available at a rate that is slower than the sampling rate of the controller. If performance degrades significantly when the controller sampling rate is reduced to equal that of the process output, a dual-rate control scheme can be implemented. The situation is depicted in Fig. 12.28, where it is assumed that the slow sampling period λT is a multiple of the fast sampling period T (i.e., λ is an integer). Thus, the ADC operates at the slower sampling rate $1/(\lambda T)$, whereas the controller and the zero-order hold operate at the faster sampling rate $1/T$.

**Figure 12.28**

Block diagram of dual-rate control. The controller and the zero-order hold (ZOH) operate with sampling period T ; the process output sampling period is λT .

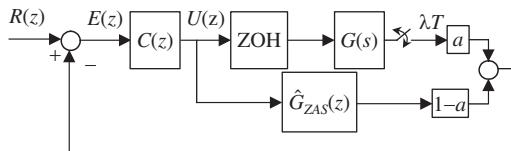


Figure 12.29
Block diagram of dual-rate inferential control.

To achieve the performance obtained when the output is sampled at the fast rate, a possible solution is to implement the so-called **dual-rate inferential control** scheme. It uses a fast-rate model of the process $\hat{G}_{ZAS}(z)$ to compute the missing output samples. The control scheme is shown in Fig. 12.29, where a is an availability parameter for the output measurement defined by

$$a = \begin{cases} 0, & t \neq k\lambda T \\ 1, & t = k\lambda T \end{cases}$$

The controller determines the values of the control variable using the measured output at $t = k\lambda T$ and using the estimated output when $t = kT$ and $t \neq k\lambda T$. In the absence of disturbances and modeling errors, the dual-rate control scheme is equivalent to the fast single-rate control scheme. Otherwise, the performance of dual-rate control can deteriorate significantly. Because disturbances and modeling errors are inevitable in practice, the results of this approach must be carefully checked.

Example 12.14

Design a dual-rate inferential control scheme with $T = 0.02$ and $\lambda = 5$ for the process (see Example 6.17)

$$G(s) = \frac{1}{(s+1)(s+10)}$$

and the controller

$$C(s) = 47.2 \frac{s+1}{s}$$

Solution

The fast rate model ($T = 0.02$) for the plant with DAC and ADC is

$$G_{ZAS}(z) = 10^{-4} \frac{1.86z + 1.729}{z^2 - 1.799z + 0.8025} = \frac{Y(z)}{U(z)} = \hat{G}_{ZAS}(z)$$

The difference equation governing the estimates of the output is

Example 12.14—cont'd

$$y(k) = 1.799y(k-1) - 0.8025y(k-2) + 1.86 \times 10^{-4}u(k-1) + 1.729 \times 10^{-4}u(k-2)$$

With $T = 0.02$, the controller transfer function obtained by bilinear transformation is

$$C(z) = \frac{47.67z - 46.73}{z - 1}$$

The controller determines the values of the control variable from the measured output at $t = 5kT$. When $t = kT$ and $t \neq 5kT$, we calculate the output estimates using the estimator difference equation.

The control scheme shown in Fig. 12.29 yields the step response shown in Fig. 12.30. The results obtained using a single-rate control scheme with $T = 0.1$ are shown in Fig. 12.31. The step response for single-rate control has a much larger first peak and settling time.

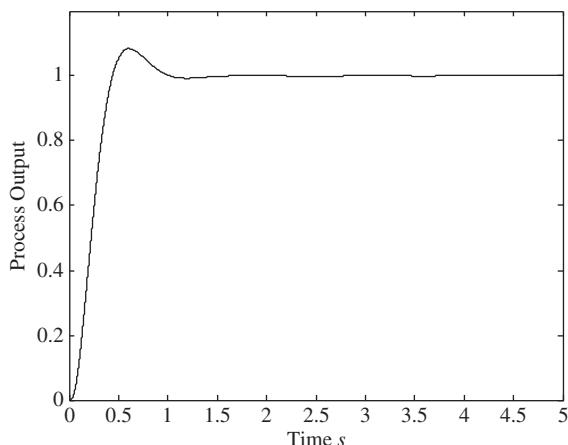


Figure 12.30

Step response described in Example 12.14 for a dual-rate controller with $T = 0.02$ and $\lambda = 5$.

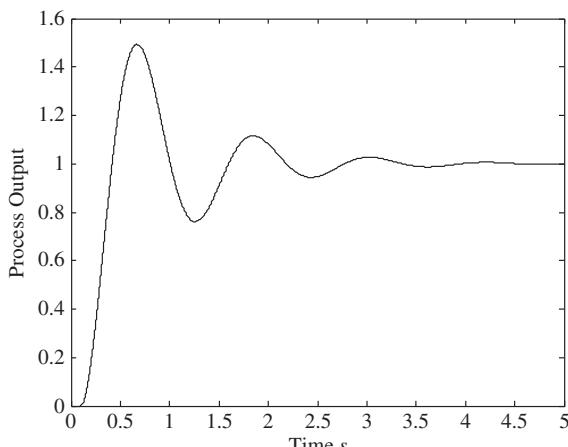


Figure 12.31

Step response described in Example 12.14 for a single-rate controller with $T = 0.1$.

Reference

Visioli, A., 2006. Practical PID Control. Springer, London, UK.

Further reading

- Albertos, P., Vallés, M., Valera, A., 2003. Controller transfer under sampling rate dynamic changes. In: Proceedings European Control Conference. Cambridge, UK.
- Åström, K.J., Hägglund, T., 2006. Advanced PID Controllers. ISA Press, Research Triangle Park, NJ.
- Cervin, A., Henriksson, D., Lincoln, B., Eker, J., Årzén, K.-E., 2003. How does control timing affect performance? *IEEE Control Syst. Mag.* 23, 16–30.
- Gambier, A., 2004. Real-time control systems: a tutorial. In: Proceedings 5th Asian Control Conference, pp. 1024–1031.
- Li, D.S., Shah, L., Chen, T., 2003. Analysis of dual-rate inferential control systems. *Automatica* 38, 1053–1059.

Problems

- 12.1 Write pseudocode that implements the following controller:

$$C(z) = \frac{U(z)}{E(z)} = \frac{2.01z - 1.99}{z - 1}$$

- 12.2 Rewrite the pseudocode for the controller described in Problem 12.1 to decrease the execution time by assigning priorities to computational tasks.
- 12.3 Design an antialiasing filter for the position control system

$$G(s) = \frac{1}{s(s + 10)}$$

with the analog controller (see Example 5.6)

$$C(s) = 50 \frac{s + 0.5}{s}$$

- Select an appropriate sampling frequency and discretize the controller.
- 12.4 Determine the mean and variance of the quantization noise when a 12-bit ADC is used to sample a variable in a range of 0–10 V for (a) rounding and (b) truncation.

- 12.5 For the system and the controller described in Problem 12.3 with a sampling interval $T = 0.02$ s, determine the decrease in the phase margin due to the presence of the ZOH.
- 12.6 Consider an oven control system (Visioli, 2006) with process transfer function

$$G(s) = \frac{1.1}{1300s + 1} e^{-25s}$$

and the PI controller

$$C(s) = 13 \frac{200s + 1}{200s}$$

Let both the actuator and the sensor signals be in the range 0–5 V, and let 1°C of the temperature variable correspond to 0.02 V. Design the hardware and software architecture of the digital control system.

- 12.7 Write the difference equations for the controller in (a) direct form, (b) parallel form, and (c) cascade form.

$$C(z) = 50 \frac{(z - 0.9879)(z - 0.9856)}{(z - 1)(z - 0.45)}$$

- 12.8 For the PID controller that results by applying the Ziegler–Nichols tuning rules to the process

$$G(s) = \frac{1}{8s + 1} e^{-2s}$$

determine the discretized PID controller transfer functions (12.11) and (12.12) with $N = 10$ and $T = 0.1$.

- 12.9 Design a bumpless manual/automatic mode scheme for the PID controller ($T = 0.1$)

$$C(z) = \frac{252z^2 - 493.4z + 241.6}{(z - 1)(z - 0.13)}$$

- 12.10 Design a bumpless manual/automatic mode scheme for the controller obtained in Example 6.19

$$C(z) = \frac{1.3932(z - 0.8187)(z - 0.9802)(z + 1)}{(z - 1)(z + 0.9293)(z - 0.96)}$$

- 12.11 Determine the digital PID controller (with $T = 0.1$) in incremental form for the analog PID controller

$$C(s) = 3 \left(1 + \frac{1}{8s} + 2s \right)$$

Computer exercises

12.12 Consider the process

$$G(s) = \frac{1}{s+1} e^{-0.2s}$$

and the PI controller with $K_p = 2$ and $T_i = 1$, discretized with sampling period $T = 0.02$, that is,

$$C(z) = \frac{2.02z - 1.98}{z - 1}$$

- (a) Using Simulink, show that by considering a range of the sensor output from -5 to 5 V and by using an 8 bit ADC converter, a limit cycle occurs in the set-point step response if the set-point value is equal to one, while there is no limit cycle if the set-point value is equal to 1.015625 .
 - (b) Repeat the simulations of (a) with unity set point for a 4-bit ADC converter.
- 12.13 Write a MATLAB script and design a Simulink diagram that implements the solution to Problem 12.8 with different filter parameter values N , and discuss the set-point step responses obtained by considering the effect of measurement noise on the process output.
- 12.14 Consider the analog process
- $$G(s) = \frac{1}{8s+1} e^{-2s}$$
- and the analog PI controller with $K_p = 3$ and $T_i = 8$. Obtain the set-point step response with a saturation limit of $u_{min} = -1.1$ and $u_{max} = 1.1$ and with a digital PI controller ($T = 0.1$) with
- a. No antiwindup
 - b. A conditional integration antiwindup strategy
 - c. A back-calculation antiwindup strategy
 - d. A digital PI controller in incremental form
- 12.15 Consider the analog process and the PI controller described in Problem 12.14. Design a scheme that provides a bumpless transfer between manual and automatic mode, and simulate it by applying a step set-point signal and by switching from manual mode, where the control variable is equal to one, to automatic mode at time $t = 60$ s. Compare the results with those obtained without bumpless transfer.
- 12.16 Design and simulate a dual-rate differential control scheme with $T = 0.01$ and $\lambda = 4$ for the plant

$$G(s) = \frac{1}{(s+1)(s+5)}$$

and the analog PI controller (see Problem 5.8). Then apply the controller to the process

$$\tilde{G}(s) = \frac{1}{(s+1)(s+5)(0.1s+1)}$$

to verify the robustness of the control system.

- 12.17 Consider the analog process and the analog PI controller described in Problem 12.16. Write a MATLAB script that simulates the step response with a digital controller when the sampling period switches at time $t = 0.52$ from $T = 0.04$ to $T = 0.01$.
- 12.18 Consider the analog process and the analog PI controller described in Problem 12.16. Write a MATLAB script that simulates the step response with a digital controller when the sampling period switches at time $t = 0.52$ from $T = 0.01$ to $T = 0.04$.
- 12.19 For the system of Example 12.4, write a MATLAB program to calculate the poles of $F(z)$ using the function **c2d** with the Tustin method and with sampling periods $T = 0.05\text{ s}$, and $T = 0.2\text{ s}$. Use both single precision (32 bit) and double precision (64 bit) for the calculations. Discuss the stability of the system in each case.

Linear matrix inequalities

Objectives

After completing this chapter, the reader will be able to

1. Formulate problems as linear matrix inequalities.
2. Transform linear inequalities to simplify them.
3. Write MATLAB functions that solve linear inequalities.

There are many problems in control system analysis and design that can be recast in the form of a linear matrix inequality (LMI). We provide a brief introduction to LMIs and their use in digital control. We also introduce MATLAB LMI commands that are part of the Robust Control Toolbox of MATLAB. Because of the introductory nature of this presentation, we avoid the complex theory underlying the numerical methods used to solve LMIs.

Chapter Outline

13.1 Linear matrix inequalities (LMI) from matrix equation 615

 13.1.1 From Linear Equations to LMIs 616

13.2 The Schur complement 617

13.3 Decision variables 620

13.4 MATLAB LMI commands 620

 13.4.1 LMI editor 626

Further reading 627

Problems 627

13.1 Linear matrix inequalities (LMI) from matrix equation

The inequality in an LMI refers to the notation $P > 0$, for a positive definite matrix, or $P < 0$, for a negative definite matrix. As an example of an LMI, consider an $n \times n$ matrix

A and we have an unknown $n \times n$ matrix P such that $A^T P A - P$ must be negative definite, we would want to solve the LMI

$$A^T P A - P < 0$$

If a solution exists for the LMI, we say that the LMI is **feasible**, otherwise it is **infeasible**. The feasibility of the LMI does not change if we perform a **congruence transformation**, i.e., for any compatible, square, invertible matrix T

$$P < 0 \Leftrightarrow T^T P T < 0$$

This follows from the fact that for any invertible matrix T

$$P < 0 \Leftrightarrow \mathbf{x}^T T^T P T \mathbf{x} = \mathbf{y}^T P \mathbf{y} < 0, \mathbf{x} \neq \mathbf{0}$$

13.1.1 From Linear Equations to LMIs

In our stability analysis and control system design for systems in state-space form, we typically encounter matrix equations. We show how such equations can be reformulated as a matrix inequality through a simple example.

Example 13.1

Rewrite the Lyapunov equation

$$A^T P A - P = -Q \quad (13.1)$$

as an LMI, where Q is a positive definite $n \times n$ symmetric matrix and A is the state matrix for discrete time linear time-invariant system.

Solution

The Lyapunov equation is known to have a symmetric positive definite solution P for any positive definite symmetric choice of the matrix Q if and only if the matrix A is Schur stable, i.e., all the eigenvalues of A are inside the unit circle. Because of our freedom to choose any matrix Q , we can rewrite the equation as two matrix inequalities

$$A^T P A - P < 0$$

$$P > 0$$

P is a symmetric positive definite matrix and hence its eigenvalues are real and can be written as

$$P = V_p \Lambda_p V_p^T$$

where V_p is the modal matrix of eigenvectors and

$$\Lambda_p = \text{diag}\{\lambda_1(P), \dots, \lambda_n(P)\}, \quad \lambda_i(P) > 0, i = 1, \dots, n$$

Multiplying P by -1 results in a negative definite matrix with negative eigenvalues and we have the equivalence

$$P > 0 \Leftrightarrow -P < 0$$

Example 13.1—cont'd

Because the eigenvalues of a block diagonal matrix are simply the union of the eigenvalues of the diagonal blocks, we can write multiple LMIs as a single LMI. We now have the matrix inequality

$$\begin{bmatrix} A^T P A - P & 0 \\ 0 & -P \end{bmatrix} < 0 \quad (13.2)$$

The solution of the LMI in Example 13.1 is nonunique since it corresponds to any matrix Q in the Lyapunov equation. The solution of the Lyapunov equation for a specific positive definite symmetric matrix Q is unique.

We now present a result that allows us to simplify many LMIs.

13.2 The Schur complement

Theorem 13.1

For matrices $Q = Q^T$, $S = S^T$, then for any compatible matrix R

$$\begin{bmatrix} Q & R \\ R^T & S \end{bmatrix} < 0 \Leftrightarrow S < 0, Q - RS^{-1}R^T < 0 \Leftrightarrow Q < 0, S - R^T Q^{-1}R < 0 \quad (13.3)$$

Proof

(\Rightarrow) Consider the quadratic form

$$[\mathbf{x}^T \quad \mathbf{y}^T] \begin{bmatrix} Q & R \\ R^T & S \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{x}^T Q \mathbf{x} + \mathbf{x}^T R \mathbf{y} + \mathbf{y}^T R^T \mathbf{x} + \mathbf{y}^T S \mathbf{y} < 0, \forall [\mathbf{x}^T \quad \mathbf{y}^T] \neq \mathbf{0}^T \quad (13.4)$$

Because the inequality holds for any choice of the vector $[\mathbf{x}^T \quad \mathbf{y}^T]$, we can prove the result by choosing specific forms. For $[\mathbf{0}^T \quad \mathbf{y}^T]$, the quadratic form (13.4) becomes $\mathbf{y}^T S \mathbf{y} < 0$, $\forall \mathbf{y} \neq \mathbf{0}$, which shows that $S < 0$ and is therefore invertible.

For $\mathbf{y} = -S^{-1}R^T \mathbf{x}$, the quadratic form (13.4) becomes

$$\mathbf{x}^T Q \mathbf{x} - 2\mathbf{x}^T R S^{-1} R^T \mathbf{x} + \mathbf{x}^T R S^{-1} S S^{-1} R^T \mathbf{x} = \mathbf{x}^T [Q - R S^{-1} R^T] \mathbf{x} < 0, \forall \mathbf{x} \neq \mathbf{0},$$

which shows that

$$Q - R S^{-1} R^T < 0$$

(\Leftarrow) Assume that

$$S < 0, \quad Q - R S^{-1} R^T < 0$$

Proof—cont'd

Differentiating with respect to \mathbf{y} gives a necessary condition for a maximum of the quadratic form

$$2S\mathbf{y} + 2R^T\mathbf{x} = \mathbf{0}$$

Solving for \mathbf{y} gives $\mathbf{y} = -S^{-1}R^T\mathbf{x}$, and substituting in the quadratic form in Eq. (13.4) gives $\mathbf{x}^T(Q - RS^{-1}R^T)\mathbf{x}$, which is negative definite by assumption. Thus, the quadratic form (13.4) has a maximum value of zero that corresponds to $[\mathbf{x}^T \quad \mathbf{y}^T] = \mathbf{0}^T$ and the quadratic form is negative definite.

The proof of $Q < 0$, $S - RQ^{-1}R^T$ is similar and is left as an exercise.

Multiplying the inequalities in Schur's complement by (-1) gives

$$-\begin{bmatrix} Q & R \\ R^T & S \end{bmatrix} > 0 \Leftrightarrow -S > 0, -Q - R(-S)^{-1}R^T > 0 \Leftrightarrow -Q > 0, -S - R^T(-Q)^{-1}R > 0$$

This shows that the result can be written as a positive inequality if the matrices Q and S are positive definite.

When using Schur's complement, we can also apply the congruence transformation to simplify the LMI, i.e., premultiply by any invertible matrix and postmultiply by its transpose. Next, we apply Schur's complement together with congruence transformation to the Lyapunov inequality of Example 13.1.

Example 13.2

Use Schur's complement to obtain an alternative form for the Lyapunov inequality.

Solution

If a matrix S is negative definite and symmetric then it can be written in terms of its modal matrix of eigenvectors V_s and matrix of eigenvalues Λ_s as

$$S = V_s \Lambda_s V_s^T$$

$$\Lambda_s = \text{diag}\{\lambda_1(S), \dots, \lambda_n(S)\}, \quad \lambda_i(S) < 0, i = 1, \dots, n$$

Inverting the matrix gives

$$S^{-1} = V_s^T \Lambda_s^{-1} V_s$$

$$\Lambda_s^{-1} = \{\lambda_1^{-1}(S), \dots, \lambda_n^{-1}(S)\}, \quad \lambda_i^{-1}(S) < 0, i = 1, \dots, n$$

The inverse of a negative definite matrix is therefore also negative definite. Hence we can write the condition $-P < 0$, as $-P^{-1} < 0$.

We rewrite the Lyapunov inequality as

$$A^T P A - P = -P - A^T(-P)A < 0$$

Example 13.2—cont'd

Applying Schur's complement to the inequalities $S < 0$, $Q - RS^{-1}R^T < 0$, with $Q = -P$, $R = A^T$, $S = -P^{-1}$, gives the form

$$\begin{bmatrix} -P & A^T \\ A & -P^{-1} \end{bmatrix} < 0$$

Perform a congruence transformation with the matrix

$$\begin{bmatrix} I & 0 \\ 0 & P \end{bmatrix}$$

to obtain

$$\begin{bmatrix} I & 0 \\ 0 & P \end{bmatrix} \begin{bmatrix} -P & A^T \\ A & -P^{-1} \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & P \end{bmatrix} = \begin{bmatrix} -P & A^T P \\ PA & -P \end{bmatrix} < 0$$

Example 13.3

Rewrite the algebraic Riccati equation (see [Section 10.3.1](#)) as the inequality

$$A^T \left\{ P - PB(R + B^T PB)^{-1} B^T P \right\} A - P + Q > 0 \quad (13.5)$$

where R is a positive definite $m \times m$ symmetric, A is the state matrix for discrete time linear time-invariant system and B is its input matrix. Obtain an equivalent LMI for [\(13.5\)](#)

Solution

Expanding the Riccati inequality gives

$$(A^T PA - P) - (A^T PB)(R + B^T PB)^{-1}(A^T PB)^T > 0$$

The solution of equation P must be positive definite and this gives the inequality $P > 0$.

Applying Schur's complement to the inequality $S > 0$, $Q - RS^{-1}R^T > 0$ with $Q = A^T PA - P + Q$, $R = -A^T PB$, $S = R + B^T PB$, gives the simpler form

$$\begin{bmatrix} A^T PA - P + Q & -A^T PB \\ -B^T PA & R + B^T PB \end{bmatrix} > 0, \quad P > 0$$

13.3 Decision variables

The unknown variables in an LMI are known as decision variables. Although this is not required to use LMI solvers, it is beneficial to examine the way that these solvers view an LMI. In general, the algorithms that these solvers use consider the form

$$F(\mathbf{p}) = F_0 + \sum_{i=1}^m F_i p_i, \quad \mathbf{p} = [p_1 \quad p_2 \quad \dots \quad p_m]^T \quad (13.6)$$

where F_i , $i = 1, \dots, m$, are square symmetric matrices and p_i , $i = 1, \dots, m$, are known as **decision variables**. The following simple example shows how an LMI can be easily written in the form of Eq. (13.6).

Example 13.4

Write the Lyapunov LMI for a 2×2 matrix in companion form in the LMI form of Eq. (13.6).

Solution

The matrices we need are

$$A = \begin{bmatrix} 0 & 1 \\ -a_0 & -a_1 \end{bmatrix}, \quad P = \begin{bmatrix} p_1 & p_2 \\ p_2 & p_3 \end{bmatrix}$$

The vector of decision variable is

$$\mathbf{p} = [p_1 \quad p_2 \quad p_3]^T$$

Expanding the Lyapunov inequality gives

$$A^T P A - P = \begin{bmatrix} a_0^2 p_3 - p_1 & -(a_0 + 1)p_2 + a_1 a_0 p_3 \\ -(a_0 + 1)p_2 + a_1 a_0 p_3 & p_1 - 2a_1 p_2 + (a_1^2 - 1)p_3 \end{bmatrix} < 0$$

We now have the expansion

$$\begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} p_1 + \begin{bmatrix} 0 & -(a_0 + 1) \\ -(a_0 + 1) & -2a_1 \end{bmatrix} p_2 + \begin{bmatrix} a_0^2 & a_1 a_0 \\ a_1 a_0 & a_1^2 - 1 \end{bmatrix} p_3 < 0$$

13.4 MATLAB LMI commands

The MATLAB Robust Control Toolbox includes commands for solving LMIs of the form

$$A < B$$

Here A, B are functions of the variables to be evaluated. The following are the steps needed to solve an LMI with MATLAB.

1. Initiate an LMI

The first step is to tell MATLAB that we are about to define an LMI with the command.

```
>> setlmis([])
```

2. Define the LMI variables.

Next, you define the LMI variables needed with the command `lmivar`. The type and structure of the variable are the input arguments to `lmivar`.

```
>> lmivar(type, struct)
```

The variable is a matrix with its type defined by:

$$type = \begin{cases} 1, & \text{symmetric block diagonal} \\ 2, & \text{rectangular} \\ 3, & \text{complex structure} \end{cases}$$

The number of blocks and their characteristics are defined in a vector through the command **struct** that has two entries for each block. The first entry gives the size of the block while the second entry is given by

$$struct(r, 2) = \begin{cases} 1, & \text{full matrix} \\ 0, & \text{scalar matrix}, \quad r = 1, 2, \dots, n_{bl} \\ -1, & \text{zero matrix} \end{cases}$$

where n_{bl} is the number of blocks. For example, with constants p_{ij} , $i,j = 1,2,3$, the 6×6 matrix

$$P = \begin{bmatrix} p_{11} & p_{12} & 0 & 0 \\ p_{12} & p_{22} & 0 & 0 \\ 0 & 0 & p_{33}I_3 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

has three blocks and **struct** is defined with three rows, one for each block. The first row denotes a 2×2 symmetric matrix, the second a 3×3 scalar matrix (i.e., a scalar times the identity matrix), and the third a 1×1 zero. The matrix has the structure.

```
>> struct = [2,1;3,0;1,-1]
```

Because the matrix is symmetric, we define it using.

```
>> p = lmivar(1,struct)
```

For a rectangular matrix (type = 2), struct gives the order of the matrix. For example, we define a rectangular three by five matrix with the command.

```
>> p = lmivar(2,[3,5])
```

3. Define the terms of the LMI

The LMIs involving this term are defined with the command.

```
>> lmitem(termID, A, B, flag)
```

Because we often define several LMIs, each LMI is referred to by a number indicating its order in the set of LMIs. The argument **termID** is a 4×1 vector whose first entry is the LMI number with a sign indicating the sign of the term defined by the command

$$\text{termID}(1) = \begin{cases} +n_{LMI} \\ -n_{LMI} \end{cases}$$

where a positive number indicates “ <0 ” and a negative number indicates “ >0 ”. The second and third entry indicate the location of the term in the LMI matrix. The fourth term gives the LMI variable in the term and is set equal to zero if the term is a constant. For example, to have the term $-P$ in the (1,1) location, i.e., first row and first column, the command is.

```
>> lmitem([1 1 1 p], -1, 1) % LMI #1, location (1,1):-1 P 1<0  
>> lmitem([-1 1 1 p], 1, 1) % LMI #1, location (1,1):1 P 1>0
```

For a constant matrix A , the command is.

```
>> lmitem([-1 1 2 0], A) % LMI #1, location (1,2):A
```

4. Identify the LMI with a label

After specifying the LMIs to be solved, we need to send the description to the LMI solvers with the command.

```
>> lmi_name = getlmis
```

The command now associates the LMIs that have been defined with **lmi_name** for later reference. Information about the LMI can be obtained with the command.

```
>> lmiinfo(lmi_name)
```

The command starts a dialog and the user is prompted to provide the information needed.

LMI ORACLE.

This is a system of 1 LMI(s) with 1 matrix variables.

Do you want information on

- (v) matrix variables (l) LMIs (q) quit

We can refer to the LMI by its name when calling the LMI solver with the command.

5. Solve the LMI

There are several LMI solvers in MATLAB.

- (a) `feasp`: To solve an LMI feasibility problem, we have the command.

```
>> [tmin,xfeas] = feasp(lmi_name, options, target)
```

The output $tmin \leq 0$ if the LMI is feasible. The input `target` is optional, and it sets a target value for `tmin` at which the program stops. Its default value is zero. **Options** has five entries and the only two we need here are the second, which specifies the number of iterations, and the third which specifies a bound on the Frobenius norm of the matrix variable of the LMI. The other entries of the option vector can be set to zero for the default values. To run 200 iterations to solve for a matrix $P = [p_{ij}]$ satisfying

$$\|P_F\| = \sqrt{\sum_{i=1}^n \sum_{j=1}^n |p_{ij}|^2} < 10$$

we need the options vector

```
>> options = [0,200,10,0,0];
```

- (b) `gevp`: A generalized eigenvalue problem is to find the solution (λ, \mathbf{x}) , $\mathbf{x} n \times 1$, of the equation

$$A\mathbf{x} = \lambda B\mathbf{x} \Leftrightarrow [A - \lambda B]\mathbf{x} = \mathbf{0}$$

where A , B , are $m \times n$ matrices. The matrix $[A - \lambda B]$ is known as a `matrix pencil`. If $m = n$, $[A - \lambda B]$ is said to be `regular`. To solve a generalized eigenvalue problem, we have the command.

```
>> [lopt, xopt] = gevp(lmi_name, nlfc, options, linit, xinit, target)
```

The parameter `nlfc` is the number of LMIs including λ . At least one such LMI must be specified and it is in the form `options` is a vector specified as for the `feasp`, `target` is an upper bound for λ . If initial estimates $(\lambda_0, \mathbf{x}_0)$ are known, then they are used as `(linit, xinit)`. The vector `xinit` is of length equal to the number of decision variables (i.e., the length of `xopt`, which is the number of entries needed to define the matrix that the command solves for). We can obtain this number with the command.

```
>> decnbr(lmis)
```

- (c) **mincx:** Linear programming is the problem of finding the solution of the optimization problem.

$\min_{\mathbf{x}} \mathbf{c}^T \mathbf{x}$ subject to $A\mathbf{x} = \mathbf{b}$.

where \mathbf{c}, \mathbf{b} are $n \times 1$, and A is $n \times n$.

The MATLAB command to solve the problem is

```
>> [copt, xopt] = mincx(lmi_name, c, options, xinit, target)
```

The parameters **option** and **target** are defined as for the other solvers.

6. Obtain the LMI variables

The solution is in $xfeas$ but is not in a useful form and to obtain the desired matrix we need the command.

```
>> P = dec2mat(lmi_name, xfeas, p)
```

Example 13.5

Write a MATLAB function to solve the Lyapunov LMI

$$\begin{bmatrix} A^T P A - P & 0 \\ 0 & -P \end{bmatrix} < 0$$

Use the function to test the stability of the matrix

$$A = \begin{bmatrix} 0.3000 & -0.2000 & 0 \\ 0.8000 & -0.3000 & -0.4000 \\ 0 & 0.4000 & -0.9000 \end{bmatrix}$$

Solution

The inequality must be solved for one symmetric matrix P , and has a symmetric matrix with 2 terms in the (1,1) location. The following function solves the inequality.

```
function [P,tmin]=dtlyaplmi(a)
setlmis([]); % Initiate LMI formulation
[n,n]=size(a);
p=lmivar(1,[n 1]); % 1=symmetric, [n 1]= n by n & 1 block
% First LMI
lmiterm([1 1 1 p],a',a) % LMI term #1:a'P a symmetric
lmiterm([1 1 1 p],-1,1) % LMI term #1:-1 P 1 symmetric
% % Second LMI
lmiterm([-2 1 1 p],1,1) % LMI term #1:P
dtlyap=getlmis;
[tmin,xfeas]=feasp(dtlyap);%Solve the feasibility problem
% feasible for tmin <=0, strictly feasible if tmin <0
% tmin positive and very small for feasible but not strictly
P=dec2mat(dtlyap,xfeas,p); % change from decision variables to matrix
```

Example 13.5—cont'd

The output for the given matrix is:

Solver for LMI feasibility problems $L(x) < R(x)$

This solver minimizes t subject to $L(x) < R(x) + t*I$

The best value of t should be negative for feasibility

Iteration : Best value of t so far

$1 -0.479,818$

Result: best value of t : $-0.479,818$

f-radius saturation: 0.000% of $R = 1.00e+09$

$P =$

$2.6615 -0.4879 -0.0401$

$-0.4879 1.9813 -0.5670$

$-0.0401 -0.5670 2.1878$

$t_{\min} =$

-0.4798

The value of t_{\min} is negative and the LMI is feasible. We verify that P is positive definite since all its eigenvalues are positive.

>> eig(P)

ans =

1.3772

2.4517

3.0017

This is the correct result in this case where the matrix is Schur stable with all its eigenvalues inside the unit circle.

>> eig(A)

ans =

-0.1000

-0.3000

-0.5000

13.4.1 LMI editor

It is often easier to create an LMI using the LMI editor. The editor is a graphical user interface that is invoked with the command.

>> lmiedit

This starts a window as in Fig. 13.1, which shows the LMIs corresponding to the discrete time Lyapunov equation. We enter a name for the LMI next to the tag “**name the LMI system**”. If “**describe the matrix variables**” is selected”, then one can specify the variables of the LMI by listing them under “**variable name**” with each line corresponding to a different variable. The type of the variable is given under “**type (S/R/G)**” where S is for symmetric, R is for rectangular, and G is for general. The **structure** of the variable is given as an array in the “**structure**” window. If “**describe the LMIs as MATLAB expressions**” is selected, we can type the LMIs as MATLAB expression including inequality signs.

We can see the command corresponding to the LMI by selecting “**view commands**” as shown in Fig. 13.2. The LMI commands can also be uploaded from a file by clicking on “**read**” or written to a file by clicking on “**write**” The commands are executed to create an LMI by clicking on “**create**”. The next step is to solve the LMI with one of the LMI solvers. For the example of Fig. 13.1, the command to use the solver feasp is

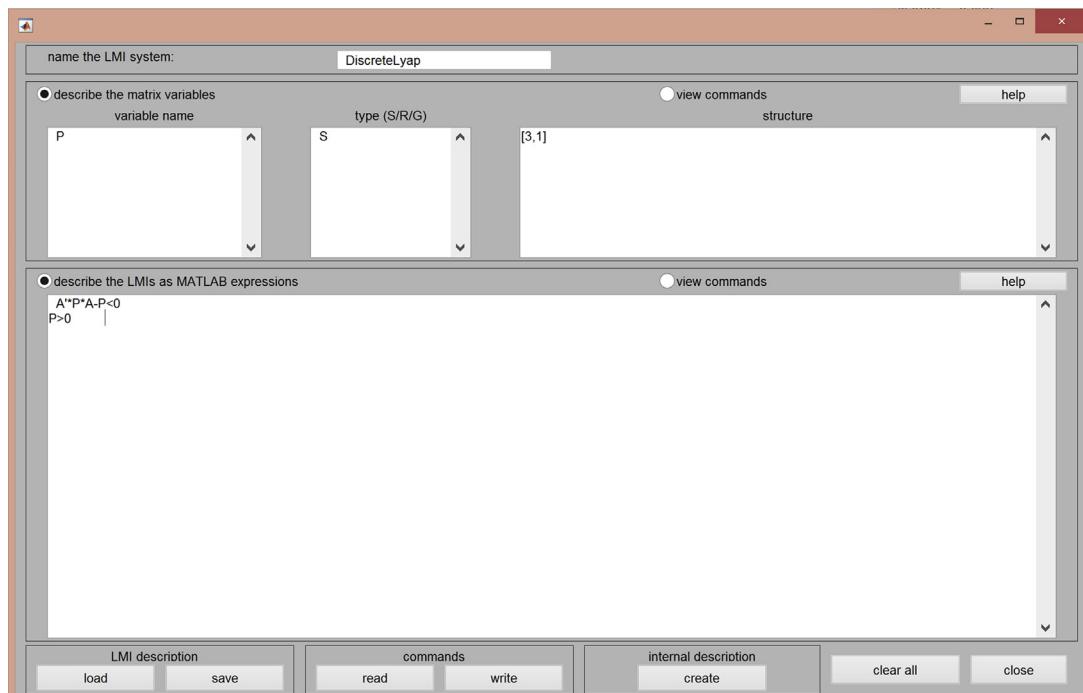


Figure 13.1

MATLAB LMI editor: displaying LMI variables and MATLAB expressions.

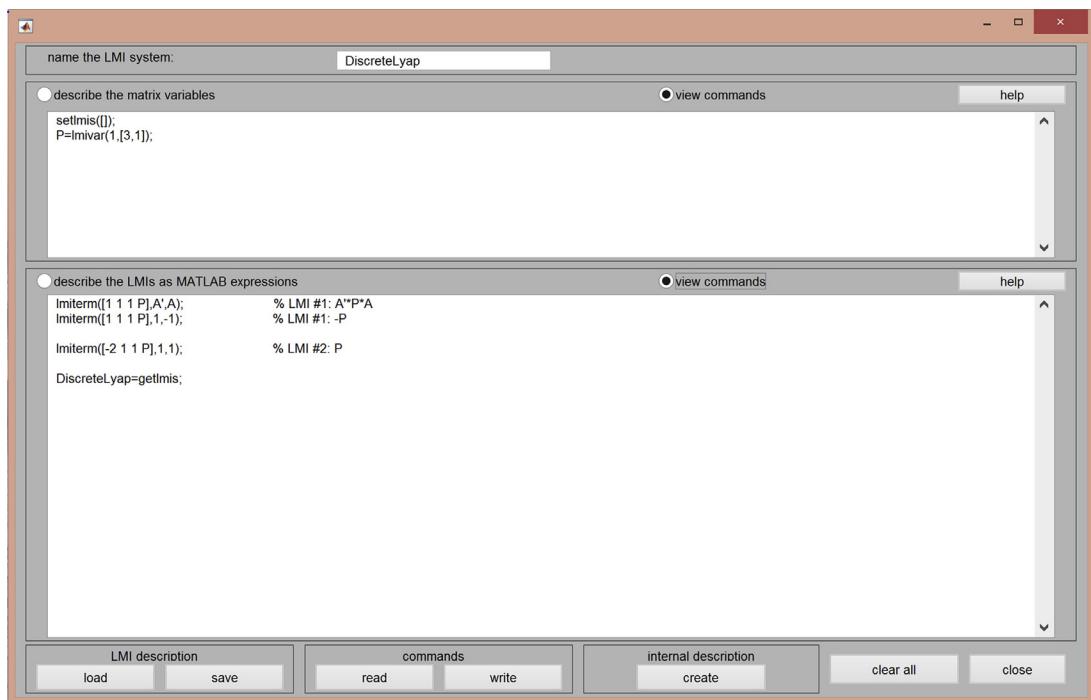


Figure 13.2
MATLAB LMI editor: displaying MATLAB commands.

```
>> [tmin, xfeas] = feasp(DiscreteLyap)
```

We obtain the variable P with the command.

```
>> P = dec2mat(DiscreteLyap, xfeas, P)
```

Further reading

Duan, G.R., Yu, H.-H., 2013. LMIs in Control Systems. CRC Press, Boca Raton, Fl.

VanAntwerp, J.G., Braatz, R.D., 2000. A tutorial on linear and bilinear matrix inequalities. *J. Process Control* 10, 363–385.

Problems

- 13.1. Complete the proof of the Schur complement result.
- 13.2. Apply Schur's complement to obtain the general conditions on the leading principal minors (minor along the diagonal) for a negative definite matrix then apply the test to the matrix with numerical entries.

$$(i) \begin{bmatrix} a_{11} & a_{12} \\ a_{12} & a_{22} \end{bmatrix} < 0, \quad \begin{bmatrix} -1 & 1 \\ 1 & -2 \end{bmatrix}$$

$$(ii) \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{12} & a_{22} & a_{23} \\ a_{13} & a_{23} & a_{33} \end{bmatrix} < 0, \quad \begin{bmatrix} -1 & -1 & -1 \\ -1 & -2 & 0 \\ -1 & 0 & -3 \end{bmatrix}$$

- 13.3. Apply Schur's complement to obtain the general conditions on the leading principal minors (minors along the diagonal) for a positive definite matrix then apply the test to the matrix with numerical entries.

$$(iii) \begin{bmatrix} a_{11} & a_{12} \\ a_{12} & a_{22} \end{bmatrix} > 0, \quad \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$$

$$(iv) \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{12} & a_{22} & a_{23} \\ a_{13} & a_{23} & a_{33} \end{bmatrix} > 0, \quad \begin{bmatrix} 1 & -1 & -1 \\ -1 & 2 & 0 \\ -1 & 0 & 3 \end{bmatrix}$$

- 13.4. Write a MATLAB function to solve the discrete-time Lyapunov inequality obtained using the Schur identity in Example 13.2 then use it to test the stability of the matrix

$$A = \begin{bmatrix} 0 & 1 \\ -0.1 & -0.02 \end{bmatrix}$$

- 13.5. Show that to verify that the eigenvalues of a matrix A satisfy $|\lambda_i| < r$, $i = 1, 2, \dots, n$, we can test the condition

$$\begin{bmatrix} -P & A^T/r \\ A/r & -P^{-1} \end{bmatrix} < 0$$

- 13.6. It can be shown that a continuous time linear system is asymptotically stable if and only if the following LMI are feasible

$$A^T P + PA < 0, \quad P > 0$$

Write a MATLAB function to solve the LMIs and use it to check the stability of the state matrix of the motor position control system

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -11 & -6 & -11 \end{bmatrix}$$

- 13.7. Show that to verify that the eigenvalues of a matrix A satisfy $\operatorname{Re}\{\lambda_i\} < -\alpha$, $i = 1, 2, \dots, n$, we can test the condition

$$\begin{bmatrix} -P & (A + \alpha I_n)^T \\ A + \alpha I_n & -P \end{bmatrix} < 0$$

- 13.8. Consider the problem of selecting a state feedback gain matrix for a discrete time system (A, B) to ensure its stability. A solution to the problem can be obtained by solving the inequality

$$(A - BK)^T P(A - BK) - P < 0$$

for the gain K with condition $P > 0$ to guarantee the stability of the closed loop state matrix $A - BK$.

- (a) Use Schur's complement to transform the inequalities to a single LMI.
- (b) Write a MATLAB function whose input is the pair (A, \mathbf{b}) with output as (i) the feedback gain matrix K that stabilizes the system, (ii) the matrix P , and (iii) the parameter **tmin**.
- (c) Obtain the results of running the function with the pair

$$A = \begin{bmatrix} 3 & -2 & 0 \\ 8 & -3 & -4 \\ 0 & 4 & -9 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Show that the closed-loop system with the gain K is Schur stable.

- 13.9. Repeat problem 13.6 with the eigenvalues of the closed-loop system required to have real parts less than $-\alpha$ and obtain numerical results for the system with $r = 0.2$
- 13.10. Consider the problem of selecting a state feedback gain matrix for a continuous-time system (A, B) to ensure its stability. A solution to the problem can be obtained by solving the inequality

$$(A - BK)P + P(A - BK)^T < 0$$

for the gain K with condition $P > 0$ to guarantee the stability of the closed loop state matrix $A - BK$.

- (a) Write a MATLAB function whose input is the pair (A, \mathbf{b}) with output as (i) the feedback gain matrix K that stabilizes the system, (ii) the matrix P , and (iii) the parameter **tmin**.
- (b) Obtain the results of running the function with the pair

$$A = \begin{bmatrix} 3 & -2 & 0 \\ 8 & -3 & -4 \\ 0 & 4 & -9 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Show that the closed-loop system with the gain K is Hurwitz stable.

- 13.11. Repeat Problem 13.10 with the eigenvalues of the closed-loop system required to have real parts less than $-\alpha$ and obtain numerical results for the system with $\alpha = -2$
- 13.12. Prove that for any compatible matrices Q_1 and Q_2 , there exists a nonnegative constant α such that

$$\exists \alpha > 0, \begin{bmatrix} \alpha Q_1^T Q_1 & \pm Q_1^T Q_2 \\ \pm Q_2^T Q_1 & \alpha^{-1} Q_2^T Q_2 \end{bmatrix} \geq 0$$

In addition, prove that if $Q_1^T Q_1 < \gamma Q_4$, then $\exists \alpha > 0$, $\begin{bmatrix} \alpha \gamma Q_4 & \pm Q_1^T Q_2 \\ \pm Q_2^T Q_1 & \alpha^{-1} Q_2^T Q_2 \end{bmatrix} \geq 0$

- 13.13. Prove that for any compatible matrices Q_1 , Q_2 , and Q_3 , there exists a nonnegative constant α such that

$$\begin{bmatrix} \alpha Q_1^T Q_1 & \pm Q_1^T Q_2 & \pm Q_1^T Q_3 \\ \pm Q_2^T Q_1 & \alpha^{-1} Q_2^T Q_2 & \pm \alpha^{-1} Q_2^T Q_3 \\ \pm Q_3^T Q_1 & \pm \alpha^{-1} Q_3^T Q_2 & \alpha^{-1} Q_3^T Q_3 \end{bmatrix} \geq 0$$

In addition, prove that if $Q_1^T Q_1 < \gamma Q_4$, then

$$\exists \alpha > 0, \begin{bmatrix} \alpha \gamma Q_4 & \pm Q_1^T Q_2 & \pm Q_1^T Q_3 \\ \pm Q_2^T Q_1 & \alpha^{-1} Q_2^T Q_2 & \pm \alpha^{-1} Q_2^T Q_3 \\ \pm Q_3^T Q_1 & \pm \alpha^{-1} Q_3^T Q_2 & \alpha^{-1} Q_3^T Q_3 \end{bmatrix} \geq 0$$

- 13.14. Show that

- (a) For Q , symmetric and $Q < 0$ we have

$$\max_{\mathbf{x}} \{2\mathbf{x}^T P \mathbf{y} + \mathbf{x}^T Q \mathbf{x}\} = -\mathbf{y}^T P^T Q^{-1} P \mathbf{y}$$

- (b) For Q , symmetric and $Q > 0$ we have

$$\min_{\mathbf{x}} \{2x^T P y + x^T Q x\} = -y^T P^T Q^{-1} P y$$

- 13.15. Show that for any matrices P , Q , R with R positive definite symmetric

$$\begin{bmatrix} P^T R P & \pm P^T Q \\ \pm Q P^T & Q^T R^{-1} Q \end{bmatrix} \geq 0$$

- 13.16. Use Schur's complement to show that for $A_{11} = A_{11}^T$, $A_{22} = A_{22}^T$, $B = B^T$

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{12}^T & A_{22} \end{bmatrix} < 0 \Leftrightarrow \begin{bmatrix} A_{11} - B & A_{12} & B \\ A_{12}^T & A_{22} & \mathbf{0} \\ B & \mathbf{0} & -B \end{bmatrix} < 0$$

- 13.17. Use Schur's complement to show that

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{12}^T & I - P \end{bmatrix} > 0, \quad I - P > 0 \Leftrightarrow \begin{bmatrix} A_{11} & A_{12} & \mathbf{0} \\ A_{12}^T & I & P \\ \mathbf{0} & P & P \end{bmatrix} > 0$$

- 13.18. Use the condition of Problem 13.17 to find the maximum value for which the matrix is positive definite

$$\begin{bmatrix} 5 & 2 & 0 \\ 2 & 1 & p \\ 0 & p & p \end{bmatrix}$$

- 13.19. Use Schur's complement to show that for negative definite $A_{ii} = A_{ii}^T, i = 1, 2, 3$, we have the equivalence

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{12}^T & A_{22} \end{bmatrix} < 0, \quad \begin{bmatrix} A_{11} & A_{13} \\ A_{13}^T & A_{33} \end{bmatrix} < 0 \Leftrightarrow \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{12}^T & A_{22} & B \\ A_{13}^T & B & A_{33} \end{bmatrix} < 0, \quad B = A_{12}^T A_{11}^{-1} A_{13}$$

- 13.20. Given real matrices, A , B , X , of order $n \times m$, $n \times m$, and $n \times n$, respectively, with $X^T X \leq \gamma I_n$. Show that for any $\delta > 0$

$$\gamma \delta A^T A + \gamma \delta^{-1} B^T B \geq A^T X^T X B + B^T X^T X A$$

Table of Laplace and z-transforms*

No.	Continuous time	Laplace transform	Discrete time	z-transform
1	$\delta(t)$	1	$\delta(k)$	1
2	$1(t)$	$\frac{1}{s}$	$1(k)$	$\frac{z}{z-1}$
3	t	$\frac{1}{s^2}$	kT^{**}	$\frac{zT}{(z-1)^2}$
4	t^2	$\frac{2!}{s^3}$	$(kT)^2$	$\frac{z(z+1)T^2}{(z-1)^3}$
5	t^3	$\frac{3!}{s^4}$	$(kT)^3$	$\frac{z(z^2+4z+1)T^3}{(z-1)^4}$
6	$e^{-\alpha t}$	$\frac{1}{s+\alpha}$	a^k^{***}	$\frac{z}{z-a}$
7	$1 - e^{-\alpha t}$	$\frac{\alpha}{s(s+\alpha)}$	$1 - a^k$	$\frac{(1-a)z}{(z-1)(z-a)}$
8	$e^{-\alpha t} - e^{-\beta t}$	$\frac{\beta-\alpha}{(s+\alpha)(s+\beta)}$	$a^k - b^k$	$\frac{(a-b)z}{(z-a)(z-b)}$
9	$te^{-\alpha t}$	$\frac{1}{(s+\alpha)^2}$	kTa^k	$\frac{azT}{(z-a)^2}$
10	$\sin(\omega_n t)$	$\frac{\omega_n}{s^2 + \omega_n^2}$	$\sin(\omega_n kT)$	$\frac{\sin(\omega_n T)_z}{z^2 - 2\cos(\omega_n T)z + 1}$
11	$\cos(\omega_n t)$	$\frac{s}{s^2 + \omega_n^2}$	$\cos(\omega_n kT)$	$\frac{z[z - \cos(\omega_n T)]}{z^2 - 2\cos(\omega_n T)z + 1}$
12	$e^{-\varsigma\omega_n t} \sin(\omega_d t)$	$\frac{\omega_d}{(s+\varsigma\omega_n)^2 + \omega_d^2}$	$e^{-\varsigma\omega_n kT} \sin(\omega_d kT)$	$\frac{e^{-\varsigma\omega_n T} \sin(\omega_d T)z}{z^2 - 2e^{-\varsigma\omega_n T} \cos(\omega_d T)z + e^{-2\varsigma\omega_n T}}$
13	$e^{-\varsigma\omega_n t} \cos(\omega_d t)$	$\frac{s + \varsigma\omega_n}{(s + \varsigma\omega_n)^2 + \omega_d^2}$	$e^{-\varsigma\omega_n kT} \cos(\omega_d kT)$	$\frac{z[z - e^{-\varsigma\omega_n T} \cos(\omega_d T)]}{z^2 - 2e^{-\varsigma\omega_n T} \cos(\omega_d T)z + e^{-2\varsigma\omega_n T}}$
14	$\sinh(\beta t)$	$\frac{\beta}{s^2 - \beta^2}$	$\sinh(\beta kT)$	$\frac{\sinh(\beta T)_z}{z^2 - 2\cosh(\beta T)z + 1}$
15	$\cosh(\beta t)$	$\frac{s}{s^2 - \beta^2}$	$\cosh(\beta kT)$	$\frac{z[z - \cosh(\beta T)]}{z^2 - 2\cosh(\beta T)z + 1}$

* The discrete time functions are generally sampled forms of the continuous time functions.

** Sampling t gives kT , whose transform is obtained by multiplying the transform of k by T .

*** The function $e^{-\alpha kT}$ is obtained by setting $a = e^{-\alpha T}$.

Properties of the z-transform

Number	Name	Formula
1	Linearity	$\mathcal{Z}\{\alpha f_1(k) + \beta f_2(k)\} = \alpha F_1(z) + \beta F_2(z)$
2	Time delay	$\mathcal{Z}\{f(k-n)\} = z^{-n}F(z)$
3	Time advance	$\mathcal{Z}\{f(k+1)\} = zF(z) - zf(0)$ $\mathcal{Z}\{f(k+n)\} = z^n F(z) - z^n f(0) - z^{n-1} f(1) \cdots - zf(n-1)$
4	Discrete time convolution	$\mathcal{Z}\{f_1(k)^* f_2(k)\} = \mathcal{Z}\left\{ \sum_{i=0}^k f_1(i) f_2(k-i) \right\} = F_1(z)F_2(z)$
5	Multiplication by exponential	$\mathcal{Z}\{a^{-k}f(k)\} = F(az)$
6	Complex differentiation	$\mathcal{Z}\{k^m f(k)\} = \left(-z \frac{d}{dz}\right)^m F(z)$
7	Final value theorem	$f(\infty) = \lim_{k \rightarrow \infty} f(k) = \lim_{z \rightarrow 1} (1 - z^{-1})F(z) = \lim_{z \rightarrow 1} (z - 1)F(z)$
8	Initial value theorem	$f(0) = \lim_{k \rightarrow 0} f(k) = \lim_{z \rightarrow \infty} F(z)$

Review of linear algebra

A.1 Matrices

An $m \times n$ matrix is an array of entries¹ denoted by

$$A = [a_{ij}] = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}$$

with m rows and n columns. The matrix is said to be of order $m \times n$.

Rectangular matrix $m \neq n$

Square matrix $m = n$

Row vector $m = 1$

Column vector $n = 1$

Example A.1 Matrix representation

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad \text{real } 2 \times 3 \text{ rectangular matrix}$$

$$A = \begin{bmatrix} 1 & 2-j \\ 4+j & 5 \end{bmatrix} \quad \begin{aligned} &\gg A = [1, 2, 3; 4, 5, 6] \\ &\quad \text{complex } 2 \times 2 \text{ square matrix} \\ &\gg A = [1, 2-j; 4+j, 5] \end{aligned}$$

A.2 Equality of matrices

Equal matrices are matrices of the same order with equal corresponding entries.

¹ Relevant MATLAB commands are given as necessary and are preceded by “ $>>$.”

$$A = B \Leftrightarrow [a_{ij}] = [b_{ij}], \quad i = 1, 2, \dots, m \\ j = 1, 2, \dots, n$$

Example A.2 Equal matrices

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad B = \begin{bmatrix} 1.1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad C = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

$A \neq B, \quad C \neq B, \quad A = C$

A.3 Matrix arithmetic**A.3.1 Addition and subtraction**

The sum (difference) of two matrices of the same order is a matrix with entries that are the sum (difference) of the corresponding entries of the two matrices.

$$C = A \pm B \Leftrightarrow [c_{ij}] = [a_{ij} \pm b_{ij}], \quad i = 1, 2, \dots, m \\ j = 1, 2, \dots, n$$

Example A.3 Matrix addition/subtraction

`>> C = A+B`
`>> C = A-B`

Note: MATLAB accepts the command

`>> C = A + b`

if **b** is a scalar. The result is the matrix $C = [a_{ij}+b]$.

A.3.2 Transposition

Interchanging the rows and columns of a matrix,

$$C = A^T \Leftrightarrow [c_{ij}] = [a_{ji}], \quad i = 1, 2, \dots, m \\ j = 1, 2, \dots, n$$

Example A.4 Matrix transpose

$$A = \begin{bmatrix} 1 & 2-j \\ 4+j & 5 \end{bmatrix}$$

$\gg \mathbf{B} = \mathbf{A}'$

$$B = \begin{bmatrix} 1 & 4-j \\ 2+j & 5 \end{bmatrix}$$

Note: The ('') command gives the complex conjugate transpose for a complex matrix.

Symmetric matrix $A = A^T$

Hermitian matrix $A = A^*$ (*denotes the complex conjugate transpose).

Example A.5 Symmetric and hermitian matrix

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 5 \end{bmatrix} \quad \text{symmetric}$$

$$A = \begin{bmatrix} 1 & 2-j \\ 2+j & 5 \end{bmatrix} \quad \text{Hermitian}$$

Notation

Column vector $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^T = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

Row vector $\mathbf{x}^T = [x_1 \ x_2 \ \dots \ x_n]$

Example A.6 Column and row vectors

$$\mathbf{x} = [1 \ 2 \ 3]^T \quad \mathbf{y}^T = [1 \ 2 \ 3]$$

Example A.6 Column and row vectors—cont'd

```
>> x = [1; 2; 3]
```

```
x = 1
```

```
2
```

```
3
```

```
>> y = [1, 2, 3]
```

```
y = 1 2 3
```

A.3.3 Matrix multiplication**Multiplication by a scalar**

Multiplication of every entry of the matrix by a scalar.

$$C = \alpha A \Leftrightarrow [c_{ij}] = [\alpha a_{ji}], \quad i = 1, 2, \dots, m \\ j = 1, 2, \dots, n$$

>> C = a*A

Note: MATLAB is a case-sensitive mode that distinguishes between uppercase and lowercase variables.

Multiplication by a matrix

The product of an $m \times n$ matrix and an $n \times l$ matrix is an $m \times l$ matrix—that is, $(m \times n) \cdot (n \times l) = (m \times l)$.

$$C = AB \Leftrightarrow [c_{ij}] = \left[\sum_{k=1}^n a_{ik} b_{kj} \right], \quad i = 1, 2, \dots, m \\ j = 1, 2, \dots, l$$

Noncommutative $AB \neq BA$ (in general).

Normal matrix $A^*A = AA^*$ (commutative multiplication with its conjugate transpose).

Clearly, any symmetric (Hermitian) matrix is also normal. But some normal matrices are not symmetric (Hermitian).

Premultiplication by a row vector $m = 1$

$$(1 \times n) \cdot (n \times l) = (1 \times l) \quad C = \mathbf{c}^T = [c_1 \ c_2 \ \dots \ c_l]$$

Postmultiplication by a column vector $l = 1$

$$(m \times n) \cdot (n \times 1) = (m \times 1) \quad C = \mathbf{c} = [c_1 \ c_2 \ \dots \ c_m]^T$$

Multiplication of a row by a column $m = l = 1$

$$(1 \times n) \cdot (n \times 1) = (1 \times 1) \quad c = \text{scalar}$$

Note that this product is the same for any two vectors regardless of which vector is transposed to give a row vector because

$$c = \mathbf{a}^T \mathbf{b} = \mathbf{b}^T \mathbf{a} = \sum_{k=1}^n a_i b_i.$$

This defines a dot product for any two real vectors and is often written in the form

$$\langle \mathbf{a}, \mathbf{b} \rangle$$

Multiplication of a column by a row $n = 1$

$$(m \times 1) \cdot (1 \times l) = (m \times l) \quad C = m \times l \text{ matrix}$$

Positive integral power of a square matrix

$$A^s = A A \dots A \quad (A \text{ repeated } s \text{ times})$$

$$A^s A^r = A^{s+r} = A^r A^s \text{ (commutative product)}$$

Example A.7 Multiplication

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad B = \begin{bmatrix} 1.1 & 2 \\ 4 & 5 \\ 1 & 2 \end{bmatrix}$$

1. Matrix by scalar

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

$$= \begin{bmatrix} 4 & 8 & 12 \\ 16 & 20 & 24 \end{bmatrix}$$

$$\gg \mathbf{C} = 4^*[\mathbf{1}, \mathbf{2}, \mathbf{3}; \mathbf{4}, \mathbf{5}, \mathbf{6}]$$

Example A.7 Multiplication—cont'd

2. Matrix by matrix

$$\begin{aligned}
 C = AB &= \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \begin{bmatrix} 1.1 & 2 \\ 4 & 5 \\ 1 & 2 \end{bmatrix} \\
 &= \begin{bmatrix} 1 \times 1.1 + 2 \times 4 + 3 \times 1 & 1 \times 2 + 2 \times 5 + 3 \times 2 \\ 4 \times 1.1 + 5 \times 4 + 6 \times 1 & 4 \times 2 + 5 \times 5 + 6 \times 2 \end{bmatrix} \\
 &= \begin{bmatrix} 12.1 & 18 \\ 30.4 & 45 \end{bmatrix} \\
 &\gg \mathbf{C} = \mathbf{A}^* \mathbf{B}
 \end{aligned}$$

3. Vector-matrix multiplication

$$\begin{aligned}
 C = \mathbf{x}^T B &= [1 \ 2 \ 3] \begin{bmatrix} 1.1 & 2 \\ 4 & 5 \\ 1 & 2 \end{bmatrix} \\
 &= [1 \times 1.1 + 2 \times 4 + 3 \times 1 \ 1 \times 2 + 2 \times 5 + 3 \times 2] \\
 &= [12.1 \ 18] \\
 D = A \mathbf{y} &= \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \begin{bmatrix} 1.1 \\ 4 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} 1 \times 1.1 + 2 \times 4 + 3 \times 1 \\ 4 \times 1.1 + 5 \times 4 + 6 \times 1 \end{bmatrix} = \begin{bmatrix} 12.1 \\ 30.4 \end{bmatrix} \\
 &\gg \mathbf{C} = [\mathbf{1}, \mathbf{2}, \mathbf{3}]^* \mathbf{B}; \\
 &\gg \mathbf{D} = \mathbf{A}^* [\mathbf{1.1}; \mathbf{4}; \mathbf{1}];
 \end{aligned}$$

4. Vector-vector

$$\begin{aligned}
 \mathbf{z} = \mathbf{x}^T \mathbf{y} &= [1 \ 2 \ 3] \begin{bmatrix} 1.1 \\ 4 \\ 1 \end{bmatrix} \\
 &= 1 \times 1.1 + 2 \times 4 + 3 \times 1 = 12.1
 \end{aligned}$$

Example A.7 Multiplication—cont'd

$$\begin{aligned}
 D = \mathbf{y}\mathbf{x}^T &= \begin{bmatrix} 1.1 \\ 4 \\ 1 \end{bmatrix} [123] \\
 &= \begin{bmatrix} 1.1 \times 1 & 1.1 \times 2 & 1.1 \times 3 \\ 4 \times 1 & 4 \times 2 & 4 \times 3 \\ 1 \times 1 & 1 \times 2 & 1 \times 3 \end{bmatrix} \\
 &= \begin{bmatrix} 1.1 & 2.2 & 3.3 \\ 4 & 8 & 12 \\ 1 & 2 & 3 \end{bmatrix}
 \end{aligned}$$

$$\begin{aligned}
 \gg \mathbf{z} &= [1, 2, 3]^* [1.1; 4; 1] \\
 \gg \mathbf{D} &= [1.1; 4; 1]^* [1, 2, 3]
 \end{aligned}$$

5. Positive integral power of a square matrix

$$S = \begin{bmatrix} 1 & 2 \\ 0 & 4 \end{bmatrix} \quad S^3 = \begin{bmatrix} 1 & 2 \\ 0 & 4 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & 4 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & 4 \end{bmatrix} = \begin{bmatrix} 1 & 42 \\ 0 & 64 \end{bmatrix}$$

$$\gg S^3$$

Diagonal of a matrix

The diagonal of a square matrix are the terms a_{ii} , $i = 1, 2, \dots, n$.

Diagonal matrix

A matrix whose off-diagonal entries are all equal to zero.

$$A = \text{diag}\{a_{11}, a_{22}, \dots, a_{nn}\}$$

$$= \begin{bmatrix} a_{11} & 0 & \cdots & 0 \\ 0 & a_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_{nn} \end{bmatrix}$$

Example A.8 Diagonal matrix

$$A = \text{diag}\{1, 5, 7\} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 7 \end{bmatrix}$$

$\gg \mathbf{A} = \text{diag}([1, 5, 7])$

Identity or unity matrix

A diagonal matrix with all diagonal entries equal to unity.

$$I = \text{diag}\{1, 1, \dots, 1\} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$

We denote an $n \times n$ identity matrix by I_n . The identity matrix is a multiplicative identity because any $m \times n$ matrix A satisfies $AI_m = I_n A = A$. By definition, we have $A^0 = I_n$.

Example A.9 Identity matrix

$$I_3 = \text{diag}\{1, 1, 1\} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$\gg \text{eye}(3)$

Zero matrix

A matrix with all entries equal to zero.

$$C = \mathbf{0}_{m \times n} \Leftrightarrow [c_{ij}] = [0], \quad i = 1, 2, \dots, m \\ j = 1, 2, \dots, n$$

For any $m \times n$ matrix A , the zero matrix has the properties

$$\begin{aligned} A \mathbf{0}_{n \times l} &= \mathbf{0}_{n \times l} \\ \mathbf{0}_{l \times m} A &= \mathbf{0}_{l \times n} \\ A \pm \mathbf{0}_{m \times n} &= A \end{aligned}$$

Example A.10 Zero matrix

$$\mathbf{0}_{2 \times 3} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

`>> zeros(2,3)`

A.4 Determinant of a matrix

The determinant of a square matrix is a scalar computed using its entries. For a 1×1 matrix, the determinant is simply the matrix itself. For a 2×2 matrix, the determinant is

$$\det(A) = |A| = a_{11}a_{22} - a_{12}a_{21}$$

For higher-order matrices, the following definitions are needed to define the determinant.

Minor

The ij th minor of an $n \times n$ matrix is the determinant of the $(n-1) \times (n-1)$ matrix obtained by removing the i th row and the j th column and is denoted M_{ij} .

Cofactor of a matrix

The ij th cofactor of an $n \times n$ matrix is a signed minor given by

$$C_{ij} = (-1)^{i+j} M_{ij}$$

The sign of the ij th cofactor can be obtained from the ij th entry of the matrix

$$\begin{bmatrix} + & - & + & \cdots \\ - & + & - & \cdots \\ + & - & + & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

Determinant

$$\det(A) = |A| = \sum_{i=1}^n a_{is} C_{is} = \sum_{j=1}^n a_{sj} C_{sj}$$

that is, the determinant can be obtained by expansion along any row or column.

Singular matrix $\det(A) = 0$

Nonsingular matrix $\det(A) \neq 0$

Properties of determinants

For an $n \times n$ matrix A ,

$$\begin{aligned}\det(A) &= \det(A^T) \\ \det(\alpha A) &= \alpha^n \det(A) \\ \det(AB) &= \det(A)\det(B)\end{aligned}$$

Example A.11 Determinant of a matrix

$$\begin{aligned}A &= \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & -1 \\ 1 & -5 & 0 \end{bmatrix} \\ |A| &= 3 \times [4 \times (-5) - 5 \times 1] - (-1) \times [1 \times (-5) - 2 \times 1] + 0 \\ &= 3 \times (-25) + (-7) = -82 \\ &\gg \det(A)\end{aligned}$$

Adjoint matrix

The transpose of the matrix of cofactors

$$\text{adj}(A) = [C_{ij}]^T$$

A.5 Inverse of a matrix

The inverse of a square matrix is a matrix satisfying

$$AA^{-1} = A^{-1}A = I_n$$

The inverse of the matrix is given by

$$A^{-1} = \frac{\text{adj}(A)}{\det(A)}$$

Example A.12 Inverse matrix

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 5 \\ 0 & 6 & 7 \end{bmatrix} \quad A^{-1} = \frac{\text{adj}(A)}{\det(A)} = \frac{1}{6} \begin{bmatrix} (28 - 30) & -(14 - 18) & (10 - 12) \\ -(14 - 0) & (7 - 0) & -(5 - 6) \\ (12 - 0) & -(6 - 0) & (4 - 4) \end{bmatrix} / 6$$

$$= \begin{bmatrix} -0.333 & 0.667 & -0.333 \\ -2.333 & 1.167 & 0.167 \\ 2 & -1 & 0 \end{bmatrix}$$

Use the command

```
>> inv(A)
>> A/B
```

to calculate $A^{-1}B$ and the command

```
>> A/B
```

to calculate AB^{-1} .

Combinations of operations

$$(ABC)^T = C^T B^T A^T$$

$$(ABC)^{-1} = C^{-1} B^{-1} A^{-1}$$

$$(A^T)^{-1} = (A^{-1})^T = A^{-T}$$

Orthogonal matrix:

A matrix whose inverse is equal to its transpose

$$A^{-1} = A^T$$

that is,

$$A^T A = A A^T = I_n$$

Using the properties of a determinant of a square matrix,

$$\det(I_n) = \det(A)\det(A^T) = \det(A)^2 = 1$$

$$\det(I_n) = \det(A)\det(A^T) = \det(A)^2 = 1$$

that is, $\det(A) = \pm 1$ for an orthogonal matrix.

Example A.13 Orthogonal matrix

The coordinate rotation matrix for a yaw angle (rotation about the z-axis) α is the orthogonal matrix.

$$R(\alpha) = \begin{bmatrix} \cos(\alpha) & -\sin(\alpha) & 0 \\ \sin(\alpha) & \cos(\alpha) & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

with $\det(R) = \cos^2(\alpha) + \sin^2(\alpha) = 1$.

Unitary matrix

A matrix whose inverse is equal to its complex conjugate transpose

$$A^{-1} = A^*$$

A.6 Trace of a matrix

The sum of the diagonal elements of a square matrix

$$\text{tr}(A) = \sum_{i=1}^n a_{ii}$$

The trace satisfies the following properties:

$$\begin{aligned} \text{tr}(A^T) &= \text{tr}(A) \\ \text{tr}(AB) &= \text{tr}(BA) \\ \text{tr}(A + B) &= \text{tr}(A) + \text{tr}(B) \end{aligned}$$

Example A.14 Trace of a matrix

Find the trace of the matrix $R(\alpha)$ shown in Example A.13.

$$\text{tr}(R) = \cos(\alpha) + \cos(\alpha) + 1 = 1 + 2 \cos(\alpha)$$

For $\alpha = \pi$, $\cos(\alpha) = -1$ and $\text{tr}(R) = -1$.

```
>> trace(R)
-1
```

A.7 Rank of a matrix

Linearly independent vectors

A set of vectors $\{\mathbf{x}_i, i = 1, 2, \dots, n\}$ is linearly independent if

$$\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \cdots + \alpha_n \mathbf{x}_n = 0 \Leftrightarrow \alpha_i = 0, i = 1, 2, \dots, n$$

Otherwise, the set is said to be **linearly dependent**.

Example A.15 Linear independence

Consider the following row vectors: $\mathbf{a}^T = [3 \ 4 \ 0]$, $\mathbf{b}^T = [1 \ 0 \ 0]$, and $\mathbf{c}^T = [0 \ 1 \ 0]$.

The set $\{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$ is linearly dependent because $\mathbf{a} = 3\mathbf{b} + 4\mathbf{c}$. But the sets $\{\mathbf{a}, \mathbf{b}\}$, $\{\mathbf{b}, \mathbf{c}\}$, and $\{\mathbf{a}, \mathbf{c}\}$ are linearly independent.

Column rank

Number of linearly independent columns.

Row rank

Number of linearly independent rows.

The rank of a matrix is equal to its row rank, which is equal to its column rank.

For an $m \times n$ (rectangular) matrix A , the rank of the matrix is

$$r(A) \leq \min\{n, m\}$$

If equality holds, the matrix is said to be **full rank**. A full rank square matrix is nonsingular.

Example A.16 Rank of a matrix

The matrix

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

has the row vectors considered in Example A.15. Hence, the matrix has two linearly independent row vectors (i.e., row rank 2). The first two columns of the matrix are also

Example A.16 Rank of a matrix—cont'd

linearly independent (i.e., it has column rank 2). The largest square matrix with nonzero determinant is the 2×2 matrix:

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

Clearly, the matrix has rank 2.

A.8 Eigenvalues and eigenvectors

The eigenvector of a matrix A are vectors that are mapped to themselves when multiplied by the matrix A :

$$\begin{aligned} A\mathbf{v} &= \lambda\mathbf{v} \\ [\lambda I_n - A]\mathbf{v} &= 0 \end{aligned}$$

The scale factor on the right-hand side of the equation is known as the eigenvalue. For a nonzero solution \mathbf{v} to the preceding equation to exist, the premultiplying matrix must be rank deficient—that is, λ must be an eigenvalue of the matrix A .

The eigenvector is defined by a direction or by a specific relationship between its entries. Multiplication by a scalar changes the length but not the direction of the vector.

The eigenvalues of an $n \times n$ matrix are the n roots of the characteristic equation:

$$\det[\lambda I_n - A] = 0$$

Distinct eigenvalues:

$$\lambda_j \neq \lambda_i, i \neq j, i, j, = 1, 2, \dots, n$$

Repeated eigenvalues:

$$\lambda_i \neq \lambda_j, \text{ for some } i \neq j$$

Multiplicity of the eigenvalue:

The number of repetitions of the repeated eigenvalue (also known as the **algebraic multiplicity**).

Spectrum of matrix A :

The set of eigenvalues $\{\lambda_i, i = 1, 2, \dots, n\}$.

Spectral radius of a matrix:

Maximum absolute value over all the eigenvalues of the matrix.

Trace in terms of eigenvalues:

$$\text{tr}(A) = \sum_{i=1}^n \lambda_i$$

Upper triangular matrix

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ 0 & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & a_{nn} \end{bmatrix}$$

Lower triangular matrix

$$A = \begin{bmatrix} a_{11} & 0 & \cdots & 0 \\ a_{21} & a_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix}$$

For lower triangular, upper triangular, and diagonal matrices,

$$\{\lambda_i, i = 1, 2, \dots, n\} = \{a_{ii}, i = 1, 2, \dots, n\}$$

Example A.17 Eigenvalues and eigenvectors

Find the eigenvalues and eigenvectors of the matrix

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

$$\lambda I_3 - A = \begin{bmatrix} \lambda - 3 & -4 & 0 \\ -1 & \lambda & 0 \\ 0 & -1 & \lambda \end{bmatrix}$$

$$\det[\lambda I_3 - A] = [(\lambda - 3)\lambda - 4]\lambda = [\lambda^2 - 3\lambda - 4]\lambda = (\lambda - 4)(\lambda + 1)\lambda$$

$$\lambda_1 = 4 \quad \lambda_2 = -1 \quad \lambda_3 = 0$$

$$AV = VA$$

$$\begin{bmatrix} 3 & 4 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} v_{11} & v_{12} & v_{13} \\ v_{21} & v_{22} & v_{23} \\ v_{31} & v_{32} & v_{33} \end{bmatrix} = \begin{bmatrix} v_{11} & v_{12} & v_{13} \\ v_{21} & v_{22} & v_{23} \\ v_{31} & v_{32} & v_{33} \end{bmatrix} \begin{bmatrix} 4 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Example A.17 Eigenvalues and eigenvectors—cont'd

1. $\lambda_3 = 0$
2. $v_{13} = v_{23} = 0$ and v_{33} free. Let $v_{23} = 1$.
3. $\lambda_2 = -1$
4. $v_{12} = -v_{22}$
5. $v_{22} = -v_{32}$. Let $v_{12} = 1$.
6. $\lambda_1 = 4$
7. $v_{11} = 4 v_{21}$
8. $v_{21} = 4 v_{31}$. Let $v_{31} = 1$.

Hence, the modal matrix of eigenvectors is

$$V = \begin{bmatrix} 16 & 1 & 0 \\ 4 & -1 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$

The lengths or 2-norms of the eigenvectors are

$$\begin{aligned}\|\mathbf{v}_1\| &= [16^2 + 4^2 + 1^2]^{1/2} \\ \|\mathbf{v}_2\| &= [1^2 + 1^2 + 1^2]^{1/2} \\ \|\mathbf{v}_3\| &= [0 + 0 + 1^2]^{1/2}\end{aligned}$$

The three eigenvectors can be normalized using the vector norms to obtain the matrix

$$\begin{aligned}V &= \begin{bmatrix} 0.9684 & 0.5774 & 0 \\ 0.2421 & -0.5774 & 0 \\ 0.0605 & 0.5774 & 1 \end{bmatrix} \\ \gg \mathbf{A} &= [3, 4, 0; 1, 0, 0; 0, 1, 0] \\ \gg [\mathbf{V}, \mathbf{L}] &= \mathbf{eig}(\mathbf{A})\end{aligned}$$

$$\begin{aligned}V = \\ 0 &\quad 0.5774 & -0.9684 \\ 0 &\quad -0.5774 & -0.2421 \\ 1.0000 &\quad 0.5774 & -0.0605\end{aligned}$$

$$\mathbf{L} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 4 \end{pmatrix}$$

The trace of the preceding matrix is

$$\text{tr}(A) = 3 + 0 + 0 = 3 = 0 + (-1) + 4 = \lambda_1 + \lambda_2 + \lambda_3$$

Normal matrix

Multiplication by its (conjugate) transpose is commutative.

$$A^T A = AA^T \quad (A^* A = AA^*)$$

This includes symmetric (Hermitian) matrices as a special case.

The matrix of eigenvectors of a normal matrix can be selected as an orthogonal (unitary) matrix:

$$A = V \Lambda V^T \quad (A = V \Lambda V^*)$$

A.9 Partitioned matrix

A matrix partitioned into smaller submatrices.

$$\left[\begin{array}{c|c|c} A_{11} & A_{12} & \cdots \\ \hline \cdots & \cdots & \cdots \\ A_{21} & A_{22} & \cdots \\ \hline \cdots & \cdots & \cdots \\ \vdots & \vdots & \ddots \end{array} \right]$$

Transpose of a partitioned matrix:

$$\left[\begin{array}{c|c|c} A_{11}^T & A_{21}^T & \cdots \\ \hline \cdots & \cdots & \cdots \\ A_{12}^T & A_{22}^T & \cdots \\ \hline \cdots & \cdots & \cdots \\ \vdots & \vdots & \ddots \end{array} \right]$$

Sum/difference of partitioned matrices:

$$C = A \pm B \Leftrightarrow C_{ij} = A_{ij} \pm B_{ij}$$

Product of partitioned matrices:

Apply the rules of matrix multiplication with the products of matrix entries replaced by the noncommutative products of submatrices.

$$C = AB \Leftrightarrow C_{ij} = \sum_{k=1}^n A_{ik}B_{kj}, \quad i = 1, 2, \dots, r \\ i = 1, 2, \dots, s$$

Determinant of a partitioned matrix:

$$\begin{vmatrix} A_1 & | & A_2 \\ \hline \hline A_3 & | & A_4 \end{vmatrix} = \begin{cases} |A_1||A_4 - A_3A_1^{-1}A_2|, & A_1^{-1} \text{ exists} \\ |A_4||A_1 - A_2A_4^{-1}A_3|, & A_4^{-1} \text{ exists} \end{cases}$$

Inverse of a partitioned matrix:

$$\begin{bmatrix} A_1 & | & A_2 \\ \hline \hline A_3 & | & A_4 \end{bmatrix}^{-1} = \begin{bmatrix} (A_1 - A_2A_4^{-1}A_3)^{-1} & | & -A_1^{-1}A_2(A_4 - A_3A_1^{-1}A_2)^{-1} \\ \hline \hline -A_4^{-1}A_3(A_1 - A_2A_4^{-1}A_3)^{-1} & | & (A_4 - A_3A_1^{-1}A_2)^{-1} \end{bmatrix}$$

Example A.18 Partitioned matrices

$$A = \begin{bmatrix} 1 & 2 & | & 5 \\ 3 & 4 & | & 6 \\ \hline \hline 7 & 8 & | & 9 \end{bmatrix} \quad B = \begin{bmatrix} -3 & 2 & | & 5 \\ 3 & 1 & | & 7 \\ \hline \hline -4 & 0 & | & 2 \end{bmatrix}$$

$$A + B = \begin{bmatrix} 1-3 & 2+2 & | & 5+5 \\ 3+3 & 4+1 & | & 6+7 \\ \hline \hline 7-4 & 8+0 & | & 9+2 \end{bmatrix} = \begin{bmatrix} -2 & 4 & | & 10 \\ 6 & 5 & | & 13 \\ \hline \hline 3 & 8 & | & 11 \end{bmatrix}$$

$$AB = \begin{bmatrix} \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} -3 & 2 \end{bmatrix} + \begin{bmatrix} 5 \end{bmatrix} \begin{bmatrix} -4 & 0 \end{bmatrix} & \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 5 \end{bmatrix} + \begin{bmatrix} 5 \end{bmatrix} \begin{bmatrix} 2 \end{bmatrix} \\ \begin{bmatrix} 7 & 8 \end{bmatrix} \begin{bmatrix} -3 & 2 \end{bmatrix} + 9 \begin{bmatrix} -4 & 0 \end{bmatrix} & \begin{bmatrix} 7 & 8 \end{bmatrix} \begin{bmatrix} 5 \end{bmatrix} + 9 \times 2 \end{bmatrix}$$

$$= \begin{bmatrix} -17 & 4 & | & 29 \\ -21 & 10 & | & 55 \\ \hline \hline -33 & 22 & | & 109 \end{bmatrix}$$

```

>> A1 = [1, 2; 3, 4];
>> a2 = [5; 6];
>> a3 = [7, 8];
>> a4 = 9;
>> A = [A1, a2; a3, a4];
>> B = [[-3, 2; 3, 1], [5; 7]; [-4, 0], 2];

```

Example A.18 Partitioned matrices—cont'd**>> A+B**

$$\begin{matrix} -2 & 4 & 10 \\ 6 & 5 & 13 \\ 3 & 8 & 11 \end{matrix}$$

>> A*B

$$\begin{matrix} -17 & 4 & 29 \\ -21 & 10 & 55 \\ -33 & 22 & 109 \end{matrix}$$

Matrix Inversion Lemma

The following identity can be used in either direction to simplify matrix expressions:

$$[A_1 + A_2 A_4^{-1} A_3]^{-1} = A_1^{-1} - A_1^{-1} A_2 [A_4 + A_3 A_1^{-1} A_2]^{-1} A_3 A_1^{-1}$$

A.10 Norm of a vector

The norm is a measure of size or length of a vector. It satisfies the following axioms, which apply to the familiar concept of length in the plane.

Norm axioms

1. $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = 0$
2. $\|\mathbf{x}\| > 0$ for $\mathbf{x} \neq 0$
3. $\|\alpha\mathbf{x}\| = |\alpha| \|\mathbf{x}\|$
4. $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ (triangle inequality)

 l_p norms

$$l_\infty \text{norm: } \|\mathbf{x}\|_\infty = \max_i |x_i|$$

$$l_2 \text{norm: } \|\mathbf{x}\|_2^2 = \sum_{i=1}^n |x_i|^2$$

$$l_1 \text{norm: } \|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$$

Equivalent norms

Norms that satisfy the inequality

$$k_1 \|\mathbf{x}\|_i \leq \|\mathbf{x}\|_j \leq k_2 \|\mathbf{x}\|_i$$

with finite constants k_1 and k_2 . All norms for $n \times 1$ real vectors are equivalent. All equivalent norms are infinite if and only if any one of them is infinite.

Example A.19 Vector norms

```

 $\mathbf{a}^T = [1, 2, -3]$ 
 $\|\mathbf{a}\|_1 = |1| + |2| + |-3| = 6$ 
 $\|\mathbf{a}\|_2 = \sqrt{1^2 + 2^2 + (-3)^2} = 3.7417$ 
 $\|\mathbf{a}\|_\infty = \max\{|1|, |2|, |-3|\} = 3$ 

>> a = [1; 2; -3]
>> norm(a, 2)      % 2-norm (square root of sum of squares)
3.7417
>> norm(a, 1)      % 1-norm (sum of absolute values)
6
>> norm(a, inf)    % infinity-norm (max element)
3

```

A.11 Matrix norms

Satisfy the norm axioms.

Induced matrix norms

Norms that are induced from vector norms using the definition

$$\|A\|_i = \max_{\mathbf{x}} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\|\mathbf{x}\|=1} \|A\mathbf{x}\|$$

where $\|\cdot\|$ is any vector norm.

Submultiplicative property

$$\begin{aligned} \|A\mathbf{x}\| &\leq \|A\| \|\mathbf{x}\| \\ \|AB\| &\leq \|A\| \|B\| \end{aligned}$$

All induced norms are submultiplicative, but only some noninduced norms are.

\mathbf{l}_1 Norm $\|A\|_1 = \max_j \sum_{i=1}^m |a_{ij}|$ (maximum absolute column sum)

\mathbf{l}_∞ Norm $\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|$ (maximum absolute row sum)

\mathbf{l}_2 Norm $\|A\|_2 = \max_i \lambda_i^{1/2}(A^T A)$ (maximum singular value = maximum eigenvalue of $A^T A$)

Frobenius norm

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\text{tr}\{A^T A\}}$$

Other matrix norms

$$\begin{aligned}\|A\| &= \max_{i,j} |a_{ij}| \\ \|A\|_F &= \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}\end{aligned}$$

The Frobenius norm is **not** an induced norm.

Example A.20 Norm of a matrix

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

$$\|A\|_1 = \max\{|1| + |3|, |2| + |-4|\} = 6$$

$$\|A\|_2 = \lambda_{\max}^{1/2} \left\{ \begin{bmatrix} 1 & 3 \\ 2 & -4 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & -4 \end{bmatrix} \right\} = \lambda_{\max}^{1/2} \left\{ \begin{bmatrix} 10 & -10 \\ -10 & 20 \end{bmatrix} \right\} = 5.1167$$

$$\|A\|_\infty = \max\{|1| + |2|, |3| + |-4|\} = 7$$

$$\|A\|_F = \sqrt{|1|^2 + |2|^2 + |3|^2 + |-4|^2} = 5.4772$$

```
>> norm(A, 1) % 1 induced norm (maximum of column sums)
```

Example A.20 Norm of a matrix—cont'd

```

>> norm(A, 2)      % 2 induced norm (maximum singular value)
5.1167
>> norm(A, inf)    % infinity induced norm (maximum of row sums)
7
>> norm(A, 'fro')   % 2 (square root of sum of squares)
5.4772

```

A.12 Quadratic forms

A quadratic form is a function of the form

$$V(\mathbf{x}) = \mathbf{x}^T P \mathbf{x} = \sum_{i=1}^n \sum_{j=1}^n p_{ij} x_i x_j$$

where \mathbf{x} is an $n \times 1$ real vector and P is an $n \times n$ matrix. The matrix P can be assumed to be symmetric without loss of generality. To show this, assume that P is not symmetric and rewrite the quadratic form in terms of the symmetric component and the skew-symmetric component of P as follows:

$$\begin{aligned} V(\mathbf{x}) &= \mathbf{x}^T \left(\frac{P + P^T}{2} \right) \mathbf{x} + \mathbf{x}^T \left(\frac{P - P^T}{2} \right) \mathbf{x} \\ &= \mathbf{x}^T \left(\frac{P + P^T}{2} \right) \mathbf{x} + \frac{1}{2} (\mathbf{x}^T P \mathbf{x} - (P \mathbf{x})^T \mathbf{x}) \end{aligned}$$

Interchanging the row and column in the last term gives

$$\begin{aligned} V(\mathbf{x}) &= \mathbf{x}^T \left(\frac{P + P^T}{2} \right) \mathbf{x} + \frac{1}{2} (\mathbf{x}^T P \mathbf{x} - \mathbf{x}^T P \mathbf{x}) \\ &= \mathbf{x}^T \left(\frac{P + P^T}{2} \right) \mathbf{x} \end{aligned}$$

Thus, if P is not symmetric, we can replace it with its symmetric component without changing the quadratic form.

The sign of a quadratic form for nonzero vectors \mathbf{x} can be invariant depending on the matrix P . In particular, the eigenvalues of the matrix P determine the sign of the quadratic form. To see this, we examine the eigenvalues-eigenvector decomposition of the matrix P in the quadratic form

$$V(\mathbf{x}) = \mathbf{x}^T P \mathbf{x}$$

We assume, without loss of generality, that P is symmetric. Hence, its eigenvalues are real and positive, and its modal matrix of eigenvectors is orthogonal. The matrix can be written as

$$\begin{aligned} P &= V_p \Lambda V_p^T \\ \Lambda &= \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\} \end{aligned}$$

Using the eigenvalues decomposition of the matrix, we have

$$\begin{aligned} V(\mathbf{x}) &= \mathbf{x}^T V_p \Lambda V_p^T \mathbf{x} \\ &= \mathbf{y}^T \Lambda \mathbf{y} \\ &= \sum_{i=1}^n \lambda_i y_i^2 > 0 \\ \mathbf{y} &= [y_1 \ y_2 \cdots y_n] \end{aligned}$$

Because the modal matrix V_p is invertible, there is a unique \mathbf{y} vector associated with each \mathbf{x} vector. The expression for the quadratic form in terms of the eigenvalues allows us to characterize it and the associated matrix as follows.

Positive definite

A quadratic form is positive definite if

$$V(\mathbf{x}) = \mathbf{x}^T P \mathbf{x} > 0, \quad \mathbf{x} \neq \mathbf{0}$$

This is true if the eigenvalues of P are all positive, in which case, we say that P is a positive definite matrix, and we denote this by $P > 0$.

Negative definite

A quadratic form is negative definite if

$$V(\mathbf{x}) = \mathbf{x}^T P \mathbf{x} < 0, \quad \mathbf{x} \neq \mathbf{0}$$

This is true if the eigenvalues of P are all negative, in which case we say that P is a negative definite matrix, and we denote this by $P < 0$.

Positive semidefinite

A quadratic form is positive semidefinite if

$$V(\mathbf{x}) = \mathbf{x}^T P \mathbf{x} \geq 0, \quad \mathbf{x} \neq \mathbf{0}$$

This is true if the eigenvalues of P are all positive or zero, in which case we say that P is a positive semidefinite matrix, and we denote this by $P \geq 0$. Note that in this case, if an eigenvalue λ_i is zero, then the nonzero vector \mathbf{y} with its i th entry equal to 1 and all other entries zero gives a zero value for V . Thus, there is a nonzero vector $\mathbf{x} = V_p \mathbf{y}$ for which V is zero.

Negative semidefinite

A quadratic form is negative semidefinite if

$$V(\mathbf{x}) = \mathbf{x}^T P \mathbf{x} \leq 0, \quad \mathbf{x} \neq \mathbf{0}$$

This is true if the eigenvalues of P are all negative or zero, in which case we say that P is a negative semidefinite matrix, and we denote this by $P \leq 0$. In this case, if an eigenvalue λ_i is zero, then V is zero for the nonzero $\mathbf{x} = V_p \mathbf{y}$, where \mathbf{y} is a vector with its i th entry equal to 1 and all other entries are zero.

Indefinite

If the matrix Q has some positive and some negative eigenvalues, then the sign of the corresponding quadratic form depends on the vector \mathbf{x} , and the matrix is called indefinite.

A.13 Singular value decomposition and pseudoinverses

Any $n \times m$ real matrix A can be decomposed into the product

$$A = U \Sigma V^T$$

where U is $n \times n$, V is $m \times m$, and Σ is $n \times m$. For a matrix of rank r , the matrices in the decomposition satisfy

$$\begin{aligned} U^{-1} &= U^T \\ V^{-1} &= V^T \\ \Sigma &= \begin{bmatrix} \Sigma_r & \mathbf{0}_r \times (m-r) \\ \mathbf{0}_{(n-r) \times r} & \mathbf{0}_{(n-r) \times (m-r)} \end{bmatrix} \\ \Sigma_r &= \text{diag}\{\sigma_1, \sigma_2, \dots, \sigma_r\} \end{aligned}$$

The terms on the diagonal are real and positive and are known as the singular values of the matrix. The columns of U are called the left singular vectors and the columns of V are

called the right singular vectors. We assume without loss of generality that the diagonal terms are arranged such that

$$\sigma_{\max} = \sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r = \sigma_{\min}$$

If the matrix A is full rank, then the matrix Σ has one of the following forms:

$$\Sigma \begin{cases} \begin{bmatrix} \Sigma_n & 0_{n \times (m-n)} \end{bmatrix}, & n < m \\ \Sigma_m, & n = m \\ \begin{bmatrix} \Sigma_m \\ 0_{(n-m) \times m} \end{bmatrix}, & n > m \end{cases}$$

The singular value decomposition is obtained using the MATLAB command

```
>> [u, s, v] = svd(A)% A = u*s*v'
```

Example A.21

Find the singular value decomposition of the matrix $A = \begin{vmatrix} 3 & 4 & 5 \\ 6 & 8 & 10 \end{vmatrix}$

Solution

```
>> A = [1, 2, 3; 3, 4, 5; 6, 8, 10; 1, 2, 3];
>> [u, s, v] = svd(A)
```

```
u =
-0.2216 -0.6715 0.6959 -0.1256
-0.4247 0.1402 0.1588 0.8802
-0.8494 0.2803 -0.0794 -0.4401
-0.2216 -0.6715 -0.6959 0.1256
```

```
s =
16.6473 0 0
0 0.9306 0
0 0 0.0000
0 0 0
```

```
v =
-0.4093 0.8160 0.4082
-0.5635 0.1259 -0.8165
-0.7176 -0.5642 0.4082
```

Example A.21—cont'd**Eigenvalues and Singular Values**

For any $n \times m$ real matrix A , we have

$$AA^T = U\Sigma\Sigma^T U^T = U\text{diag}\left\{\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2, \mathbf{0}_{1 \times (n-r)}\right\}U^T$$

$$A^T A = V\Sigma^T\Sigma V^T = V\text{diag}\left\{\sigma_1^2, \sigma_2^2, \dots, \sigma_r^2, \mathbf{0}_{1 \times (m-r)}\right\}V^T$$

Thus, the nonzero eigenvalues of either of the above products are equal to the singular values of the matrix A .

The determinant of an $n \times n$ matrix A is

$$|\det[A]| = \left| \prod_{i=1}^n \lambda_i(A) \right| = \prod_{i=1}^n \sigma_i(A)$$

The eigenvalues of a normal $n \times n$ matrix A and its singular values are related by

$$|\lambda_i(A)| = \sigma_i, i = 1, 2, \dots, n$$

Singular Value Inequalities

For any invertible square matrix,

- $\sigma_{\max}(A) = 1/\sigma_{\min}(A^{-1})$
- $\sigma_{\min}(A) = 1/\sigma_{\max}(A^{-1})$

For any two compatible matrices A and B ,

- $\sigma_{\max}(A+B) \leq \sigma_{\max}(A) + \sigma_{\max}(B)$
- $\sigma_{\max}(AB) \leq \sigma_{\max}(A) \sigma_{\max}(B)$

Pseudoinverse

Using the singular value decomposition, we can define a pseudoinverse for any $n \times m$ real matrix A as

$$A^\# = V\Sigma^\#U^T$$

with the pseudoinverse of the matrix of singular values given by

$$\Sigma = \begin{bmatrix} \Sigma_r^{-1} & \mathbf{0}_{r \times (n-r)} \\ \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (n-r)} \end{bmatrix}$$

$$\Sigma_r^{-1} = \text{diag}\{1/\sigma_1, 1/\sigma_2, \dots, 1/\sigma_r\}$$

Example A.21—cont'd

Note that the inversion is similar to the inversion of the product of three matrices, with the inverse of the middle matrix being a pseudoinverse. For a full rank matrix, we have

$$\Sigma^\# = \begin{cases} \begin{bmatrix} \Sigma_r^{-1} \\ \mathbf{0}_{(m-n) \times n} \end{bmatrix}, & n < m \\ \Sigma_m, & n = m \\ \begin{bmatrix} \Sigma_m & \mathbf{0}_{m \times (n-m)} \end{bmatrix} & n > m \end{cases}$$

Clearly, the pseudoinverse of a nonsingular square matrix is simply its inverse.

The following MATLAB command computes the pseudoinverse:

`>> pinv(A)`

Example A.22

Find the pseudoinverse of the matrix $A = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 5 \\ 6 & 8 & 10 \\ 1 & 2 & 3 \end{bmatrix}$

`>> A = [1,2,3;3,4,5;6,8,10;1,2,3];`

`>> pinv(A)`

`ans =`

$$\begin{bmatrix} -0.5833 & 0.1333 & 0.2667 & -0.5833 \\ -0.0833 & 0.0333 & 0.0667 & -0.0833 \\ 0.4167 & -0.0667 & -0.1333 & 0.4167 \end{bmatrix}$$

Pseudoinverse of a Full-Rank Matrix

For a full-rank matrix, the pseudoinverse of an m by n matrix reduces to the following:

$$A^\# = \begin{cases} (A^T A)^{-1} A^T, & m > n \\ A^{-1}, & m = n \\ A^T (A A^T)^{-1}, & m < n \end{cases}$$

The first is a **left inverse**, the second is the usual matrix inverse, and the third is a **right inverse** of the matrix. The terms *right inverse* and *left inverse* are due to the products

$$(A^T A)^{-1} A^T A = I_n$$

$$A A^T (A A^T)^{-1} = I_m$$

$$A^\# = V \Sigma^\# U^T$$

Example A.23

>>A = [1, 2; 3, 4; 6, 8; 1, 2]

A =

1 2

3 4

6 8

1 2

>>Apinv = (A'*A)/A' % Left inverse

Apinv =

-1.0000 0.2000 0.4000 -1.0000

0.7500 -0.1000 -0.2000 0.7500

>>Apinv*A

ans =

1.0000 0.0000

-0.0000 1.0000

>>B = A'

B =

1 3 6 1

2 4 8 2

>>Bpinv = B'/(B*B') % Right inverse

Bpinv =

-1.0000 0.7500

0.2000 -0.1000

0.4000 -0.2000

-1.0000 0.7500

>>B*Bpinv

ans =

1.0000 -0.0000

0.0000 1.0000

A.14 Matrix differentiation/integration

The derivative (integral) of a matrix is a matrix whose entries are the derivatives (integrals) of the entries of the matrix.

Example A.24 Matrix differentiation and integration

$$A(t) = \begin{bmatrix} 1 & t & \sin(2t) \\ t & 0 & 4+t \end{bmatrix}$$

$$\int_0^t A(\tau)d\tau = \begin{bmatrix} \int_0^t 1d\tau & \int_0^t \tau d\tau & \int_0^t \sin(2\tau)d\tau \\ \int_0^t \tau d\tau & 0 & \int_0^t (4+\tau)d\tau \end{bmatrix}$$

$$= \begin{bmatrix} t & t^2/2 & \{1 - \cos(t)\}/2 \\ t^2/2 & 0 & 4t + t^2/2 \end{bmatrix}$$

$$\frac{dA(t)}{dt} = \begin{bmatrix} \frac{d1}{dt} & \frac{dt}{dt} & \frac{d \sin(2t)}{dt} \\ \frac{dt}{dt} & 0 & \frac{d(4+t)}{dt} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 2 \cos(2t) \\ 1 & 0 & 1 \end{bmatrix}$$

Derivative of a product $\frac{dAB}{dt} = A \frac{dB}{dt} + \frac{dA}{dt} B$

Derivative of the inverse matrix $\frac{d(A^{-1})}{dt} = -A^{-1} \frac{dA}{dt} A^{-1}$

Gradient vector

The derivative of a scalar function $f(\mathbf{x})$ with respect to the vector \mathbf{x} is known as the gradient vector and is given by the n by 1 vector.

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \left[\frac{\partial f(\mathbf{x})}{\partial x_i} \right]$$

Some authors define the gradient as a row vector.

Example A.24 Matrix differentiation and integration—cont'd**Jacobian matrix**

The derivative of an $n \times 1$ vector function $\mathbf{f}(\mathbf{x})$ with respect to the vector \mathbf{x} is known as the Jacobian matrix and is given by the $n \times n$ matrix.

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \left[\frac{\partial f_i(\mathbf{x})}{\partial x_j} \right]$$

Gradient of inner product

$$\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \left[\frac{\partial \sum_{i=1}^n a_i x_i}{\partial x_i} \right] = [a_i] = \mathbf{a}$$

Gradient matrix of a quadratic form

$$\begin{aligned} \frac{\partial \mathbf{x}^T P \mathbf{x}}{\partial \mathbf{x}} &= \mathbf{x}^T \frac{\partial P \mathbf{x}}{\partial \mathbf{x}} + \frac{\partial (P^T \mathbf{x})^T}{\partial \mathbf{x}} \mathbf{x} \\ &= (P + P^T) \mathbf{x} \end{aligned}$$

Because P can be assumed to be symmetric with no loss of generality, we write

$$\frac{\partial \mathbf{x}^T P \mathbf{x}}{\partial \mathbf{x}} = 2P\mathbf{x}$$

Hessian matrix of a quadratic form

The Hessian or second-derivative matrix is given by

$$\frac{\partial^2 \mathbf{x}^T P \mathbf{x}}{\partial \mathbf{x}^2} = \frac{\partial 2P\mathbf{x}}{\partial \mathbf{x}} = 2 \left[\frac{\partial \mathbf{p}_i^T \mathbf{x}}{\partial x_j} \right] = 2[p_{ij}]$$

where the i th entry of the vector $P\mathbf{x}$ is

$$\mathbf{p}_i^T \mathbf{x}$$

$$P = [p_{ij}] = \begin{bmatrix} \mathbf{p}_1^T \\ \mathbf{p}_2^T \\ \vdots \\ \mathbf{p}_n^T \end{bmatrix}$$

$$\frac{\partial^2 \mathbf{x}^T P \mathbf{x}}{\partial \mathbf{x}^2} = 2P$$

A.15 Kronecker product

The Kronecker product of two matrices A of order $m \times n$ and B of order $p \times q$ is denoted by \otimes and is defined as

$$A \otimes B = \begin{bmatrix} a_{11}B & a_{12}B & \cdots & a_{1n}B \\ a_{21}B & a_{22}B & \cdots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mn}B \end{bmatrix}$$

The resulting matrix is of order $m.p \times n.q$.

Example A.25 Kronecker matrix product

The Kronecker product of the two matrices:

$$\begin{aligned} A &= \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} & B &= \begin{bmatrix} 1.1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \\ A \otimes B &= \begin{bmatrix} 1 \begin{bmatrix} 1.1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} & 2 \begin{bmatrix} 1.1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} & 3 \begin{bmatrix} 1.1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \\ 4 \begin{bmatrix} 1.1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} & 5 \begin{bmatrix} 1.1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} & 6 \begin{bmatrix} 1.1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \end{bmatrix} \\ &= \begin{bmatrix} 1.1 & 2 & 3 & | & 2.2 & 4 & 6 & | & 3.3 & 6 & 9 \\ 4 & 5 & 6 & | & 8 & 10 & 12 & | & 12 & 15 & 18 \\ \hline 4.4 & 8 & 12 & | & 5.5 & 10 & 15 & | & 6.6 & 12 & 18 \\ 16 & 20 & 24 & | & 20 & 25 & 30 & | & 24 & 30 & 36 \end{bmatrix} \end{aligned}$$

```
>> kron(a,b)
ans =
1.1000    2.0000    3.0000    2.2000    4.0000    6.0000    3.3000    6.0000    9.0000
4.0000    5.0000    6.0000    8.0000   10.0000   12.0000   12.0000   15.0000   18.0000
4.4000    8.0000   12.0000   5.5000   10.0000   15.0000   6.6000   12.0000   18.0000
16.0000   20.0000   24.0000   20.0000   25.0000   30.0000   24.0000   30.0000   36.0000
```

Further reading

- Barnett, S., 1984. Matrices in Control Theory. R.E. Krieger, Malabar, FL.
- Brogan, W.L., 1985. Modern Control Theory. Prentice Hall, Englewood Cliffs, NJ.
- Fadeeva, V.N., 1959. Computational Methods of Linear Algebra. Dover, New York.
- Gantmacher, F.R., 1959. The Theory of Matrices. Chelsea, New York.
- Noble, B., Daniel, J.W., 1988. Linear Algebra. Prentice Hall, Englewood Cliffs, NJ.

Index

Note: ‘Page numbers followed by “*f*” indicate figures and “*t*” indicate tables’.

A

- abs** command, 153–157, 602–608
- Absolute stability, 556–560
 - closed-loop system with linear and nonlinear blocks, 556f
 - nonlinear system for saturation nonlinearity, 558f
 - nonlinearity and associated disc, 557f
 - Nyquist plot of furnace and actuator, 557f
- Absolutely summable system, 325–328
- Accuracy, in digital control, 319
- acker** command, 381
- Ackermann’s formula, 395–397, 402
- Actuator
 - control signals to, 3
 - discrete-time first-order, 89–90
 - noise, 80
 - nonlinearity, 590
 - Nyquist plot, 78f, 557f
 - saturation, 66, 580f, 590, 590f, 593f
 - types, 7
- ADC. *See* Analog-to-digital converter (ADC)
- Adjoint matrix, 272–273
- Aircraft turbojet engine, computer control of, 4, 5f
- Algebraic loop, 594–596
- Algebraic Lyapunov equation, 467
- Aliasing, 572, 577
- Allocation, pole, 389–390
- Analog control system design
 - design specifications and effect of gain variation, 147–149

- digital implementation of, 195–216
- root locus, 142–146
- Analog controller, 2
- Analog disturbances
 - in digital system, 80–82
 - digital system with, 81f
- Analog filter, 196–197
 - bilinear transformation, 201–214
 - discrete approximation, 199
 - pole-zero matched digital filter approximation, 200
 - with transfer function, 200
- Analog plant
 - discrete-time approximation, 537
 - stability condition for discretization, 538–539
- Analog system
 - cascade of, 65f
 - fluid level control system, 11f
 - with piecewise constant inputs, 10–12, 10f
 - root locus with PID control, 212f–213f
 - with sampled input, 66f
 - stability with digital control, 537–539
 - step response with digital control, 92f
- Analog-to-digital converter (ADC), 3, 62–63, 62f, 68–75, 119, 397, 570–572. *See also* Digital-to-analog converter (DAC)
- automotive vehicle, 70f
- cascade of, 69f
- digital control system configuration, 62f
- quantization characteristics of, 575f
- rounding, 574–575
- schematic of furnaces, 73f
- truncation, 574–575
- Analytic continuation, 6f
- angle** command, 153–157
- Angle of arrival, root locus method, 143–146
- Angle of departure, on root locus method, 143
- Anthropomorphic manipulator, 259–261
- Antialiasing filters, 572–574, 574f
- Antiwindup techniques, 590–591
- Assignment, pole, 389–390, 429–434, 429f
- Asymptote angle, 321
- Asymptotic stability, 106–107, 320–324
- Asymptotically stable system, 321
- Automatic mode, 593–596
- Autonomous underwater vehicles (AUVs), 312–313

B

- Back-calculation, 591–593
- Backward differencing methods, 197–198
- Backward iteration, 273
- Band-limited signal, 51–52
- Bandwidth
 - closed-loop system, 572–573
 - finite, 51–52
 - infinite, 53
- Bessel filter, 572
- BIBO stability. *See* Bounded-input–bounded-output stability (BIBO stability)

- Bilinear transformation, 201–214
digital filter by applying,
203–204
frequency response of digital,
202f
relationship between analog
filter and associated digital
filter frequencies, 203f
root locus
for analog speed control
system, 211f
of analog system with PID
control, 212f–213f
for PD design, 209f
time step response for digital
PID design, 214f
- bode** command, 49
Bode plots, 221–222, 225f
of analog filter, 204f
of digital filter, 204f
- Bounded sequence, 104, 105f
- Bounded-input–bounded-output
stability (BIBO stability),
106–110, 320,
325–329
- Break-in points, 143
- Breakaway points, 143
- Bumpless transfer, 594
bumpy manual(M)/automatic
(A) transfer, 593f–596f
between manual and automatic
mode, 593–596
- Butterworth filter, 572
- C**
- c2d** command, 353–355
CAD. *See* Computer-aided design
(CAD)
- canon** command, 306–307
- Cascade
of analog systems, 65f
of DAC, 69f
sampler on transfer function
effect of, 65–68
- Case statement, 569–570
- Causal signals, 13
- Cayley–Hamilton theorem,
317–318, 330–332, 395
- Chaotic behavior, 316
- Characteristic equations,
invariance of, 306–307
- Closed-loop drug delivery
system, 3–4, 4f
- Closed-loop eigenvalues, choice
of, 397–402
analog velocity of motor with
deadbeat control, 401f
zero-input state response and
control variable, 399f–401f
- Closed-loop matrix, 403–404, 456
- Closed-loop state equation, 388
- Closed-loop system, 409f, 411f,
412–413, 413f, 572–573
- Closed-loop transfer function,
78–80, 114, 229
single-loop digital control
system, 78f
system with sampling in
feedback path, 79f
- Co-state equation, 450
- Completely controllable system.
See Controllable system
- Completely observable system.
See Observable system
- Complex conjugate eigenvalues
in discrete-time state–space
equations, 292–293
real form for, 284–285
- Computed torque method, 510
- Computer-aided design (CAD),
33–34, 119, 182–184,
316
- Conditional integration, 590–591
- Constituent matrices, 277
properties, 282–283
- Constraints, 490–491
- Continuous-time systems,
234–235
realizations for, 377–378
- Control jitter, 568–569
- Control law, 490
- Controllability, 329–343
MATLAB commands for
controllability testing, 337
matrix, 333–334
rank condition, 332–333
of systems in normal form,
337–338
tests, 334–336
- Controllable canonical
realization, 358–363
controllable form in MATLAB,
363
- parallel realization, 363–368
simulation diagram, 362f
- systems
with input differencing,
360–363
with no input differencing,
358–360
- Controllable form in MATLAB,
363
- Controllable subspaces, 330f
- Controllable system, 330
- Controller form, 358, 530
- Controller structure, 585–588
- Convergence rate, 526–527
- Convolution
property of Laplace transforms,
267
summation, 36–39, 38f, 38t
causal LTI discrete-time
system, 36f
theorem, 39–41, 43f
- Cost
in digital control, 320–321
function, 442–443, 489–490
- D**
- DAC.** *See* Digital-to-analog
converter (DAC)
- Damping ratio, 114, 147
- dare** command, 468
- Data acquisition system,
571–572
- ddamp** command, 114
- Deadbeat controller, 235
- Decision variables, 620
- Degree of freedom (DOF), 4–5
- Delay
computational, 568–569
conversion, 570
intrinsic, 234–235
phase, 573–574, 584–585
sensor, 53
time, 17–18, 63–64, 171–172
transport, 76
- den** command, 49
- den** command, 114

- denp** command, 34
Detectability, 348–350
Diagonal state matrix, state-transition matrix for, 278–283
Differencing methods, 196–199
Differential equation, 12–13, 34–36, 515
Digital control involves systems, 9
Digital control system, 2–3, 568
 closed-loop drug delivery system, 3–4, 4f
 computer control of aircraft turbojet engine, 4, 5f
 design
 digital implementation of analog controller design, 195–216
 direct control design, 229–234
 direct z-domain digital controller design, 216–221
 finite settling time design, 234–247
 frequency response design, 221–229
 Z-domain digital control system design, 184–195
 Z-domain root locus, 182–184
 modeling, 68
 ADC model, 62–63, 68–75
 analog disturbances in digital system, 80–82
 analog system with sampled input, 66f
 cascade of two analog systems, 65f
 closed-loop transfer function, 78–80
 DAC model, 63, 68–75
 MATLAB commands, 86–90
 sampled ramp input, 84–86
 sampled step input, 84
 sampler effect on transfer function of cascade, 65–68
 sensitivity analysis, 93–97
 steady-state error and error constants, 82–86
 systems with transport lag, 76–78
 transfer function of ZOH, 63–64
 robotic manipulator control, 4–5, 6f
 system structure, 3
Digital control system stability, 105f
 conditions, 106–113
 asymptotic, 106–107
 BIBO, 107–110
 internal, 110–113
 determination, 114–117, 126f
 MATLAB, 114
 Routh–Hurwitz criterion, 115–117
 Jury test, 117–122
 Nyquist criterion, 122–137
 stable z-domain pole locations, 105–106
Digital controllers, practical issues of
 choice of sampling period, 572–585
 controller structure, 585–588
 design of hardware and software architecture, 568–572
 proportional–integral–derivative control, 588–598
 sampling period switching, 598–610
Digital filter, 201
Digital implementation of analog controller design, 195–216
 bilinear transformation, 201–214
 differencing methods, 196–199
 empirical digital PID controller tuning, 214–216
 pole-zero matching, 199–201
Digital signal processing chips (DSP chips), 4–5
Digital-to-analog converter (DAC), 3, 63, 68–75, 69f, 119, 570–572. *See also* Analog-to-digital converter (ADC)
 automotive vehicle, 70f
 cascade of, 69f
 model, 63f
 schematic of furnace, 73f
Diophantine equation, 430
Dirac delta, 44–45
Direct control design, 229–234
 ad hoc procedure, 231–234
 step response for, 232f–234f
Direct z-domain digital controller design, 216–221
 PD compensator zero, 217f
 process output with the digital PID controller, 216f
 root locus for digital PI control, 220f
 time step response for digital PID control, 221f
Discrete-time linear system, 322–323
Discrete-time models, 508
Discrete-time nonlinear controller design, 543–548
 saddle point, 540f
 stable focus, 541f
Discrete-time state-transition matrix, 296
Discrete-time state–space equations, 289–293
 complex conjugate eigenvalues, 292–293
 MATLAB commands for, 292
 solution, 293–300
 z-transform solution, 295–300
Discrete-time systems, 13, 65
 analog systems with piecewise constant inputs, 10–12
CAD, 33–34
difference equations, 12–13
frequency response of, 44–50
sampling theorem, 51–55
time response of, 36–41
z-transform, 13–33
Discretization of nonlinear systems, 508–517
dlqr command, 468–470
dlqry command, 468–470
dlyap command, 531

- DMC. *See* Dynamic Matrix Control (DMC)
- DOF. *See* Degree of freedom (DOF)
- Domain of attraction, 534–536
stable node, 540f
- dsort** command, 114
- DSP chips. *See* Digital signal processing chips (DSP chips)
- Dual-rate control, 608–610, 608f–609f
step response described, 610f
- Duality, 370–372
- Duals, 370–372
- Dynamic Matrix Control (DMC), 492–498
closed-loop unit step response, 495f–498f
- E**
- Eigenstructure, 389
- Empirical digital PID controller tuning, 214–216
- Equilibrium of nonlinear discrete-time systems, 518–521
phase portrait for nonlinear system, 520f
- Equilibrium point, 263
- Equilibrium state system, 321
- Error constants, 82–86
- evalfr** command, 153–157
- Extended linearization, 508–509
controller design using, 544f, 546–548
simulation diagram for system, 542f
step response for linear design, 545f
unstable focus, 541f–542f
by input and state redefinition, 511–512
by input redefinition, 509–510
using matching conditions, 514–517
- F**
- Feasible solution, 616
- Feedback control law, 388
- Feedforward action, 408–411
- Filtering derivative action, 588–590
- Final value theorem, 31–33, 230
- Finite impulse response (FIR systems), 108–110
- Finite settling time design, 234–247
block diagram for, 237f
eliminating intersample oscillation, 239–247
control variable for deadbeat control, 239f, 242f, 244f, 247f
ripple-free deadbeat controller, 240–244
sampled and analog step response for deadbeat control, 247f
sampled and analog step response for deadbeat control, 236f
- FIR systems. *See* Finite impulse response (FIR systems)
- First-order approximation, 262
- First-order hold, 63
- Fitzhugh–Nagumo model, 315–316
- Flexibility, in digital control, 319
- Folding frequency, 49
- Forward differencing, 197
- Free final state, 455–465
inertial control system
phase plane trajectory for, 462f, 464f
position trajectory for, 461f, 463f
velocity trajectory for, 461f, 464f
plot of feedback gains vs. time, 462f, 465f, 504f
- Frequency
of oscillations, 190
response design, 221–229
Bode plots, 225f
steps for controller design, 223
- response of discrete-time systems, 44–50
digital system, 49f
- MATLAB commands for, 49–50
- Full-order observer, 414–416, 414f
- Furnace
discrete-time first-order, 89–90
Nyquist plot, 78f, 557f
- G**
- gf** feedback transfer function, 114
- Global linearization, 508
- Globally positive definite function, 522
- H**
- Hamiltonian system, 473–481
eigenstructure of, 477–481
- Hankel matrix, 373–374
- Hankel realization, 372–377
- Hardware architecture design, 568–572
- Hessian, 443–445
- Homogeneous function, 12–13
- I**
- If-then-else statement, 569–570
- Implementation errors, in digital control, 319
- Incremental form, 596–598
- Indefinite function, 522
- Inertial control system
phase plane trajectory for, 462f, 464f
position trajectory for, 461f, 463f
velocity trajectory for, 461f, 464f
- Infeasible solution, 616
- Initcond, 602–608
- Input zero direction, 353
- Input-decoupling zeros, 301, 321
- Input-output-decoupling zeros, 301, 321
- Input–output differential equation, 254
- Input–output stability, 507
- Input–output stability and small gain theorem, 548–560

-
- absolute stability, 556–560
 closed-loop system with linear and nonlinear blocks, 556f
 nonlinear system for saturation nonlinearity, 558f
 nonlinearity and associated disc, 557f
 Nyquist plot of furnace and actuator, 557f
 closed-loop system with initial conditions, 555f–556f
 with saturation nonlinearity, 555f
 nonlinear closed-loop system, 550f
 unit circle in complex plane, 553f
 Instability theorems, 533–534
 Integral control, 408–411
 Integral controller, 162
 Integral time constant, 172
 Integrator windup, 590–593, 592f–593f
 control loop with actuator saturation, 590f, 592f
 Internal dynamics, 513–514
 Internal stability, 110–113
 digital control system with disturbance, 111f
 stable pole locations in z-plane, 110f
 Interpolator, 600
 Invariance
 of system zeros, 411–413
 of transfer functions and characteristic equations, 306–307
 Invariant zeros, 356
 Inverse z -transform, 361
 Invertible matrix ($A - In$), 323
 Irreducible realization, 301
iztrans command, 34
- J**
 Jacobians, 264–265
 Joseph form of Riccati equation, 457
 Jury test, 117–122, 118t–119t
- L**
 Lagrange multipliers, 445, 447
 Lagrangian, 445
 Laplace transform, 17–18, 23, 34–36, 44
 of state equation, 265–266
 Laplace transformation, 63–65, 76–77
 Left half plane (LHP), 115
 Leverrier algorithm, 272–277
 LHP. *See* Left half plane (LHP)
 Linear continuous-time SISO system, 254
 Linear equations, 616–617
 Linear matrix inequalities (LMI)
 from matrix equation, 615–617
 decision variables, 620
 editor, 626–627, 626f–627f
 MATLAB LMI commands, 620–627
 Schur complement, 617–619
 Linear quadratic regulator, 453–465
 free final state, 455–465
 Linear quadratic tracking controller, 470–473
 step response, 473f
 Linear state–space equations, 259–261
 solution of, 265–285
 Leverrier algorithm, 272–277
 real form for complex conjugate eigenvalues, 284–285
 state-transition matrix for diagonal state matrix, 278–283
 Sylvester’s expansion, 277–278
 Linear time-invariant system (LTI system), 66, 268
 case, 36–37
 digital systems, 103
 Linearity equation, 17
 Linearization of nonlinear state equations, 262–265
 lmiedit, 626
 Local maximum, 442–443
 Local minimum, 442–443
- Locally positive definite function, 522
 Logarithmic transformation, 517–518
 Long division, 21–22
 LQR control, 483–484
 LQR yields, 484
 LTI system. *See* Linear time-invariant system (LTI system)
 Lyapunov stability of linear systems, 527–530
 Lyapunov stability theory, 507, 522–536
 controller design based on, 546–548
 saturation nonlinearity, 546f
 velocity plot for step response for linear design, 545f
 estimation of domain of attraction, 534–536
 instabilitys theorems, 533–534
 of linear systems, 527–530
 linearization method, 532–533
 Lyapunov difference equation, 457
 Lyapunov equation, 531
 Lyapunov functions, 522–523
 MATLAB, 531
 rate of convergence, 526–527
 stability theorems, 524–526
- M**
 Manual mode, 593–596
 Marginally stable system, 321
 Markov parameter matrices, 372
 MATLAB commands, 114, 491, 601–608, 605f–607f
 bilinear transformation, 115f
 control error interpolation $\tilde{e}(t)$, 603f
 control signal, 603f
 for controllability testing, 337
 of digital control systems, 86–90
 for discrete-time state–space equations, 292
 for discrete-time systems, 49–50

MATLAB commands (*Continued*)
LMI commands, 620–627
for observability, 347
pole, 351
for pole placement, 402–403
Simulink, 87–90
solution of steady-state regulator problem, 468–470
time response of output regulator discussed, 470f
state–space representation in, 259
step response, 608f
in transfer function matrix, 287–289
tzero, 356–357
Z-transfer function in, 303
Matrix exponential, 267
Matrix norms, 325
Matrix Riccati equation, 455
MIMO system. *See* Multiinput –multioutput system (MIMO system)
Mincx, 624
Minimal polynomial, 351
Minimal realization, 301
Model predictive control (MPC), 441, 488–491
computation of control law, 490
constraints, 490–491
cost function, 489–490
MATLAB commands, 491
model, 489
Modes of system, 268–272
MPC. *See* Model predictive control (MPC)
Multiinput–multioutput system (MIMO system), 254, 350–351, 402
control, 571–572
parallel realization for, 366–368
Multiplexer (MUX), 571–572
Multiplication by exponential, 19–20
Multivariable systems
poles of, 350–357
from transfer function matrix, 351–355
zeros of, 350–357

from state–space models, 355–357
from transfer function matrix, 351–355
Multivariable zero, 353
MUX. *See* Multiplexer (MUX)

N

Negative definite function, 522
Negative semidefinite function, 522
Newton’s method, 114
nichols command, 49
Nonlinear control system design, 543
Nonlinear difference equations, 508, 517–518
logarithmic transformation, 517–518
Nonlinear digital control systems
discrete-time nonlinear controller design, 543–548
discretization of nonlinear systems, 508–517
equilibrium of nonlinear discrete-time systems, 518–521
input–output stability and small gain theorem, 548–560
Lyapunov stability theory, 522–536
nonlinear difference equations, 517–518
stability of analog systems with digital control, 537–539
state–plane analysis, 539–540
Nonlinear state equations, linearization of, 262–265
Nonlinear state–space equations, 259–261
Nonlinear systems, discretization of, 508–517
extended linearization
by input and state redefinition, 511–512
by input redefinition, 509–510
using matching conditions, 514–517

by output differentiation, 512–514
Nonminimal pole-zero cancellation, 301
Norm of vector, 325
num command, 49
Numerical analysis, 196–197
Numerical computation of matrix exponential, 277–278
Nyquist command, 49
Nyquist criterion, 122–137, 125f, 484–485
closed contours, 123f
contour for stability determination, 124f
furnace transfer function, 127f
modified contour for stability determination, 126f
phase margin and gain margin, 129–137, 134f–135f, 137f
model perturbation, 129f
Nyquist plot with, 130f
of system, 128f
Nyquist plot, 122, 130, 134f–135f
for furnace and actuator, 132f
furnace and actuator in vicinity, 133f
with negative gain margin, 131f
phase margin and gain margin, 130f
for position control system, 136f

O

Observability, 343–350
MATLAB commands, 347
rank condition, 345–346
of systems in normal form, 347
Observable realization, 369
Observable system, 325–328
Observer form, 369, 530
Observer state feedback, 421–429
observer eigenvalues, 423–428
step response, 428f, 434f
Zero-input response, 425f, 427f
obsv command, 347
Optimal control, 441, 447–453, 463f, 465f
Hamiltonian system, 473–481

- linear quadratic regulator, 453–465
 modification of reference signal, 491–498
 MPC, 488–491
 optimization, 442–447
 return difference equality and stability margins, 481–488
 steady-state quadratic regulator, 466–473
O
 Optimization, 442–447
 constrained, 445–447
 optimal trajectory for scalar system, 453f
 optimality conditions, 450t
 unconstrained, 442–445
 minimization and maximization, 443f
O
 Output equation, 257–258
 Output quadratic regulator, 467–468
 Output zero direction, 353
 Output-decoupling zeros, 301, 321
- P**
- Parallel realization, 363–368
 block diagram, 364f
 for MIMO systems, 366–368
 observable form, 369–370
 simulation diagram, 364f
 Partial differential equation, 515
 Partial fraction expansion, 23–31, 367
 PD control. *See* Proportional-derivative control (PD control)
 Percentage overshoot (PO), 147
 Perturbation, 263
 Phase and Gain margin, 129–137, 130f
 Phase margin (PM), 485
 Phase plane, 254–255
 Phase portrait, 254–255
 Phase variables, 254–255, 358
 PI control. *See* Proportional-integral control (PI control)
- PID control. *See* Proportional-integral-derivative control (PID control)
 PID controller. *See* Proportional-integral-derivative controller (PID controller)
p
place command, 381
 Planning horizon, 466
 Plant, 10–12
 PM. *See* Phase margin (PM)
 PO. *See* Percentage overshoot (PO)
 Pole allocation, 389–390
 Pole assignment, 389–390, 429–434, 429f
pole command, 114, 351
 Pole placement, 389–407
 closed-loop eigenvalues, 397–402
 MATLAB commands for pole placement, 402–403
 using matrix polynomial, 395–397
 for multiinput systems, 403–406
 by output feedback, 406–407
 by transformation to controllable form, 393–395
 Pole polynomial, 351
 Pole sensitivity, 95–97
 Pole-zero matched digital filters, 200
 Pole-zero matching, 199–201
 Poles of multivariable systems, 350–357
 from transfer function matrix, 351–355
poly command, 317
polyvalm command, 317
 Pontryagin's minimum principle, 450, 454
 Positive semidefinite, 522
 Positive system, 312
 Practical implementation of digital controllers, 567
 Prediction horizon, 488
 Prewarping equality, 203
 Primary strip, 186
 Proportional (P), 150
- Proportional control, 151–152
 design in z-domain, 190–195
 Proportional-derivative control (PD control), 150, 152–162, 156f–157f
 root locus plot of PD-compensated systems, 161f
 Proportional-integral control (PI control), 150, 162–168
 plot of controller angle (ϕ), 164f
 pole-zero diagram of, 163f
 root locus of, 166f
 Proportional-integral-derivative control (PID control), 150, 168–171, 171f, 174f
 empirical tuning of, 171–176
 time response for design, 170f
 with Ziegler-Nichols method, 175f
 Proportional-integral-derivative controller (PID controller), 588–598, 598f
 bumpless transfer between manual and automatic mode, 593–596
 filtering derivative action, 588–590
 incremental form, 596–598
 integrator windup, 590–593
- Q**
- Quadratic form $xTPx$, 522–523
 Quadratic programming problem, 491
 Quantization effects, 570
 Quantization error effects, 574–583
 characteristics of ADC, 575f
 coefficients of $F(z)$ computed, 582t
 control loop with actuator saturation, 580f
 controller output, 579f–580f
 poles of continuous time transfer function, 583t
 process output, 578f–579f

R

rank command, 347
Reachability, 330
Real form for complex conjugate eigenvalues, 284–285
Real-time system, 569
Reduced-order observer, 417–420, 418f
system with observer state feedback, 421f
Reducible pole-zero cancellation, 301
Reference signal, modification of, 491–498, 492f
 DMC, 492–498
Relative stability, 129
Resolvent matrix, 267
Return difference equality and stability margins, 481–488
plot of closed-loop pole location, 486f
stability region for frequency response, 485f
RHP. *See* Right half plane (RHP)
Riccati equation, 456, 467
 Joseph form of, 457
Right half plane (RHP), 115
rlocus command, 146, 182–184
Robotic manipulator control, 4–5, 6f
Robust Control Toolbox, 620
Robustness, 402
Root locus, 142–146
 design, 149–171
 control configurations, 151f
 PD control, 152–162
 PI control, 162–168
 PID control, 168–171
 pole locations and associated time responses, 150f
 proportional control, 151–152
 in design of second-order system, 148f
using MATLAB, 146
of PI-compensated system, 166f
plot of PD-compensated systems, 161f
plot of uncompensated systems, 160f

of second-and third-order systems, 144f
of system with integrator, 165f
Rosenbrock's system matrix, 356
Routh–Hurwitz criterion, 115–117

S

s-degree-of-freedom (*s*-D.O.F.), 259–261
Sampled parabolic, 82
Sampled step input, 84
Sampling frequency selection, 52–55
Sampling period
 antialiasing filters, 572–574
 choice of, 572–585
 effects of quantization errors, 574–583
 phase delay introduced by ZOH, 584–585
 switching, 598–610
 controller sampling from sampling period, 600f–601f
 dual-rate control, 608–610
 MATLAB commands, 601–608
Sampling theorem, 51–55, 52f
 selection of sampling frequency, 52–55
 waveforms with identical samples, 50f
Schur complement, 617–619
Schur stable system. *See*
 Asymptotically stable system
Schur–Cohn test, 117–118
Second-order hold, 63
Sensitivity analysis, 92f, 93–97
 pole sensitivity, 95–97
Sensor, 3
 delay, 53
Separation theorem, 422
Servo problem, 407–411
 closed-loop system, 409f, 411f
 control scheme with integral control, 409f
Settling time, 190
SI system. *See* Single-input system (SI system)

Similarity transformation, 303–307
invariance of transfer functions and characteristic equations, 306–307

Similarity transformation, 349–350
Simplex algorithm, 601
Simulink, 87–90, 89f
 discrete transfer function, 90f
 Library Browser, 88–90
parameters of discrete transfer function, 89f

selecting scope parameters, 91f
simulation diagram for step input and scope, 88f, 91f

step response of analog system with digital control, 92f

Simultaneous sample-and-hold system (SSH system), 571–572

Sinc function, 64
Single-input system (SI system), 388
Single-input–single-output system (SISO system), 254, 350–351, 393, 402, 429
 transfer functions, 357

Singular value, 535
SISO system. *See* Single-input –single-output system (SISO system)

Software architecture design, 568–572
selection of ADC and DAC, 570–572
 data acquisition system, 571f
software requirements, 568–570, 568f

Software verification, 569–570
Speed of computer hardware, 320
SSH system. *See* Simultaneous sample-and-hold system (SSH system)

Stability theorems, 524–526
Stabilizability, 338–343
Stable *z*-domain pole locations, 105–106

- State and output feedback, 388–389, 389f
 State equations, 257, 450, 457, 508
 State estimation, 413–420
 full-order observer, 414–416
 reduced-order observer, 417–420
 State feedback control, 388f, 411–412
 invariance of system zeros, 411–413
 observer state feedback, 421–428
 pole assignment using transfer functions, 429–434
 pole placement, 389–407
 servo problem, 407–411
 state and output feedback, 388–389
 state estimation, 413–420
 State plane, 254–255
 State portrait, 254–255
 State trajectories, 254–255
 State variables, 254–257
 State vector, 254–255
 State-transition matrix, 268
 for diagonal state matrix, 278–283
 for discrete-time system, 294
 State–plane analysis, 539–540
 equilibrium point classification, 540t
 unstable node, 540f
 State–space, 254–255
 State–space realizations, 357–370
 for continuous-time systems, 377–378
 controllable canonical realization, 358–363
 controllable form in MATLAB, 363
 parallel realization, 363–368
 simulation diagram, 362f
 systems with input differencing, 360–363
 systems with no input differencing, 358–360
 duality, 370–372
 Hankel realization, 372–377
 stability of, 320–329
 asymptotic stability, 320–324
 BIBO stability, 325–329
 State–space representation, 257–262
 discrete-time state–space equations, 289–293
 linear vs. nonlinear state–space equations, 259–261
 linearization of nonlinear state equations, 262–265
 in MATLAB, 259
 similarity transformation, 303–307
 solution of discrete-time state–space equations, 293–300
 solution of linear state–space equations, 265–285
 state variables, 254–257
 transfer function matrix, 285–289
 Z-transfer function from state–space equations, 300–303
 Steady-state error, 82–86
 Steady-state quadratic regulator, 466–473
 linear quadratic tracking controller, 470–473
 MATLAB solution of, 468–470
 output quadratic regulator, 467–468
 Steady-state regulator problem, 466
struct commands, 621
 Structure window, 626
 Submultiplicative norm property, 325–328
 Sylvester’s expansion, 277–278
 System state, 254–255
- T**
- tan(theta)** command, 153–157
- Tangent method, 172–173
 application of, 173f–174f
 Ziegler-Nichols tuning rules, 173t
- Terminal penalty, 447
- tf** command, 288
- theta** command, 153–157
- Time advance equation, 18–19
- Time constant, 190
- Time delay
 equation, 17–18
 theorem, 63–64
- Time response of discrete-time system, 36–41
 convolution summation, 36–39
 convolution theorem, 39–41
- Time-limited function, 53
- Tracking
 error, 82
 problem, 470
 time constant, 591–593
- Transfer function matrix, 285–289, 372
 MATLAB commands, 287–289
- Transfer functions, 67–68
 invariance of, 306–307
- Transmission zeros, 356
- Transport lag, systems with, 76–78
 continuous time function, 78f
- Two degree-of-freedom control scheme, 407
- tzero** command, 356–357
- U**
- Uncertainty equivalence principle, 422
- Uncontrollable modes, 330
- Uncontrollable subspaces, 330f
- Unobservable modes, 343
- Unobservable systems, standard form for, 348–349
- Unstable system, 321
- W**
- W-plane, 222–223
- Windup, integrator, 590–593, 592f–593f
- Z**
- Z-domain digital control system design, 184–195
 observation, 186
 remarks, 186

- Z-domain digital control system
design (*Continued*)
proportional control design in
z-domain, 190–195
Z-domain contours, 187–190
z-domain pole locations and
associated temporal
sequences, 185f
Z-domain root locus, 182–184
Z-domain transfer function,
67–68
z-transfer functions, 329–330,
357, 424
matrix, 301
from state-space equations,
300–303
in MATLAB, 303
z-transform function, 13–33,
40–44
final value theorem, 31–33
inversion of, 21–31
long division, 21–22
- partial fraction expansion,
23–31
properties of, 17–21
complex differentiation,
20–21
linearity equation, 17
multiplication by exponential,
19–20
time advance equation, 18–19
time delay equation, 17–18
solution
of difference equations,
34–36
of discrete-time state
equations, 295–300
of standard discrete-time signals,
14–17, 16f
discrete-time impulse, 15f
Zero dynamics, 513–514
Zero-input response, 267–268,
294
Zero-order hold (ZOH), 63, 575
frequency response of, 65f
- model of DAC, 63f
phase delay introduced by,
584–585
transfer function of, 63–64
Zero-state response, 268, 294
Zeros of multivariable systems,
350–357
from state-space models,
355–357
from transfer function matrix,
351–355
zgrid command, 189–190
Ziegler–Nichols method,
173–175
PID controller tuned with, 175f
tuning rules for closed-loop
method, 175t
Ziegler–Nichols tuning rules,
215–216
ZOH. *See* Zero-order hold (ZOH)
zpk command, 114
ztrans command, 34

DIGITAL CONTROL ENGINEERING

Analysis and Design

Third Edition

M. SAMI FADALI ANTONIO VISIONI

Provides concise and accessible coverage of core concepts in digital controls.

Digital controllers are part of nearly all modern personal, industrial, and transportation systems. Every senior or graduate student of electrical, chemical, or mechanical engineering should therefore be familiar with the basic theory of digital controllers. This new text covers the fundamental principles and applications of digital control engineering, with emphasis on engineering design.

Fadali and Visioli cover analysis and design of digitally controlled systems and describe applications of digital control in a wide range of fields. With worked examples and Matlab applications in every chapter and many end-of-chapter assignments, this text provides both theory and practice for those coming to digital control engineering for the first time, whether as a student or practicing engineer.

New to this edition:

- This new edition covers new topics such as Model Predictive Control and Linear Matrix Inequalities.
- To engage students, it has more illustrations and simple examples; the mathematical notation is reduced where possible, and it also includes intermediate mathematical steps in derivations.
- Companion website features resources for instructors, including Powerpoint slides and solutions.
- Extensive use of CAD Packages: Matlab and Simulink sections at the end of each chapter show how to implement concepts from the chapter.
- Contains review material to aid understanding of digital control analysis and design.
- Includes some advanced material to make it suitable for an introductory graduate level class or for two quarters at the senior/graduate level.
- The mathematics background required for understanding most of the book is based on what can be reasonably expected from the average electrical, chemical, or mechanical engineering senior.

About the authors:

M. Sami Fadali

Professor and Chair of Department of Electrical & Biomedical Engineering, College of Engineering, University of Nevada, Reno, NV, USA.

Antonio Visioli

Professor in Control Systems at the Department of Mechanical and Industrial Engineering of the University of Brescia, Brescia, Italy.

Test, Instrumentation, and Control



ACADEMIC PRESS

An imprint of Elsevier

elsevier.com/books-and-journals

ISBN 978-0-12-814433-6



9 780128 144336